#### ABSTRACT

Title of Dissertation:	STRATEGIES AND RESOURCES FOR RATIONAL VACCINE DESIGN AND ANTIBODY-ANTIGEN DOCKING AND AFFINITY PREDICTION
	Johnathan Guest, Doctor of Philosophy, 2022
Dissertation directed by:	Assistant Professor, Brian G. Pierce, Department of Cell Biology and Molecular Genetics

Antibody recognition of antigens is a unique class of protein-protein interactions, and increased knowledge regarding the determinants of these interactions has advanced fields such as computational vaccine design and protein docking. However, the diversity and flexibility of antibodies and antigens can hinder generation of potent vaccine immunogens or prediction of correct antibody-antigen interfaces, slowing progress in the design of vaccines and antibody therapeutics. In this thesis, we present strategies to design vaccine candidates for a difficult viral target and describe expanded resources for benchmarking and training antibody-antigen docking and affinity prediction algorithms.

We utilized rational design to develop candidate immunogens for a vaccine against hepatitis C virus (HCV), which represents a global disease burden despite recent advances in antiviral treatments. This design strategy produced a soluble and secreted E1E2 glycoprotein heterodimer with native-like antigenicity and immunogenicity by fusing ectodomains with a leucine zipper scaffold and a furin cleavage site. We developed additional constructs that incorporated synthetic or non-eukaryotic scaffolds or alternative ectodomains that included consensus sequences designed using a large reference database. Finally, we utilized previously published data on HCV antibody neutralization and E1E2 mutagenesis to predict residues that impact antibody neutralization and E1E2 heterodimerization, offering potential insights that can aid vaccine design.

To improve our knowledge of and accuracy in modeling antibody-antigen recognition, we assembled a set of antibody-antigen complex structures from the Protein Data Bank (PDB) that expanded Docking Benchmark 5, a widely used benchmark for protein docking. These complexes more than doubled the number of antibody-antigen structures in the benchmark and, based on tests of current algorithms, highlight significant challenges for docking and affinity prediction. Building on this resource, we assembled and curated a dataset of ~400 antibodyantigen affinities and corresponding structures, forming an expanded and updated benchmark to guide  $\Delta G$  prediction of antibody-antigen interactions. Using this dataset, we retrained combinations of terms from existing scoring functions and potentials, demonstrating that this resource can be used to improve antibody-antigen  $\Delta G$  prediction. Overall, these findings can advance HCV vaccine design and antibody-antigen interactions and to better display vaccine immunogens for induction of neutralizing antibodies.

# STRATEGIES AND RESOURCES FOR RATIONAL VACCINE DESIGN AND ANTIBODY-ANTIGEN DOCKING AND AFFINITY PREDICTION

by

Johnathan Guest

### Dissertation submitted to the Faculty of the Graduate School of the University of Maryland, College Park, in partial fulfillment of the requirements for the degree of Doctor of Philosophy 2022

Advisory Committee: Assistant Professor Brian G. Pierce, Chair Professor Roy A. Mariuzza Professor John Moult Assistant Professor Margaret A. Scull Professor Lai-Xi Wang © Copyright by Johnathan Guest 2022

## Dedication

To my wife Christine, my mother Kathryn, my father Daniel, and my sister Monica: thank you

for your constant love and support.

### Acknowledgements

I feel incredibly fortunate just to get to this point, and I have so many people to thank for helping me on my journey. Too many to mention here, in fact, so I will acknowledge those who had the largest impact.

I could not have conducted this research without the wisdom and guidance of my advisor, Professor Brian Pierce. He gave me the opportunity to pursue a variety of projects, even when they were a mix of wet lab and dry lab, and shared many insightful comments along the way. I would also like to thank my thesis committee for their useful comments and suggestions: Professor Roy Mariuzza, Professor John Moult, Professor Lai-Xi Wang, and Professor Margaret Scull, both as a committee member and a co-advisor for several years.

All members of the Pierce lab and our colleagues at the Institute of Bioscience and Biotechnology Research have provided invaluable advice and feedback. In the Pierce lab, thank you to Ragul Gowthaman, Rui Yin, Dongxiu Zhang Spiering, Ghazaleh Taherzadeh, and Stefan Ivanov for sharing your expertise and helpful suggestions in lab meeting and beyond. At IBBR, thank you to Ruixue Wang, Eric Toth, Andrezza Chagas, Thomas Fuerst, Khadija Elkholy, Liudmila Kulakova, Kyle Garagusi, Yunus Abdul, Yuxing Li, and Andrey Galkin for helping to guide and greatly advance computational and experimental projects. Special thanks to Ruixue and Dongxiu, whose research provided several key figures in chapters 2 and 3. Finally, I would like to thank our collaborators; in particular, Zhen-Yong Keck and Steven Foung at Stanford University for their contributions to chapters 2 and 3, along with Jing Zhou and Jeffrey Gray at Johns Hopkins for their contributions to chapter 5. Specific contributions to all figures and tables are listed below:

Figure 2.13: Western blot data provided by Liudmila Kulakova (IBBR) and Eric Toth (IBBR)

Figure 2.14: AUC data provided by Kinlin L. Chao (IBBR); SEC-MALS data provided by Thomas E. Cleveland IV (IBBR/NIST)

- Figure 2.15: mbE1E2 native gel western blot data provided by Andrezza Chagas (IBBR)
- Figure 2.16: ELISA data provided by Ruixue Wang (IBBR)
- Table 2.1: Quantitative ELISA data provided by Young Chang Kim (Stanford) and Zhen-Yong Keck (Stanford)
- Figure 2.17: SPR data provided by Eric Toth (IBBR)
- Figure 2.18: Endpoint titer and ID50 data provided by Ruixue Wang (IBBR)
- Figure 2.19: Neutralization data provided by Ruixue Wang (IBBR)
- Figure 3.3: ELISA data provided by Ruixue Wang (IBBR)
- Figure 3.4: ELISA data provided by Ruixue Wang (IBBR)
- Table 3.1: Dose-dependent ELISA data provided by Zhen-Yong Keck (Stanford)
- Figure 3.9: Western blot data provided by Dongxiu Zhang Spiering (IBBR)
- Figure 3.10: ELISA data provided by Dongxiu Zhang Spiering (IBBR)
- Table 5.8: SnugDock performance provided by Jing Zhou (Johns Hopkins)

This work was supported by NIH grants R21 AI154100, R21 AI126582, R01 AI102766,

R01 AI132213, and T32 AI125186.

# Table of Contents

Dedication	ii
Acknowledgements	iii
Table of Contents	v
List of Tables	ix
List of Figures	X
List of Abbreviations	xii
Chapter 1: Introduction 1.1 Hepatitis C virus 1.1.1 Discovery and health impacts 1.1.2 Current treatments 1.2 Vaccines and vaccine design 1.2.1 History and techniques	1 1 2 2 2 2
1.2.2 Structure-based vaccine design 1.2.3 Making an HCV vaccine	3
1.3 HCV immune evasion	4
1.3.2 Hypervariable regions 1.3.3 Resistance to antibodies	7 7 8
1.4.1 Antibody structure and classification 1.4.2 Recognition of antigens	8
<ul> <li>1.4.3 Structural characterization and modeling</li> <li>1.5 Protein docking</li> <li>1.5.1 Strategies and challenges</li> </ul>	9 11 11
<ul> <li>1.5.2 CAPRI and docking benchmarks</li> <li>1.6 Protein affinity prediction</li> </ul>	12
1.6.2 Community resources 1.7 Dissertation overview	14
Chapter 2: Design of a native, secreted hepatitis C virus E1E2 heterodimer Abstract	19 19
2.1 Introduction	20
2.2.1 Protein expression 2.2.2 Antibodies 2.2.3 Protein purification and size exclusion chromatography	22 23 23
2.2.4 Computational design of coiled coil assemblies 2.2.5 Peptide synthesis and characterization	24
2.2.6 SEC-MALS 2.2.7 SDS-PAGE and western blot	26 26

2.2.8 Analytical ultracentrifugation (AUC)	27
2.2.9 Enzyme-linked immunosorbent assay (ELISA)	28
2.2.10 Determination of antibody affinity by quantitative ELISA	28
2.2.11 Surface plasmon resonance	29
2.2.12 Animal immunization	29
2.2.13 HCV pseudoparticle generation	30
2.2.14 HCVpp neutralization assavs	30
2.2.15 Statistical comparisons	31
2.3 Results	31
2.3.1 Design of sE1E2 constructs	31
2.3.2 sE1E2.LZ forms an intact E1E2 complex	34
2.3.3 Purification of sE1E2.LZ.	38
2.3.4 Analytical characterization of heterogeneity in solution	46
2.3.5 sE1E2.LZ exhibits native-like E1E2 antigenicity and robust immunogenicity	49
2.4 Discussion	55
Chapter 3: Design of soluble hepatitis C virus E1E2 assemblies with alternative scaffolds of	or
ectodomains	59
Abstract	59
3.1 Introduction	60
3.2 Methods	63
3.2.1 E1E2 consensus and alternative isolate sequences	63
3.2.2 Selection of alternative scaffolds	64
3.2.3 Protein expression	65
3.2.4 Antibodies	66
3.2.5 Protein purification and size exclusion chromatography	66
3.2.6 SDS-PAGE and western blot	67
3.2.7 Enzyme-linked immunosorbent assay (ELISA)	67
3.2.8 Determination of antibody affinity by quantitative ELISA	68
3.3 Results	69
3.3.1 Design of sE1E2 constructs with synthetic scaffolds	69
3.3.2 sE1E2.SZ and sE1E2.HH mimic sE1E2.LZ secretion and antigenicity	71
3.3.3 Design of sE1E2 constructs with alternative scaffolds	74
3.3.4 Design of sE1E2 constructs with alternative ectodomains	77
3.3.5 Characterization of alternative sE1E2 constructs	81
3.4 Discussion	85
Chapter 4: Prediction of hepatitis C virus polymorphisms impacting antibody neutralization	n and
residues critical for E1E2 heterodimeric assembly	90
Abstract	90
4.1 Introduction	91
4.2 Methods	94
4.2.1 Collection of antibody neutralization data	94
4.2.2 Prediction of polymorphisms contributing to neutralization changes with SNAPR.	94
4.2.3 Pairwise comparisons of antibody neutralization data	95
4.2.4 Computational mutagenesis of polymorphisms predicted to contribute to neutralize	ation
changes	95

4.2.5 Collection of E1E2 mutagenesis data	97
4.2.6 Clustering of mutagenesis data	97
4.3 Results	98
4.3.1 Neutralization datasets used for predictions	98
4.3.2 SNAPR predicted E1E2 polymorphisms contributing to neutralization changes	99
4.3.3 Pairwise sequence comparisons predicted E1E2 polymorphisms contributing to	
neutralization changes	103
4.3.4 Polymorphisms predicted to impact antibody neutralization modeled using	
computational mutagenesis	105
4.3.5 Hierarchical clustering of E1E2 mutagenesis datasets	110
4.3.5 Critical E1E2 interface residues predicted through clustering by residue	114
4.3.6 Predicted E1E2 contacts found in E1E2 heterodimer structure	117
4.4 Discussion	118
Chapter 5: An expanded benchmark for antibody-antigen docking and affinity prediction r	eveals
insights into antibody recognition determinants	122
Abstract	122
5.1 Introduction	122
5.2 Methods	125
5.2.1 Benchmark assembly	125
5.2.2 Protein-protein docking	127
5.2.3 Interface analysis and affinity prediction	131
5.2.4 Analysis of conformational changes	132
5.2.5 Quantification and statistical analysis	133
5.3 Results	134
5.3.1 Benchmark assembly and composition	134
5.3.2 Binding conformational changes	144
5.3.3 Global docking prediction	151
5.3.4 Local docking perturbations	158
5.3.5 Binding affinity prediction	164
5.4 Discussion	170
Chapter 6: A curated dataset of antibody-antigen affinities and structures to facilitate	
development of affinity prediction algorithms	173
Abstract	173
6.1 Introduction	174
6.2 Methods	177
6.2.1 Collection of cases in the antibody-antigen affinity dataset	177
6.2.2 Curation of antibody-antigen affinity dataset	178
6.2.3 Analysis of affinity predictors	180
6.2.4 Comparison of correlations in affinity subsets	181
6.2.5 Additional case information in affinity dataset	182
6.2.5 Individual terms in REF15 and beta16 scoring functions	183
6.2.6 Individual terms in IRAD and ZRANK scoring functions	184
6.2.7 Regression analysis and cross-validation	187
6.2.8 REF15/beta16-based sets	187
6.2.9 IRAD/ZRANK-based sets	187

6.2.9 Composite REF15/IRAD sets	188
6.2.10 Data for independent test set	188
6.2.11 Modeling and scoring of antibody-antigen complexes from independent test set	189
6.2.12 Predictions of optimized scoring functions on independent test set	190
6.3 Results	190
6.3.1 Dataset assembly and diversity	190
6.3.2 Performance of existing scoring functions as affinity predictors	194
6.3.3 Correlations of predictors by annotation subset	196
6.3.4 Correlations of individual terms with $\Delta G$ values	200
6.3.5 Selection and retraining of input terms for antibody-antigen affinity prediction	203
6.3.6 Performance of top retrained models on independent test set	211
6.4 Discussion	214
Chapter 7: Summary and future directions	219
Publication Information	222
Bibliography	224

## List of Tables

Table 2.1 Binding affinity of mbE1E2, sE1E2.LZ, and sE2 to a panel of HMAbs	53
Table 3.1 Antigenic analysis of mbE1E2, sE1E2.LZ, and sE1E2.SZ by quantitative ELI	SA.
	74
Table 4.1 Summary of neutralization data from previously published datasets	99
Table 4.2 Summary of residue positions with a polymorphism found as SNAPR hit	102
Table 4.3 Summary of residue positions predicted to contribute to neutralization change	ges
through pairwise comparisons.	
Table 4.4 $\Delta\Delta G$ predictions of changes in E1 monomer stability	108
Table 4.5 $\Delta\Delta G$ predictions of changes in E2 monomer stability.	108
Table 4.6 Computational mutagenesis summary and classification for each structure	110
Table 4.7 Relative binding averages for residue clusters from merged E2 mutagenesis d	lata.
	115
Table 4.8 Relative binding averages for residue clusters from E1E2 mutagenesis data.	116
Table 4.9 Summary of predicted E1E2 interface residues.	117
Table 5.1 New antibody-antigen benchmark cases organized by difficulty category	136
Table 5.2 Additional details for new antibody-antigen test cases.	138
Table 5.3 Additional details and references for new antibody-antigen affinity cases	140
Table 5.4 Antibody CDR loop sequences of benchmark cases.	141
Table 5.5 sdAb CDR3 average RMSDs for subsets with or without interloop disulfide.	151
Table 5.6 Global docking ranks of top Acceptable and Medium models.	155
Table 5.7 Comparison of ZDOCK results with or without glycans removed in unbound	L
antigen.	158
Table 5.8 SnugDock local perturbation performance by test case.	160
Table 5.9 Pearson correlation, and correlation p-value, of functions/terms with	
experimentally determined $\Delta Gs$	168
Table 5.10 Correlations with experimental ΔG values for ΔASA and Rosetta REF15	
stratified by I-RMSD.	170
Table 6.1 Correlations of REF15, beta16, IRAD, and ZRANK scores with ΔG values	196
Table 6.2 Scoring function correlations with affinity values by measurement method	199
Table 6.3 Scoring function correlations with affinity values by structure resolution	200
Table 6.4 Correlations of REF15 and beta16 scoring terms with affinity values	202
Table 6.5 Correlations of IRAD and ZRANK scoring terms with affinity values	202
Table 6.6 Terms and weights of retrained models selected by stepwise regression	208
Table 6.7 $\Delta G$ prediction of models retrained through stepwise regression.	209
Table 6.8 Affinity prediction of models retrained with Ridge, LASSO, Elastic net	
regression following 5-fold cross-validation.	211
Table 6.9 Correlations of top retrained models with affinities in independent test set	213
Table 6.10 Correlations of top retrained models with neutralization data of 45_01dG5	
isolate in independent test set	213

# List of Figures

Figure 2.1 Design of sE1E2 constructs.	33
Figure 2.2 Characterization of the peptide complex CC1+CC2.	34
Figure 2.3 E1 and E2 western blots of sE1E2 supernatant	36
Figure 2.4 Western blots of supernatant from E1-Jun/E2-Fos co-expression.	36
Figure 2.5 E1 and E2 western blots of sE1E2 cell lysate.	37
Figure 2.6 Quantitative western blots comparing sE1E2.LZ supernatant and cell lysate.	37
Figure 2.7 Size exclusion chromatography of sE1E2.LZ, sE1E2GS3, and mbE1E2.	40
Figure 2.8 Size exclusion chromatograph of sE1E2GS3.	41
Figure 2.9 Yield and purity of mbE1E2, sE1E2.LZ, and sE1E2GS3 in SDS-PAGE.	42
Figure 2.10 sE1E2GS3 fractions from SEC analyzed by SDS-PAGE and western blot	43
Figure 2.11 sE1E2.LZ fractions from SEC analyzed by SDS-PAGE and western blot	44
Figure 2.12 mbE1E2 elution fractions from SEC analyzed by western blot.	45
Figure 2.13 Deglycosylation of mbE1E2, sE1E2.LZ, and sE2.	45
Figure 2.14 Analytical characterization of sE1E2.LZ and mbE1E2 size and heterogeneit	y.
· · ·	48
Figure 2.15 mbE1E2 and sE1E2.LZ size and heterogeneity in native gel	49
Figure 2.16 Initial antigenicity screening of sE1E2 designs in ELISA.	52
Figure 2.17 Measurement of binding to the CD81 receptor by SPR.	54
Figure 2.18 Immunogenicity assessment of sE2, mbE1E2, and sE1E2.LZ.	54
Figure 2.19 Calculated curves for H77C HCVpp neutralization by immunized (Day 56)	
murine sera.	55
Figure 3.1 Design of sE1E2 constructs with synthetic scaffolds	70
Figure 3.2 Evaluation of sE1E2 secretion to supernatant in western blot.	72
Figure 3.3 Binding of sE1E2 constructs and mbE1E2 to HCV HMAbs in ELISA	73
Figure 3.4 Binding of sE1E2 constructs and mbE1E2 to HCV HMAbs in ELISA at eleva	ated
temperatures.	73
Figure 3.5 Design of sE1E2 constructs with alternative scaffolds	76
Figure 3.6 Phylogenetic tree of cons.80 with sequences from genotypes 1-7.	79
Figure 3.7 Phylogenetic tree of cons1.92.5 with sequences from genotype 1 subtypes	80
Figure 3.8 Comparison of H77 and consensus sequences at residue positions of key	
epitopes	81
Figure 3.9 Detection of alternative sE1E2 constructs in western blot	84
Figure 3.10 Antibody binding to alternative sE1E2 constructs in ELISA.	85
Figure 4.1 Example of SNAPR predictions of E1E2 polymorphism contributions.	.101
Figure 4.2 Visualization of antibody groups using hierarchical clustering.	.112
Figure 4.3 Heatmap of merged mutagenesis dataset clustered by residue.	.113
Figure 5.1 Docking and affinity benchmark composition.	.143
Figure 5.2 Binding conformational changes of antibody-antigen benchmark cases	.147
Figure 5.3 Comparison of residue-level conformational changes by antibody chain type.	148
Figure 5.4 Structural diversity of benchmark cases	.149
Figure 5.5 Binding conformational changes of antibody residues near the antigen interfa	ace
by amino acid.	.150
Figure 5.6 Docking performance on the antibody-antigen benchmark.	.154
Figure 5.7 Comparison of docking success rates in ZDOCK models.	.157
Figure 5.8 SnugDock binding funnels for two benchmark cases.	.162

Figure 5.9 SnugDock binding funnels for Rigid benchmark cases.	163
Figure 5.10 SnugDock binding funnels for Medium and Difficult benchmark cases	164
Figure 5.11 Affinity predictions on benchmark cases.	167
Figure 5.12 I-RMSD, ΔASA, and Rosetta REF15 scores versus experimentally determ	ined
ΔĞs	169
Figure 6.1 Summary of diversity in the antibody-antigen affinity dataset.	193
Figure 6.2 Ranges of $\Delta G$ values and structural resolution in affinity dataset	194
Figure 6.3 Predictive performance of existing scoring functions.	199
Figure 6.4 Heatmap of correlations between terms output by IRAD	203
Figure 6.5 Predictive performance of top retrained models	210
Figure 6.6 Predictions of retrained models with significant correlations to 45_01dG5	
neutralization data	213

# List of Abbreviations

aa	amino acid
AB-Bind	Antibody-Bind database
AF4	Asymmetric flow field flow fractionation
ANARCI	Antigen receptor Numbering And Receptor Classification
АроЕ	apolipoprotein E
AR	antigenic region
AUC	analytical ultracentrifugation
auc	area under the curve
beta16	Rosetta "beta_nov16" energy function
β-OG	n-Octyl-β-D-Glucopyranoside
BLAST	Basic Local Alignment Search Tool
BLI	bio-layer interferometry
BM5	Docking Benchmark 5.0
BM5.5	Docking Benchmark 5.5
bnAb	broadly neutralizing antibody
Ca	alpha Carbon atom
CAPRI	Critical Assessment of Predicted Interactions
CD81-LEL	large extracellular loop of CD81
CDR	complementarity determining region
CDR1	first complementarity determining region loop (mAb or sdAb chain)
CDR2	second complementarity determining region loop (mAb or sdAb chain)

CDR3	third complementarity determining region loop (mAb or sdAb chain)
CDRH1	first complementarity determining region loop on mAb heavy chain
CDRH2	second complementarity determining region loop on mAb heavy chain
CDRH3	third complementarity determining region loop on mAb heavy chain
CDRL1	first complementarity determining region loop on mAb light chain
CDRL2	second complementarity determining region loop on mAb light chain
CDRL3	third complementarity determining region loop on mAb light chain
cons.80	HCV E1E2 consensus sequence, genotypes 1-7
cons1.92.5	HCV E1E2 consensus sequence, genotype 1
Cryo-EM	Cryogenic electron microscopy
ΔASA	change in accessible surface area
$\Delta G$	change in Gibbs free energy
$\Delta\Delta G$	change in the change in Gibbs free energy
DAA	direct acting antiviral
DLS	dynamic light scattering
DMEM	Dulbecco's modified Eagle medium
ECL	enhanced chemiluminescence
ELISA	Enzyme-linked immunosorbent assay
Env	HIV envelope glycoprotein
FFT	Fast Fourier Transform
f <sub>nat</sub>	fraction of native contacts
f <sub>non-nat</sub>	fraction of non-native contacts
Fu	fraction unaffected

Fv	antibody variable domain
GNA	Galanthus Nivalis Agglutinin
GNL	Galanthus Nivalis Lectin
gp120	HIV glycoprotein 120
gp41	HIV glycoprotein 41
HAstV	human astrovirus
HCV	hepatitis C virus
HCVcc	HCV cell culture virus
HCVpp	HCV pseudoparticle
HETATM	hetero atom
HIV	human immunodeficiency virus
НМАЬ	Human monoclonal antibody
HMW	high molecular weight
HRP	Horseradish peroxidase
HVR1	hypervariable region 1
HVR2	hypervariable region 2
IC50	half-maximal inhibitory concentration
ID50	half-maximal inhibitory dose
IgG	immunoglobulin G
IgNAR	immunoglobulin new antigen receptor
IMAC	Immobilized metal affinity chromatography
IP	Intraperitoneal
IRAD	Integration of Residue- and Atom-based potentials for Docking

I-RMSD	interface root-mean-square distance
ITC	isothermal titration calorimetry
L-RMSD	ligand root-mean-square distance
K <sub>A</sub>	association constant
K <sub>D</sub>	equilibrium dissociation constant
kDa	kilodalton
KinExA	kinetic exclusion assay
k <sub>off</sub>	dissociation rate constant
kon	association rate constant
LASSO	Least Absolute Shrinkage and Selection Operator
LVP	lipoviroparticle
mAb	monoclonal antibody
mAU	milli absorbance units
mbE1E2	membrane-bound E1E2
MDS	motif dock score
μΜ	micromolar
MLV	murine leukemia virus
MSA	multiple sequence alignment
nM	nanomolar
OD	optical density
PBS	Phosphate-buffered saline
PCPP-R	Poly[di(carboxylatophenoxy)phosphazene] formulated with resiquimod
PDB	Protein Data Bank

pМ	picomolar
PRODIGY	protein binding energy prediction
RBD	receptor binding domain
REF15	Rosetta Energy Function 2015
RMSD	root-mean-square distance
RMSE	root-mean-square error
RSV	respiratory syncytial virus
RU	resonance unit
S	sedimentation coefficient
S2	SARS-CoV-2 spike S2 subunit
SAbDab	Structural Antibody Database
SARS-CoV	severe acute respiratory syndrome coronavirus
SARS-CoV-2	severe acute respiratory syndrome coronavirus 2
scFv	single-chain antibody fragment
sdAb	single domain antibody
SDS-PAGE	sodium dodecyl sulfate-polyacrylamide gel electrophoresis
sE1E2	secreted E1E2
sE1E2.cons.80	secreted E1E2 with cons.80 sequences as ectodomains
sE1E2.cons1.92.5	secreted E1E2 with cons1.92.5 sequences as ectodomains
sE1E2.1U0I	secreted E1E2 scaffolded with IAAL-E3/IAAL-K3
sE1E2.3CFI	secreted E1E2 scaffolded with EpsI/EpsJ
sE1E2.1.11.6	secreted E1E2 with 1.11.6 isolate sequences as ectodomains
sE1E2.1a38	secreted E1E2 with 1a38 isolate sequences as ectodomains

- sE1E2.CC secreted E1E2 scaffolded with CC1+CC2 hexamer
- sE1E2.FD secreted E1E2 scaffolded with foldon
- sE1E2.HH secreted E1E2 scaffolded with synthetic hetero-hexamer
- sE1E2GS3 secreted E1E2 fused with a glycine-serine linker
- sE1E2.LZ secreted E1E2 scaffolded with Fos-Jun leucine zipper
- sE1E2.R6 secreted E1E2 with furin cleavage site but no scaffold
- sE1E2RevGS3 secreted E1E2 fused with a glycine-serine linker, ectodomains reversed
- sE1E2.SpyC secreted E1E2 scaffolded with SpyCatcher on E1 and SpyTag on E2
- sE1E2.SpyT secreted E1E2 scaffolded with SpyTag on E1 and SpyCatcher on E2
- sE1E2.SZ secreted E1E2 scaffolded with SYNZIP1/SYNZIP2
- sE2 soluble E2
- SEC size exclusion chromatography
- SEC-MALS size exclusion chromatography with multi angle light scattering
- SiPMAB Single-Point Mutant Antibody Binding database
- SKEMPI Structural database of Kinetics and Energetics of Mutant Protein Interactions
- SNAPR Subject-adjusted Neutralization Antibody Prediction of Resistance
- SOSIP designed HIV gp120-gp41 hexamer with key stabilizing mutations
- SPR surface plasmon resonance
- SR-BI scavenger receptor class B type I
- ssRNA single-stranded RNA
- SV sedimentation velocity
- TBS Tris-buffered saline
- TMB 3,3',5,5'-Tetramethylbenzidine

TMD	transmembrane domain
tPA	tissue plasminogen activator
VACV	vaccinia virus
ZAPP	Zlab Affinity for Protein-Protein interaction
ZRANK	Zlab Rerank

## Chapter 1: Introduction

#### **1.1 Hepatitis C virus**

#### 1.1.1 Discovery and health impacts

Hepatitis C virus (HCV) is an enveloped, positive-sense ssRNA virus in the Flaviviridae family. It was initially known as "non-A, non-B Hepatitis" before being described as HCV by several labs in 1989 (1). Harvey Alter, Michael Houghton, and Charles Rice were awarded the 2020 Nobel Prize in Medicine for their discovery and characterization of HCV (2), demonstrating its clinical impact (https://www.nobelprize.org/prizes/medicine/2020/summary/). HCV continues to be a worldwide health burden, with a recent estimate of 58 million infected according to the WHO (https://www.who.int/news-room/fact-sheets/detail/hepatitis-c). In addition, an estimated 1.5 million new infections occur every year, with an estimated 290,000 deaths from HCV in 2019. HCV is a bloodborne virus largely restricted to humans and chimpanzees, with infection often through intravenous drug use, sexual transmission, or other practices leading to blood exposure (https://www.who.int/news-room/fact-sheets/detail/hepatitis-c). HCV infections can be cleared in some cases following an acute phase, but approximately 75% of cases become chronic, potentially causing severe inflammation of the liver (3-6). Deaths from HCV are often attributed to the worst outcomes caused by these chronic infections, including cirrhosis, liver failure, or hepatocellular carcinoma (4, 5). To make matters worse, a substantial number of acute and early chronic HCV infections are asymptomatic, requiring broad surveillance and testing to determine the extent of viral spread (7, 8). With increasing case numbers and deaths, meeting the WHO's 2030 global targets to reduce new HCV infections by 90% and deaths by 65% remains an enormous undertaking (9, 10).

#### 1.1.2 Current treatments

Despite the widespread and increasing health burden of HCV, rapid progress has been made to treat and cure infections. Direct acting antivirals (DAAs) with high cure rates (>90%) introduced in the last 10 years were a paradigm shift in treatment of HCV infection, and were safer and more effective than pegylated-interferon with ribavirin, the previous standard of care (11, 12). DAAs target NS3/NS4A protease, NS5A, or NS5B polymerase to inhibit HCV replication, and are often used in combination (12-14). However, the use of DAA treatments as the sole tool for HCV control and eradication comes with severe limitations. To start, the WHO estimates that only 21% of HCV-infected individuals are tested, and only 62% of those tested have been treated with DAAs (https://www.who.int/news-room/fact-sheets/detail/hepatitis-c). These metrics showing suboptimal treatment levels reflect barriers of access to DAAs, financial or otherwise (15). Though DAA cure rates are very high, resistance has been identified in numerous cases (16). Furthermore, curing HCV infection with DAAs does not prevent later reinfection with a different HCV genotype or subtype (17, 18). There is also limited evidence to suggest that the risk of hepatocellular carcinoma progression decreases following DAA treatment, showing that DAAs may not alleviate some of the worst impacts of infection (19). DAAs are underiably valuable tools to help control and eradicate HCV, but there is an ongoing and urgent need to pair DAAs with a low-cost preventative treatment, especially for at-risk populations; that treatment is an HCV vaccine (20-22).

#### **1.2 Vaccines and vaccine design**

#### 1.2.1 History and techniques

The development of vaccines has a long and storied history of protecting against oncedevastating diseases, beginning with inoculations against smallpox pioneered by Edward Jenner in 1796 and later including rabies, measles, and polio (23). Part of the success of vaccines comes from the diversity of formulations that can stimulate a protective immune response. The field is identified by several major types of vaccines: whole virus (attenuated or inactivated), subunit (protein or polysaccharide), and genetic material (DNA, viral vector) (23, 24). Each type has led to a safe and effective vaccine against a pathogen, giving researchers many avenues to develop a product that can induce protective responses (23, 25). The techniques of vaccine development have continued to advance, exemplified by cutting-edge mRNA technology that helped to produce highly effective vaccines against SARS-CoV-2 within a year (26-28). Vaccines can also be improved or optimized through rational vaccine design, or a loosely defined set of strategies to select vaccine antigens or adjuvants that best focus or boost immune responses (29, 30). This process can be interconnected with knowledge of broadly neutralizing antibody (bnAb) responses against a particular antigen and reverse vaccinology, which uses bioinformatics techniques to search the genome of a pathogen for potent immunogens (31, 32). Computational methods and advances in deep sequencing also facilitated reverse vaccinology 2.0, or structural vaccinology, where B cell repertoires and structures of antibody-antigen complexes are often used to inform the selection of optimal antigens (33).

#### 1.2.2 Structure-based vaccine design

As with rational vaccine design, structure-based vaccine design is not a single strategy but a collection of methods to improve vaccine immunogens, often by introducing mutations that stabilize a protein antigen, boost immune responses, or both (34, 35). These designs are directly informed by the structural characterization of promising vaccine antigens, both alone and in complex with bnAbs. Known structures can provide a promising starting point to boost bnAb responses, even though the capacity of an antigen to be recognized by antibodies, or antigenicity, does not guarantee that the antigen has a similar capacity to stimulate a robust response from the immune system, or immunogenicity (36, 37). Structure-based vaccine design has resulted in some high-profile advancements, namely SARS-CoV-2 vaccines that incorporate two or six targeted proline mutations to stabilize spike trimers in a pre-fusion conformation (38, 39). These SARS-CoV-2 designs also include foldon, a self-assembling trimer, as a C-terminal scaffold, stabilizing the entire glycoprotein assembly as a vaccine antigen and facilitating structural characterization (39). Using foldon or other self-assembling proteins, the incorporation of scaffolding has become a key component of several design efforts for viral glycoproteins, including vaccine candidates for RSV and influenza (40). Additional modifications to vaccine antigens are tailored to the specific requirements for stabilization or bnAb responses, but two examples are instructive. A designed RSV vaccine candidate named DS-Cav1 incorporated foldon and several mutations informed by the glycoprotein F structure, including an added disulfide and cavity-filling hydrophobic residues, to favor a pre-fusion conformation of the antigen that presented an epitope recognized by a key bnAb (41). HIV SOSIP designs have been instrumental in structurally characterizing the trimeric assembly of gp120-gp41 glycoproteins alone (42) and in complex with bnAbs such as VRC01 (43-45). These designs are defined by a host of modifications, including cysteine mutations to add a disulfide bond, a proline mutation to restrict flexibility, truncation of gp41, and insertion of a furin cleavage site (6xArg) between gp120 and gp41 (46).

#### 1.2.3 Making an HCV vaccine

Despite 30 years of attempts with a myriad of vaccine technologies, an effective HCV vaccine has not yet been developed (47, 48). Recently, a T cell based HCV vaccine completed phase 1 human clinical trials, but could not prevent chronic infection despite inducing HCV-specific responses (49). In these attempts, proteins have been targeted across much of the HCV

genome, which includes capsid protein Core, glycoproteins E1 and E2, ion channel protein p7, and six non-structural proteins. E1E2 glycoproteins have formed the primary target for most HCV vaccines, as these proteins associate as heterodimers on the surface of the virion (50) and are thought to form a larger assembly of a trimer of heterodimers mediated by contacts between C-terminal transmembrane domains and residues in the E1 and E2 ectodomains (51). E2 is an especially important target of B cell based HCV vaccines due to its critical interactions with co-receptors at several steps of attachment and viral entry, including sequential steps of binding to SR-BI and the large extracellular loop of CD81 (CD81-LEL), later triggering HCV membrane fusion and endocytosis (3, 52). As shown in a structure of E2 in complex with CD81-LEL, both proteins undergo conformational changes that appears to facilitate HCV attachment to the host membrane (53).

Naturally, blocking E2 interaction at one or more of these steps with antibodies is a longstanding vaccine strategy, which is supported by the characterization of several bnAbs directly inhibiting CD81 binding (54, 55). Some efforts have incorporated rational or structure-based design approaches to test vaccine candidates of E1 or E2 epitopes, the E2 ectodomain, or the E1E2 heterodimer (56). The design of individual epitopes, either through mutations or scaffolding, is more tractable considering the well-characterized bnAb epitopes and corresponding antibody-antigen complex structures (57). However, the design of scaffolded and homogenous E1E2 heterodimer has been elusive, with several candidates that produced secreted E1E2 but showed little binding to anti-E1E2 bnAbs AR4A and AR5A (58-60). An E1E2 heterodimer structure was reported recently (61), but has not yet been released and only represents a small fraction of HCV genotype diversity, with some E1E2 residues unresolved. Severe limitations in knowledge of the E1E2 structure have left research efforts without a suitable design template, undeniably hindering

structure-based vaccine design specifically and our understanding of assembly and entry mechanisms generally.

#### **1.3 HCV immune evasion**

#### 1.3.1 Genetic diversity and glycans

Despite the characterization of multiple glycoprotein epitopes that induce bnAbs (62, 63), HCV has several overlapping mechanisms of immune evasion that have made vaccine development extraordinarily difficult. Analogous to other difficult targets such as HIV (64, 65), genetic diversity and glycan shielding are key challenges to HCV vaccine development. There are eight documented HCV genotypes, with multiple subtypes reported in several genotypes (66). The genomes of these genotypes are approximately 70% identical at the nucleotide level (67), showing remarkable interspecies diversity that is largely driven by error-prone replication (68, 69) and recombination in rare cases (70). Additionally, individual HCV isolates usually develop quasispecies during a single infection, leading to similar but genetically distinct populations that may make clearing the virus more difficult (71, 72). No genotype is dominant worldwide, with genotype 1 most prevalent at ~46%, forcing any comprehensive vaccine strategy to address multiple genotypes (73). E1E2 glycoproteins include four conserved glycosylation sites on E1 and eleven mostly conserved glycosylation sites on E2, helping to shield epitopes recognized by bnAbs (74-76). HCV glycans also show substantial heterogeneity and varying site occupancies (77), making efforts to modulate or remove glycans more important for vaccine development or for therapeutics directly targeting viral glycans (78). The removal of some glycans on E1 and E2 also leads to abrogation of HCV assembly and infection, showing the importance of the glycan shield in viral function and immune evasion (79, 80).

#### 1.3.2 Hypervariable regions

Other more unusual mechanisms of HCV immune evasion have been described, making vaccine design even more challenging. E2 ectodomain sequences contain three hypervariable regions that are highly diverse and flexible, providing another avenue to shield epitopes recognized by bnAbs (62, 81). The most prominent example is hypervariable region 1 (HVR1), a 27 aa long N-terminal region that is thought to shield bnAb epitopes in the CD81 binding site and antigenic domain E (82-84). HVR1 is also immunodominant, inducing antibody responses that are strainspecific and not broadly neutralizing (85). Another immunodominant epitope on E2 that induces non-neutralizing antibodies is antigenic domain A, located on the back layer that does not interact with CD81 (86-88). During viral attachment and entry, HVR1 is also a crucial stabilizer in reported "viral breathing" mechanisms, specifically in maintaining a closed state of the E2 front layer by concealing the CD81 binding site and keeping domain E in a compact  $\beta$ -hairpin conformation (89, 90). With HVR1 removed, E2 preferentially adopts an open or receptor-bound state, where the CD81 binding site is exposed for engagement by bnAbs and domain E is in an elongated conformation. However, HVR1 removal in vaccine designs has failed to boost bnAb responses and may be counterproductive (91). Removal of domain A through glycan masking also did not have a large effect on immunogenicity (92).

#### *1.3.3 Resistance to antibodies*

Even when epitopes on E1E2 recognized by bnAbs are exposed, HCV has additional layers of immune evasion that may be context dependent and more difficult to measure through existing methods. In previous research of bnAb responses to HCV isolates, sequence polymorphisms that increase resistance to antibody neutralization have been observed, primarily in E2. While this effect may be restricted to a particular antibody epitope or to the context of one isolate, the high genetic diversity and emergence of quasispecies during infection gives HCV plenty of opportunities to escape bnAb responses. The most dramatic example of these polymorphisms has been observed in antigenic domain E, an epitope that is highly conserved across genotypes (93, 94). In rare cases, an asparagine at position 417 in E2 mutates to serine or threonine, causing a glycan shift from N417 to N415 (95, 96). This small shift ablates binding of bnAbs that specifically recognize the  $\beta$ -hairpin conformation of domain E (57, 96), in one case leading to viral rebound in a clinical trial involving domain E bnAb HCV1 (97). Other research has implicated polymorphisms of extra-epitopic residues, or residues not directly bound by an antibody, as contributors to neutralization resistance (98-104). These polymorphisms have often been found in HVR1 or near the CD81 binding site, though the mechanisms of these effects can be unclear and dependent on the sequence context of a particular isolate (102). In several studies, specific resistance-associated polymorphisms in HVR1 were shown to affect dependency on SR-BI for viral entry, suggesting modulation of the complex entry pathway as one mechanism for antibody resistance (89, 103). Other direct or indirect effects on antibody neutralization likely stem from interactions with ApoE, which interacts with E2 in the context of lipoviral particles (LVPs), or HCV viral particles contained within low-density lipoproteins (105-107). Both ApoE interactions and assembly of HCV virions as LVPs can also impede bnAb responses by shielding key epitopes during infection (108, 109). The heterogeneity of LVPs and their absence from some methods of in vitro models has made these impacts on antibody neutralization more difficult to quantify (110, 111).

#### 1.4 Antibody-antigen interactions

1.4.1 Antibody structure and classification

Antibodies recognize and specifically bind foreign antigens, making them a key component of the adaptive immune system. This class of protein has increasingly been used as a therapeutic, with dozens of antibody treatments approved in recent years (112). Antibody sequence and structure is a fascinating combination of conservation and tremendous diversity. In humans and many mammals, the Fab region of Immunoglobulin G antibodies contains a heavy chain and a light chain, each with a structurally conserved constant domain and framework regions in the variable domain (113). However, complementarity determining region (CDR) loops in both chains, along with mechanisms such as V(D)J recombination and affinity maturation, allow antibodies to recognize a nearly infinite pool of antigens (113-115). The CDRH3 loop is most important for dictating antibody recognition and specificity, as it is the most diverse loop and makes the most contacts with antigens (116, 117). Camelid nanobodies, which are present in llama and alpaca immune systems, perform a similar function to human antibodies but are structurally distinct (118). Most importantly, camelid nanobodies include a single chain, or VHH, that approximately resembles an antibody heavy chain. Despite fewer chains, nanobodies have shown remarkable potency to diverse antigens, leading to their recent development as therapeutics (119). Other differences between antibodies and nanobodies have been noted, including altered prevalence of some amino acids in CDRs and longer CDR3 sequences for nanobodies than antibody heavy chains (120-122).

#### 1.4.2 Recognition of antigens

Antibody-antigen interactions involve recognition of an epitope, or a set of residues on the antigen contacting the antibody, by a paratope, or the antibody residues used to contact the epitope (114, 123). Paratopes often include residues from CDR loops exclusively, though residues from framework regions may also contact the antigen (114). Antibodies may recognize continuous

epitopes, which are contained in a stretch of antigen residues and tend to be described as linear, or discontinuous epitopes, which are not contained in the same residue stretch but are spatially proximal and tend to be described as conformational (124, 125). This knowledge has stimulated efforts to predict B cell epitopes computationally, helping to guide antibody docking and design to an antigen of interest (126, 127). The identification and characterization of epitopes recognized by bnAbs is critical to recent efforts in rational or structure-based vaccine design, with the largest efforts surrounding HIV (128-131). Though the criteria for defining a bnAb may depend on the antigen targeted, these antibodies are typically found to recognize an epitope that is highly conserved by sequence, structure, or both (132, 133), and can neutralize a diverse set of virus isolates. Key examples include the CD4 binding site and MPER region in HIV Env (134), the stem region of influenza hemagglutinin (135), and several recently described epitopes in the S2 region of SARS-CoV-2 spike (136, 137). Though the discovery of bnAbs against a target antigen is incredibly valuable, a variety of factors can make inducing bnAbs with a vaccine immunogen very difficult, with some factors mentioned in the previous section. Within the context of epitope recognition, even some conserved epitopes can also be cryptic epitopes, which are difficult to access or are displayed in limited antigen conformations (138-141). In addition, bnAb responses can be impeded if neoepitopes, or immunogenic epitopes that are not often present in a native antigen, induce responses from the immune system with little to no capacity to neutralize (142, 143).

#### 1.4.3 Structural characterization and modeling

Structural characterization of free and antigen-bound antibody structures has been crucial for understanding antibody-antigen recognition, but time and resource-intensive experiments alone cannot hope to characterize the vast size and diversity of antibody repertoires (144). The structures

of CDR loops and their diversity have been studied closely, with structurally similar CDRs in each loop identified and reported by PyIgClassify (145, 146). In combination with several numbering schemes based on residue position (147), these tools allow for a comparative analysis of antibody structures and their recognition of antigens. Despite these classifications, CDR loops still show remarkable structural diversity in combination with conformational flexibility (148, 149). Paratope flexibility and its role in antigen recognition suggests that antibodies may utilize one or more proposed mechanisms of recognition, including induced fit and conformational selection (150, 151). These dynamics within antibody-antigen interactions make antibody CDR loops difficult to model, especially CDRH3 (152). Modeling accuracy has steadily improved in recent years, with databases of antibody structures (153, 154), modeling algorithms (155-158), design protocols (159), and modeling assessments (160, 161) playing a role. However, these improvements largely focus on unbound antibody states, which may be less reliable in predicting antibody-antigen interactions due to conformational changes.

#### 1.5 Protein docking

#### 1.5.1 Strategies and challenges

Protein docking is the computational prediction of protein-protein interactions using the unbound structures or models as input (162, 163). The binding of the two proteins is simulated with an algorithm to generate a set of modeled complexes, which are then ranked based on a score calculated by the algorithm. The goal of docking algorithms is to predict the native binding interface in top ranked models. These simulations rarely reach this ideal scenario, though knowledge of binding or non-binding residues in a complex can improve predictions (164). Strategies of developed docking algorithms fall into two major categories: global docking and local docking. Global docking performs an exhaustive search of the binding interface between two

proteins using methods such as Fast Fourier Transform (165-168), geometric hashing (169, 170), and Monte Carlo searches (171). Some algorithms incorporate additional steps to improve predictions, including clustering of models and potentials tailored for unique types of complexes such as antibody-antigen interactions (166, 172). Local docking involves refinement of a starting model interface, often simulating protein backbone flexibility or interface perturbations to improve the free energy of a complex. These algorithms can work in tandem with global docking algorithms, starting with a coarse-grained representation of a model complex then refining at a higher resolution. While most local docking generally introduces flexibility through Monte Carlo searches, molecular dynamics, or normal modes (173-177), SnugDock explicitly adds flexibility to the paratope in models of antibody-antigen interactions (178). Instead of local refinement, other algorithms perform reranking of global docking models based on complex scoring functions (179-183) that can also help to improve predictions. Though careful applications of protein docking can find or recapitulate a native binding interface, there are several major challenges in this field. Protein flexibility between unbound and bound states makes docking more difficult, especially for rigid-body global docking algorithms that do not model conformational changes (163, 184, 185). Flexible protein docking may help resolve this issue despite being more intensive computationally (163), but predictions of rigid-body complexes by rigid-body algorithms may not find a near-native hit, suggesting that the underlying algorithms used to score models may be suboptimal or inadequate for predicting specific types of complexes (186). Recent advances in utilizing coevolutionary information for protein docking cannot apply to antibody-antigen complexes, reducing the options for making predictions of this type of interaction (187, 188).

### 1.5.2 CAPRI and docking benchmarks

With many protein docking algorithms available to use, it can be difficult to know which method may perform the best in practical applications, and whether those success rates are high enough to warrant utilization of these methods. For two decades, the Critical Assessment of Predicted Interactions (CAPRI) process has tested the performance of docking algorithms by comparing docking predictions of solved but not yet released structures with the native interface, then reporting a ranking of participating groups that often use distinct docking algorithms (189, 190). Docking models are compared to the native complex and assessed for accuracy based on the criteria of interface RMSD (I-RMSD), ligand RMSD (L-RMSD), and fraction of native contacts (f<sub>nat</sub>). I-RMSD and L-RMSD are respective measurements in Å of positional fit in the interface and ligand, or the protein that was docked to its binding partner, once each model and the native complex are superposed. f<sub>nat</sub> indicates the fraction of residue-residue contacts within 5 Å of the native interface that are observed in a given docking model. Models within or above specific thresholds of each metric are classified as "Acceptable", "Medium", or "High"; any model that fails to meet these criteria is classified as "Incorrect" (189, 190).

This evaluation of docking model accuracy pioneered by CAPRI has helped to standardize assessments of predictive performance that can compare the success rates of docking algorithms. This type of resource is typically known as a docking benchmark, and can be used both to evaluate current docking algorithms and to validate newly developed docking algorithms. In some docking benchmarks, each complex structure is matched with structures of the unbound state of each component, allowing for impartial assessments of docking success rates by comparing docking predictions with the unbound structures to the native complex. Frequently used benchmarks such as iterations of the Docking Benchmark and DOCKGROUND (186, 191-196) have been developed for this specific purpose. Published in 2015, Docking Benchmark 5 presented an

updated and diverse set of 230 cases (186) that included enzyme-inhibitor interactions, antibodyantigen interactions, and other categories of complexes, with each case assigned a docking difficulty based on the degree of conformational change and fraction of non-native contacts (f<sub>nonnat</sub>), or the fraction of residue-residue contacts between unbound structures superposed onto bound structures that are not present in the native interface (197). Though this update increased the number of cases available for benchmarking, only 28 represented antibody-antigen interactions, a small portion of the overall benchmark.

#### **1.6 Protein affinity prediction**

#### 1.6.1 Affinity measurements and predictors

Prediction of protein-protein affinity, or the strength of binding between interacting proteins, is a related yet distinct computational problem. In this prediction scheme, structures or high-quality models of protein complexes are matched with a corresponding equilibrium dissociation constant ( $K_D$ ) that is measured experimentally (198, 199), often with techniques such as surface plasmon resonance, bio-layer interferometry, or isothermal titration calorimetry. These affinities can be calculated as the change in Gibbs free energy ( $\Delta G$ ) in kcal/mol or kJ/mol if there is a documented temperature for the measurement. When a wild-type affinity and an affinity of the same complex with a mutation are compared, the effect of the mutation on affinity can be calculated as the change in Gibbs free energy ( $\Delta\Delta G$ ). Both  $\Delta G$  and  $\Delta\Delta G$ predictions compare experimentally determined affinity values with the scores of corresponding bound complexes generated by algorithms. These scores can be derived from calculations of certain interface characteristics such as  $\Delta ASA$  (200) or from a weighted set of linear terms that forms a statistical potential or scoring function (201, 202). Most of these scoring functions have similar terms to calculate interface properties, including van der Waals forces, desolvation, and electrostatics that are also used by docking algorithms (179, 181, 203). Other more complex scoring functions have also added potentials for hydrogen bonding and residue or atom-based potentials to improve performance in free energy calculations or docking model rankings, which may help to improve affinity predictions (182, 204). Improving antibody-antigen affinity prediction would aid the design of antibodies as therapeutics, and is one of many design applications under investigation using machine learning algorithms (205), especially  $\Delta\Delta G$ prediction with algorithms such as TopNetTree (206).

#### 1.6.2 Community resources

As with protein docking, protein affinity prediction requires dedicated databases for training and testing affinity predictors that can be applied to any algorithm in development. Experimentally determined  $\Delta\Delta G$  values are matched with complex structures in several databases that are available to the community for  $\Delta\Delta G$  prediction, including SKEMPI (207, 208), SiPMAB (209), and AB-Bind, a resource specifically for  $\Delta\Delta G$  prediction in antibody-antigen interactions (210). These resources have been used to train and evaluate  $\Delta\Delta G$  predictors, demonstrating the utility of curated datasets for researchers (206, 211-215). Although databases such as PDBbind and SAbDab contain protein affinities that can be used for  $\Delta G$  prediction (153, 216), there is not a dedicated database for  $\Delta G$  prediction in antibody-antigen interactions. This lack of focus on resources for antibody-antigen  $\Delta G$  prediction persists even though a recently described affinity predictor utilized antibody-antigen affinities from PDBbind (214). Docking Benchmark 5 and a previous iteration also contain an Affinity Benchmark, where  $\Delta G$  values and corresponding bound and unbound structures can be used to facilitate predictions (186, 217). This benchmark has been used to train  $\Delta G$  prediction models, including PRODIGY (218), a model of interfacial contacts (219), and a minimal model of  $\triangle$ ASA and I-RMSD (220), integrating a metric of conformational
change between bound and unbound states to represent the thermodynamics of these interactions. Though these predictors performed reasonably well on the entire Affinity Benchmark, correlations between scores and  $\Delta G$  values of antibody-antigen complexes were noticeably lower. A host of predictors were also tested for correlations with  $\Delta G$  values in Docking Benchmark 5, showing a range of modest Pearson correlation coefficients with correlations by complex type not reported (186).

## 1.7 Dissertation overview

This thesis describes research that I conducted under the guidance of Dr. Brian Pierce over the period of five years. In this dissertation, I will detail computational and experimental research to address a set of questions that broadly encompass vaccine design and antibody-antigen prediction. Throughout this research, there is a unifying theme of utilizing antibody responses to antigens, leading to advancements in both therapeutically relevant immunogen designs and improved resources for antibody-antigen docking and affinity prediction. The research involving hepatitis C virus (HCV) vaccine design and analysis of antibody responses was started before broad computational analyses of antibody-antigen interactions, and the chapters here approximately match this order. Vaccine design (chapters 2 and 3) and antibody-antigen interaction (chapters 5 and 6) sections were laid out chronologically, while chapter 4 was strategically placed to bridge these two main topics. All projects were intended to facilitate vaccine and/or antibody design in different ways, and we hope that a complete description of these projects will be useful.

In **Chapter 2**, we describe the design and characterization of a novel and secreted HCV E1E2 (sE1E2) glycoprotein construct. Through rational design that implemented strategies used for other viruses, we found that a heterodimer of E1E2 ectodomains with a C-terminal coiled-coil scaffold

could be expressed and purified more easily than E1E2 with native transmembrane domains. The utility of this construct as a vaccine candidate is tested through antigenicity and immunogenicity studies, as well as analytical characterization.

In **Chapter 3**, we designed sE1E2 constructs with various scaffolds of synthetic or noneukaryotic origin or with glycoprotein ectodomains from consensus sequences or alternative isolates. This research expands on the proof-of-concept sE1E2 design to test additional constructs than can facilitate soluble E1E2 assembly while avoiding potential development issues and inducing broadly neutralizing antibody responses. In most cases, these constructs showed promise as vaccine candidates based on successful expression, native-like antigenicity, and analytical characterization.

In **Chapter 4**, we utilized existing datasets to predict E1E2 residues with important implications for immune evasion and viral assembly, including sequence polymorphisms that contribute to changes in antibody neutralization and E1E2 residues crucial for heterodimeric assembly. A variety of computational methods were used to examine experimental datasets, including scripts implemented in R, sequence-based comparisons, hierarchical clustering, and computational mutagenesis in Rosetta. These analyses produced a series of predictions about important E1E2 residues that can be tested experimentally and inform HCV vaccine design.

In **Chapter 5**, we curated and analyzed an expanded set of antibody-antigen structures for a benchmark that can be used for docking and affinity prediction. Through automated and manual searches of the Protein Data Bank (PDB), we identified a diverse set of antibodies that more than doubled the number of antibody-antigen structures and affinities in Docking Benchmark 5. Docking algorithms and scoring functions were evaluated on this expanded set for their success in

docking and affinity predictions, providing examples of how this benchmark update can aid future algorithm development.

In **Chapter 6**, we compiled a large dataset of antibody-antigen affinity values with corresponding PDB structures as a benchmark for antibody-antigen affinity prediction that can facilitate future algorithm development. Scores from existing functions showed modest correlations to affinities in this diverse dataset, showing that there is room to improve the performance of affinity predictors. We demonstrate the utility of this dataset by using its affinities and structures to retrain individual and combined scoring functions, which showed higher correlations following training and cross-validation but limited improvements on an independent test set.

Though general mechanisms of antibody-antigen recognition are well understood, induction of broadly neutralizing antibody responses through vaccine design and predictions of antibodyantigen docking or affinity remain difficult problems. The following thesis describes efforts to advance vaccine design for one virus and to provide resources for benchmarking docking and affinity prediction of all antibody-antigen interactions, presenting a broad and interdisciplinary perspective that should be useful in both areas of development.

# Chapter 2: Design of a native, secreted hepatitis C virus E1E2 heterodimer

## Abstract

Hepatitis C virus (HCV) is a major worldwide health burden, and a preventive vaccine is needed for global control or eradication of this virus. A substantial hurdle to an effective HCV vaccine is the high variability of the virus, leading to immune escape. The E1E2 glycoprotein complex contains conserved epitopes and elicits neutralizing antibody responses, making it a primary target for HCV vaccine development. However, the E1E2 transmembrane domains that are critical for native assembly make it challenging to produce this complex in a homogenous soluble form that is reflective of its state on the viral envelope. To enable rational design of an E1E2 vaccine, as well as structural characterization efforts, we have designed a soluble, secreted form of E1E2 (sE1E2). As with soluble glycoprotein designs for other viruses, it incorporates a scaffold to enforce assembly in the absence of the transmembrane domains, along with a furin cleavage site to permit native-like heterodimerization. This sE1E2 was found to assemble into a form closer to its expected size than full-length E1E2. Preservation of native structural elements was confirmed both by high-affinity binding to a panel of conformationally specific monoclonal antibodies, including two neutralizing antibodies specific to native E1E2, and by binding to its primary receptor, CD81. Finally, sE1E2 was found to elicit robust neutralizing antibodies in vivo. This designed sE1E2 can both provide insights into the determinants of native E1E2 assembly and serve as a platform for production of E1E2 for future structural characterization and vaccine studies, enabling rational optimization of an E1E2-based antigen.

## **2.1 Introduction**

Hepatitis C virus (HCV) is a global disease burden, with an estimated 71 million people infected worldwide (10, 221). Roughly 75% of HCV infections become chronic (4-6), and in severe cases can result in cirrhosis or hepatocellular carcinoma (222). Viral infection can be cured at high rates by direct acting antivirals (DAAs), but multiple public health and financial barriers (15, 223), along with the possibility of reinfection or continued disease progression (19, 223, 224), have resulted in a continued rise in HCV infections. An HCV vaccine remains essential to proactively protect against viral spread, yet vaccine developments against the virus have been unsuccessful to date (111, 225). The challenges posed by HCV sequence diversity (67, 225), glycan shielding (74, 79), immunodominant non-neutralizing epitopes (62, 82, 85, 226), and preparation of a homogeneous E1E2 antigen all contribute to the difficulty in generating protective B cell immune responses. Though multiple studies in chimpanzees and humans have used E1E2 formulations to induce a humoral immune response, their success in generating high titers of broadly neutralizing antibody (bnAb) responses has been limited (227). Optimization of E1E2 to improve its immunogenicity and elicitation of bnAbs through rational design may lead to an effective B cell based vaccine (228).

HCV envelope glycoproteins E1 and E2 form a heterodimer on the surface of the virion (50, 229, 230). Furthermore, E1E2 assembly has been proposed to form a trimer of heterodimers (51) mediated by hydrophobic C-terminal transmembrane domains (TMDs) (50, 231, 232) and interactions between E1 and E2 ectodomains (233-235). These glycoproteins are necessary for viral entry and infection, as E2 attaches to the CD81 and SR-B1 co-receptors as part of a multistep entry process on the surface of hepatocytes (236-239). Neutralizing antibody responses to HCV infection target epitopes in E1, E2, or the E1E2 heterodimer (62, 88, 240-244). Structural

knowledge of bnAb antibody-antigen interactions, which often target E2 epitopes in distinct antigenic domains B, D, or E (55, 62, 86), can inform vaccine design efforts to induce bnAb responses against flexible HCV epitopes (96, 245, 246). E1E2 bnAbs, including AR4A, AR5A (247), and others recently identified (243), are not only among the most broadly neutralizing (240), but also represent E1E2 quaternary epitopes unique to antibody recognition of HCV.

Though much is known about bnAb responses to E1E2 glycoproteins, induction of B cell based immunity with a E1E2-based vaccine immunogen (81, 248, 249) has remained difficult. The inherent hydrophobicity of E1 and E2 transmembrane domains (TMDs) (50, 250) may impede uniform production of an immunogenic E1E2 heterodimer that could be utilized for both vaccine development and E1E2 structural studies. Although partial E1 and E2 structures have been determined (244, 251-254), many other enveloped viruses have structures of a complete and nearnative glycoprotein assembly (41, 46, 255-257), providing a basis for rational vaccine design (34, 258, 259). Viral glycoproteins of Influenza hemagglutinin (260), respiratory syncytial virus (RSV) (41), SARS-CoV-2 (261), and others (262, 263) have been stabilized in soluble form using a Cterminal attached foldon trimerization domain to facilitate assembly. HIV gp120-gp41 proteins have been designed as soluble SOSIP trimers by introducing a furin cleavage site, along with a key proline mutation and an added disulfide between gp120 and gp41, to mediate native-like assembly when cleaved by the enzyme (46, 264). Previously described E1E2 glycoprotein designs include covalently-linked E1 and E2 ectodomains (58, 265), E1E2 with transmembrane domains intact and an IgG Fc tag for purification (266), as well as E1 and E2 ectodomains with a cleavage site (58), which presented challenges for purification either due to intracellular expression or to high heterogeneity. Two recently described scaffolded E1E2 designs, while promising, have not been shown to engage monoclonal antibodies (mAbs) that recognize the native E1E2 assembly,

though they were engaged by E1-specific and E2-specific mAbs, as well as co-receptors that recognize E2 (59). Therefore, these presentations of E1E2 glycoproteins may not represent a native and immunogenic heterodimeric assembly, and thus their potential as vaccine candidates remains unclear.

Here, we describe the design of a secreted E1E2 glycoprotein (sE1E2) that mimics both the antigenicity *in vitro*, and the immunogenicity *in vivo*, of the native heterodimer through the scaffolding of E1E2 ectodomains. In testing our designs, we found that both replacing E1E2 TMDs with a leucine zipper scaffold and inserting a furin cleavage site between E1 and E2 enabled secretion and native-like sE1E2 assembly. We assessed the size, heterogeneity, antigenicity, and immunogenicity of this construct (identified as sE1E2.LZ) in comparison with full-length membrane-bound E1E2 (mbE1E2). sE1E2.LZ binds a broad panel of bnAbs to E2 and E1E2, as well as co-receptor CD81, providing evidence of assembly into a native-like heterodimer. An immunogenicity study indicated that sera of mice injected with sE1E2.LZ neutralize HCV pseudoparticles (HCVpp) at levels comparable to sera from mice immunized with mbE1E2. This sE1E2 design is a novel form of the native E1E2 heterodimer that both improves upon current designs and represents a platform for structural characterization and engineering of additional HCV vaccine candidates.

#### 2.2 Methods

## 2.2.1 Protein expression

For expression of recombinant soluble HCV E2 (sE2), the sequence from isolate H77C (GenBank accession number AF011751; residues 384–661) was cloned into the pSecTag2 vector

(Invitrogen)\*, and expressed in mammalian (Expi293F) cells as described previously (92). The mbE1E2 and sE1E2 DNA coding sequences were synthesized with a modified tPA signal peptide (267) at the N-terminus. All E1E2 sequences were cloned into the vector pcDNA3.1+ at the cloning sites of KpnI/NotI (GenScript). Furin sequence DNA was cloned into the vector pcDNA3.1 and was a gift from Dr. Yuxing Li (University of Maryland IBBR). All sE1E2 constructs and mbE1E2 were transfected with ExpiFectamine 293 into Expi293F cells for expression (Invitrogen). Cleavable polyprotein constructs were co-transfected with the furin construct at a 2:1 ratio. A clone for mammalian expression of CD81 large extracellular loop (CD81-LEL), containing N-terminal tPA signal sequence and C-terminal twin Strep tag, was provided by Dr. Joe Grove (University College London). CD81-LEL was expressed through transient transfection in Expi293F cells (Thermo Fisher Scientific).

#### 2.2.2 Antibodies

Monoclonal antibodies used in ELISA assays and binding studies were produced as previously described (268-270), except for AR4A and AR5A, which were kindly provided by Dr. Mansun Law (Scripps Research Institute).

# 2.2.3 Protein purification and size exclusion chromatography

sE2 glycoprotein was purified from cell supernatant as described previously (92). Culture supernatant of sE1E2.LZ and E1E2 ectodomains fused with a Gly-Ser linker (sE1E2GS3) was purified by immobilized metal affinity chromatography (IMAC) with separate HiTrap chelating HP Ni<sup>2+</sup>-NTA columns (Cytiva). Expressed mbE1E2 was extracted from cell membranes using

<sup>&</sup>lt;sup>\*</sup> Certain commercial equipment, instruments, or materials are identified in this paper in order to specify the experimental procedure adequately. Such identification is not intended to imply recommendation or endorsement by the National Institute of Standards and Technology, nor is it intended to imply that the materials or equipment identified are necessarily the best available for the purpose.

1% NP-9 and purified via sequential Fractogel EMD TMAE (Millipore), Fractogel EMD SO3<sup>-</sup> (Millipore), immunoaffinity with HC84.26.WH.5DL antibody (271), and Galanthus Nivalis Lectin (GNL, Vector Laboratories) affinity chromatography. Sample concentration prior to size exclusion chromatography (SEC) was conducted with 15 ml Amicon Ultra 3 kDa centrifugal filters (Millipore Sigma). sE1E2.LZ, sE1E2GS3, and mbE1E2 were fractionated using a Superdex 200 Increase 10/300 column (Cytiva). sE1E2.LZ and sE1E2GS3 were equilibrated with 1x Phosphatebuffered saline (PBS; 10 mM sodium phosphate + 150 mM NaCl) pH 7, while mbE1E2 was equilibrated in Tris-buffered saline (TBS; 25 mM Tris-HCl + 150 mM NaCl) pH 7.5 + 0.5% n-Octyl-β-D-Glucopyranoside (Anatrace). Size exclusion fractions of 500 µl were collected on AKTA FPLC (Cytiva). Molecular weight standards from the high molecular weight (HMW) calibration kit (Cytiva) were compared to purified sE1E2.LZ, sE1E2GS3, and mbE1E2. CD81-LEL was purified using a 5 ml prepacked Streptactin XT column (IBA Lifesciences), following dialysis of expression supernatant in buffer W (100 mM Tris-HCl pH8, 150 mM NaCl, 1 mM EDTA) overnight in 4°C. CD81-LEL eluate was fractionated with a Superdex 200 column (Cytiva) on an AKTA FPLC (Cytiva) equilibrated with Tris buffered saline (TBS) (20 mM Tris + 150 mM NaCl) pH8.

# 2.2.4 Computational design of coiled coil assemblies

Coiled coil assemblies were designed using the HBNet protocol in Rosetta (272). This protocol accepts coiled coil architectures as input, performing modular hydrogen bond network generation and subsequent design to optimize packing and stability, resulting in models of designed assemblies (272). Two architectures were selected for parametric generation of coiled coil bundles for Rosetta input: supercoiled and no supercoil (parallel coil). The supercoil parameters were selected based on the GCN4 leucine zipper structure (PDB code 1ZIK) (273).

Backbones were generated with these two architectures using a Python program described previously and available in Rosetta (274), with each helix 30 amino acids in length. By varying helix phases in 18° increments for the inner and outer helices in the Python program, 400 backbones were generated per global architecture (supercoil and parallel coil). As the design subunits in this system were heterodimeric rather than monomeric, we added a minor modification to the published HBNet Rosetta Script protocol (272) to account for the chain break between heterodimeric subunits ("<Span begin="30" end="31" bb="0" chi="1"/>). HBNet design was performed with each of the 800 input backbone structures, resulting in approximately 335 output designs. Some backbone structures resulted in no output designs due to lack of candidate hydrogen bond networks identified by HBNet, while others resulted in multiple designs based on multiple candidate hydrogen bond networks and packing designs. Design models were assessed for lack of buried unsatisfied polar groups, which has been found to be associated with successful designed assemblies (272), followed by manual inspection, to select the top five candidates for experimental characterization.

# 2.2.5 Peptide synthesis and characterization

Peptides for coiled coil designs CC1+CC2, HEX-1, HEX-2, HEX-3, and HEX-4 were synthesized (Genscript) and resuspended in Milli-Q water. Pairs of peptides corresponding to each coiled coil design were mixed at a 1:1 ratio and incubated overnight in 4°C. 10X PBS was then added at 1/10<sup>th</sup> the volume of the mixture, which was centrifuged to separate any precipitate. Each peptide mixture was purified using a Superdex 75 Increase 10/300 column (Cytiva). Elution peak positions of gel filtration standards (Bio-Rad #1511901) were used to calculate molecular weights of designs CC1+CC2 and HEX-1-4 based on their observed peak positions.

## 2.2.6 SEC-MALS

For size exclusion chromatography coupled to multiple angle light scattering (SEC-MALS), a UHPLC system (Vanquish Flex, Thermo Fisher) was coupled to MALS (DAWN HELEOS-II, Wyatt) and Refractive Index (Optilab T-rEX, Wyatt) detectors. Separations were performed using a WTC-050N5 column (Wyatt) equilibrated in PBS for sE1E2.LZ or in TBS +  $0.5\% \beta$ -OG for mbE1E2, with a flow rate of 0.3 mL/min and sample injection volumes of 25 µL. Molar mass analysis was performed using the software ASTRA 7.1.3 (Wyatt) using refractive index as a concentration source.

# 2.2.7 SDS-PAGE and western blot

SDS-PAGE and western blot experiments were conducted with 12-well stain-free gels (Bio-Rad), with total protein detected using a stain-free imager (Bio-Rad). For SDS-PAGE, Precision Plus Unstained Protein Standards (Bio-Rad) were used as a molecular weight marker. E2 was detected in western blot with HCV1 (275) as the primary antibody. E1 was detected in western blot with H-111 as the primary antibody (276). In reducing conditions, each sample was incubated with loading dye (4x Laemmli buffer + 10% β-mercaptoethanol) (Bio-Rad) and heated to 95°C, except for mbE1E2, which was heated to 37°C. In non-reducing conditions, each sample was incubated with Laemmli buffer and heated to 37°C. For western blots, stain-free gels were transferred to a turbo mini 0.2 µm nitrocellulose membrane (Bio-Rad) using the trans-blot turbo transfer system (Bio-Rad). Supersignal Molecular Weight Protein Ladder (Thermo Fisher Scientific) was used as a marker for western blots. 10X concentration of supernatant for E1 western blots was conducted in 0.5 mL Amicon Ultra 3 kDa centrifugal filters (Millipore Sigma). Cell lysates of sE1E2.LZ and mbE1E2 were collected by centrifugation of 1 ml transfected cell suspension and extraction from cell membranes with 1% NP-9. For native western blots, 15-well

NativePAGE Novex 4-16% Bis-Tris protein gels (Thermo Fisher Scientific) were transferred to a turbo mini 0.2 µm PVDF membrane (Bio-Rad) using the same transfer system. NativeMark unstained protein standard (Invitrogen) was used as a molecular weight marker for native gels. To deglycosylate sE1E2.LZ, mbE1E2, and sE2 in non-denaturing conditions, 3 µg of each protein was mixed with 2 µl PNGase F enzyme (New England Biolabs), then incubated at 37°C for 24 hours before western blot preparation. Proteins were detected with goat anti-human IgG HRP conjugate (Invitrogen) and clarity western ECL substrate (Bio-Rad). All gels were imaged using the ChemiDoc system (Bio-Rad).

## 2.2.8 Analytical ultracentrifugation (AUC)

Sedimentation velocity (SV) experiments were performed at 20°C using a ProteomeLab Beckman XL-A with absorbance optical system and a 4-hole An60-Ti rotor (Beckman Coulter). For sE1E2.LZ, the sample and reference sectors of the dual-sector charcoal-filled epon centerpieces were loaded with 390 µL protein in PBS, pH 7.4 with or without 0.5%  $\beta$ -OG, and 400 µL buffer. For mbE1E2, the sample and reference sectors of the dual-sector charcoal-filled epon centerpieces were loaded with 390 µL protein in TBS + 0.5%  $\beta$ -OG, and 400 µL buffer. The cells were centrifuged at 40 krpm and the absorbance data were collected at 280 nm in a continuous mode with a step size of 0.003 cm and a single reading per step to obtain linear signals of <1.25 absorbance units. Sedimentation coefficients were calculated from SV profiles using the program SEDFIT (277). The continuous *c*(*s*) distributions were calculated assuming a direct sedimentation boundary model with maximum entry regularization at a confidence level of 1 standard deviation. The density and viscosity of buffers at 20°C and 4°C were calculated using SEDNTERP (278). The *c*(*s*) distribution profiles were prepared with the program GUSSI (C.A. Brautigam, Univ. of Texas Southwestern Medical Center).

## 2.2.9 Enzyme-linked immunosorbent assay (ELISA)

HCV human monoclonal antibody (HMAb) binding to mbE1E2, sE1E2.LZ, sE1E2GS3, and sE2 was evaluated and quantitated by ELISA. 96-well microplates (MaxiSorp, Thermo Fisher, Waltham, MA) were coated with 5 µg/mL Galanthus Nivalis Lectin (Vector Laboratories, Burlingame, CA) overnight, and purified mbE1E2, sE1E2.LZ, sE1E2GS3 and sE2 was then added to the plates at 2 ug/ml. After the plates were washed with PBS and 0.05% Tween 20, and blocked by Pierce<sup>TM</sup> Protein-Free (PBS) Blocking Buffer (Thermo Fisher, Waltham, MA), the HMAbs were tested in duplicate at 3-fold serial dilution starting at 100 ug/ml. The binding was detected by 1:5000 dilutions of HRP-conjugated anti-human IgG secondary antibody (Invitrogen, Carlsbad, CA) with TMB substrate (Bio-Rad Laboratories, Hercules, CA). The absorbance was read at 450 nm using a SpectraMax MS microplate reader (Molecular Devices, San Jose, CA). For ELISA measurements of immunized murine sera, endpoint titers were calculated by curve fitting in GraphPad Prism software, with endpoint OD defined as four times the mean absorbance value of Day 0 sera.

# 2.2.10 Determination of antibody affinity by quantitative ELISA

ELISA assays were performed as described (270) to compare antibody affinity to sE1E2.LZ, mbE1E2, and sE2. Briefly, plates were developed by coating wells with 500 ng of Galanthus nivalis agglutinin (GNA) and blocking with 2.5% non-fat dry milk and 2.5% normal goat serum. Purified sE1E2.LZ, mbE1E2, and sE2 at 5  $\mu$ g/ml were captured by GNA onto the plate and later bound by a range of 0.01-200  $\mu$ g/ml of antibody. Bound antibodies were detected by incubation with alkaline phosphatase-conjugated goat anti-human IgG (Promega), followed by incubation with p-nitrophenyl phosphate for color development. Absorbance was measured at 405 nm and 570 nm. The assay was carried out in triplicate in three independent assays for each HMAb.

The data were analyzed by nonlinear regression to measure antibody dissociation constants (K<sub>D</sub>) and binding potential (optical density at 405 nm) using Graphpad Prism software, and standard deviation values were calculated using the three independent affinity measurements.

## 2.2.11 Surface plasmon resonance

Surface plasmon resonance (SPR) analysis was performed using a Biacore<sup>™</sup> T200 system (Cytiva) and HBS-EP+ buffer was used as sample and running buffer. The analysis temperature and sample compartment were set to 25°C. mbE1E2, sE2, and sE1E2.LZ were immobilized on Series S CM5 chips using the Amine Coupling Kit per the manufacturer's instructions. Antigen capture levels were adjusted to yield approximately 2000 RU for the kinetic experiments. Purified CD81-LEL was injected over reference and active flow cells, applying a single cycle kinetics procedure using twelve concentrations. Data were fitted to a 1:1 binding model using Biacore<sup>™</sup> T200 Evaluation Software 2.0. As one concentration series was used to calculate binding parameters, no standard errors were calculated for those values.

#### 2.2.12 Animal immunization

CD-1 mice were purchased from Charles River Laboratories. Prior to immunization, sE2 and E1E2 antigens were formulated with polyphosphazene PCPP-R adjuvant (279). Poly[di(carboxylatophenoxy)phosphazene], PCPP (50  $\mu$ g, molecular weight 800,000 Da) (280) was formulated with resiquimod, R848 (25  $\mu$ g) in PBS (pH 7.4) to prepare PCPP-R as described previously (279). The resulting formulation was mixed with E1E2 antigen (70  $\mu$ g for prime or 15  $\mu$ g for boost immunization). The absence of aggregation in adjuvanted formulations was confirmed by dynamic light scattering (DLS): single peak, z-average hydrodynamic diameter – 60 nm. The formation of antigen–PCPP-R complex was confirmed by asymmetric flow field flow fractionation (AF4) as described previously (281). On scheduled vaccination days, groups of 6 female mice, age 7-9 weeks, were injected via the intraperitoneal (IP) route with a 50  $\mu$ g E1E2 prime (day 0) and boosted with 10  $\mu$ g E1E2 on days 7, 14, 28, and 42. Blood samples were collected prior to each injection with a terminal bleed on day 56. The collected samples were processed for serum by centrifugation and stored at -80°C until analysis was performed.

## 2.2.13 HCV pseudoparticle generation

HCV pseudoparticles (HCVpp) were generated as described previously (77), by cotransfection of HEK293T cells with the murine leukemia virus (MLV) Gag-Pol packaging vector, luciferase reporter plasmid, and plasmid expressing HCV E1E2 using Lipofectamine 3000 (Thermo Fisher Scientific). Envelope-free control (empty plasmid) was used as negative control in all experiments. Supernatants containing HCVpp were harvested at 48 h and 72 h posttransfection and filtered through 0.45 µm pore-sized membranes. For measurements of serum binding to HCVpp in ELISA, concentrated HCVpp were obtained by ultracentrifugation of 33 ml of filtered supernatants through a 7 ml 20% sucrose cushion using an SW 28 Beckman Coulter rotor at 25,000 rpm for 2.5 hours at 4°C, following a previously reported protocol (86).

# 2.2.14 HCVpp neutralization assays

Huh7 cells were maintained in the Dulbecco's modified Eagle medium (DMEM) supplemented with 10 % FBS.  $1.5 \times 10^4$  Huh7 cells per well, plated in white 96-well tissue culture plates (Corning), and incubated overnight at 37°C. The following day, HCVpp was mixed with serially diluted murine serum samples at 37°C. After one-hour incubation, the HCVpp-serum mixture was added to the Huh7 cells (kindly provided by Jonathan K. Ball, University of Nottingham, UK) in 96-well plates and incubated at 37°C for 5 h. After removing the inoculum, the cells were further incubated for 72 h with DMEM containing 10% fetal bovine serum (Thermo

Fisher, Waltham, MA) and the luciferase activities were measured using Bright-Glo<sup>™</sup> luciferase assay system as indicated by the manufacturer (Promega, Madison, WI).

## 2.2.15 Statistical comparisons

P-values between group endpoint titers and group half-maximal inhibitory dose (ID50) values were calculated in Graphpad Prism software, using non-parametric Kruskal-Wallis analysis of variance with Dunn's multiple comparisons test.

## 2.3 Results

# 2.3.1 Design of sE1E2 constructs

We designed and screened a set of sE1E2 constructs to determine which type of scaffold might be suitable for development of a novel secreted heterodimer (**Figure 2.1A**). Scaffolded sE1E2 constructs were synthesized as cleavable polyproteins and contain a six-arginine furin cleavage site, which was incorporated to facilitate E1E2 assembly with a method also used for HIV SOSIP constructs (46). Each cleavable polyprotein replaces E1 and E2 TMDs with a self-assembling heterodimeric, homotrimeric, or hetero-hexameric scaffold designed to enforce E1E2 ectodomain assembly in the absence of a membrane anchor. In addition, all constructs replace the N-terminal wild-type signal peptide sequence with a modified version of the signal sequence from tissue plasminogen activator (tPA) (267) and include a C-terminal 6xHis tag for purification.

sE1E2.LZ used the human Fos-Jun leucine zipper, a coiled coil obligate heterodimer with a known structure (PDB code 1FOS; **Figure 2.1B**) (282), as a scaffold. The heterodimeric Fos-Jun leucine zipper has been used as a scaffold for expression of T cell receptors (283), making it a possible candidate for maintaining heterodimeric E1E2 in secreted form. sE1E2.FD replaced the E1 TMD with a foldon domain (**Figure 2.1C**; PDB code: 4NCU) (284), a self-trimerizing protein

that has been previously used to stabilize soluble assemblies of viral glycoprotein trimers (41, 285). This construct was designed to test whether enforcing E1 trimerization (51) would be sufficient to enable E1E2 ectodomain assembly. sE1E2.CC used a scaffold that was designed to self-assemble into a hetero-hexameric peptide complex, which would reflect the previously described model of the E1E2 TMD architecture (51) in a soluble form. The corresponding scaffold, CC1+CC2 (Figure 2.1D), was designed de novo using the HBNet protocol of Rosetta protein modeling software (272). Though we were unable to confirm the structure of CC1+CC2 with experimental structural determination, it was included as a candidate scaffold given its putative hexameric assembly (Figure 2.2). To examine the importance of including scaffolds in the absence of TMDs, a separate construct with a furin cleavage site but no scaffold was generated (sE1E2.R6). Two sE1E2 constructs with a covalent linker between ectodomains were also included. In sE1E2GS3, E1 and E2 ectodomains are linked by a 15 amino acid glycine-serine sequence, resembling a previously described sE1E2 construct (58). The construct sE1E2RevGS3 reverses the order of E1 and E2 ectodomains, testing whether altering the order of ectodomains in the context of a covalent fusion may improve E1E2 assembly, which could be affected by the currently unknown proximity of the N- and C-termini of the ectodomains in native E1E2.



**Figure 2.1 Design of sE1E2 constructs.** (A) Schematic of mbE1E2, covalent linker sE1E2, and cleavable polyprotein constructs. Regions shown include tPA signal sequence (green boxes), E1 ectodomain (yellow boxes), E2 ectodomain (red boxes), wild-type TMDs (gray boxes), Gly-Ser linker (orange boxes), and various scaffolds replacing TMDs. E1E2 residue ranges for each region are noted according to H77 numbering. C-terminal His tags and furin cleavage sites are shown in boxes and labeled. The expected molecular weight of each construct is indicated, and molecular weights of expected oligomers for sE1E2.FD and sE1E2.CC are in parentheses. For molecular weight estimations, each N-glycan is approximated to be 2 kDa at each NxS/NxT sequon, a value within the molecular weight range of typical N-linked glycans (286). (B) X-ray structure of human Fos-Jun heterodimer (PDB code: 1FOS); only the coiled coil region that was used for the sE1E2.LZ scaffold is shown. Fos and Jun chains were colored to match the diagram of sE1E2.LZ. (C) X-ray structure of foldon domain (PDB code: 4NCU). All chains colored light blue to match the diagram for sE1E2.FD. (D) Model of CC1+CC2 heterohexameric peptide assembly. CC1 and CC2 chains colored to match the diagram for sE1E2.CC. All structures were visualized in PyMOL (Schrödinger).



**Figure 2.2 Characterization of the peptide complex CC1+CC2.** Shown are the chromatographic traces for the CC1+CC2 complex (blue line) and other tested designs (labeled HEX1-4) following elution from a Superdex 75 size exclusion chromatography column (Cytiva). The CC1+CC2 complex elutes at a volume consistent with hexameric assembly. Indicated on the chromatograph is the estimated molecular weight for CC1+CC2, calculated based on the retention volumes of molecular size standards (Bio-Rad).

## 2.3.2 sE1E2.LZ forms an intact E1E2 complex

Each sE1E2 construct was expressed in mammalian cells, with cleavable polyproteins coexpressed with furin. To test for successful secretion of sE1E2, we probed for the presence of E1 and E2 ectodomains in the supernatant, using the E1 HMAb H-111 (276) and the E2 HMAb HCV1 (275) in western blots. These antibodies bind to linear epitopes at or near the N-terminus of the E1 or E2 ectodomain, respectively. sE1E2.LZ was the only cleavable polyprotein design to show clear detection of both E1 and E2 in the supernatant (**Figure 2.3**), though sE1E2.FD exhibited some secretion of E2. The sE1E2.R6 construct without a scaffold showed no secretion of sE1E2, consistent with previous results that E1 and E2 ectodomains alone do not form a stable complex (59). Expression of E1-Jun and E2-Fos constructs in trans without a furin cleavage site found secretion of E1-Jun, but minimal secretion of E2-Fos (Figure 2.4). Collectively, these results determine that the combination of a furin cleavage site and leucine zipper scaffold enables secretion of the E1E2 complex. sE1E2GS3 and sE1E2RevGS3 showed high levels of E1 and E2 in supernatant, corroborating previous findings with a covalently linked sE1E2 design that is similar to sE1E2GS3 and was likewise detected in the supernatant (58). In addition, we examined if protein was expressed but not secreted by probing for the presence of E1 and E2 in lysed cells (Figure 2.5). sE1E2GS3 and sE1E2RevGS3 that was retained in cells migrated at smaller molecular weights than the corresponding secreted proteins, while sE1E2.FD and sE1E2.LZ exhibited multiple bands in E2 detection; both results may be indicative of incomplete processing or degradation of protein that was not secreted. Though some sE1E2.LZ was detected intracellularly, approximately 90% of expressed sE1E2.LZ was secreted to the supernatant, as determined by a quantitative analysis comparing supernatant and cell lysate western blots probed with the anti-E2 HMAb HCV1 (Figure 2.6). Based on these results, we selected sE1E2.LZ, as a cleaved scaffolded sE1E2 candidate, and sE1E2GS3, as a covalently linked sE1E2 candidate, for further characterization.



## Figure 2.3 E1 and E2 western

E2 western blot. H-111 antibod

samples were loaded under reducing conditions. Supermatants were concentrated 10 times prior to E1 western blot. Molecular weights, in kDa, of the western blot markers closest to observed bands are indicated on the left. Expected band positions of E1, E2, and E1E2 are indicated with black triangles on right and labeled.



Figure 2.4 Western blots of supernatant from E1-Jun/E2-Fos co-expression. sE1E2.LZ components E1-Jun and E2-Fos, both with tPA signal sequence, were co-expressed in trans and then probed with HCV1 (anti-E2) at 5 µg/ml or H-111 (anti-E1) at 10 µg/ml under reducing conditions. E1 and E2

detection was compared to expression levels of the full sE1E2.LZ construct. Supernatants were concentrated 10X prior to E1 western blot. Molecular weights, in kDa, of the marker closest to observed bands are labeled. In both western blots, E1-Jun/E2-Fos and sE1E2.LZ were loaded in non-adjacent wells but were placed together to aid viewing.



Figure 2.5 E1 and E2 western blots of sE1E2 cell lysate. HCV1 antibody at 5  $\mu$ g/ml was used for the E2 western blot. H-111 antibody at 10  $\mu$ g/ml was used for the E1 western blot. All sE1E2 lysate samples were loaded under reducing conditions. Supernatants were concentrated 10X prior to E1 western blot. Molecular weights, in kDa, closest to observed bands are labeled. Expected band sizes of E1, E2, and E1E2 are indicated with black triangles and labeled accordingly. E1 detection of sE1E2.R6 and sE1E2.CC were loaded in non-adjacent wells but are grouped together in this figure to aid comparisons.



HCV1 (E2)

Figure 2.6 Quantitative western blots comparing sE1E2.LZ supernatant and cell lysate. One  $\mu$ l of each sample was used for E2 probing with anti-E2 antibody HCV1 at a concentration of 5  $\mu$ g/ml in separate western blots. A standard curve with defined values of sE2 purified protein (50, 100, or 200 ng) was included in each western blot, with a representative standard curve shown. Band intensities of

supernatant and cell lysate samples were compared with the standard curve via ImageQuant software (Cytiva) to estimate protein amount and the proportion of expressed sE1E2.LZ that was secreted in supernatant. Molecular weights, in kDa, closest to observed bands are labeled. E2 detection of sE1E2.LZ supernatant and cell lysate was aligned by molecular weight range of markers from separate western blots. sE1E2.LZ cell lysate and 50 ng of sE2 were loaded in non-adjacent wells but are grouped together in this figure to aid viewing.

# 2.3.3 Purification of sE1E2.LZ

We purified both sE1E2.LZ and sE1E2GS3 using IMAC, and then examined the molecular weight and heterogeneity of each construct with SEC (Figure 2.7A; Figure 2.8). Expression and purification of all three constructs produced sufficiently pure protein for characterization, with sE1E2.LZ providing the highest yield at 480 µg per 100 ml of transfected cells (Figure 2.9). Both constructs eluted in SEC across a broad molecular weight range, with the peak for each estimated at approximately 400 kDa. The resultant SEC peaks were directly compared with the peak SEC fractions of purified mbE1E2 (Figure 2.7D). Though sE1E2.LZ, along with sE1E2GS3, exhibited a broad peak in SEC, it eluted at a volume consistent with a molecular weight that is both smaller than mbE1E2, which eluted as a peak in void volume (approximately 700 kDa), and closer to the expected size of the heterodimeric assembly (94 kDa; Figure 2.1). To further investigate the size distribution and heterogeneity of purified constructs, we examined fractions eluted from SEC under non-reducing conditions, using western blot for sE1E2.LZ (Figure 2.7B-C), mbE1E2 (Figure 2.7E-F), and sE1E2GS3 (Figure 2.10D,F), and SDS-PAGE for sE1E2GS3 (Figure **2.10B**) and sE1E2.LZ (Figure 2.11B). Both sE1E2.LZ and sE1E2GS3 SEC fractions showed two predominant species migrating in the range between 150 and 250 kDa when probed for E1 and E2 under non-reducing conditions, which is smaller than expected based on the SEC chromatographs but confirms the heterogeneity of each protein. mbE1E2 SEC fractions probed by western blot under non-reducing conditions showed several species, including prominent bands corresponding to free E1 and E2 along with higher molecular weight aggregates. In addition, the anti-E1 nonreducing western blot shows discrete bands corresponding to self-associating E1 dimers and trimers as observed previously (51), suggesting that, while the purified protein is a heterogenous mixture, the mixture contains a significant population of E1E2 assembled natively. In contrast, under reducing conditions the E1 and E2 components migrated at the expected molecular weight for both sE1E2.LZ (Figure 2.11) and mbE1E2 (Figure 2.12) fractions, and at a molecular weight corresponding to covalently linked E1E2 in sE1E2GS3 (Figure 2.10) fractions. The spread of the bands in SDS-PAGE and western blot is likely due to in part heterogeneity in glycoforms, as observed previously (77, 287). To examine the contribution of glycosylation to observed size distributions, we subjected the purified proteins to PNGase F cleavage to remove the glycans. An examination of the deglycosylated proteins on a non-reducing western blot showed more species (Figure 2.13), indicating that the heterogeneity in solution we observed for all constructs is dominated by another factor, possibly disulfide crosslinking or exchange. Although these results suggest that sE1E2.LZ is closer to expected size of a heterodimer than mbE1E2, the ranges of observed sizes led us to utilize more sensitive methods of characterization to examine molecular size and heterogeneity.



**Figure 2.7 Size exclusion chromatography of sE1E2.LZ, sE1E2GS3, and mbE1E2.** Chromatographic traces for (A) sE1E2.LZ and (D) mbE1E2 shown in blue lines plotted with molecular weight standards shown in grey lines after elution from a Superdex 200 SEC column (Cytiva). Molecular weight estimates for the center of each peak are labeled based on comparisons with elution of HMW standards (Cytiva), with molecular masses of 670, 440, 158, 73, and 44 kDa. The range for elution fractions F1-F10 used for analysis is shown as a red line. Western blots of sE1E2.LZ for E2 (B), sE1E2.LZ for E1 (C), mbE1E2 for E2 (E), and mbE1E2 for E1 (F) under non-reducing conditions. HCV1 antibody at 5  $\mu$ g/ml was used to probe for E2, while H-111 antibody at 10  $\mu$ g/ml was used to probe for E1. Molecular weights, in kDa, of the western blot markers closest to observed bands are indicated on the left of each panel. All fractions had 250 ng loaded for improved visualization of size. For E1 western blots, all fractions were concentrated 10X prior to loading. Putative E1 monomer, dimer, and trimer populations shown in panel (F) are highlighted with red initials.



**Figure 2.8 Size exclusion chromatograph of sE1E2GS3.** Chromatograph of sE1E2GS3 shown as a blue line plotted with molecular weight standards shown as a grey line after elution from a Superdex 200 SEC column (Cytiva). The elution fractions F1-F9 used for subsequent analysis is shown as a red line. A molecular weight estimate for the center of the peak is labeled based on comparisons with elution of HMW standards (Cytiva), with values of 670, 440, 158, 73, and 44 kDa.



Figure 2.9 Yield and purity of mbE1E2, sE1E2.LZ, and sE1E2GS3 in SDS-PAGE. Yield of each protein in  $\mu$ g per 100 ml of transfected cells is shown underneath the corresponding sample. 3.75  $\mu$ g of protein was loaded for each purified protein. Expected band sizes of E1, E2, and E1E2 are indicated with black triangles and labeled accordingly. Molecular weight markers, in kDa, closest to observed bands are also indicated.



**Figure 2.10 sE1E2GS3 fractions from SEC analyzed by SDS-PAGE and western blot.** Fractions F1-F9 show a gradient of molecular weights following elution. SDS-PAGE results for sE1E2GS3 fractions under (A) reducing and (B) non-reducing conditions with molecular weights, in kDa, of the marker labeled. Western blot of sE1E2GS3 fractions under (C) reducing and (D) non-reducing conditions probed with HCV1 (anti-E2) antibody at 5  $\mu$ g/ml. In panel (C), the fraction with the highest concentration had 250 ng loaded, with other fractions scaled accordingly. In panel (D), 250 ng of sE1E2GS3 fractions were loaded to improve visualization of size. Western blot of sE1E2GS3 fractions under (E) reducing and (F) non-reducing conditions probed with H-111 (anti-E1) antibody at 10  $\mu$ g/ml. All fractions were concentrated 10X prior to E1 western blots. In panel (E), the fraction with the highest concentration had 250 ng loaded, with other fractions scaled accordingly. In panel (F), 250 ng of sE1E2.LZ fractions were concentrated 10X prior to E1 western blots. In panel (C), the fraction with the highest concentration had 250 ng loaded, with other fractions scaled accordingly. In panel (F), 250 ng of sE1E2.LZ fractions were loaded to improve visualization of size. In panels (C-F), molecular weight, in kDa, of the western blot marker closest to observed bands is indicated.



Figure 2.11 sE1E2.LZ fractions from SEC analyzed by SDS-PAGE and western blot. Elution fractions F1-F10 show both E1 and E2 in SDS-PAGE under reducing conditions (A) and a molecular weight gradient in SDS-PAGE under non-reducing conditions (B). Molecular weights, in kDa, for the SDS-PAGE protein ladder (Bio-Rad) are indicated. Western blots of sE1E2.LZ fractions under reducing conditions when probed with HCV1 (anti-E2) antibody at 5  $\mu$ g/ml (C) or H-111 (anti-E1) antibody at 10  $\mu$ g/ml (D). Molecular weight, in kDa, of the western blot marker closest to observed bands is indicated. In both western blots, the fraction with the highest concentration had 250 ng loaded, with other fractions scaled accordingly. For the E1 western blot, all fractions were concentrated 10X prior to loading.



Figure 2.12 mbE1E2 elution fractions from SEC analyzed by western blot. Elution fractions were probed under reducing conditions with HCV1 (anti-E2) antibody at 5  $\mu$ g/ml (A) or H-111 (anti-E1) antibody at 10  $\mu$ g/ml (B). Molecular weight, in kDa, of the western blot marker closest to observed bands is indicated.



HCV1 (E2)

**Figure 2.13 Deglycosylation of mbE1E2, sE1E2.LZ, and sE2.** Purified proteins were analyzed by western blot with HCV1 (anti-E2) antibody at 5  $\mu$ g/ml under reducing (left) and non-reducing (right) conditions with molecular weights, in kDa, of the marker labeled. 800 ng of each deglycosylated sample, along with a paired sample with intact glycans, were loaded in each lane of the reducing western blot. Some degradation of deglycosylated sE2 is apparent as the band intensity is markedly reduced. To aid detection of the full range of species present in the non-reducing western, additional sample was added as needed. It is apparent that deglycosylation either allows separation or induces formation of additional species in the non-reducing western blot. Figure provided by Liudmila Kulakova and Eric Toth.

## 2.3.4 Analytical characterization of heterogeneity in solution

sE1E2.LZ and mbE1E2 purified constructs were also characterized using AUC, which can separate a mixture of protein populations more precisely than SEC (288). A comparison of AUC results offers further support that sE1E2.LZ is less heterogeneous than mbE1E2. AUC for sE1E2.LZ showed two prominent peaks between sedimentation coefficient (S) values 4.9 and 7.5, which are approximately consistent with a monomer and dimer of the sE1E2.LZ heterodimer, respectively, and resemble what we observed in the non-reducing western blot. To control for potential effects of 0.5% n-Octyl-β-D-Glucopyranoside (β-OG), a detergent required for mbE1E2 purification, we performed a parallel AUC experiment with sE1E2.LZ in the presence of 0.5%  $\beta$ -OG (Figure 2.14A). The size distribution in that experiment closely matched that of the sample without  $\beta$ -OG, indicating that the detergent itself does not contribute to heterogeneity. mbE1E2 showed three large peaks between S values 4 and 9.1, suggesting that mbE1E2 exhibits more heterogeneity than sE1E2.LZ (Figure 2.14B). Furthermore, the peak with the highest intensity for mbE1E2 closely resembles the S value found for free E2. sE1E2.LZ by contrast shows no peak at that S value. Though sE1E2.LZ is not a uniform single species, it is a less complex mixture of E1E2 assemblies than mbE1E2.

SEC-MALS was used as another analytical technique to examine the heterogeneity and size of sE1E2.LZ. Since the presence of  $\beta$ -OG detergent had little to no effect on sE1E2.LZ in AUC, we expected that an absence of  $\beta$ -OG would not affect analytical characterization of sE1E2.LZ in SEC-MALS. When compared with standards and analyzed by light scattering, sE1E2.LZ exhibited a single peak in SEC-MALS with an estimated molecular weight at peak center of 173 kDa, corresponding approximately to a dimer of the sE1E2.LZ heterodimer (**Figure 2.14C**). This estimated size is generally consistent with the observed AUC peak around 7.5 S,

though the breadth of the peak in SEC-MALS still suggests that sE1E2.LZ displays some heterogeneity in size, corresponding to 1-2 sE1E2.LZ heterodimers, in accordance with the two major peaks from AUC measurements. In SEC-MALS, mbE1E2 was characterized as a single, very broad peak with an estimated molecular weight of 1.1 MDa at peak center (**Figure 2.14D**). The broad range of this peak identified mbE1E2 as a mixture containing a broad range of species, with approximately 5 to over 20 E1E2 heterodimers. Additionally, sE1E2.LZ was directly compared to mbE1E2 in a native western blot, showing differences in overall size (**Figure 2.15**). In assessments by multiple analytical techniques, sE1E2.LZ forms a moderately heterogeneous mixture that is nonetheless smaller and closer to expected size than mbE1E2, representing a potentially improved immunogen for HCV vaccine development and a candidate for structural characterization. In addition, sE1E2.LZ does not require detergents for solubility, allowing for simpler formulations than mbE1E2.



Figure 2.14 Analytical characterization of sE1E2.LZ and mbE1E2 size and heterogeneity. AUC profiles of (A) purified sE1E2.LZ with or without detergent  $\beta$ -OG and (B) purified mbE1E2. The distribution of Lamm equation solutions c(s) for the two proteins (blue or black lines) is shown. Calculated sedimentation coefficients for the peaks are labeled. Observed species for sE1E2.LZ approximately correspond to a heterodimer at 4.9 S, a dimer of heterodimers at 7.7 S, and higher-order aggregates at 10.3 S. Observed species for mbE1E2 approximately correspond to free E2 at 4.0 S, a dimer of heterodimers at 6.6 S, a trimer of heterodimers at 9.1 S, and a tetramer of heterodimers and higherorder aggregates at >10 S. Figures provided by Kinlin L. Chao. (C) sE1E2.LZ and (D) mbE1E2 characterization with SEC-MALS. The chromatographs of each protein are shown as blue lines. For reference, chromatographs of molecular weight standards are shown as grey lines in panels (C) and (D), corresponding to molecular masses of 670, 158, 44, 17, and 1.35 kDa. The MALS scattering sizes between the peak half-maxima are shown as red lines, with the estimated molecular weight at the center of each peak labeled, and size distribution of each range in parentheses. Based on calculated molecular weights of each heterodimer and SEC-MALS molecular size ranges, these peaks predominantly contain oligomers of (C) 1-2 sE1E2.LZ heterodimers and (D) 5-27 mbE1E2 heterodimers. Figures provided by Thomas E. Cleveland IV.



Figure 2.15 mbE1E2 and sE1E2.LZ size and heterogeneity in native gel. Purified proteins were compared through blue native gel electrophoresis followed by western blot probed with either HCV1 (anti-E2) at 5  $\mu$ g/ml or H-111 (anti-E1) antibody at 10  $\mu$ g/ml. E2 detection of mbE1E2 and sE1E2.LZ originated from different gels, which were then aligned to make the range of molecular weights equivalent. Molecular weights, in kDa, closest to observed bands are labeled. E1 detection of mbE1E2 and sE1E2.LZ was also conducted on separate gels, then aligned by molecular weight range. Figure provided, in part, by Andrezza Chagas.

# 2.3.5 sE1E2.LZ exhibits native-like E1E2 antigenicity and robust immunogenicity

We next examined the native-like properties of sE1E2.LZ by measuring the binding affinities to a panel of bnAbs in comparison with sE2 and mbE1E2. Unlike the antibodies used in western blot, most bnAbs used for this analysis recognize conformational epitopes on E2 (55, 270, 289), and E1E2 (247). We first performed ELISA at one antibody concentration to compare mbE1E2 and sE1E2.LZ antibody reactivity, along with purified sE1E2GS3 and sE2. This screening was used to assess lack of reactivity by any of the constructs to conformationally

sensitive antibodies, versus quantitative comparisons of affinities, which were undertaken later. The antibodies utilized were a representative panel of bnAbs to antigenic domain B, D, and E epitopes in E2 and the E1E2 bnAbs AR4A and AR5A (**Figure 2.16**). At the tested antibody concentration (0.185  $\mu$ g/ml), mbE1E2 and sE1E2.LZ exhibited similar binding levels for all antibodies. Importantly, sE1E2.LZ maintained reactivity to E1E2 bnAbs, providing evidence that this sE1E2 construct contains a soluble, native-like form of the E1E2 heterodimer. In contrast, sE1E2GS3 and sE2 showed little to no reactivity to AR4A and AR5A; this was not unexpected for sE2, which lacks key residues comprising the E1E2 bnAb epitopes (87). Based on the AR4A and AR5A binding results, the lack of E1-E2 cleavage or scaffold in sE1E2GS3 appears to lead to a severe disruption of native-like assembly, thus we focused on sE1E2.LZ for subsequent characterization.

To confirm more precisely our initial measurements of bnAb reactivity, we tested the affinity of sE1E2.LZ to a larger panel of HCV antibodies (**Table 2.1**) and CD81 (**Figure 2.17**). K<sub>D</sub>s were measured by dose-dependent ELISA to antibodies that recognize discrete epitopes of E2 (62) and E1E2 bnAbs. For comparison, we performed the same analysis for purified mbE1E2 and sE2. sE1E2.LZ and mbE1E2 showed similar affinities to almost all tested HCV HMAbs, within a 2-3 fold difference. One notable exception was an 8-fold lower affinity of AR4A for sE1E2.LZ relative to mbE1E2. Although sE1E2.LZ maintained affinity to AR5A, a decrease in affinity to AR4A may stem from subtle differences in heterodimer assembly or dynamics when compared to mbE1E2, which may be difficult to elucidate without detailed structural characterization of the epitope. Regardless, AR4A binds sE1E2.LZ with nanomolar affinity (16 nM), indicating that the overall structure of the AR4A epitope and the E1E2 interface in that region are intact. In addition to measurements of binding to conformationally sensitive E2 and E1E2 HMAbs, we also tested

binding to the CD81 receptor, which recognizes a region on the E2 ectodomain overlapping with epitopes for many bnAbs (87). sE1E2.LZ showed robust binding to CD81-LEL in SPR (10.8 nM; **Figure 2.17**), establishing that this sE1E2 construct displays receptor binding critical for native HCV infection. While measured CD81-LEL  $K_D$  values show comparable or higher affinity than corresponding glycoprotein affinities for antibodies in **Table 2.1**, due to the different experimental measurement methods, these results provide a comparison between antigens rather than a comparison between absolute glycoprotein affinities of receptor versus antibodies.

After confirming the native-like antigenicity of sE1E2.LZ, we tested the native-like properties of sE1E2.LZ in vivo to determine whether it will elicit antibodies that effectively recognize HCV and inhibit infection. Mice were immunized with either mbE1E2, sE1E2.LZ, or sE2 and tested for the presence of antibodies that target E1E2 and neutralize the virus (Figure **2.18**). sE1E2.LZ elicited anti-mbE1E2 antibody responses that mimicked responses from mbE1E2-immunized mice, while serum binding of mbE1E2 from sE2-immunized mice was lower, particularly when compared with the mbE1E2-immunized group (p < 0.01) (Figure 2.18A). Binding of immunized sera to H77-pseudotyped HCV pseudoparticles (HCVpp) was also tested for all groups (Figure 2.18B), and while mean serum titer was highest for the sE1E2.LZ group, there were no significant differences found between immunized group titers based on nonparametric (Kruskal-Wallis) assessment. Serum neutralization of H77C HCVpp was tested for all groups to assess for elicitation of neutralizing antibodies that target the homologous virus (Figure **2.18C**). Testing of pre-immune sera for background neutralization showed no detectable HCVpp neutralization (Figure 2.19). sE1E2.LZ-immunized sera showed robust neutralization of HCVpp, with neutralization titers (ID50s) that showed no significant difference from mbE1E2-immunized and sE2-immunized groups. This initial test of sE1E2.LZ immunogenicity shows that this secreted
E1E2 construct induces an antibody response comparable to mbE1E2 and sE2 that can recognize homologous E1E2 on the surface of HCVpp and neutralize the virus.



**Antigenicity screening** 



A máile a dar	Domain <sup>1</sup>	K <sub>D</sub> (nM) <sup>2</sup>			Standard Error (nM)		
Anudody		mbE1E2	sE1E2.LZ	sE2	mbE1E2	sE1E2.LZ	sE2
CBH-4D	А	28	26	1	3.2	3.4	0.2
CBH-4G	А	7.8	18	0.5	2.3	3.1	0.3
HC-1 AM <sup>3</sup>	В	1.5	2.9	3.6	0.06	0.5	0.4
HC-11	В	1.8	3.2	11	0.09	0.4	0.6
CBH-7	С	1	1.7	0.3	0.1	0.1	0.04
HC84.24	D	0.5	1.3	0.7	0.07	0.1	0.1
HC84.26	D	1.2	2.6	0.4	0.03	0.4	0.1
HC33.1	Е	3.8	0.9	1.9	0.3	0.09	0.2
HCV1	Е	9.8	3.5	6.2	0.3	0.2	0.3
AR4A	E1E2	2.3	16	-	0.2	1.5	-
AR5A	E1E2	1.5	1.7	-	0.2	0.2	-

Table 2.1 Binding affinity of mbE1E2, sE1E2.LZ, and sE2 to a panel of HMAbs.

"-" denotes no binding detected

<sup>1</sup>Antigenic domain on E2 targeted by antibody (A-E), as previously described (63). "E1E2" denotes antibodies that target the E1E2 heterodimer.

<sup>2</sup>Measured by dose-dependent ELISA, with standard error values shown for each affinity measurement. Figure provided by Young Chang Kim and Zhen-Yong Keck.

<sup>3</sup>Affinity-matured HC-1 antibody, as previously described (290).



Figure 2.17 Measurement of binding to the CD81 receptor by SPR. CD81 binding kinetic curves to (A) mbE1E2, (B) sE1E2.LZ, and (C) sE2 are shown. Kinetic ( $k_{on}$ ,  $k_{off}$ ) and steady-state ( $K_D$ ; calculated as  $k_{off}/k_{off}$ ) binding parameters were calculated based on 1:1 model and are shown in each panel. Figure provided by Eric Toth.



**Figure 2.18 Immunogenicity assessment of sE2, mbE1E2, and sE1E2.LZ.** Six mice per group were immunized with sE2, mbE1E2 or sE1E2.LZ, and sera were tested for binding to (A) mbE1E2 and (B) H77C-pseudotyped HCVpp in ELISA. One mouse in the sE2-immunized group died prior to final bleed, thus responses for five mice are shown for that group. Endpoint titers were calculated using Graphpad Prism, and geometric mean titers are shown for each group as black lines. (C) Neutralization of H77C HCVpp by immunized murine sera. ID50 values were calculated in Graphpad Prism for individual mice, and average ID50 titers for each immunized group are shown as black lines. The minimal serum dilution used for ID50 measurement (1:64) is shown as a horizontal dashed line, for reference. P-values between group endpoint titer or ID50 values were calculated using Kruskal-Wallis analysis of variance with Dunn's multiple comparison test (*ns, not significant*: p > 0.05; \*:  $p \le 0.05$ ; \*\*:  $p \le 0.01$ ). Figure provided by Ruixue Wang.



**Figure 2.19 Calculated curves for H77C HCVpp neutralization by immunized (Day 56) murine sera.** Data are shown for individual mice, and names (key on right) correspond to immunized groups (G1: mbE1E2, G2: sE1E2.LZ, G3: sE2), with six mice per group. Pooled pre-immune sera from each group were tested as controls. One mouse from G3 died prior to Day 56, thus had no serum available for testing. Serum dilutions (x-axis) are two-fold serial dilutions, starting at 1:64 (Reciprocal Serum Dilution = 64). Figure provided by Ruixue Wang.

## **2.4 Discussion**

The development and characterization of a native-like E1E2 antigen containing a leucine zipper scaffold offers a proof of principle platform for designing E1E2 vaccine antigens within a soluble and secreted backbone. Exploration of this scaffold approach for producing E1E2 from other HCV genotypes is warranted, as sE1E2.LZ was only designed using the H77C sequence. E2 ectodomains from other strains have been characterized structurally (244, 254, 291), and the E1E2 sequences of those strains could be targets for sE1E2.LZ backbone expression and characterization. However, strain-specific sequence changes may affect sE1E2.LZ secretion, as differences in E1 and E2 stalk regions could modulate assembly and export from cellular components (292, 293). In addition, further studies of sE1E2 secretion may shed light on cellular factors that facilitate efficient sE1E2 assembly, which could then be used to either improve production levels or examine mechanisms of viral assembly and secretion.

There are several avenues for subsequent design and optimization of the sE1E2.LZ platform. As a potential vaccine immunogen, a human leucine zipper in sE1E2.LZ poses potential problems related to immunizing humans with human protein sequences (294, 295). As the Fos-Jun leucine zipper is structurally defined at high resolution, this can be used as a template for identification of heterodimeric leucine zipper structures from non-human proteins or de novo designs of synthetic leucine zipper scaffolds. Furthermore, although the CC1+CC2 sE1E2 design (sE1E2.CC) did not yield appreciable secretion, it is possible that alternative hetero-hexameric scaffolds, possibly generated using the Fos-Jun leucine zipper structure as a subunit, could promote stable E1E2 assembly. Finally, recent studies have shown that cage-like protein nanoparticles can provide scaffolds for viral glycoproteins such as RSV F (296, 297) and Influenza hemagglutinin (255). A nanoparticle recapitulating the Fos-Jun leucine zipper structure as attachment points could be identified or designed to present sE1E2 in a similar nanoparticle format. Binding to E1E2specific antibodies, such as AR4A and AR5A, is particularly important for validation of scaffolded E1E2 antigens. Since sE1E2.LZ exhibited slightly impaired binding to AR4A, new designed or synthetic scaffolds may provide an opportunity to improve upon a human leucine zipper scaffold by matching or exceeding wild-type binding to E1E2-specific antibodies. High-resolution structural characterization of sE1E2.LZ or subsequent designs, enabled by effective secretion and purification of this native-like assembly, can permit an improved view of the determinants of E1E2 assembly and support structure-based modifications to enhance assembly and stability.

Although sE1E2.LZ was observed as closer to expected size of a heterodimer than mbE1E2, our extensive analytical characterization indicated a likely mix of heterodimers and higher-order oligomers. This degree of sample heterogeneity has been found during purification of previous soluble construct designs, both with a covalent linker (58) and a designed

heterodimeric scaffold (59). Although glycoform heterogeneity is apparent in both constructs, our results suggest that it is not the primary source of observed oligomerization. Instead, these constructs demonstrate that removing the heterodimer from its natural membrane-attached environment does not preclude formation of large assemblies. The E2 ectodomain likely plays a large role in aggregation via additional hydrophobic interactions or disulfide crosslinking, as its ectodomain contains conserved and surface-exposed tyrosines, tryptophans, and cysteines (62). These residues are critical for co-receptor interactions (241, 298), proper ectodomain folding, and assembly (87, 292), but could readily mediate E1E2 aggregation without TMDs present. Self-association of E2 ectodomains has also been noted previously (299), offering additional support for the propensity of soluble E2 to exhibit crosslinking. Future studies will examine specific determinants of sE1E2.LZ heterogeneity and methods to mitigate it, building on recent efforts to obtain homogenous secreted glycoprotein (300). Experimental structural characterization, as noted above, would help to delineate the stoichiometry and oligomerization modes of sE1E2 designs.

In summary, replacing the native TMDs of E1 and E2 with a leucine zipper scaffold provides proof of concept that this approach can be used to develop a native-like, antigenically and immunogenically intact E1E2 complex without requiring a membrane or detergent environment. The design and validation of additional scaffolds that adopt dimeric, trimeric, or hetero-hexameric quaternary structures could elucidate key determinants of E1E2 complex assembly, another area of research that has been hindered by membrane association of E1E2. In addition, this scaffold approach could serve as a platform to study how the substantial genetic diversity of HCV translates to structural diversity and envelope glycoprotein dynamics, and how structural and dynamic changes, including "open" and "closed" envelope glycoprotein states, may promote immune evasion, as noted by recent work (89). Finally, in addition to their use in

structural characterization, designed soluble E1E2 complexes with functional TMD replacements that retain all essential structural properties can serve as an integral component of rational vaccine design.

# Chapter 3: Design of soluble hepatitis C virus E1E2 assemblies with alternative scaffolds or ectodomains

# Abstract

Hepatitis C virus (HCV) represents a worldwide disease burden, and a vaccine is urgently needed to control and eradicate this virus. HCV vaccine design has been difficult in part because of substantial genetic diversity and of challenges in isolating and producing the E1E2 glycoprotein heterodimer, the primary target of broadly neutralizing antibodies (bnAbs). Recently, a scaffolded E1E2 heterodimeric assembly was shown to be secreted, antigenically native, and immunogenic *in vivo*. However, the use of a human sequence to scaffold E1E2 makes the use of this design as a vaccine candidate problematic. Previous studies with other viruses have utilized protein scaffolds that are synthetic or non-eukaryotic in origin or glycoprotein sequences that are a consensus of a reference panel, but the impact of these strategies in improving on the initial sE1E2 design has not been examined. In this study, we design and characterize sE1E2 constructs that incorporate synthetic or non-eukaryotic heterodimers as scaffolds or consensus sequences or alternative isolates as replacement E1E2 ectodomains. Alternative scaffolds were selected either by identifying structures in the Protein Data Bank (PDB) with similarity to sE1E2 scaffold with a human leucine zipper or by noting its incorporation in other vaccine candidates. Consensus sequences for genotypes 1-7 and genotype 1 were designed using a reference panel of full-length E1E2 sequences from multiple sources, while alternative isolates were selected from an antigenically diverse panel of E1E2 sequences functional as HCV pseudoparticles (HCVpp). Several new sE1E2 constructs secreted and displayed native-like antigenicity with a panel of HCV

antibodies, providing promising options for HCV vaccine design and additional strategies for effectively scaffolding E1E2 glycoproteins.

## **3.1 Introduction**

In Chapter 2, we detailed the global impact of Hepatitis C virus (HCV) infection despite the success of direct acting antivirals (DAAs), showing the ongoing and urgent need for a vaccine. The primary target for a B cell based HCV vaccine is the complex of E1E2 glycoproteins, which are essential for viral entry and infection through interactions with co-receptors (236, 237) and are thought to form a trimer of heterodimers at the surface of the virion (51). As discussed in Chapter 2, various challenges to vaccine development posed by HCV have led to numerous efforts to delineate the effective immune response to HCV, including characterization of bnAb responses to this target, especially to E2 regions classified as antigenic domains (247, 269, 270). Though this characterization has led to valuable efforts in structure-based vaccine design (92, 301, 302), the lack of an accessible E1E2 heterodimer structure has hindered these efforts. An experimentally determined structure of the E1E2 heterodimer was recently reported in a preprint, offering an exciting glimpse of this elusive complex that can lead to rapid advances in structure-based vaccine design against this target (61). However, this structure does not show the putative "trimer of dimers" hetero-hexameric assembly that may best capture E1E2 presentation in native virus (51, 303). As described in Chapter 2, this slow progress of E1E2 structural characterization is in stark contrast with the structural characterization of multiple viruses using various methods of modification to stabilize full glycoprotein assemblies (38, 41, 46, 255, 304) and to proceed with these immunogens as vaccine candidates (305). Without structural determination of a complete and native E1E2 glycoprotein assembly, the goal of a B cell based vaccine that can prevent HCV infection may never be reached.

Chapter 2 presented a proof of concept for a scaffolded native-like E1E2 heterodimer that can help HCV vaccine design and structural characterization efforts, but important questions about the feasibility of this approach remain. Importantly, this scaffolded E1E2 design uses a human leucine zipper sequence, and the possibility of generating autoreactive antibodies to a critical human transcription factor (294, 295) following vaccine administration makes the inclusion of this sequence potentially problematic for HCV vaccine development. Previous research has harnessed knowledge of leucine zipper domains to produce coiled coil peptide heterodimers of synthetic origin for a variety of potential functions (306-308), roughly resembling the structure of the designed heterodimeric scaffold shown to assemble secreted E1E2 (59). Scaffolded immunogens for vaccine design have also been combined with methods of multimeric display, boosting immunogenicity of this design while maintaining secretion and proper assembly. Self-assembling proteins such as ferritin (255, 309, 310) and other non-eukaryotic or synthetic assemblies (297, 302, 311-314) have enabled multimeric nanoparticle assembly of viral glycoproteins. Other designs have utilized the SpyCatcher/SpyTag system to facilitate multimeric assembly through an irreversible isopeptide linkage that enforces attachment of a nanoparticle and a target protein (315, 316). This system has been used for mosaic nanoparticles of viral glycoproteins that show promise as vaccine candidates (317-320). Another strategy to advance HCV vaccine design would be alteration of the E1E2 ectodomains used in the scaffolded constructs instead of the scaffold itself, as the reference isolate H77 used for a previous design is relatively sensitive to antibody neutralization (321, 322) and often induces limited neutralization of other genotypes as an immunogen (56), though the cause of these limitations is unclear. Consensus sequences of viral glycoproteins of HIV (323, 324) and Influenza (325, 326) were found to be highly stable while retaining conserved epitopes targeted by neutralizing antibodies. In addition, a recent panel of HCV sequences isolated from patients has characterized neutralization phenotypes more representative of antibody responses to the panel overall (321), which may identify HCV isolates that better reflect typical neutralization sensitivity. By exploring multiple options for native-like presentation and ectodomain representation of E1E2 glycoproteins, the secreted and scaffolded design that eliminates risk of inducing autoreactive antibodies, facilitates multimeric display, and better induces bnAb responses may become clearer.

Here, we characterize a variety of scaffolded E1E2 constructs that either replace the previous Fos-Jun leucine zipper scaffold with a synthetic or non-eukaryotic heterodimeric assembly or replace the previous E1E2 ectodomain sequences with a consensus sequence or different isolate. Alternate scaffolds were selected through structural similarity with the Fos-Jun leucine zipper scaffold, stemming from a broad search of the PDB for heterodimeric or heterohexameric assemblies that matched certain criteria. E1E2 consensus sequences were generated for genotype 1 only and genotypes 1-7 through a multiple sequence alignment of E1E2 reference sequences, and alternative HCV isolates were tested based on previous evaluations of their representative neutralization profiles. All constructs were tested for secretion of E1E2 complex into supernatant following transfection and expression, with a range of secretion levels detected. Promising sE1E2 candidates were tested for binding to a diverse panel of HCV antibodies in ELISA that included anti-E1, anti-E2, and anti-E1E2 antibodies. Several new designs showed antibody binding levels that were comparable to E1E2 scaffolded with a leucine zipper (sE1E2.LZ) and membrane-bound E1E2 (mbE1E2), including for anti-E1E2 antibodies AR4A and AR5A, suggesting that these constructs produce properly folded E1E2 heterodimer. One new construct also showed comparable binding affinities to HCV bnAbs as measured in quantitative ELISA. By using the Fos-Jun leucine zipper scaffold as a reference, we have identified additional scaffolds and ectodomains that enforce secreted and native-like E1E2 assembly while revealing that a coiled coil scaffold or a genotype 1a sequence is not determinative. This characterization of additional vaccine candidates defines new strategies of E1E2 scaffolding and sequence selection, directly advancing HCV vaccine development.

# 3.2 Methods

## 3.2.1 E1E2 consensus and alternative isolate sequences

Approximately 3600 full-length E1E2 genotype 1-7 aa sequences from NCBI (www.ncbi.nlm.nih.gov) LANL HCV (327, 328), and euHCVdb (329) were downloaded and aligned using MAFFT (330). To determine an E1E2 genotype 1-7 consensus sequence using this reference database, the database was first filtered by an in-house Perl script so that the remaining sequences (N = 43) were no more than 80% identical to each other. This filtered set of sequences was then used to generate the consensus sequence cons.80. To determine an E1E2 genotype 1 consensus sequence, the portion of the reference alignment with approximately 2700 genotype 1 sequences was filtered by an in-house Perl script so that the remaining sequences were no more than 92.5% identical to each other. This filtered set of sequences was then used to generate the consensus sequence cons.1.92.5. Phylogenetic trees comparing cons.80 and cons1.92.5 to E1E2 sequences used to design consensus sequences were visualized in R using the APE software package (331). Distances between sequences for each tree were exported from ClustalX2 (332) in PHYLIP format. To aid visualization of cons1.92.5 placement, the phylogenetic tree only includes a reduced but representative set of genotype 1 sequences used to design cons1.92.5. This reduced set was defined by filtering with an in-house Perl script so that the remaining sequences (N = 44) were no more than 88% identical to each other. The alignment of H77, cons.80, and cons1.92.5 sequences in key epitopes were visualized and captured in SeaView version 5 (333). Genotype 1a

HCV isolates with neutralization profiles distinct from reference sequence H77 (GenBank accession number AF011751) (334) were selected from a genetically diverse panel of HCVpp (321, 335). Isolate 1.11.6 (GenBank accession number ALV85487) was chosen as a Tier 1 representative more sensitive to antibody neutralization than H77, while isolate 1a38 (GenBank accession number AHV90280) was chosen because of its high HCVpp infectivity within the context of the entire panel.

#### 3.2.2 Selection of alternative scaffolds

In order to find additional scaffolds for a soluble and secreted E1E2 candidate, the PDB (336) was searched for structures based on the following criteria: heterodimeric or heterohexameric, synthetic or non-eukaryotic in origin, parallel, N-termini in chains proximal to each other, and short sequences to minimize the possibility of immunogenicity to the scaffold. Searching for these structures was conducted in part by running an automated feature on PDB pages to find similar proteins by 3D structure using individual chains from several coiled coil scaffolds as input. These structures include the Fos-Jun human leucine zipper heterodimer structure (PDB code 1FOS), a portion of which is used in sE1E2.LZ, a synthetic coiled coil heterodimer SYNZIP1/SYNZIP2 (PDB code 3HE5), and a synthetic coiled coil hetero-hexamer (PDB code 3R48). These structures with synthetic sequences had been found in a previous search for self-assembling peptides that involved manual inspection of candidates in the PDB and corresponding literature. Both structures were included to aid searches of additional structures even though we also planned to test these structures as alternative scaffolds. Any structures identified through these structural similarity searches that also matched the above criteria were visually inspected in PyMOL (Schrödinger). Structure similarity searches with 1FOS as input could not be focused to the portions of chains used in sE1E2.LZ (chain E, residues 161-200 and chain F, residues 285-324), but only structural similarity to these portions following alignment in PyMOL was considered for candidate structures identified during these searches.

In addition, candidate structures in the PDB were identified using a text search for "heterodimer", followed by the filtering of results to structures of non-eukaryotic origin; these structures were also manually inspected in PyMOL. No specific distance cutoff was used to define proximity of N-termini, but in selected structures, the maximum distance at N-termini between Cα atoms was approximately 10 Å. The structure of the SpyTag/SpyCatcher complex (PDB code 4MLS) was identified in an explicit search for the complex following its inclusion in vaccine designs. Each selection was included as a C-terminal scaffold in alternative secreted E1E2 (sE1E2) designs. 3HE5 and 3R48 structures were superposed on the component of the 1FOS structure used in sE1E2.LZ, ensuring that the N-terminal fusion points of the scaffolds for sE1E2.SZ and sE1E2.HH were structurally similar to sE1E2.LZ. The residues from 3CFI and 4MLS used as N-terminal fusion points were selected though structural analysis in PyMOL to estimate which N-terminal residues would most likely be on the same horizontal plane. All fusions between ectodomain and scaffold are separated by a small "PGG" linker.

# 3.2.3 Protein expression

The mbE1E2 and sE1E2 DNA coding sequences were synthesized with a modified tPA signal peptide (267) at the N-terminus. All E1E2 sequences were cloned into the vector pcDNA3.1+ at the cloning sites of KpnI/NotI (GenScript). Human furin sequence DNA in the vector pcDNA3.1 was a gift from Dr. Yuxing Li (University of Maryland IBBR). mbE1E2 was transfected with ExpiFectamine 293 into Expi293F cells for expression (Invitrogen). sE1E2 constructs were co-transfected with furin at a 2:1 ratio with ExpiFectamine 293 into Expi293F cells for expression (Invitrogen), using standard protocols outlined by the manufacturer. Harvested

supernatants were immediately filtered with a bottle top vacuum 0.22 µm filter and supplemented with protease inhibitor cocktail (Thermo Fisher Scientific).

### *3.2.4 Antibodies*

To conduct western blot detection, HCV1 and H-111 antibodies were purified in-house with a protein A column. Antibodies AR4A, AR5A, AR3A, HC33.1, and IGH526 were used for screening of sE1E2 design antigenicity by ELISA. These antibodies were also purified in-house except for AR4A and AR5A, which were kindly provided by Mansun Law.

## 3.2.5 Protein purification and size exclusion chromatography

sE1E2 designs were purified using a 5 ml HiTrap Chelating HP IMAC column (Cytiva) with a stepwise imidazole gradient from 5 mM to 1 M. Fractions eluted by imidazole were loaded onto SDS-PAGE to confirm protein purification, then fractionated with a Superdex 200 increase 10/300 GL column (Cytiva) equilibrated with 1x Phosphate buffered saline (PBS) pH7 on an AKTA FPLC (Amersham). Expressed mbE1E2 was extracted from cell membranes using 1% NP-9 and purified via sequential Fractogel EMD TMAE (Millipore), Fractogel EMD SO<sub>3</sub><sup>-</sup> (Millipore), HC84.26 immunoaffinity (271), and Galanthus Nivalis Lectin (GNL, Vector Laboratories) affinity chromatography. mbE1E2 was equilibrated in Tris-buffered saline (TBS; 25 mM Tris-HCl + 150 mM NaCl) pH 7.5 + 0.5% n-Octyl- $\beta$ -D-Glucopyranoside (Anatrace) and fractionated with a Superdex 200 column (Cytiva) on an AKTA FPLC (Amersham). All samples were concentrated prior to size exclusion chromatography with 15 ml Amicon Ultra 3 kDa centrifugal filters (Millipore Sigma). All fractions from AKTA FPLC were collected in 500 µl increments. Molecular weight standards from the high molecular weight (HMW) calibration kit (Cytiva) were compared to purified sE1E2 designs.

# 3.2.6 SDS-PAGE and western blot

SDS-PAGE and western blot experiments were conducted with 12-well stain-free gels (Bio-Rad), with total protein detected using a stain-free imager (Bio-Rad). For SDS-PAGE, Precision Plus Unstained Protein Standards (Bio-Rad) were used as a molecular weight marker. E2 was detected in western blot with HCV1 (275) at 5 µg/ml as the primary antibody. E1 was detected in western blot with H-111 at 10 µg/ml as the primary antibody (276). Each sample was incubated with loading dye (4x Laemmli buffer + 10% β-mercaptoethanol) (Bio-Rad) and heated to 95°C, except for mbE1E2, which was heated to 37°C. For western blots, stain-free gels were transferred to a turbo mini 0.2 µm nitrocellulose membrane (Bio-Rad) using the trans-blot turbo transfer system (Bio-Rad). Supersignal Molecular Weight Protein Ladder (Thermo Fisher Scientific) was used as a marker for western blots. 10X concentration of supernatant for E1 western blots was conducted in 0.5 mL Amicon Ultra 3 kDa centrifugal filters (Millipore Sigma). Proteins were detected with 1:10,000 dilution of goat anti-human IgG HRP conjugate (Invitrogen) and clarity western ECL substrate (Bio-Rad). All gel pictures were captured using the ChemiDoc imaging system (Bio-Rad).

## 3.2.7 Enzyme-linked immunosorbent assay (ELISA)

HCV human monoclonal antibody (HMAb) binding to mbE1E2 and sE1E2 designs was evaluated and quantitated by ELISA. 96-well microplates (MaxiSorp, Thermo Fisher) were coated with 5 µg/mL GNL (Vector Laboratories) overnight. Purified mbE1E2 was added to the plates at 1 ug/ml. Expi293 cell culture supernatant of sE1E2 constructs was added to the ELISA plate. After the plates were washed with PBS and 0.05% Tween 20, and blocked by Pierce<sup>TM</sup> Protein-Free (PBS) Blocking Buffer (Thermo Fisher Scientific), the HMAbs were tested in duplicate at 0.5 µg/ml. HMAb binding to supernatant of alternative sE1E2 designs was tested in ELISA using the same procedure, except that HMAbs were tested in duplicate at 1 µg/ml. The binding was detected by 1:5000 dilutions of HRP-conjugated anti-human IgG secondary antibody (Invitrogen) with TMB substrate (Bio-Rad). The absorbance was read at 450 nm using a SpectraMax MS microplate reader (Molecular Devices). Prior to temperature-dependent ELISA, samples of mbE1E2, sE1E2.LZ, sE1E2.SZ, and sE1E2.HH at elevated temperatures 37°C or 56°C were heated in a water bath for one hour.

## 3.2.8 Determination of antibody affinity by quantitative ELISA

ELISA assays were performed as described (270) to compare antibody affinity to sE1E2.LZ, mbE1E2, and sE1E2.SZ. Briefly, plates were developed by coating wells with 500 ng of Galanthus Nivalis Agglutinin (GNA) and blocking with 2.5% non-fat dry milk and 2.5% normal goat serum in Tris-buffered saline (TBS; 20 mM Tris-HCl, pH 7.5, 150 mM NaCl) + 0.1% Tween 20. Purified sE1E2.LZ, mbE1E2, and sE1E2.SZ at 5 µg/ml were captured by GNA onto the plate and later bound by a range of 0.01-200 µg/ml of antibody. Bound antibodies were detected by incubation with alkaline phosphatase-conjugated goat anti-human IgG (Promega), followed by incubation with p-nitrophenyl phosphate for color development. Absorbance was measured at 405 nm and 570 nm. The assay was carried out in triplicate in three independent assays for each HMAb. The data were analyzed by nonlinear regression to measure antibody dissociation constants (K<sub>D</sub>) and binding potential (optical density at 405 nm) using Graphpad Prism software, and standard deviation values were calculated using the three independent affinity measurements.

# **3.3 Results**

## 3.3.1 Design of sE1E2 constructs with synthetic scaffolds

We sought to design additional sE1E2 constructs that built on the proof of concept established by secretion and native-like antigenicity of sE1E2.LZ (337). Specifically, we aimed to avoid potential issues of incorporating a human sequence into a vaccine immunogen, as the Fos-Jun sequence included in sE1E2.LZ could induce an autoimmune response to this common protein. As shown in Figure 3.1, two sE1E2 constructs were designed with a synthetic scaffold that replaced the Fos-Jun human leucine zipper sequence in sE1E2.LZ. Both designs were built from the framework of sE1E2.LZ, in which 1) the E1E2 transmembrane domains (TMDs) are replaced by a heterodimeric scaffold to enforce assembly, 2) a furin cleavage site (6xArg) is inserted between the first scaffold and the E2 ectodomain to facilitate proper folding of the construct prior to secretion, 3) a tissue plasminogen activator (tPA) signal sequence is added at the N-terminus to boost expression, and 4) a His tag (6xHis) at the C-terminus to aid protein purification. The synthetic scaffolds used in these designs were found to be structurally similar to the human Fos-Jun leucine zipper, allowing for possible characterization of sE1E2 constructs without the potential of an autoimmune response. sE1E2.SZ was constructed with SYNZIP1/SYNZIP2, a designed coiled coil heterodimer (306), as a scaffold to facilitate sE1E2 assembly. sE1E2.HH was constructed with synthetic coiled coil peptides known to self-assemble as a hetero-hexamer (307), leading to scaffolded sE1E2 designed as a trimer of heterodimers rather than a single heterodimer, potentially mimicking predicted E1E2 display on the viral surface (51).



**Figure 3.1 Design of sE1E2 constructs with synthetic scaffolds.** (A) Schematic of mbE1E2, sE1E2.LZ, sE1E2.SZ, and sE1E2.HH. Regions shown include tPA signal sequence (green box), E1 ectodomain (yellow boxes), E2 ectodomain (red boxes), wild-type TMDs (gray boxes), and various scaffolds replacing TMDs. E1E2 residue ranges for each region are noted according to H77 numbering. C-terminal His tags and furin cleavage sites are shown in boxes and labeled. The expected molecular weight of each construct is indicated, and molecular weight of expected oligomers for sE1E2.HH is in parentheses. For molecular weight estimations, each N-glycan is approximated to be 2 kDa at each NxS/NxT sequon, a value within the molecular weight range of typical N-linked glycans (286). (B) X-ray structure of human Fos-Jun heterodimer (PDB code: 1FOS); only the coiled coil region that was used for the sE1E2.LZ scaffold is shown. c-Fos and c-Jun chains were colored to match the diagram of sE1E2.LZ. (C) X-ray structure of SYNZIP1/SYNZIP2 assembly (PDB code: 3HE5). SYNZIP1 and SYNZIP2 were colored to match the diagram for sE1E2.SZ. (D) X-ray structure of W22-L24H/Y15-L24D hetero-hexameric assembly (PDB code: 3R48). All chains were colored to match the diagram of sE1E2.HH. All structures were visualized in PyMOL (Schrödinger).

Α

# 3.3.2 sE1E2.SZ and sE1E2.HH mimic sE1E2.LZ secretion and antigenicity

Following transfection and expression, levels of sE1E2.SZ and sE1E2.HH secretion to supernatant were directly compared to sE1E2.LZ in western blot (**Figure 3.2**), with the anti-E2 antibody HCV1 (275) used to test E2 secretion and the anti-E1 antibody H-111 (276) used to test E1 secretion. Both sE1E2.SZ and sE1E2.HH showed robust detection of E1 and E2 protein in separate western blots, showing that these constructs secrete sE1E2 in amounts comparable to or higher than the sE1E2.LZ construct. The antigenicity of both designs was also compared in the supernatant of sE1E2.LZ and mbE1E2 using a panel of HCV antibodies in ELISA, including anti-E1, anti-E2, and anti-E1E2 antibodies (**Figure 3.3**). All binding levels of sE1E2.SZ and sE1E2.HH to these antibodies were comparable to sE1E2.LZ, suggesting that both constructs mimic the native-like antigenicity of the original sE1E2 design. Crucially, binding levels of anti-E1E2 antibodies AR4A and AR5A were maintained, showing that sE1E2.SZ and sE1E2.HH adequately present epitopes sensitive to proper heterodimeric assembly, despite switching to a synthetic scaffold, designing a larger oligomeric form, or both.

We also compared the effects of temperature on antibody binding to sE1E2.SZ, sE1E2.HH, sE1E2.LZ, and mbE1E2 as a method of assessing the relative stability of E1E2 assembly in each construct. In this experiment, all constructs were heated to 37°C or 56°C prior to a test of AR4A binding in ELISA (**Figure 3.4**). As temperatures increased, all constructs showed reduced binding to AR4A, suggesting that higher temperatures destabilized the E1E2 heterodimer and subsequently disrupted the conformational and heterodimer-dependent antibody epitope. Of the synthetic constructs, sE1E2.SZ showed a smaller reduction of AR4A binding than sE1E2.HH, providing evidence that sE1E2.SZ stability may be higher than sE1E2.HH and comparable to sE1E2.LZ. The combination of a synthetic scaffold with relatively higher stability at increased temperatures makes

sE1E2.SZ a promising candidate for vaccine design. Additionally, affinity of HCV antibodies to purified sE1E2.SZ protein was measured by quantitative ELISA and compared with purified sE1E2.LZ and mbE1E2 (**Table 3.1**). sE1E2.SZ antibody affinities were comparable to or higher than sE1E2.LZ antibody affinities, including for anti-E1E2 antibodies AR4A and AR5A, demonstrating that sE1E2.SZ also displays characteristics of a secreted and native-like heterodimer. These results also support the strategy of selecting a structurally similar scaffold based on a characterized template, despite the synthetic origin of the new scaffold.



Figure 3.2 Evaluation of sE1E2 secretion to supernatant in western blot. HCV1 antibody at 5  $\mu$ g/ml was used for the E2 western blot. H-111 antibody at 10  $\mu$ g/ml was used for the E1 western blot. All sE1E2 supernatant samples were loaded under reducing conditions. Supernatants were concentrated 10X prior to E1 western blot. Molecular weights, in kDa, of the western blot markers closest to observed bands are indicated on the left. Expected band positions of E1 and E2 are indicated with black triangles on right and labeled.



Figure 3.3 Binding of sE1E2 constructs and mbE1E2 to HCV HMAbs in ELISA. Supernatant from expressed sE1E2.LZ, sE1E2.SZ, and sE1E2.HH was added to ELISA plates and tested for binding to a panel of E1, E2, and E1E2 HMAbs, representing E1 N-terminus (H-111), E2 antigenic domains E (HCV1 and HC33.1), B (AR3A), and D (HC84.26.WH.5DL; abbreviated to HC84.26), as well as E1E2 domains AR4 (AR4A) and AR5 (AR5A). mbE1E2 protein was coated on ELISA plates at a concentration of 1  $\mu$ g/ml. Binding was measured at 450 nm with an antibody concentration of 0.5  $\mu$ g/ml. Negative control shown is an unrelated antibody (CA45). Figure provided by Ruixue Wang.



**Figure 3.4 Binding of sE1E2 constructs and mbE1E2 to HCV HMAbs in ELISA at elevated temperatures.** Prior to ELISA assay, samples were heated to 37°C, 56°C, or kept at room temperature for

one hour. Heated or room temperature supernatant from expressed sE1E2.LZ, sE1E2.SZ, and sE1E2.HH was added to ELISA plates and tested for binding to a panel of E1, E2, and E1E2 HMAbs, representing E1 N-terminus (H-111), E2 antigenic domains E (HCV1 and HC33.1), B (AR3A), and D (HC84.26.WH.5DL; abbreviated to HC84.26 in panels), as well as E1E2 domains AR4 (AR4A) and AR5 (AR5A). Heated or room temperature mbE1E2 protein was coated on ELISA plates at a concentration of 1  $\mu$ g/ml. Binding was measured at 450 nm with an antibody concentration of 0.5  $\mu$ g/ml. The % decreases shown in each panel refer to the reduction in O.D. for supernatant heated to 56°C compared to O.D. at room temperature. Asterisks above the bars for AR4A/AR5A binding at 56°C are only meant to highlight the labeled reductions, and do not indicate statistical significance. Negative control shown is an unrelated antibody (CA45). Figure provided by Ruixue Wang.

Antibody	Domain	$K_{\rm D} (nM)^1$				
		mbE1E2	sE1E2.LZ	sE1E2.SZ		
CBH-4D	А	$17 \pm 1$	$28 \pm 3$	$22 \pm 3$		
CBH-4G	А	$13 \pm 1$	$28 \pm 2$	$23 \pm 2$		
AR3A	В	$1.2 \pm 0.3$	$4.6\pm0.2$	$1.9\pm0.2$		
HEPC74	В	$0.57\pm0.04$	$3.0\pm0.1$	$0.88\pm0.03$		
HC84.26.WH.5DL	D	$0.94\pm0.06$	$1.5 \pm 0.1$	$0.94\pm0.02$		
HC84.1	D	$0.81\pm0.11$	$1.3 \pm 0.1$	$0.43\pm0.08$		
HC33.1	Е	$5.7 \pm 0.2$	$0.35\pm0.03$	$0.27\pm0.02$		
HCV1	Е	$5.7 \pm 0.2$	$0.35\pm0.01$	$0.35\pm0.01$		
AR4A	E1E2	$2.3 \pm 0.2$	$16 \pm 1$	$1.5 \pm 0.1^2$		
AR5A	E1E2	$2.2 \pm 0.2$	$7.4 \pm 2.2$	$3.5 \pm 0.2$		

Table 3.1 Antigenic analysis of mbE1E2, sE1E2.LZ, and sE1E2.SZ by quantitative ELISA.

<sup>1</sup>Figure provided by Zhen-Yong Keck

<sup>2</sup>Affinity of sE1E2.SZ to a given antibody showed an increase of more than 10-fold when compared to sE1E2.LZ

# 3.3.3 Design of sE1E2 constructs with alternative scaffolds

Encouraged by the initial characterization of sE1E2 designs with synthetic scaffolds, we expanded our search and design process to test alternative scaffolds and ectodomain sequences for assembly, secretion, and native-like antigenicity. These searches were conducted to find an optimal sE1E2 scaffold that may better enforce heterodimer or hetero-hexamer assembly or present a smaller immunogenic target and to find E1E2 ectodomains that can better induce bnAb responses, characteristics that could all lead to an improved sE1E2 immunogen for a vaccine. Constructs in **Figure 3.5** include alternative scaffolds that were designed following searches of the PDB (336) for heterodimeric or hetero-hexameric structures that were synthetic or non-eukaryotic in origin.

These searches led to the design of sE1E2 scaffolded from a variety of sources, with only some alternative scaffolds showing structural similarity with the coiled coil leucine zipper scaffold in sE1E2.LZ. The synthetic heterodimer IAAL-E3/IAAL-K3 was designed *de novo* as a coiled coil heterodimer and characterized as highly stable despite its smaller size (308); this heterodimer was used as an alternative scaffold for the construct sE1E2.1U0I due to its structural similarity to the Fos-Jun leucine zipper scaffold and to its smaller size that would present a minimized immunogenic target in a vaccine candidate.

Other structures selected as alternative scaffolds were non-eukaryotic in origin and, though matching the principle of heterodimeric scaffolding in sE1E2.LZ, showed less structural similarity to the Fos-Jun leucine zipper scaffold. A structure of the EpsI/EpsJ heterodimer from a bacterial type 2 secretion system (338) was incorporated as a scaffold for sE1E2.3CFI because it contained a N-terminal and parallel coiled coil assembly that resembled the coiled coil association of the Fos-Jun leucine zipper scaffold. The SpyTag/SpyCatcher complex, which forms a heterodimer when the SpyTag peptide irreversibly links to the single domain SpyCatcher protein (315), was also tested as a non-eukaryotic scaffold despite little structural similarity with Fos-Jun. This complex as structurally characterized (316) was used as an alternative scaffold in two designs, one with the SpyTag after E1 and SpyCatcher after E2 (sE1E2.SpyC) and one with the placement of the scaffolds switched (sE1E2.SpyT). This selection of two designs was based on the documented strength of SpyTag/SpyCatcher association, its utility in linking glycoprotein antigens of other viruses to nanoparticles for vaccine designs (317-320), and its successful stabilization of a large heterodimer after fusion to the C-terminus that aided structural characterization (339). In addition, the constructs that include 3CFI or SpyTag/SpyCatcher help to test the determinants for

scaffolding E1E2; specifically, whether a helical and coiled coil scaffold is necessary for assembly, or if spatial proximity of N-termini is sufficient.



**Figure 3.5 Design of sE1E2 constructs with alternative scaffolds.** (A) Schematic of sE1E2 constructs sE1E2.1U0I, sE1E2.3CFI, sE1E2.SpyC, and sE1E2.SpyT. Regions shown include tPA signal sequence (green box), E1 ectodomain (yellow boxes), E2 ectodomain (red boxes), wild-type TMDs (gray boxes), and various scaffolds replacing TMDs. E1E2 residue ranges for each region are noted according to H77 numbering. C-terminal His tags and furin cleavage sites are shown in boxes and labeled. The expected molecular weight of each construct is indicated. For molecular weight estimations, each N-glycan is approximated to be 2 kDa at each NxS/NxT sequen, a value within the molecular weight range of typical N-linked glycans (286). (B) NMR structure of synthetic IAAL-E3/IAAL-K3 heterodimer (PDB code: 1U0I). IAAL-E3 and IAAL-K3 chains were colored to match the diagram of sE1E2.1U0I. (C) X-ray structure of EpsI/EpsJ assembly (PDB code: 3CFI). EpsI and EpsJ chains were colored to match the diagram for sE1E2.1U0I. (D) X-ray structure of SpyTag/SpyCatcher complex (PDB code: 4MLS). SpyTag and SpyCatcher chains were colored to match the diagrams of sE1E2.SpyC and sE1E2.SpyT. All structures were visualized in PyMOL (Schrödinger).

#### 3.3.4 Design of sE1E2 constructs with alternative ectodomains

Constructs containing alternative E1E2 ectodomain sequences were selected either from designed consensus sequences or from an existing panel of functional HCV isolates with known neutralization profiles. To best match the original Fos-Jun leucine zipper assembly while still testing ectodomains with a non-human scaffold, the SYNZIP1/SYNZIP2 heterodimer was used as a scaffold for all alternative ectodomain constructs. Two designs incorporated consensus ectodomains, each utilizing a different scope of HCV genetic diversity. One consensus sequence that formed the ectodomains of sE1E2.cons.80 is pan-genotypic, as it was designed using a non-redundant set of aligned E1E2 sequences from genotypes 1-7. The other consensus sequence, contained in sE1E2.cons1.92.5, was designed with a non-redundant set of E1E2 sequences only from genotype 1, making this sequence more restricted genetically but focused on the most prevalent genotype in the utilized E1E2 reference dataset. In separate phylogenetic trees, both consensus sequences do not form a branch with reference sequences used to generate each consensus sequence, showing that these sequences are not biased toward any genotype (Figure 3.6) or genotype 1 subtype (Figure 3.7).

Though the generation of these consensus sequences is unbiased by genotype or subtype, it should be noted that genotype 6 sequences in **Figure 3.6** and genotype 1b sequences in **Figure 3.7** appear overrepresented in their respective trees. However, this observation may largely be explained by a higher level of genetic diversity within these groups, leading to more selections within a sequence identity threshold. This contribution is most evident in genotype 6, as it contains the most confirmed subtypes of any HCV genotype (67). An increase in genetic diversity within genotype 1b is less clear, but recent research has found higher levels of diversity in hypervariable region 1 (HVR1) of genotype 1b than genotype 3 and higher mutational flexibility than genotype

1a, which could contribute to the increased number of unique sequences within genotype 1b (340, 341). At the same time, both consensus sequences mostly retained common bnAb epitopes in E1 and E2, suggesting that these consensus ectodomains can readily bind antibodies that recognize highly conserved epitopes that are relevant for HCV neutralization (Figure 3.8). Sequence differences between H77 and either consensus sequence were found in H-111 (276), domain D (271), and domain B (254) antibody epitopes (Figure 3.8A, Figure 3.8D-E). These differences were the least concerning in H-111, which is a non-neutralizing antibody and likely does not impact bnAb responses directly. Though sequence differences were found in domain B and D epitopes that contain bnAbs, key residues in those domains were nearly or completely conserved in both consensus sequences (62). Alternative ectodomains from HCV isolates were selected from a diverse panel of E1E2 HCVpp (321) that described a reduced set of isolates with distinct patterns of sensitivity to antibody neutralization, allowing us to pick isolates with promising properties. Isolate 1.11.6 in the reduced panel was more sensitive to antibody neutralization than the reference sequence H77 used for sE1E2.LZ. Isolate 1a38 showed high HCVpp infectivity that matched levels of isolate 1b09, which has been used for structural characterization of the E2 ectodomain (244). Since their characteristics resembled those of well characterized genotype 1a isolates, 1a38 and 1.11.6 were incorporated into separate sE1E2 constructs to test the utility of these sequences in proper assembly and secretion of E1E2.



**Figure 3.6 Phylogenetic tree of cons.80 with sequences from genotypes 1-7.** Branches containing E1E2 sequences from genotypes 1-7 are highlighted with arced or straight lines and labeled. The branch corresponding to the cons.80 E1E2 consensus sequence is indicated with an asterisk. The scale for the branch lengths within this tree is identified with red text.



Figure 3.7 Phylogenetic tree of cons1.92.5 with sequences from genotype 1 subtypes. Branches containing E1E2 sequences from genotype 1 subtypes are highlighted with arced or straight lines and labeled. Sequences in this tree represent subtypes a-e, g-j, l-n, and unassigned subtypes (as of March 2022 (66)). The scale for the branch lengths within this tree is identified with red text.



**Figure 3.8 Comparison of H77 and consensus sequences at residue positions of key epitopes.** Any sequence difference between H77 and cons1.92.5 or cons.80 was highlighted with an arrow above the sequence alignment. (A) H-111 epitope comparison, with five sequence differences indicated. (B) IGH526 epitope comparison showed complete sequence conservation. (C) Domain E epitope comparison showed complete sequence conservation, with four sequence differences indicated. (F) AR4A epitope comparison showed complete sequence conservation. Sequence ranges of H77, cons1.92.5, and cons.80 are labeled by sequence name on left-side panels A, C, and E, but these labels also apply to right-side panels B, D, and F. Sequence ranges or individual residues in E1 or E2 include a label above the corresponding residue position in the alignment. Residues are colored using default settings in visualization of sequences with SeaView.

# 3.3.5 Characterization of alternative sE1E2 constructs

Following expression, each construct was tested for secretion of E1E2 ectodomains to the supernatant. The ectodomains were detected in separate western blots, probing with the same antibodies used to separately assess secretion of E1 and E2 in sE1E2.LZ, sE1E2.SZ, and

sE1E2.HH. E2 was detected in the supernatant for all constructs, with constructs that have larger scaffolds found at higher molecular weights than sE1E2.LZ (**Figure 3.9A**). E1 was also detected by western blot for all constructs except sE1E2.1U0I and sE1E2.cons.80 (**Figure 3.9B**), largely confirming that these constructs secrete as designed and often at levels comparable to sE1E2.LZ. sE1E2.SpyC and sE1E2.SpyT showed single bands with an estimated molecular weight of about 100 kDa even under reducing conditions. This result suggests that E1E2 scaffolded by the SpyTag/SpyCatcher complex is still covalently linked in the presence of SDS, consistent with previous characterizations of the SpyTag/SpyCatcher complex (315). The lack of E1 detection for sE1E2.cons.80 may be explained by multiple sequence differences in the H-111 epitope between the H77 and cons.80 sequences, as shown in **Figure 3.8A**.

Supernatant of designed sE1E2 constructs was also tested for binding to a panel of five HCV bnAbs using ELISA, including anti-E1, anti-E2, and anti-E1E2 antibodies (**Figure 3.10**). Binding to antibodies HCV1 and HC84.26.WH.5DL (271) was relatively high for all constructs. This result was especially encouraging for the tested consensus sequences, which showed sequence differences in domain D (**Figure 3.8D**) but still bound to domain D antibody HC84.26.WH.5DL at levels comparable to sE1E2.SZ. sE1E2.1U0I, sE1E2.3CFI, sE1E2.1a38, and sE1E2.1.11.6 showed decreased binding to IGH526, AR4A, and AR5A, suggesting that these constructs may not form native-like E1E2 heterodimers. However, the constructs sE1E2.SpyC, and sE1E2.SpyT, sE1E2.cons.80, and sE1E2.cons1.92.5 showed antibody binding that was comparable to sE1E2.SZ and sE1E2.HH, showing that each design exhibits native-like antigenicity and represents a promising candidate of secreted E1E2 heterodimer with a synthetic or non-eukaryotic scaffold. Results for sE1E2.SpyC and sE1E2.SpyT suggest that methods of E1E2 scaffolding independent of coiled coil assembly are also feasible, as the SpyTag/SpyCatcher complex uses an entirely

different mechanism of self-assembly and is structurally dissimilar to coiled coil peptide structures. Alternative ectodomain constructs sE1E2.cons.80 and sE1E2.cons1.92.5 exhibited native-like antigenicity while sE1E2.1a38 and sE1E2.1.11.6 showed decreased binding to AR4A and AR5A, suggesting that the utilization of consensus sequences may be more promising than selecting alternative isolates. Consensus sequence immunogens have also induced cross-reactive antibodies against diverse viral glycoproteins (342, 343), supporting the consideration of a E1E2 consensus sequence for vaccine design. The reduction in AR4A and AR5A binding to alternative ectodomain constructs was unexpected, as there was no clear reason why these genotype 1a ectodomains would show differences in assembly with H77, a sequence in the same subtype. There may be sequence determinants within 1a38 and 1.11.6 that modulate heterodimer conformation or dynamics and lead to reduced AR4A and AR5A recognition in an sE1E2 context, but any investigation of these contributions was beyond the scope of this study.



Figure 3.9 Detection of alternative sE1E2 constructs in western blot. (A) Western blot probing for E2 with HCV1 antibody at 5  $\mu$ g/ml. (B) Western blot probing for E1 with H-111 antibody at 10  $\mu$ g/ml. All sE1E2 supernatant samples were loaded under reducing conditions. Molecular weights, in kDa, of the western blot markers closest to observed bands are indicated on the left. Expected band positions of E1, E2, or E1E2 heterodimer are indicated with black triangles on right and labeled. The "sE1E2" part of sE1E2 construct names was omitted from labels for clarity. Figure provided by Dongxiu Zhang Spiering.



**Figure 3.10 Antibody binding to alternative sE1E2 constructs in ELISA.** Supernatant from expressed sE1E2 constructs was added to ELISA plates and tested for binding to a panel of E1, E2, and E1E2 bnAbs, representing E1 (IGH526), E2 antigenic domains E (HCV1) and D (HC84.26.WH.5DL), as well as E1E2 domains AR4 (AR4A) and AR5 (AR5A). Binding was measured at 450 nm with an antibody concentration of 1  $\mu$ g/ml. The "sE1E2" part of sE1E2 construct names was omitted from labels for clarity. Figure provided by Dongxiu Zhang Spiering.

#### **3.4 Discussion**

We have designed and characterized various sE1E2 constructs that contain either an alternative scaffold (synthetic or non-eukaryotic) or alternative ectodomains (consensus or different HCV isolate). Many of these constructs showed effective secretion of E1E2 ectodomains in western blot and native-like antigenicity in ELISA, with sE1E2.SZ found as a synthetic analog to sE1E2.LZ. Other constructs showed little structural similarity with sE1E2.LZ, but still produced secreted and native-like sE1E2 in the same assays. sE1E2 designs with consensus sequences of genotype 1-7 and genotype 1 E1E2 were especially promising, demonstrating that a sequence designed from a pool of naturally occurring HCV sequences preserved both epitopes critical for bnAb recognition and ectodomain residues required for proper assembly of the heterodimer.

Assessing alternatives for sE1E2 design found possible improvements in the scaffold and ectodomain sequence components, providing further support for the proof of concept established by sE1E2.LZ while also broadening the options for future iterations.

Although all sE1E2 constructs were scrutinized for proper secretion and assembly via interactions with HCV bnAbs, the determinants of and potential differences in construct formation are largely unclear. For sE1E2.LZ, analytical ultracentrifugation (AUC) and size exclusion chromatography with multi-angle light scattering (SEC-MALS) were performed to fully characterize the formation of purified construct, finding both sE1E2.LZ dimers and monomers in solution (337). The sE1E2 constructs that incorporate alternative scaffolds or ectodomains have not been tested using the same techniques, making the heterogeneity and true molecular weight of these constructs unknown. This characterization is especially important for sE1E2.HH, given its design as a hetero-hexameric construct and expected difference in molecular weight from sE1E2 constructs designed as heterodimers. Gathering data on sample heterogeneity for each construct will not only determine how much the predicted sE1E2 design differs from its behavior in solution, but will help elucidate the impact of selected scaffolds on the prevalence of higher molecular weight populations, which may inform development of an HCV vaccine candidate. Additionally, further research can test the native-like properties of subsequent sE1E2 designs beyond recognition by bnAbs. As with sE1E2.LZ characterization, binding of sE1E2 constructs to the large extracellular loop of CD81 (CD81-LEL) should be tested, as this interaction is crucial for infection and entry in native virus (344, 345) and would provide more confirmation that E1E2 is properly assembled in these constructs.

sE1E2 constructs with synthetic or non-eukaryotic scaffolds showed a range of secretion levels and antigenicity, helping to assess the boundaries of scaffold usage that can result in secreted and native-like sE1E2 heterodimer. Coiled coil scaffolds of synthetic origin showed the most success, with sE1E2.SZ and sE1E2.HH suggested to form native-like complexes. However, sE1E2.1U0I did not show E1 secreted in western blot or native-like antigenicity in ELISA, despite including a parallel and synthetic coiled coil heterodimer as a scaffold. One possible explanation for diminished secretion and assembly could be the length of the scaffold; coiled coil chains in 1U0I are only 21 amino acids long (308), less than half the length of chains used for scaffolds in sE1E2.LZ and sE1E2.SZ. Though this reduced scaffold length presents a smaller immunogenic target as a prospective vaccine candidate, the results for sE1E2.1U0I suggest that this scaffold may be too small to produce an E1E2 heterodimer with native-like properties. Even the smallest scaffold by individual chain is 32 amino acids in length (sE1E2.HH), and reducing this length further may make any coiled coil scaffold unable to facilitate E1E2 assembly effectively, regardless of the ability of the scaffold to form a stable heterodimer on its own. Another possibility is that the current orientation and placement of 1U0I components is not optimal for facilitating sE1E2 assembly and secretion. With the native orientation of E1E2 ectodomain assembly unknown, it is unclear if heterodimeric coiled coil scaffolds need to be attached to specific Cterminal positions or if proper assembly would occur regardless of the order of scaffold components. To investigate these possibilities, the scaffold components for sE1E2 designs with coiled coil scaffolds could be switched and retested for secretion and antigenicity. In identifying the optimal placement of coiled coil scaffolds, this future work could help determine structural restraints in sE1E2 scaffolding that may provide insights for native E1E2 assembly and rational vaccine design.

sE1E2 constructs with non-eukaryotic scaffolds also showed a range of secretion and antigenicity, with the most promising results from scaffolds with no structural similarity to coiled
coil heterodimers. The constructs utilizing the SpyTag/SpyCatcher complex (sE1E2.SpyC, sE1E2.SpyT) showed more native-like antigenicity than a scaffold with N-terminal coiled coil regions (sE1E2.3CFI), which may stem from the unique mechanism of SpyTag/SpyCatcher self-assembly through covalent isopeptide linkage (315). This strong and irreversible linkage between scaffold components in sE1E2.SpyC and sE1E2.SpyT produced secreted and native-like E1E2 regardless of where SpyTag and SpyCatcher were placed, suggesting that this scaffolding method could be a powerful tool for facilitating E1E2 assembly. Furthermore, this complex could also lead to multimeric display of sE1E2, either alone or in combination with other scaffolds. This approach may utilize self-assembling nanoparticles in a method analogous to previous studies of vaccine candidates using SpyTag/SpyCatcher (319, 320), which could facilitate coupling to nanoparticles following expression and allow production of nanoparticle-displayed sE1E2 with increased immunogenicity.

For the tested alternative ectodomains, consensus sequence antigenicity was much closer to native E1E2 than ectodomains from isolates 1.11.6 or 1a38. Although both isolates were previously characterized as functional and infectious in HCVpp (321), it is still possible that differences with the H77 sequence may contribute to decreased stability of a scaffolded E1E2 heterodimer and decreased antigenicity. In contrast, the designed consensus sequences may represent ectodomains that are highly stable versions of a heterodimer for genotype 1 or genotypes 1-7, producing robust immunogens without having to account for the sequence of one isolate that could contain unique characteristics affecting sensitivity to bnAb neutralization. Future work should assess potential changes in sE1E2 construct yield and stability when changing ectodomain sequences, helping to elucidate sequence determinants for native-like antigenicity of sE1E2 and improved protein yield for a prospective immunogen. This research could incorporate previously developed E1E2 consensus sequences, which often were designed from genotype 1 sequences (346, 347). Also, the alternative isolates may not have been ideal candidates for sE1E2 designs, as they may have modulated E1E2 heterodimeric assembly in a way that makes sE1E2 non-native such as shifts in conformational dynamics or changes in the relative stability of the complex. In future work, ectodomains from other isolates could be selected that have low or medium levels of HCVpp infectivity – high for isolates H77 and 1b09 (321) that have been structurally characterized (244, 251) – to better determine how this metric could influence native-like presentation in an sE1E2 context.

Chapter 4: Prediction of hepatitis C virus polymorphisms impacting antibody neutralization and residues critical for E1E2 heterodimeric assembly

### Abstract

Hepatitis C virus (HCV) is a worldwide disease burden, and an effective vaccine is needed to facilitate reduced infection and global eradication. Surface glycoproteins E1 and E2 form a heterodimer at the surface of the virion and are the primary target for broadly neutralizing antibodies (bnAbs), making these proteins the primary target as immunogens for vaccine development. However, factors such as HCV sequence variability and the lack of an accessible E1E2 heterodimer structure have hindered HCV vaccine design, as the sequence determinants of increased resistance to bnAbs and of proper E1E2 assembly remain largely unclear. Here, we analyze previously published datasets on HCV antibody neutralization and E1E2 mutagenesis to predict polymorphisms contributing to changes in antibody neutralization and residues required for proper assembly of the E1E2 heterodimer. Polymorphisms in 80 E1E2 residue positions were predicted to affect antibody neutralization through the analysis of neutralization data with a predictive algorithm and comparisons of neutralization results between highly similar sequences, with some positions not examined in previous studies. Hierarchical clustering of E1E2 mutagenesis datasets found residues that primarily affected anti-E1E2 antibody binding, suggesting that these residues are critical for E1E2 assembly. Both sets of predictions were evaluated for possible structural effects, either on partial E1 or E2 structures with computational mutagenesis or on a recently reported structure of the E1E2 heterodimer. Overall, these in-depth analyses of existing datasets detected E1E2 residues with predicted impacts on antibody neutralization and heterodimeric assembly with unclear mechanisms that warrant further investigation.

### 4.1 Introduction

As discussed in Chapters 2 and 3, a vaccine against hepatitis C virus (HCV) is urgently needed to combat infection globally, even though direct acting antivirals (DAAs) can cure infection at high rates. Characterization of broadly neutralizing antibody (bnAb) responses to HCV E1E2 glycoproteins have helped to inform vaccine design approaches. Both chapters also mentioned that antibody responses to E1E2 can be defined by patterns of key residues of antibody epitopes, especially for regions of antibody binding on E2 that have been classified as antigenic domains (A-E) (62, 63), antigenic regions (AR1-AR3) (55), or epitopes (I-III) (348). Additional antigenic regions have been defined for antibody epitopes that contact residues on both E1 and E2 (AR4-AR5) (87, 243), along with characterized anti-E1 antibodies without previously defined classifications (276, 349, 350). Some of these antigenic domains are highly conserved (301, 351, 352) or overlap with the binding site of critical co-receptor CD81 (54, 353), making these antibody epitopes promising targets for generating bnAb responses (243, 244, 354). Several of these bnAbs have been characterized extensively, both in the breadth of HCV isolate neutralization (55, 240, 355) and the structural basis of antibody interactions with bnAb epitopes (96, 252, 254).

However, E1E2 genetic diversity and sequence polymorphisms have restricted the breadth of neutralizing antibodies and contributed to viral immune evasion and escape, including in clinical trials (97) and *in vitro* despite substantial fitness costs in some cases (99, 100, 102). Though some bnAb epitopes are highly conserved, sequence polymorphisms in E1E2 have been found to reduce antibody neutralization even when those polymorphisms are not directly contacted by the antibody (101-103). These extra-epitopic residues have been found primarily in regions of E2, including hypervariable region 1 (HVR1), hypervariable region 2 (HVR2), and the E2 back layer. To examine contributions of extra-epitopic residues to neutralization resistance more closely, El-Diwany et al. utilized neutralization data of domain E antibody HC33.1 (268) and anti-E1E2 antibody AR4A (247) against 113 genotype 1a and 1b HCV isolates to predict and experimentally validate sequence polymorphisms that increase resistance to antibody neutralization (103). Though this analysis suggested possible mechanisms for antibody resistance via polymorphisms, including a shift in SR-BI co-receptor dependency, the full breadth and possible mechanisms of extra-epitopic contributions to antibody resistance remain unclear. In other studies, the breadth of antibody neutralization to a diverse set of HCV sequences has been evaluated (240, 356), but the contributions of extra-epitopic residues were not specifically investigated, despite the possibility of substantial impacts on antibody neutralization.

Recently, mechanisms of neutralization resistance have been explored through the effects of HCV E1E2 flexibility and dynamics (89, 98). For instance, sequence differences in HVR1 were associated with a shift in SR-BI dependency and altered dynamics of viral breathing, specifically between putative open and closed E2 states that control accessibility of the CD81 binding site to co-receptors. While these sequence differences and others may play a role in resistance to antibody neutralization through shifts in glycoprotein flexibility, the potential structural effects of polymorphisms have not been examined within the context of E1E2 assembly. E1E2 is thought to form a trimer of heterodimers on the viral surface (51), and several studies have performed mutagenesis on E1E2 residues to determine critical residues for heterodimeric assembly. Most residues suspected of contributing to E1E2 heterodimerization were in the hydrophobic E2 transmembrane domain (TMD), E2 stem domain, or the E1 N-terminal domain (234, 298, 357), but it has been difficult to confirm which E1E2 residues are critical for heterodimerization without a characterized and accessible heterodimer structure. Elucidation of critical E1E2 residues continues to be confounded by this knowledge gap, despite high-throughput mutagenesis data of E1E2 that have helped to characterize the epitopes of neutralizing antibodies in recent years (87, 88). Although mutations in determinants of E1E2 assembly may contribute to changes in binding of neutralizing antibodies to a variety of epitopes, this information has not been harnessed to help elucidate the mechanisms of E1E2 heterodimerization, leaving the overall picture of E1E2 assembly determinants unclear.

Here, we utilize previously published neutralization and mutagenesis datasets to predict both polymorphisms associated with changes to antibody neutralization and E1E2 residues critical for heterodimerization. For polymorphism predictions, we used neutralization data from two different sources (103, 356), where several antibodies to different epitopes were tested for neutralization on a diverse set of HCV isolates. An algorithm called Subject-adjusted Neutralization Antibody Prediction of Resistance (SNAPR) was modified to predict polymorphisms contributing to neutralization changes in a set of genotype 1a sequences with IC50 data, mimicking the strategy of predicting polymorphisms from a previous study (103). Highly similar sequences in both neutralization datasets were also compared to find the largest fold changes for a given antibody between E1E2 sequences with one or two mutations. Following these strategies, polymorphisms predicted to contribute to neutralization changes for a larger number of antibodies were noted. Predicted polymorphisms were also modeled on partial E1 and E2 structures using computational mutagenesis to examine the predicted effects of these polymorphisms on glycoprotein stability. To predict E1E2 residues that are crucial for heterodimerization, mutagenesis data from two different sources (87, 88) were analyzed with hierarchical clustering, which visualized distinct patterns of antibody binding disruption that largely matched pre-defined annotations of the tested antibodies. Mutagenesis data were also clustered by E1E2 residues, and two of those clusters showed a disruption of anti-E1E2 antibody binding, even though binding to other antigenic domains and CD81 was largely intact. Some of the residues in these clusters corresponded to residues that affected E1E2 heterodimerization in previous studies, and several predicted E1E2 contacts were found in figures presenting a structure of the E1E2 heterodimer. Overall, these predictions broadly investigate polymorphism contributions to antibody neutralization and residue determinants of E1E2 heterodimerization, potentially providing a better understanding of E1E2 immune escape and assembly that can inform HCV vaccine design.

### 4.2 Methods

### 4.2.1 Collection of antibody neutralization data

Data on antibody neutralization were obtained for binding of two antibodies to 113 HCV isolates from genotype 1 (103), and for binding of five antibodies to 69 HCV isolates from genotypes 1-6 (356). The Urbanowicz et al. dataset measured antibody neutralization as IC50 values in  $\mu$ g/ml. The El-Diwany et al. dataset measured antibody neutralization as fraction unaffected (F<sub>u</sub>) values by comparing HCV pseudoparticle (HCVpp) infection in the presence and absence of neutralizing antibody.

### 4.2.2 Prediction of polymorphisms contributing to neutralization changes with SNAPR

The Subject-adjusted Neutralization Antibody Prediction of Resistance (SNAPR) R script (103) was kindly shared by Justin Bailey. This script was modified to calculate differences in halfmaximal inhibitory concentration (IC50) values rather than F<sub>u</sub> values to predict polymorphisms contributing to changes in antibody neutralization for 39 genotype 1a sequences in the Urbanowicz et al. dataset. In brief, the SNAPR script reads a multiple sequence alignment (MSA) of HCV E1E2 sequences with corresponding neutralization data, separates sequences by polymorphism at a given position of the alignment, then predicts the impact of a polymorphism on neutralization by assessing differences using Wilcoxon rank sum tests. The SNAPR script summarizes its predictions by outputting a list of the ten most impactful polymorphisms based on these statistical tests. In separate runs after modification of the script, SNAPR was used to predict polymorphisms with the highest contributions to antibody resistance and with the highest contributions to antibody sensitivity. This analysis was conducted by antibody in the Urbanowicz et al. dataset, which included neutralization data for the antibodies AP33 (358), 1:7 (359), D03 (360), L1 (359), and XTL68 (361).

### 4.2.3 Pairwise comparisons of antibody neutralization data

For both neutralization datasets, E1E2 sequences with no more than two sequence changes between them (>99.5% sequence identity) were noted and the neutralization differences for each antibody in the dataset were analyzed. Highly similar sequences from MSAs of isolates in each dataset were found using BLAST (362). The largest neutralization differences between highly similar sequences were quantified as fold changes of neutralization and identified for all seven antibodies within the two datasets. Sequence differences in a comparison were identified as predicted contributors to changes in antibody neutralization if a neutralization fold change in that comparison was greater than five for any antibody.

# 4.2.4 Computational mutagenesis of polymorphisms predicted to contribute to neutralization changes

Polymorphisms predicted to contribute to neutralization changes, either with SNAPR or a pairwise comparison, were modeled onto E1 and E2 structures to assess the predicted effects on

protein stability. 31 polymorphisms could not be modeled on any E1 or E2 structure, including 14 E1 polymorphisms and 17 E2 polymorphisms. These polymorphisms could not be introduced because the residue position was either disordered in every structure or absent from the ectodomain in E1 or E2 used for structural characterization (251, 363). Polymorphisms were modeled with the protein modeling software Rosetta version 2.3, specifically applications within Rosetta for computational mutagenesis (364, 365). In Rosetta version 2.3, computational mutagenesis was conducted on E1 N-terminal domain (PDB code: 4UOI; (363)) and antibody-bound E2 ectodomain structures from three different subtypes in two genotypes: genotype 1a (PDB code: 4MWF; (251)), genotype 1b (PDB code: 6MEJ; (244)), and genotype 2a (PDB code: 7MWW; (53)). The - min\_interface and -int\_chi flags were added to each command for mutagenesis to refine residues nearby the mutation prior to  $\Delta\Delta G$  prediction. The command for a typical run with 4MWF is below, with 12 selected as a constant seed:

/piercehome/pierce/rosetta/rosetta2.3/rosetta++/rosetta.gcc34 aa 4MWF
\_ -interface -output\_structure -intout 4MWF.ddg.ros.out ignore\_unrecognized\_res -safety\_check -skip\_missing\_residues -mutlist
4MWF.muts.txt -min\_interface -int\_chi -extrachi\_cutoff 1 -ex1 -ex2 ex3 -constant\_seed -jran 12 -yap -s 4MWF

Polymorphisms with predicted  $\Delta\Delta G$  values, in kcal/mol, that were less than or equal to -0.7 were considered stabilizing mutations. Polymorphisms with predicted  $\Delta\Delta G$  values, in kcal/mol, that were greater than or equal to 0.7 were considered destabilizing mutations. All other  $\Delta\Delta G$  values were classified as neutral.

### 4.2.5 Collection of E1E2 mutagenesis data

Experimental mutagenesis data were obtained from previously published datasets showing the impact of E2 alanine mutants on the binding of 21 antibodies (88) and the impact of E1E2 alanine mutants on the binding of 13 antibodies and the co-receptor CD81 (87). These datasets were combined for the subsequent analysis. Merged data on relative binding of all antibodies and CD81 to E2 alanine mutants were used for generation of the mutagenesis heatmap and residue clusters based on E2 mutagenesis data alone. The dataset from Gopal and Jackson et al. was also examined independently for clustering E1E2 residues because only this dataset contained antibody binding to mutants in both E1 and E2.

### 4.2.6 Clustering of mutagenesis data

Clustering of antibodies and E1E2 residues by patterns in the experimental mutagenesis data was performed using hierarchical clustering in R (366). An unrooted tree of antibodies based on E2 merged mutagenesis data was generated with the ggplot2 package in R (367). A heatmap of residue and antibody clustering was generated using the ggplot2, grid, reshape, and brackets packages in R (367-369). Once antibody clusters were visualized, an in-house Perl script was used to output averages of relative antibody binding values for each pre-defined antibody group in E1E2 mutagenesis and merged E2 mutagenesis datasets. In the merged dataset, these groups were antigenic domain A, antigenic domain D/antigenic domain B, antigenic domain C, antigenic domain E, AR2, AR3/domain B, AR4-5, distinct E1 antibodies A4 and IGH526, and CD81-LEL. Two groups were identified with multiple classifications to reflect how antibodies with similar classifications clustered together based on mutagenesis data. AR3 and domain B essentially define the same E2 epitope under different classification schemes (62). Another group was labeled as domain D/domain B because two antibodies in this cluster are classified as domain B but are also

known to compete with domain D antibodies for binding (88). In the E1E2 dataset, these groups were AR1, AR2, AR3, AR4-5, distinct E1 antibodies A4 and IGH526, antigenic domain E, and CD81-LEL. The number of residue clusters for merged E2 data (N = 16) and E1E2 data (N = 30) were set following examination of several cluster numbers. These values were selected through evaluations of average relative antibody binding in clusters, finding that the chosen numbers led to clusters that best isolated residues with the largest impacts on E1E2 antibody binding, as indicated by average relative binding values.

### 4.3 Results

### 4.3.1 Neutralization datasets used for predictions

Several groups have published antibody neutralization datasets with diverse HCV E1E2 isolates in recent years (103, 356). The sequences included in these studies were isolated from HCV-infected subjects and were generated as functional HCV pseudoparticles (HCVpp) and HCV cell culture virus (HCVcc). In one dataset, hereafter known as the Bailey dataset (103), 113 genotype 1a and 1b isolates were tested for antibody neutralization with HC33.4 (84), a domain E antibody, and AR4A (247), an anti-E1E2 antibody. The other dataset, hereafter known as the Ball dataset (356), included 69 sequences of HCV isolates that included genotypes 1-6 and were tested for neutralization against five antibodies: domain E antibody AP33 (370), 1:7 (359), D03 (360), L1 (359), and XTL68 (361). These studies used different methods to measure neutralization, yet each found a broad range of neutralization values; the Bailey dataset was measured using fraction unaffected (F<sub>u</sub>), or an estimated percentage of isolate that was not neutralized (103), while the Ball dataset was measured as IC50s (356). The contents and characteristics of each dataset are summarized in **Table 4.1**.

Dataset	Neutralization	Total	Genotype	Genotype	Genotype	Antibodies
	measurements	sequences	1a	1b	2-6	tested
			sequences	sequences	sequences	
Bailey	Fraction	113	71	42	0	HC33.4,
(103)	unaffected					AR4A
	$(F_u)$					
Ball (356)	IC50 ( $\mu$ g/ml)	69	39	17	13	AP33, 1:7,
						D03, L1,
						XTL68

Table 4.1 Summary of neutralization data from previously published datasets.

### 4.3.2 SNAPR predicted E1E2 polymorphisms contributing to neutralization changes

In the study presenting the Bailey dataset, the Subject-adjusted Neutralization Antibody Prediction of Resistance (SNAPR) was used to separate HCV isolates by polymorphisms at a given position in an MSA, then to make statistical comparisons between differences in median F<sub>u</sub> values of these subsets (103). The comparisons found to have the best SNAPR values, or subject-adjusted p-value, were predicted to make the largest contributions to resistance or sensitivity to antibody neutralization, depending on type of search in which the polymorphism was found. The separation by polymorphism and statistical comparisons made by SNAPR were also plotted to visualize neutralization differences (example in Figure 4.1). Using SNAPR, we analyzed neutralization data in the Ball dataset to predict sequence polymorphisms that contribute to changes in antibody neutralization, expanding on previous research. Though the Ball dataset encompassed genotypes 1-6, genotype 1 sequences were highly overrepresented (356), making predictions with HCV isolates from multiple genotypes difficult despite their potential value. Instead, SNAPR predictions were conducted using a similar genotypic scope as when the algorithm was used previously (103); more specifically, SNAPR primarily analyzed genotype 1a isolates from the Ball dataset. Though a separate analysis with 1a and 1b isolates was also performed, El-Diwany et al. noted the potential for inherent sequence differences between subtypes to dominate predictions, leading us to mimic

their approach while already working with a smaller set of sequences. Polymorphisms that were predicted as contributors by SNAPR for one or more antibodies are summarized in **Table 4.2**.

SNAPR predicted that polymorphisms in 36 E1E2 residue positions had the largest contributions to changes in neutralization by antibodies from the Ball dataset. Since a residue position may have been identified in an "antibody resistance" or "antibody sensitivity" search, the polymorphism or set of polymorphisms predicted to be more resistant to antibody neutralization is indicated for each residue position. Polymorphisms in 12 residue positions, all within E2 (residue range 384-746), were predicted contributors to antibody neutralization in three or more antibodies in this dataset, suggesting that these polymorphisms may have a global effect on HCV antibody responses. Some polymorphisms in these predictions have been implicated in resistance to antibody neutralization in previous studies. Residue positions 242 and 438 were predicted by SNAPR for antibody L1; these polymorphisms were also predicted by SNAPR for antibodies AR4A and HC33.4 in El-Diwany et al., and the V438 polymorphism was found to significantly change antibody neutralization when introduced by site-directed mutagenesis into HCVpp (103). This polymorphism was also found to alter binding of the SR-BI and CD81 co-receptors, a mechanism that could extend to predictions for residue positions 442 and 531, which are nearby residues critical for binding to CD81 (298, 371), were found in the interface of a recent E2-CD81 complex structure (53), and were previously identified as polymorphic in resistant E1E2 sequences (102). Other polymorphisms in residue positions listed in Table 4.2 have been identified in the same analysis of resistant E1E2 sequences (positions 416, 446, 461, 475, 520, 524) (102), in the polymorphisms of an escape pathway from a broadly neutralizing antibody (position 610) (93), or in polymorphisms that may enhance viral fitness in combination with resistance-associated polymorphisms at positions 415/417 (positions 399, 463, 603) (101). However, other polymorphisms predicted by SNAPR have not been studied for possible effects on antibody neutralization, primarily at residue positions in HVR1 and HVR2 along with six residue positions in E1. Interestingly, some of these residue positions were predicted to have the broadest effects on antibody neutralization changes, as they were the only polymorphisms to be predicted by SNAPR in all antibodies of the Ball dataset (position 474), or in all but one antibody (positions 397, 460, 709, 742). Most of the polymorphisms in **Table 4.2** occur in residue positions outside of documented antibody epitopes (359-361, 370), suggesting that these polymorphisms may modulate antibody neutralization in uncharacterized ways.



Figure 4.1 Example of SNAPR predictions of E1E2 polymorphism contributions. IC50 values of sequences in Ball dataset were separated into groups by residue at a given E1E2 position. Each black dot in the plot represents the neutralization value of antibody AP33 to one genotype 1a isolate in the dataset. In this comparison, SNAPR identified A742 as significantly more sensitive to AP33 neutralization than S742. A Wilcoxon rank sum test was used to test statistically significant differences between IC50 values of the most sensitive residue with the combined IC50 values of every other residue. This comparison found a high SNAPR value and a significant p-value (p < 0.001) when IC50s were not subject adjusted, or copied to ensure that values from each HCV-infected subject were represented equally (103). Median IC50 values are shown as black bars in each boxplot.

Residue	AP33	1:7	D03	L1	XTL68	Resistant
position <sup>1</sup>						residue(s)
198				Х		S
219	Х					S, T
234			Х		Х	N
241	Х					Р
242				Х		V
308		Χ				Ι
330		Χ				Т
386				Х		H, Y
395	Х	Х			Х	S, T
396				Х		T, V
397	Х		Х	Х	Х	A, S
401	Х				Х	S
405		Χ	Х		Х	P, R
411	Х	Х	Х			Ι
416				Х		Т
438				Х		V
442		Х	Х	Х		F
446		Х				R
460	Х	Χ	Х	Х		K
461	Х					P, R
463				Х		A
466		Х		Х		A, N
474	Х	Х	Х	Х	Х	Y
475	Х				Х	Т
477	Х	Х				G
478	Х	Х	Х			S
500	Х				Х	K, L, Q, R
520					Х	D
524					Х	V
531					Х	E
580			Х			R, Y
603			Х			L, V
610	Х		Х		Х	Н
653	Х	Х			Х	D
709	Х	Х	Х		Х	Ι
742	Х		Х	Х	Х	S

Table 4.2 Summary of residue positions with a polymorphism found as SNAPR hit.

<sup>1</sup>E1E2 residue positions according to H77 numbering. Residue positions found by SNAPR in three or more antibodies are highlighted in bold.

## 4.3.3 Pairwise sequence comparisons predicted E1E2 polymorphisms contributing to neutralization changes

We utilized SNAPR to make dataset-level predictions of contributions to changes in antibody neutralization, which was useful in examining the potential breadth of polymorphism effects. However, this algorithm does not consider identity between two aligned sequences as a potential indicator of polymorphisms contributing to neutralization changes, meaning that SNAPR may miss highly impactful polymorphisms within a given dataset simply due to decreased prevalence or an insufficient number of sequences in the analysis. In both datasets, we found highly similar E1E2 sequences and examined pairwise differences in neutralization for each antibody, complementing the predictive approach of SNAPR. Sequence changes in comparisons with neutralization differences greater than fivefold were identified, offering additional predictions of polymorphisms that contribute to neutralization changes, with residue positions listed in Table **4.3**. Through pairwise comparisons of sequences in both datasets, we found 52 residue positions that were the only sequence change, or one of two sequence changes, when differences in neutralization values were greater than fivefold for any antibody. Nineteen of these residue positions represent the only sequence change in highly similar E1E2 sequences with neutralization fold changes above the set threshold, with some of these positions also predicted by SNAPR (positions 653, 709, 742). As with SNAPR predictions, several residue positions in Table 4.3 have been examined as possible contributors to changes in antibody neutralization in previous studies (positions 399, 408, 434, 456, 492, 520, 538, 558, 570, 629) (93, 101-103) or repeatedly documented to help escape neutralization by AP33 and similar antibodies through a glycan shift (position 417) (62, 97, 101, 372). Single sequence changes at residue positions such as 260 and 496 revealed high fold changes for multiple antibodies, suggesting that these polymorphisms have

broad effects on antibody neutralization and represent novel predictions. Residue positions identified by pairwise comparisons also include more E1E2 TMD polymorphisms than predicted by SNAPR, which indicates that the pairwise comparison method could be finding sequence changes that are rarer yet no less impactful to antibody neutralization.

Residue	AR4A	HC33.4	AP33	1:7	D03	L1	XTL68	Resistant
position <sup>1</sup>								residue(s)
195	1							Н
199	1	1						R
208	2	1						S
234	1							D
254			1	1				А
255						1	1	А
260			2	2	2	1	2	Н
275							1	L
284		1						А
290	1	1						L
334						1		V
346						1	1	G
354	1	1						А
361			1				1	H, Y
363	1	4						S, F
364		5						М
368						2	1	W, G
383		1						V
384	1							Т
387		1						А
389		1						G
396			1	1	1	1		М
397	1		1	1	1	1		I, F
399	1							L
408			1		2			R
410		1						G
417				1			3	Ν
432	3	3						G
434			1		1			D
437		1						R
456		2						М
473	1	1			1			G

Table 4.3 Summary of residue positions predicted to contribute to neutralization changes through pairwise comparisons.

478				1		1	1	G
492							1	R
496			2	2	2	1	2	Ι
520				1			1	D, N
538	1	1						V
555	1	1						М
558					1		1	S, T
570		1						G
576					1			D
591			1				1	E
593	1	3						Т
629			1				1	Ι
653			1	1	1	1		D, G
687	1	1						S
709	1	1						V
710	1							А
733	1	4						S
742	1	2						F, V
743		1						Q
744	1	2						Α

<sup>1</sup>Counts indicate the number of pairwise comparisons containing a sequence change at this residue position with fold change >5 for a given antibody. Residue positions that were the only sequence change in one or more pairwise comparisons are highlighted in bold.

### 4.3.4 Polymorphisms predicted to impact antibody neutralization modeled using computational

### mutagenesis

Direct effects to antibody neutralization, either through altered co-receptor binding or mutation of an antibody epitope, are suggested for some of the 80 residue positions predicted across SNAPR and pairwise comparisons such as residue positions in domain B, D, or E epitopes. However, this diverse set of polymorphisms may be more likely to impact antibody neutralization through indirect effects, with polymorphisms that could change glycoprotein conformations, dynamics, or stability, which would cause a change in free energy. To examine the potential structural effects of predicted polymorphisms, several HCV glycoprotein structures were used to introduce polymorphisms computationally with the protein modeling software Rosetta (364, 365). One structure of E1 (363) and structures of antibody-bound E2 with sequences from three different subtypes (53, 244, 251) were used to model E1E2 polymorphisms and predict the change in free energy ( $\Delta\Delta$ G) to the monomeric protein. About 26% of predicted polymorphisms in regions of E1 (residue ranges 246-257, 277-383) and E2 (residue ranges 384-404, 646-746) could not be introduced computationally because they were either not included or disordered in the crystal structures used for modeling. Polymorphisms introduced by computational mutagenesis were classified as destabilizing, stabilizing, or neutral based on the predicted  $\Delta\Delta$ G value, investigating whether a particular polymorphism is predicted to shift glycoprotein stability in a way that may impact antibody neutralization.

Most modeled polymorphisms were classified as neutral, with minimal predicted effect on glycoprotein stability (~61%; **Table 4.4**, **Table 4.5**). If a predicted resistant residue was already present, that residue was mutated to one or more residues predicted to be more sensitive to neutralization. Overall, 12 mutations were classified as stabilizing and 31 mutations were classified as destabilizing, with some overlap of these classifications between structures (**Table 4.6**). One polymorphism in E1 was stabilizing (A241P) and another two were destabilizing (R195H, R260H), forming a small group of exposed residues (**Table 4.4**). The predicted destabilizing effect of R260H was the most intriguing, given the predicted contributions to neutralization changes for a broad set of antibodies. In E2, just 22 residue positions showed polymorphisms that were predicted to be destabilizing or stabilizing (**Table 4.5**). Some polymorphisms showed similar  $\Delta\Delta G$  values in several structures, suggesting that these polymorphisms might shift the stability of E2 regardless of genotype.

Polymorphisms with the most consistent stabilizing or destabilizing effects were within exposed residues of domain B or D epitopes targeted by bnAbs (A531E, W437R) or in the vicinity

of those same epitopes (1538V, M555T). E531 was predicted both as resistant and stabilizing, or destabilizing when mutated to a sensitive polymorphism, an effect that may be localized to the domain B epitope but affect dynamics of both bnAb binding and CD81 co-receptor binding, in which the domain B epitope switches to an elongated and open state (53). In a similar fashion, R437 may shift domain D conformational dynamics in ways that directly or indirectly impact antibody neutralization and co-receptor binding, an effect that could mimic the effects of polymorphisms found in residue 438 (103). Polymorphisms V538 and T555 were both classified as destabilizing, but V538 was predicted as resistant while T555 was predicted as sensitive, suggesting that lowering glycoprotein stability may have dramatically different effects on antibody neutralization or viral fitness depending on residue location. Residue position 555 is closer to the hydrophobic core of the E2 ectodomain than 538, and destabilization of this core may be more likely to disrupt glycoprotein stability in a way that decreases viral fitness and increases sensitivity to antibody neutralization.

Through these predictions of protein stability changes, the predicted effects of some polymorphisms on antibody neutralization could have mechanisms of indirectly altering the conformational dynamics of E1E2 glycoproteins, warranting further investigation. However, major caveats for this analysis should be noted. Even if  $\Delta\Delta G$  predictions made by Rosetta were consistently accurate, changes in stability only apply to partial E1 and E2 ectodomains that may contain structural inaccuracies, revealing little about the native E1E2 assembly that interacts with antibodies. Furthermore, no explicit modeling of glycoprotein flexibility or dynamics was conducted after mutations were introduced, nor were glycans included during  $\Delta\Delta G$  prediction, drastically limiting the ability of this analysis to assess possible mechanisms of neutralization changes.

Residue	$4UOI mut(s)^2$	$\Delta\Delta G$
position <sup>1</sup>		(kcal/mol)
195	R195H	1.6
198	S198T	0
199	G199R	0.6
208	P208S	0.6
219	A219S	0.2
	A219T	-0.1
234	Q234N	-0.1
	Q234D	0.1
241	A241P	-1.1
242	V242L	0.6
	V242M	0.2
	V242I	0.3
260	R260H	1.7

Table 4.4  $\Delta\Delta G$  predictions of changes in E1 monomer stability.

<sup>1</sup>E1 residue position according to H77 numbering <sup>2</sup>Mutation introduced in E1 N-terminal domain structure (363)

Table 4.5  $\Delta\Delta G$  predictions of changes in E2 monomer stability.

Residue	4MWF	$\Delta\Delta G$	6MEJ	$\Delta\Delta G$	7MWW	$\Delta\Delta G$
position <sup>1</sup>	$mut(s)^2$	(kcal/mol)	$mut(s)^3$	(kcal/mol)	$mut(s)^4$	(kcal/mol)
405			P405R	-0.9		
408			K408R	-0.4		
410			N410G	2.3		
411			I411V	0.1		
416			T416S	0.1		
			T416K	-0.1		
417			N417S	-0.1		
432	S432G	2	S432G	0.6	S432G	0.4
434	N434D	0.4	N434D	0.4	H434D	-0.4
437	W437R	2.2	W437R	2.4	F437R	0.8
438	L438V	0.9	L438V	1.3	I438V	0.2
442	F442I	0.6	F442I	1.5	F442I	0
446	K446R	-0.3	K446R	-0.8	S446R	-2.4
456			M456V	1.1		
460			R460K	0.6		
461			P461R	0		
			P461L	0.9		
463			T463A	0.1		
466			D466A	1		
			D466N	0		
473			S473G	0		

474			H474Y	0.4		
475			A475T	-0.6		
477			G477T	5.2		
478			S478G	0.8		
492	R492K	0.3	R492K	0.3	R494K	1.3
	R492Q	0.2	R492Q	0.1	R494Q	1
496	I496V	0.3	I496V	-0.6	V498I	0.1
500	K500L	1.4	K500L	0.6	K502L	1
	K500Q	0.3	K500Q	-0.3	K502Q	0
	K500R	0.3	K500R	-1.3	K502R	-0.7
520	D520N	0.3	D520N	-0.4	D522N	-0.1
524	A524V	-0.2	A524V	-0.3	A526V	-0.6
531	A531E	-1.4	E531A	1.4	E533A	1.5
			D531D	0.4	E533D	-0.1
538	V538I	-1	V538I	-0.9	L540V	1.4
555	M555T	2	M555T	2.1	M557T	2.7
558	T558S	0	T558S	0.2	S560T	-0.4
570	V570G	1.1	A570G	1		
576	D576A	1.4				
580	L580R	-0.4	Y580L	0.2		
	L580Y	0.3	Y580H	0.6		
591			E591G	1.6		
593			T593S	0.4		
603	I603L	0.4	I603L	0.2	L607I	0.9
	I603V	0.8	I603V	0.6		
610	D610H	0.2	D610H	0.2	D614H	-0.7
629	V629I	-0.5	I629V	0.6	I633V	0.2
029	V 0291	-0.5	1029 V	0.0	1033 V	0.2

<sup>1</sup>E2 residue position according to H77 numbering
<sup>2</sup>Mutation introduced in E2 ectodomain structure, genotype 1a (251)
<sup>3</sup>Mutation introduced in E2 ectodomain structure, genotype 1b (244)
<sup>4</sup>Mutation introduced in E2 ectodomain structure, genotype 2a (53)

Structure	Total	Stabilizing	Residue	Destabilizing	Residue
	mutations	mutations	positions <sup>1</sup>	mutations	positions
4UOI (E1)	13	1	241	2	195, 260
4MWF (E2 1a)	26	2	531, 538	8	432, 437, 438,
					500, 555, 570,
					576, 603
6MEJ (E2 1b)	47	4	405, 446,	13	410, 437, 438,
			500, 538		442, 456, 461,
					466, 477, 478,
					531, 555, <mark>570</mark> ,
					591
7MWW (E2 2a)	22	3	446, 500,	8	437, 492, 500,
			610		531, <mark>538</mark> , 555,
					603

Table 4.6 Computational mutagenesis summary and classification for each structure.

<sup>1</sup>Residue positions colored blue had a resistant residue present in the structure, and one or more polymorphisms predicted to be sensitive were modeled. Residue positions colored magenta had a sensitive residue present in the structure, and one or more polymorphisms predicted to be resistant were modeled. All other mutations were classified as neutral.

### 4.3.5 Hierarchical clustering of E1E2 mutagenesis datasets

As new HCV bnAbs have been isolated and characterized, additional research has been conducted to thoroughly delineate and compare antibody epitopes. A step beyond broad classifications of antibody epitopes, this research uses alanine scanning mutagenesis to reveal how a mutation at each E1E2 residue affects a diverse set of antibodies. Two studies have obtained this data in recent years, offering an in-depth look into disruption of bnAb responses. Gopal and Jackson et al. performed high-throughput mutagenesis on E1E2 from the H77 reference isolate to assess binding changes of 13 antibodies and CD81-LEL, a dataset that includes anti-E1, anti-E2, and anti-E1E2 antibodies (87). Another study provided mutagenesis data with the same isolate to a panel of 21 antibodies from antigenic domains B-E that does not overlap with antibodies tested in Gopal and Jackson et al. (88). Since the mutagenesis datasets from different studies both reported relative antibody binding of mutants, these datasets could be combined and analyzed as a

larger set of mutagenesis experiments. The combined datasets were analyzed for similar patterns of disruption to antibody binding through hierarchical clustering implemented in R. Mutagenesis data were clustered both by antibody and by mutated residue to examine how they separated into discernible groups. Hierarchical clustering by antibody helped to visualize 10 groups that also matched previously defined annotations of antibodies by antigenic domain and antigenic region (AR), corresponding to unique epitopes (**Figure 4.2**). Hierarchical clustering found groups of antibodies in domain A (e.g. CBH-4B), domain B (e.g. AR3A), domain C (e.g. CBH-7), domain D (e.g. HC84.26.WH.5DL), domain E (e.g. AP33), AR2 (AR2A), AR4-5 (AR4A, AR5A), distinct E1 antibodies (A4, IGH526), and the CD81 co-receptor (CD81-LEL). Clustering these antibodies by mutated residue and visualizing in a heatmap also showed distinct patterns of relative antibody binding that reflected unique antibody epitopes, allowing this analysis to identify how residues were separated (**Figure 4.3**).



**Figure 4.2 Visualization of antibody groups using hierarchical clustering.** Antibodies in merged mutagenesis dataset were clustered through data on relative binding of E1E2 alanine mutants. Antibodies generally clustered by antigenic domain or region, and groups of antibodies that clustered together were given distinct colors for clarity. Antibody groups are colored as follows: antigenic domain A antibodies are in gray, AR2 antibody AR2A is in purple, AR4-5 antibodies are in dark green, antigenic domain C antibodies are in cyan, antigenic domain E antibodies are in blue, E1 antibody A4 is in red, E1 antibody IGH526 is in light pink, AR3/antigenic domain B antibodies are in magenta, antigenic domain B/D antibodies are in orange, and co-receptor CD81-LEL is in light green. HC84.26.WH.5DL is abbreviated to HC84.26 in this figure. CD81-LEL was placed in its own group and colored accordingly to reflect its unique role as a co-receptor, despite clustering with an antibody group.



**Figure 4.3 Heatmap of merged mutagenesis dataset clustered by residue.** Antibodies are labeled at the bottom of the heatmap and colored according to antibody groups visualized through hierarchical clustering. Antibody groups are colored as follows: antigenic domain A antibodies are in gray, AR2 antibody AR2A is in purple, AR4-5 antibodies are in dark green, antigenic domain C antibodies are in cyan, antigenic domain E antibodies are in blue, E1 antibody A4 is in red, E1 antibody IGH526 is in light pink, AR3/antigenic domain B antibodies are in magenta, antigenic domain B/D antibodies are in orange, and co-receptor CD81-LEL is in light green. Relative binding values for each mutant to wild-type antibody binding are colored as follows: 0-20% in red, 21-40% in orange, 41-60% in yellow, 61-90% in white, 91-150% in green, and >150% in blue.

### 4.3.5 Critical E1E2 interface residues predicted through clustering by residue

Following hierarchical clustering of mutagenesis data, we determined clusters of residues separated by patterns of disruption to antibody binding and calculated the average relative binding of each pre-defined group for every cluster. This analysis was performed on the merged dataset, which includes data on E2 residues from both datasets and was separated into 16 clusters (Table 4.7), with some clusters overlapping with the clusters found in a previous analysis of E2 mutagenesis data (63). The Gopal and Jackson et al. dataset that included residues from both E1 and E2 was separated into 30 clusters (**Table 4.8**). In examining these results, we found that two defined clusters of mutated residues in E2 merged data (clusters 10 and 14; Table 4.7) showed low relative binding to anti-E1E2 antibodies AR4A and AR5A, and medium to high relative binding to virtually all other antibody groups and CD81. Since AR4A and AR5A epitopes involve E1 and E2 residues, and the native conformation of the E1E2 heterodimer is required, the residues in these clusters may disrupt E1E2 heterodimerization while leaving other antigenic domains largely unaffected. Residues in cluster 14 also showed low relative binding to domain A antibodies, which are non-neutralizing and recognize an epitope on the back layer of E2 (86), but this effect was deemed not to change the predictions of these residues as putative E1E2 heterodimer interface residues. Other clusters in this combined dataset showed low relative binding to anti-E1E2 antibodies, but were not considered as putative E1E2 heterodimer interface residues due to low relative binding of bnAbs to conformational epitopes and CD81, suggesting that these mutants disrupt native E2 conformations. Clusters from the E1E2 dataset had a similar pattern of relative binding, with two clusters showing low relative binding to anti-E1E2 antibodies and relative binding comparable to wild-type E1E2 for all other groups (clusters 7 and 8; Table 4.8). From these pairs of clusters, 17 E1E2 residues were predicted as E1E2 heterodimer interface residues based on mutagenesis data (**Table 4.9**), including seven residues in the E1 N-terminal domain and ten E2 residues ranging from HVR2 to the E2 stem region. Twelve cysteines from E1 and E2 were also included in the identified clusters, but were excluded from this list of predictions, as it appeared unlikely for these residues to be directly involved in the E1E2 interface. Unsurprisingly, some of the predicted E1E2 interface residues are part of the AR4A or AR5A antibody epitope (247), but a majority of predictions have not been identified as anti-E1E2 binding residues, though it should be noted that some predictions have been previously identified as probable binding residues for anti-E1E2 antibodies (243). At the same time, other residues in **Table 4.9** have been identified as determinants of native E1E2 assembly or infectivity in previous studies, offering support for these predictions as residues critical for the E1E2 heterodimeric interface.

Cluster	CD81-	AR1/		AR3/				Dom	Dom	Dom
number <sup>1</sup>	LEL <sup>2</sup>	Dom C <sup>3</sup>	AR2	Dom B	AR4-5	IGH526	A4	E	B/D	Α
1	71	98	94	86	91	90	92	87	95	100
2	78	101	113	116	111	96	85	7	132	105
3	18	115	95	44	106	98	93	89	65	130
4	3	116	96	10	108	102	95	88	50	106
5	2	126	93	1	121	107	88	118	5	77
6	1	121	104	10	144	91	96	93	50	110
7	61	241	98	100	98	88	95	91	112	202
8	39	73	65	54	55	85	82	76	73	46
9	8	56	48	21	40	79	90	88	38	27
10	115	90	130	89	15	74	97	75	77	143
11	2	2	21	1	2	76	98	85	1	1
12	3	18	16	3	10	76	90	80	10	5
13	3	22	21	2	5	94	104	84	6	125
14	80	69	113	74	12	73	87	76	72	15
15	72	87	107	71	52	83	85	76	78	6
16	16	32	45	15	13	85	81	72	30	3

Table 4.7 Relative binding averages for residue clusters from merged E2 mutagenesis data.

<sup>1</sup>Cluster numbers containing putative E1E2 interface residues are in bold.

<sup>2</sup>Average relative binding values for each mutant to wild-type antibody binding are colored as follows: 0-20% in red, 21-40% in orange, 41-60% in yellow, 61-90% in white, 91-150% in green, and >150% in blue.

<sup>3</sup>Antigenic domains are identified as abbreviations "Dom" A-E.

Table 4.8 Relative binding averages for residue clusters from E1E2 mutagenesis data.

Cluster	CD81-							
number <sup>1</sup>	LEL <sup>2</sup>	AR1	AR2	AR3	AR4-5	IGH526	A4	Dom E <sup>3</sup>
1	92	99	105	104	118	104	101	91
2	64	85	91	89	88	98	77	53
3	73	92	97	91	93	89	99	89
4	79	97	105	89	43	77	90	79
5	53	84	82	78	72	81	85	73
6	42	66	60	67	65	76	58	45
7	91	99	114	94	14	74	95	80
8	89	92	115	94	3	70	99	82
9	71	87	96	92	101	16	82	68
10	34	79	73	57	86	85	91	82
11	78	95	113	116	111	96	85	3
12	3	94	97	72	81	88	82	75
13	12	70	79	28	79	100	105	96
14	3	104	101	6	111	107	95	85
15	1	114	99	1	135	98	93	101
16	5	102	100	23	124	97	96	82
17	3	90	97	16	98	96	104	98
18	34	64	75	42	18	72	88	73
19	105	85	96	77	62	88	102	104
20	2	2	27	1	2	80	99	85
21	24	41	48	38	49	82	79	67
22	3	12	24	3	15	77	88	88
23	4	15	16	6	13	76	83	76
24	12	23	37	22	24	82	89	75
25	25	29	67	41	71	83	85	85
26	73	14	129	80	23	87	83	65
27	55	6	116	83	48	66	103	95
28	2	8	3	2	6	73	77	73
29	0	1	0	1	0	73	96	98
30	3	38	0	14	38	74	126	94

<sup>1</sup>Cluster numbers containing putative E1E2 interface residues are in bold.

<sup>2</sup>Average relative binding values for each mutant to wild-type antibody binding are colored as follows: 0-20% in red, 21-40% in orange, 41-60% in yellow, 61-90% in white, and 91-150% in green. <sup>3</sup>Antigenic domain E.

Residue	Location	AR4-5 hotspot residue? <sup>1</sup>	Noted effects on assembly
			or infectivity? <sup>2</sup>
Y201	E1	Yes	No
N205	E1	Yes	No
I212	E1	No	Yes (234)
I220	E1	No	No
H222	E1	No	Yes (234)
P228	E1	No	No
W239	E1	No	Yes (234)
W487	E2	Yes	Yes (298)
R543	E2	No	Yes (298)
D584	E2	No	No
F586	E2	No	No
Y594	E2	No	No
R657	E2	Yes	No
D658	E2	Yes	No
F679	E2	No	No
L692	E2	Yes	Yes (357)
D698	E2	Yes	No

### Table 4.9 Summary of predicted E1E2 interface residues.

<sup>1</sup>Epitope residues based on documented critical residues for AR4A and AR5A binding based on alanine scanning (247).

<sup>2</sup>Specific references for previous research on a given residue are listed in the table.

### 4.3.6 Predicted E1E2 contacts found in E1E2 heterodimer structure

Until recently, the E1E2 heterodimer structure had only been predicted in several models (373, 374) that had been compared in a previous review (57). With a set of residues predicted to be critical for the E1E2 heterodimeric interface, we examined the report of an E1E2 heterodimer structure, which provided some insights about residues in the E1E2 interface (61). Though the coordinates of this cryo-EM structure have not yet been released, images of the full heterodimer and portions of the interface in figures were compared to our set of predictions. As described in this preprint, the E1 N-terminus and stem domain contact the E2 stem domain and back layer, forming a discontinuous and conformationally sensitive interface. The E1 N-terminus and E2 back layer form multiple contacts, with two predicted contacts from E1 (residues 201, 205) and from

E2 (residues 594, 679) highlighted as interacting residues in this portion of the interface. This overlap with predicted E1E2 interface residues suggests that the analysis of mutagenesis data led to some accurate determinations of critical heterodimer residues. However, other residue regions within the reported E1E2 interface such as the E1 stem domain were not predicted to be critical based on mutagenesis data, showing that aspects of the E1E2 interface may not have been captured by data on relative antibody binding. Other predicted E1E2 interface residues may also be near the heterodimeric interface, but it is difficult to assess the position of these residues solely through the figures in the initial report. Interestingly, residue positions with polymorphisms predicted to contribute to neutralization changes are also highlighted in figures of the E1E2 interface, with residues in both E1 (positions 199, 308) and E2 (position 709). Though we are currently unable to assess the impact of predicted polymorphisms on E1E2 heterodimerization, the presence of these residues in the interface suggests that modulation of E1E2 assembly could contribute to changes in antibody neutralization. Once the coordinates of the E1E2 structure are fully accessible (61), possible impacts on heterodimerization and antibody neutralization can be investigated in earnest.

### 4.4 Discussion

In this study, we utilized multiple published datasets to analyze HCV E1E2 sequences and predict contributors to changes in antibody neutralization and interface residues in E1E2 heterodimeric assembly. These predictions were generated from diverse sets of antibody neutralization data, which included E1E2 sequences from multiple genotypes, and from mutagenesis data, which included a comprehensive and diverse set of HCV antibodies. We predicted polymorphisms contributing to neutralization changes using two methods: a modified SNAPR script that assessed neutralization differences of aligned E1E2 sequences, and pairwise comparisons of highly similar sequences from the Bailey or Ball datasets. Predicted

polymorphisms from both methods were modeled on existing E1 and E2 structures, examining the possible effects of these polymorphisms on glycoprotein stability. At the same time, we predicted determinants of E1E2 assembly through hierarchical clustering of E1E2 and merged E2 mutagenesis datasets. Two clusters were identified as possible E1E2 interface residues, as the residues in these clusters primarily disrupted anti-E1E2 antibody binding while largely preserving antibody binding to other antigenic domains and CD81-LEL. In a recently reported E1E2 structure, four residues predicted by this analysis were highlighted as residues in the heterodimeric interface. Both analyses provide a broader set of predictions for key elements of E1E2 characterization and could be useful for finding crucial residue determinants in E1E2 that influence HCV vaccine design.

Although the predictions presented and assessed in this study were generated using large and diverse datasets, both types of data have their limitations. The genotypic diversity in the Bailey dataset is only contained within two genotype 1 subtypes, and the genotypic diversity in the Ball dataset is dramatically overrepresented by genotype 1. In both cases, only genotype 1 sequences had enough data to conduct SNAPR predictions. Restricting these predictions by genotype could make claims of breadth for contributions to antibody neutralization changes less reliable, as the observed effects on antibody neutralization could theoretically be genotype specific. However, this risk of extrapolating the neutralization effects of polymorphisms to other HCV genotypes could be reduced for extra-epitopic residues that are conserved across genotypes. This possibility is supported by computational mutagenesis that classified polymorphisms as stabilizing or destabilizing using E2 structures from multiple genotypes, suggesting that polymorphisms in key epitopes or conserved regions may have similar effects across genotypes. However, these predicted effects on protein stability would only be applicable to E2, severely limiting the value of this analysis in predicting changes in stability or flexibility of the E1E2 assembly. Lastly, the prediction of polymorphisms contributing to changes in antibody neutralization may also be hindered by the underlying neutralization data. Both datasets tested antibody neutralization using HCVpp of isolates, a pseudovirus system shown to be more susceptible to antibody neutralization than HCVcc (110). Although both systems are suitable for assessing antibody neutralization, the HCVpp system does not incorporate lipoproteins, potentially leaving a mechanism of neutralization resistance untested in combination with observed polymorphisms (108, 109). Future research must validate resistance-associated polymorphisms experimentally as in previous studies (102, 103) by introducing predictions into an array of HCV isolates, ideally from multiple HCV genotypes. This validation will help determine which polymorphisms have the broadest or most potent effect on antibody neutralization, confirming contributions to neutralization change that could aid vaccine design. At the same time, more work is needed to assess mechanisms of increased resistance or sensitivity to antibody neutralization, especially for polymorphisms far from documented antibody epitopes. Future work on these polymorphisms can validate predicted changes in ectodomain stability, and possible changes to flexibility or viral breathing could be examined with molecular dynamics simulations.

Though the combined mutagenesis dataset represents the most comprehensive understanding of the impacts of individual E1E2 residues on antibody binding, this data also has limitations that may complicate predictions. The dataset only included mutated residues that were tested with every antibody, which missed some residues in HVR1 and most of the E1 and E2 TMDs. E1E2 TMDs have been implicated in E1E2 assembly in multiple studies (51, 375), but the absence of these residues from mutagenesis data made us unable to predict the involvement of TMD residues in E1E2 assembly. In addition, residues that were predicted through disruptions to antibody binding naturally overlap with known epitope residues for anti-E1E2 antibodies without being able to fully distinguish between residues that only affect antibody binding and residues that also affect E1E2 assembly. While it is possible that solvent exposed E1E2 interface residues are critical both for an anti-E1E2 epitope and for E1E2 heterodimerization, this dual role is difficult to elucidate through mutagenesis data alone and would require additional research, starting with an in-depth analysis of the heterodimer in complex with an anti-E1E2 antibody. Lastly, mutagenesis data may be influenced by the sequence of, and antibody responses to, genotype 1 sequence H77. Since the H77 isolate is more sensitive to antibody neutralization than most other isolates (321), it is unclear if the degree of decreased antibody binding to alanine mutants in a more resistant isolate would be more pronounced or display different patterns for certain classifications of antibodies. Mutagenesis data on a different HCV isolate could offer clarity on conserved shifts in antibody binding induced by alanine scanning, elucidating which antibody binding impacts may be pan-genotypic and which may be genotype-specific. Additional experimental work should be conducted on predicted E1E2 interface residues, including a study of how alanine mutations introduced experimentally may affect E1E2 assembly.

Chapter 5: An expanded benchmark for antibody-antigen docking and affinity prediction reveals insights into antibody recognition determinants

### Abstract

Accurate predictive modeling of antibody-antigen complex structures and structure-based antibody design remain major challenges in computational biology, with implications for biotherapeutics, immunity, and vaccines. Through a systematic search for high-resolution structures of antibody-antigen complexes and unbound antibody and antigen structures, in conjunction with identification of experimentally determined binding affinities, we have assembled a non-redundant set of test cases for antibody-antigen docking and affinity prediction. This benchmark more than doubles the number of antibody-antigen complexes and corresponding affinities available in our previous benchmarks, providing an unprecedented view of the determinants of antibody recognition and insights into molecular flexibility. Initial assessments of docking and affinity prediction tools highlight the challenges posed by this diverse set of cases, which includes camelid nanobodies, therapeutic monoclonal antibodies, and broadly neutralizing antibodies targeting viral glycoproteins. This dataset will enable development of advanced predictive modeling and design methods for this therapeutically relevant class of protein-protein interactions.

### **5.1 Introduction**

Protein-protein interactions are crucial for many biological processes, and structural characterization of those interactions has provided valuable insights into their binding

mechanisms. Immune recognition of antigens by antibodies represents an important and wellcharacterized class of these interactions, with over 3000 antibody structures deposited in the Protein Data Bank (PDB) (153, 336). However, experimental structural characterization is not possible for all interactions due to resource and experimental limitations, particularly in the case of antibody-antigen interactions due to the vast size of immune repertoires and amount of antigen targets. To address these constraints, numerous computational techniques for predictions of protein-protein interactions have been developed.

Computational predictions of protein-protein interactions can address either the likely mode or relative strength of protein binding. Docking algorithms to model protein-protein complexes utilize a variety of search strategies (163), including Fast Fourier Transform (FFT) (165-168) and Monte Carlo searches (171). Additional algorithms conduct refinement of docking models, utilizing increased conformational sampling (376), rescoring (179, 377), or clustering (378, 379) to improve prediction accuracy. A variety of affinity prediction methods have been developed that employ physics-based and knowledge-based potentials (202, 218, 380). Evaluations of algorithm accuracy have included both community-wide prediction experiments with CAPRI (189) and benchmarking with databases of non-redundant protein complexes.

Structures of protein complexes have been curated and assembled in multiple benchmarks for docking and affinity prediction. These include Docking Benchmark 5.0 (BM5) (186), Affinity Benchmark versions 1 and 2 (186, 217), DOCKGROUND (191, 192), protein model benchmarks (381, 382), and a benchmark of homology docking templates (383). These benchmarks include numerous protein-protein interactions validated for analysis, facilitating impartial assessments of docking and affinity predictions to develop, refine, and compare algorithms. Structures of antibody-antigen interactions have been assembled in public repositories (153, 154, 384, 385), and
have been used for interface analysis (386, 387) and predictions of binding energy changes due to mutations (210). A relatively small subset of BM5 cases consists of antibody-antigen interactions (186), making assessments specific to antibody-antigen docking and affinity predictions somewhat difficult.

Antibody-antigen interactions pose a unique challenge for predictive computational modeling, even though antibody recognition of antigens is largely limited to the complementarity determining region (CDR) loops. The therapeutic and disease relevance of antibody-antigen recognition has led to the development of numerous predictive docking algorithms (172, 178, 388, 389) and affinity predictors (209, 213, 215) for this class of complexes. However, the structural flexibility of CDR loops (116, 390, 391) can confound docking and affinity prediction despite the development of methods that integrate flexible backbone modeling (173, 184, 392). The limited success of current algorithms may stem from an incomplete view of antibody structural flexibility. Improving predictions of antibody-antigen interactions will benefit from a larger benchmark that better represents the scope of sequence diversity and structural flexibility used for antigen recognition.

Here we describe an expanded dataset of antibody-antigen complex structures to enable improved docking and affinity predictions for these distinctive protein-protein interactions. After a comprehensive search for antibody-antigen complex structures in the PDB, we identified 41 nonredundant cases in which an antibody-antigen complex and both of its unbound components were experimentally determined. These cases were added to 26 antibody-antigen cases from BM5, increasing the total number of cases to 67. An analysis of conformational changes of antibodyantigen complexes revealed tremendous diversity in the structural flexibility of both antibodies and antigens. Success rates of docking algorithms were examined, highlighting the challenges presented by this diverse benchmark. A variety of affinity predictors were also assessed, resulting in a broad range of correlations with experimental affinity values. These analyses highlight our benchmark as a useful resource for understanding, predicting, and designing antibody-antigen recognition. This benchmark is available set at https://piercelab.ibbr.umd.edu/antibody benchmark/, included in an updated version of a proteindocking benchmark (version 5.5) protein at https://zlab.umassmed.edu/benchmark/benchmark5.5.html, downloadable from and https://github.com/piercelab/antibody benchmark.

#### 5.2 Methods

# 5.2.1 Benchmark assembly

New benchmark cases were identified through searches of experimentally determined structures in the Protein Data Bank (PDB) (336). A list of all antibody-containing structures in the PDB was downloaded from the Structural Antibody Database (SAbDab) database (153) in May 2019. An automated script was then used to perform BLAST (362) searches against all PDB amino acid sequences ("pdb\_seqres.txt", downloaded from the PDB site, May 2019) to identify structures of unbound antibodies and unbound antigens for each antibody-antigen complex from the SAbDab set. Unbound structures with 80% sequence coverage and 98% or greater sequence identity (more stringent than the 96% identity threshold used for BM5) were considered as matches at this stage. As with BM5 (186), only structures within the resolution cutoff ( $\leq 3.25$  Å) were considered, and structures with peptide antigens (< 30 aa) were excluded. As a complementary approach to ensure coverage, we also performed an advanced search through the PDB website, with search terms to filter structures by resolution ( $\leq 3.25$  Å), release date (since 1/1/2015, corresponding to BM5 release (186)), collection method (X-ray diffraction), and the term 'antibody' in the text of

structure descriptions. Cases identified from these automated and manual searches were combined and filtered for non-redundancy with the unbound-unbound antibody-antigen structures from BM5 (186). Manual and automated inspection was then used to assess cases for anomalies warranting exclusion, including missing or mutated interface residues, resulting in elimination of a small set of cases. BM5 cases that were removed from this set include twelve bound-unbound antibody cases to ensure the inclusion of only unbound-unbound test cases in the current set. Cases 1BGX and 1BVK were removed based on manual inspection for structural quality and non-redundancy. Scripts and example input files used to search for and identify cases are in the Github site for the antibody benchmark (https://github.com/piercelab/antibody benchmark). We have also uploaded the output files from the search of the PDB, which include unprocessed redundant sets of complexes and unbound structures ("bound unbound pdbs.txt", "bound unbound pdbs.camelids.txt"). Separate files are available for antibodies and single domain antibodies (sdAbs), each containing bound complex PDB codes along with, when identified, one or more PDB codes corresponding to unbound antibody and antigen structures.

Once verified, new benchmark cases were processed to remove extraneous atoms (e.g. water HETATMs) and to resolve double occupancy atoms. As with BM5 (186), all cases were classified as 'Rigid', 'Medium', or 'Difficult' based on binding conformational changes, as described previously (195). Superposition of unbound onto bound interface residue  $C_{\alpha}$  atoms was used to calculate interface root mean squared distances (I-RMSDs) for antibody interface, antigen interface, and all interface residues, with interface residues defined as those within 10 Å of any non-hydrogen protein atom in the binding partner. Fraction of non-native contacts (f<sub>non-nat</sub>) values were calculated as previously defined (197), using the same superposed unbound structures. Following previously used difficulty classification criteria (186, 193-195), complexes with I-

RMSD > 2.2 Å were classified as Difficult, complexes with I-RMSD < 1.5 Å and  $f_{non-nat} < 0.40$  were classified as Rigid, and all other complexes were classified as Medium.

For experimentally determined affinities, free energy ( $\Delta G$ ) was calculated using the equation  $\Delta G = RTlnK_D$ , where R is the gas constant and T is the temperature in degrees Kelvin. For some cases, equilibrium dissociation constant ( $K_D$ ) values were reported in the literature without a listed experimental temperature or  $\Delta G$  value, and corresponding authors of the respective studies were contacted to request temperature information. Those who provided binding affinity temperatures permitted inclusion of  $\Delta G$  values for those cases. In the event of binding affinity measurements reported for multiple antibody or antigen formats (e.g. antibody Fab, IgG, or scFv), the measured affinity was selected that corresponded to the molecular format present in the structure.

## 5.2.2 Protein-protein docking

ZDOCK is a docking algorithm that uses Fast Fourier Transform (FFT) for rapid rigidbody docking through six-dimensional sampling in translational and rotational space (165, 393). ZDOCK version 3.0.2 with 15° sampling was executed through the command line for docking predictions on all cases. With these parameters, 3600 docking predictions were generated for each complex and ranked by ZDOCK score. As in previous ZDOCK studies (393, 394), antibody framework residues in heavy and light chains were "blocked" to avoid their presence in predicted interfaces. These blocked residue ranges were assigned based on Chothia-numbered antibody structures from SAbDab, and correspond to residue ranges in Chothia numbering (395): Heavy chain: 6-22, 39-46, 81-91, 106-end

Light chain: 7-20, 38-44, 76-85, 101-end

sdAb structures were not subjected to blocking.

ClusPro is a protein-protein docking algorithm that uses FFT-based sampling with the PIPER algorithm (166), followed by clustering of low energy conformations and minimization to avoid steric clashes (396). Unbound structures for each benchmark case were submitted to the online ClusPro server for docking predictions. Docking predictions used antibody mode, which incorporates an asymmetrical statistical potential for improved predictions of antibody-antigen complexes (172). As in the published method, non-CDR regions in antibody structures were masked during docking through selection of that option during ClusPro submission.

SnugDock is an antibody-antigen docking protocol in Rosetta (178) that enables induced fit modeling by modeling CDR loops and interfacial side chains during docking. Additionally, ensembles of structures can be used to allow conformer selection.

We used the FastRelax protocol in Rosetta to relieve potential clashes and find a low energy conformation for the initial unbound crystal structures of antibody and antigen for each complex (397). To take backbone flexibility into consideration, we also used FastRelax to generate an ensemble of 10 decoys, both for the unbound antibody and the unbound antigen. To ensure low-energy starting side-chain conformations, the 10 relaxed antibody and antigen structures for each complex were prepacked apart from each other using the Prepack protocol (171).

The SnugDock protocol was used to dock the antibody ensemble with the antigen ensemble for each case (178). All cases were tested with SnugDock except for 2I25, as the architecture of the IgNAR antibody of that case was not compatible with the SnugDock protocol. To explore the local docking energy funnel, we superimposed the starting structures onto the bound crystal structures. For each docking target, SnugDock started with a random perturbation of 8° rotation and 3 Å translation in each Cartesian direction. We used the new Motif Dock Score (MDS) during the low-resolution docking phase (173). 1000 decoys were generated for each complex. Example command line used in FastRelax protocol: relax.mpi.linuxiccrelease -s 5JMO\_l\_u\_cleaned.pdb -relax:constrain\_relax\_to\_start\_coords -relax:ramp\_constraints true -ex1 -ex2 -use\_input\_sc -flip\_HNQ -no\_optH false -min\_cycles 5000 -nstruct 10

# Example command line used in Prepack protocol:

docking\_prepack\_protocol.mpi.linuxiccrelease -in:file:s 5JMO\_u\_cpx.pdb -ex1 -ex2 -partners H\_A -ensemble1 antibody\_ensemble.list -ensemble2 antigen\_ensemble.list -docking:dock\_rtmin

Example command line used in SnugDock protocol:

```
129
```

```
snugdock.mpi.linuxiccrelease
-s 5JMO u cpx 0001.pdb
-antibody:auto generate kink constraint
-antibody:all atom mode kink constraint
-spin
-dock pert 3 8
-loops:refine outer cycles 2
-loops:max inner cycles 20
-out:path:pdb out pdb
-pdb gz
-detect disulf false
-partners H A
-out:file:scorefile score-snugdock.sf
-nstruct 1000
-docking low res score motif_dock_score
-mh:path:scores BB BB
/work/06525/tg858246/stampede2/motif_dock/score_data_
-mh:score:use ss1 false
-mh:score:use ss2 false
-mh:score:use aa1 true
-mh:score:use aa2 true
```

Docking models were evaluated using Critical Assessment of Predicted Interactions (CAPRI) criteria, which ranks predictions as 'High', 'Medium', 'Acceptable', or incorrect based

on their degree of similarity to the bound complex. This similarity criteria includes I-RMSD, ligand RMSD (L-RMSD), and fraction of native interface residue contacts  $(f_{nat})$  (197).

#### 5.2.3 Interface analysis and affinity prediction

All non-protein HETATMs were removed prior to surface area analysis and scoring. Change in accessible surface area ( $\Delta$ ASA) upon complex formation was calculated using NACCESS v2.1.1 (398), with default probe size (1.4 Å).  $\Delta$ ASA values were negated during  $\Delta$ G prediction to facilitate comparison of correlation with affinity prediction terms. Prior to affinity prediction calculations, complex structures were processed using the "score" application in Rosetta (399) to add missing side chain atoms, remove double occupancies, and ensure consistent atom naming. Affinity prediction methods ZRANK (179) and ZRANK2 (180) were executed from downloadable command line programs. PRODIGY binding affinity predictions were calculated for each complex by the PRODIGY web server (218). ZAPP scores were computed as published, using command-line scripts and programs (380). PYDOCK\_TOT, dDFIRE, PISA, FIREDOCK\_AB, ROSETTADOCK, AP\_T2, CP\_TB, LK\_SOLV, ELE, DDG\_W, HBOND2, VDW, DCOMPLEX scores were computed by the CCharPPI server (400). Function descriptions from CCharPPI are as follows:

PYDOCK\_TOT: Total pyDock energy (377)

dDFIRE: dDFIRE interaction energy (401)

PISA: PISA score (402)

FIREDOCK\_AB: FireDock antibody-antigen function total energy (181)

ROSETTADOCK: Total RosettaDock energy, calculated by PyRosetta (403)

T2: The second atomic two-step potential described by Tobi (404)

CP\_TB: The residue-level interface contact potential from Tobi and Bahar (405)

LK\_SOLV: The effective solvation energy function from Lazaridis and Karplus (406) ELE: Total electrostatic energy, as calculated by pyDock (377) DDG\_W: A weighted atomic potential derived from mutation data (201) HBOND2: Hydrogen bonding potential energy from PyRosetta (403) VDW: van der Waals energy calculated by pyDock (377)

DCOMPLEX: The DCOMPLEX potential (407)

Calculation of Rosetta "REF15" and "beta\_nov16" binding scores was performed using the "score" application in Rosetta (weekly release 2017.52), subtracting the separately scored antibody and antigen components from the score of the bound complex. Prior to REF15 and beta\_nov16 affinity calculation, structures were pre-processed using Rosetta FastRelax ("relax" executable) (397) to perform constrained minimization of the bound complex structures.

#### 5.2.4 Analysis of conformational changes

To assess and compare residue-level and CDR-level binding conformational change, antibody chains were re-numbered in the AHo numbering scheme (408) using the ANARCI software tool (409). This processing only retained re-numbered atoms in the heavy and light chain variable domains, and discarded residues outside of the scope of AHo. The antibody in BM5 case 2125 was not recognized by ANARCI for re-numbering and was excluded from analysis of residuelevel and CDR-level binding conformational changes as it contains a shark IgNAR, which has a non-canonical antibody structure (410). CDR boundaries followed definitions, proposed by North et al. (145) and implemented by PyIgClassify (146), that maximize structural similarities between sdAb, heavy, and light chain loops. Specific CDR definitions for assessing conformational changes are as follows in AHo numbering: CDRH1/sdAb CDR1: 24-42 CDRH2/sdAb CDR2: 57-69 CDRH3/sdAb CDR3: 107-138 CDRL1: 24-42 CDRL2: 57-72 CDRL3: 107-138

The structural alignment software ProFit V3.1 (411) was used to fit unbound and bound antibody variable domains, and calculate RMSDs of C $\alpha$  backbone atoms for individual CDRs and entire variable domains. All antibody-antigen interface contacts within 5 Å were identified to determine the frequency at which each residue forms an interface contact in benchmark complexes. The RMSDs of antibody residues identified as interface contacts were grouped by type of amino acid, allowing for specific comparisons of average conformational changes. Significant differences in conformational changes between one amino acid and all other amino acids were calculated using two-sided Wilcoxon rank sum tests.

# 5.2.5 Quantification and statistical analysis

All statistical analysis was performed using the R program. For statistical analysis of Pearson correlations between scores and binding affinity measurements, p-values were computed using the cor.test() function in R (366), which is based on the t distribution, and "n" corresponds to the total number of affinity values (n = 51) (**Table 5.9**) or a subset thereof as indicated in **Table 5.10**. For Wilcoxon rank sum tests between one amino acid and all other amino acids, p-values were computed using the wilcox.test() function in R.

# 5.3 Results

#### 5.3.1 Benchmark assembly and composition

A comprehensive search of PDB structures identified 41 antibody-antigen complexes with corresponding structures of unbound antibody and antigen chains available that are not present in BM5 (186), increasing the benchmark to 67 cases. The antibody-antigen structures verified as new benchmark cases are summarized in **Table 5.1** (additional details of new cases in **Table 5.2**). Affinity values were found in the literature for 28 new cases (details in **Table 5.3**), increasing the number of benchmark cases with affinities to 51. These structures were released in the PDB as early as 2004 and as recently as May 2019, consisting of complexes that were released after the publication of BM5, newly complemented with unbound structures, or overlooked in previous searches.

This benchmark update expands and diversifies the set of antibody-antigen cases for docking and affinity prediction. **Figure 5.1A** highlights notable additions for antibody and antigen classifications, as well as an increase of cases designated 'Medium' or 'Difficult' for docking algorithms. Twelve new complexes include camelid nanobodies, giving the antibody-antigen benchmark a small, but meaningful, subset of 13 single domain antibodies (sdAbs) that allows for comparisons by antibody type within the benchmark. Underscoring the growing therapeutic relevance of antibodies, several cases in this benchmark include therapeutic monoclonal antibodies (mAbs) in complex with their targets, with 6 new cases in this category, and 12 in total (**Table 5.1**, **Table 5.2**). Viral antigens in new cases include SARS-CoV spike receptor binding domain, Influenza hemagglutinin, human astrovirus 2 (HAstV) spike, vaccinia virus (VACV) D8, and HIV gp120.

Nearly half of the new complexes are classified as Medium or Difficult, based on binding conformational change criteria represented by I-RMSD and fraction of non-native contacts  $(f_{non_nat})$ , which were used for previous docking benchmarks (186), for unbound structures superposed on bound complexes. This update dramatically increases the number of antibody-antigen cases with large conformational changes. In BM5, only six out of 28 unbound antibody-antigen cases (21%) were categorized as Medium or Difficult, a lower percentage than BM5 cases from all other categories of interactions (36%) (186). In this antibody-antigen benchmark, 23 of 67 (34%) cases are classified as Medium or Difficult, resembling the proportion of challenging complexes found among all classes of interactions in BM5 (186).

The addition of 41 new antibody-antigen complexes also expands the breadth and depth of CDR3 lengths, interface sizes, and binding affinities available in the benchmark. Sequence lengths for the CDR3 region vary greatly among these cases, ranging from 6 to 26 residues for mAb cases and 6 to 23 residues for sdAbs (**Table 5.4**). Benchmark sdAbs were found to have longer CDRH3 loops than mAbs, with some significance (p = 0.02) (**Figure 5.1B**). This expanded benchmark encompasses a broad range of binding interface sizes (**Table 5.1**, **Figure 5.1C**) from approximately 1000 Å<sup>2</sup> for an sdAb case (5VNW) to over 2500 Å<sup>2</sup> for several cases, with an average interface size of 1650 Å<sup>2</sup>. sdAbs (open circles in **Figure 5.1C**), while lower on average in interface size, have several cases with buried surface area at or above 1650 Å<sup>2</sup>. The wide range of experimentally determined binding affinities spans many orders of magnitude (**Figure 5.1D**), with an average  $\Delta G$  of approximately -11.7 kcal/mol, corresponding to approximately 2 nM K<sub>D</sub> (25°C), which is commensurate with typical antibody-antigen interaction affinities. The range of  $\mu$ M to pM affinities in this set allows testing of affinity prediction methods for discrimination of high affinity antibody interfaces based on structure.

Complex PDB <sup>1</sup>	Antibody PDB <sup>1</sup>	Antibody <sup>2</sup>	Antigen PDB <sup>1</sup>	Antigen	I-RMSD (Å) <sup>3</sup>	fnon-nat <sup>3</sup>	ΔASA (Ų)	KD (nM)	ΔG (kcal/mol)
Rigid									
1S78_DC: A	1L7I_HL	pertuzumab (Perjeta)	2A91_A	ErbB2	1.13	0.38	2175.1	500	-8.45
3MJ9_HL: A	3MJ8_HL	HL4E10	3MJ6_A	JAML	1.48	0.36	2456.6	8	-11.05
3SE8_HL: G	5JXA_HL	VRC03	3TGT_A	HIV 93TH057 gp120	1.22	0.34	2690.3		
3U7Y_HL: G	3U7W_H L	NIH45-46	3TGT_A	HIV 93TH057 gp120	0.84	0.26	2543.9	160	-9.27
3WD5_HL: ABC	4NYL_H L	adalimumab (Humira)	1TNF_AB C	TNFalpha	0.93	0.31	2328.2	0.115	-13.56
4FP8_HL: A	4FNL_HL	C05	4FNK_A BCDEF	Influenza H3 HA	0.34	0.07	1321.5	430	-8.83
4M5Z_HL: A	4M5Y_H L	5J8	3UYX_A	Influenza H1N1 HA1	0.73	0.21	1591.1	10	-11.10
4Y7M_A:C	4QGY_A	nb25	4Y7L_A	E coli TssM CTD	0.84	0.17	1102.5	1.61	-11.79
5GRJ_HL: A	4NKI_HL	avelumab scFv (Bavencio)	4Z18_A	PD-L1	1.14	0.26	1752.5	0.0421	-14.15
5JMO_D:B	5JMR_A	Nb14	5JXI_A	Furin	0.29	0.08	1393.9		
5014_HL: A	5UR8_AB	1A12	3KVD_D	Neisseria meningitidis fHbp	0.63	0.21	1523.7	0.019	-14.62
501R_HL: A	5NYX_H L	5H2	6CUJ_A	Neisseria meningitidis NHBA	0.9	0.3	1851.9		
5SV3_C:D	5SV4_A	A3C8	1IFT_A	Ricin	0.57	0	1293.6	0.627	-12.55
5WK3_WV :D	5WK2_H L	M116	1NR4_A	CCL17	0.57	0.18	1222.3		
5WUX_HL :EFG	5WUV_H L	certolizumab (Cimzia)	1TNF_CB A	TNFalpha	0.8	0.31	2072.8	0.0274	-14.41
5X0T_AB: E	5X4G_AB	6H8	3B5H_A	CD147	1.25	0.34	1316.2		
5Y9J_HL: ABC	5Y9K_HL	belimumab (Benlysta)	1KXG_A BC	B lymphocyte stimulator	0.96	0.34	1976.3	0.995	-12.28
6A77_HL: A	6A76_HL	B5209B	505I_A	ROBO1 Ig5	0.95	0.22	1532.8		
6B0S_HL: C	6B0W_H L	1710 Fab fragment	3VDJ_A	circumsporozo ite protein aTSR domain	0.72	0.31	1353.2	178	-9.11

Table 5.1 New antibody-antigen benchmark cases organized by difficulty category.

6BPC_EF: D	6BPB_AB	4F7	5W53_A	Plasmodium vivax reticulocyte- binding	0.47	0.19	1644.5		
				protein 2b					
6CWG_B: A	6CWK_A	A9	1IFT_A	Ricin	0.76	0.14	1151.2	0.1	-13.64
6DBG_C:B	6DBA_A	R3-03	1H6T_A	Listeria monocytogene s internalin B	0.46	0.15	1525.2	0.15	-13.4
6OC3_AB: F	60BZ_HL	FluA-20	6CHX_A	H1 hemagglutinin head	0.59	0.29	1536.9		
Medium									
2DD8_HL: S	2G75_AB	m396	2GHV_E	SARS-CoV spike	2.19	0.2	1709.7	20	-10.5
3RJQ_B:A	3R0M_A	A12	3TGR_A	C1086 HIV gp120	0.79	0.81	1734.4		
4ETQ_HL: C	4EBQ_HL	LA5	4E9O_X	vaccinia D8L IMV	0.47	0.41	2277.9	0.18	-13.29
4M3K_B:A	4M3J_A	cAb-H7S	4BLM_A	B. licheniformis beta-lactamase	1.77	0.32	1588		
4POU_B:A	4POY_A	VHHmetal	6ETL_A	bovine RNase A	1.83	0.41	1313.3	157	-9.28
5CBA_AB: E	5C2B_HL	3B4 scFv	4ZAI_A	CXCL13	1.49	0.76	1790.2	0.3715	-13.38
5E5M_B:A	5E03_A	H11	5E56_A	mouse CTLA- 4	1.56	0.43	1341.4		
5HGG_T:A	5HDO_A	Nb4	4FUD_A	uPA	0.84	0.42	1969	0.054	-14.01
5HYS_CD: JK	2XA8_HL	omalizumab (Xolair)	4GT7_AB	IgE-Fc3-4	1	0.47	1331.6	10	-10.91
5KOV_C:A B	5I30_HL	PL-2 scFv	5KOU_A B	astrovirus 2 capsid protein spike domain	1.69	0.65	1735	1.87	-11.91
5VNW_D: A	5VNV_A	Nb.b201	1E78_A	human serum albumin	1.49	0.43	966.8	430	-8.68
5WHK_HL :AB	5WHJ_H L	DX-2507	5BXF_AB	FcRn-B2M	1.88	0.35	1849.6	1.81	-11.93
6A0Z_HL: A	6A0X_AB	13D4	2FK0_A	H5N1 hemagglutinin head	1.28	0.45	1660.7	53	-9.92
6AL0_HL: A	4YNY_A B	NZ-1	6AKQ_A	A. aeolicus site-2 protease homolog with PA tag insertion	1.89	0.5	1622.8		

6EY6_I:AB	5FWO_A	nb130	6EY5_AB	P. gingivalis PorM	1.9	0.49	1806.8	8	-11.05
Difficult									
2FJG_HL: VW	2FJF_HL	G6	4KZN_A B	VEGF	2.51	0.56	1678.2	20	-10.92
4DW2_HL: U	4DVB_A B	mAb-112	4DVA_U	uPA	2.27	0.76	2037.6		
5C7X_HL: A	5D7S_HL	MOR04357	2GMF_A	GM-CSF	2.31	0.65	2523.3	0.007	-15.22

<sup>1</sup>PDB code is followed by chain IDs for antibody and antigen chains; for complexes, antibody chains are

shown first. See also Tables 5.2, 5.3, and 5.4. <sup>2</sup>Trade names for therapeutic antibodies in new benchmark cases are shown in parentheses. <sup>3</sup>Binding interface RMSD (I-RMSD) and fraction of non-native contacts ( $f_{non-nat}$ ), which were used to assign docking difficulty level, were calculated by superposition of unbound antibody and antigen structures onto the bound complex structure using root-mean-square fit of interface residues.

Table 5.2 Additional details for new antibody-antigen test cases.

Complex	Antibody	Antigen	Antigen Source	Antibody Source	Antibody Type	Therapeutic <sup>1</sup>	Associated Disease/Target
1S78	pertuzumab	ErbB2	Human	Human	mAb	Yes (Perjeta)	Cancer
2DD8	m396	SARS spike	Virus	Human	mAb	No	Viral infection (SARS- CoV)
2FJG	G6	VEGF	Human	Human	mAb	No	Cancer
3MJ9	HL4E10	JAML	Mouse	Hamster	mAb	No	None (T-cell stimulation)
3SE8	VRC03	HIV 93TH057 gp120	Virus	Human	mAb	No	Viral infection (HIV)
3U7Y	NIH45-46	HIV 93TH057 gp120	Virus	Human	mAb	No	Viral infection (HIV)
3WD5	adalimumab	TNFalpha	Human	Human	mAb	Yes (Humira)	Arthritis/Skin disorders
4DW2	mAb-112	uPA	Human	Mouse	mAb	No	Thrombosis/Bleeding disorder
4ETQ	LA5	vaccinia D8L IMV	Virus	Mouse	mAb	No	Viral infection (VACV)
4FP8	C05	Influenza H3 HA	Virus	Human	mAb	No	Viral infection (Influenza)
4M5Z	5J8	Influenza HA1	Virus	Human	mAb	No	Viral infection (Influenza)
5C7X	MOR04357	GM-CSF	Human	Human	mAb	No	Arthritis/Autoimmunity
5CBA	3B4 scFv	CXCL13	Human	Human	scFv	No	None (B lymphocyte signaling)
5GRJ	avelumab scFv	PD-L1	Human	Human	scFv	Yes (Bavencio)	Cancer
5HGG	Nb4	uPA	Human	Camel	Camelid/ VHH	No	Thrombosis/Bleeding disorder
5HYS	omalizumab	IgE-Fc3-4	Human	Human (humanized)	mAb	Yes (Xolair)	Asthma
5JMO	Nb14	Furin	Human	Camel	Camelid/ VHH	No	Various (Cancer, Infectious disease)
5KOV	PL-2	astrovirus 2 capsid protein spike domain	Virus	Mouse	scFv	No	Viral infection (HAstV)
5014	1A12	Neisseria meningitidis fHbp	Bacterium	Human	mAb	No	Bacterial infection (Meningitis)

501R	5H2	Neisseria meningitidis NHBA	Bacterium	Human	mAb	No	Bacterial infection (Meningitis)
5SV3	A3C8	Ricin	Plant	Llama	Camelid/ VHH	No	Protein toxicity
5WHK	DX-2507	FcRn-B2M	Human	Human	mAb	No	Autoimmune disorders
5WK3	M116	CCL17	Human	Human	mAb	No	Asthma
5WUX	certolizumab	TNFalpha	Human	Human (humanized)	mAb	Yes (Cimzia)	Autoimmune disorders
5X0T	6H8	CD147	Human	Human (humanized)	mAb	No	Parasitic infection (malaria)
5Ү9Ј	belimumab	B lymphocyte stimulator	Human	Human	mAb	Yes (Benlysta)	Autoimmune disorders
6A77	B5209B	ROBO1	Human	Mouse	mAb	No	Cancer
6A0Z	13D4	hemagglutinin head A. aeolicus site-	Virus	Mouse	mAb	No	Viral infection (Influenza)
6AL0	NZ-1	2 protease homolog with PA tag insertion	Bacterium	Rat	mAb	No	None (protein tagging/structural determination)
6B0S	1710	circumsporozoit e protein aTSR domain Plasmodium	Protozoa	Human	mAb	No	Parasitic infection (malaria)
6BPC	4F7	vivax reticulocyte- binding protein 2b	Protozoa	Mouse	mAb	No	Parasitic infection (malaria)
3RJQ	A12	C1086 HIV gp120	Virus	Llama	Camelid/ VHH	No	Viral infection (HIV)
4M3K	cAb-H7S	B. licheniformis beta-lactamase	Bacterium	Llama	Camelid/ VHH	No	Antibiotic resistance
4POU	VHHmetal	bovine RNase A	Animal	Llama	Camelid/ VHH	No	None (dual-specific metalloproteins)
4Y7M	nb25	E coli TssM CTD	Bacterium	Llama	Camelid/ VHH	No	None (Type VI secretion system)
5E5M	H11	mouse CTLA-4	Animal	Llama	Camelid/ VHH	No	Cancer
5VNW	Nb.b201	human serum albumin	Human	Llama	Camelid/ VHH	No	None (nanobody production)
6CWG	A9	Ricin	Plant	Llama	Camelid/ VHH	No	Protein toxicity
6DBG	R303	Listeria monocytogenes internalin B	Bacterium	Llama	Camelid/ VHH	No	Bacterial infection (Listeria)
6EY6	nb130	P. gingivalis PorM	Bacterium	Llama	Camelid/ VHH	No	Periodonitis
6OC3	FluA-20	HINI Hemagglutinin head	Virus	Human	mAb	No	Viral infection (Influenza)
1AHW	Fab 5g9	Tissue factor	Human	Mouse	mAb	No	Thrombosis/Bleeding disorder
1DQJ	Fab Hyhel63	HEW lysozyme	Chicken	Human	mAb	No	None (antibody-antigen recognition)
1E6J	Fab 13B5	HIV-1 capsid protein p24	Virus	Mouse	mAb	No	Viral infection (HIV)
1JPS	Fab D3H44	Tissue factor	Human	Human (humanized)	mAb	No	Thrombosis/Bleeding disorder
1MLC	Fab44.1	HEW lysozyme	Chicken	Mouse	mAb	No	None (antibody-antigen recognition)
1VFB	Fv D1.3	HEW lysozyme	Chicken	Mouse	scFv	No	None (antibody-antigen recognition)
1WEJ	Fab E8	Cytochrome C	Horse	Mouse	mAb	No	None (antibody-antigen recognition)
2FD6	Fab ATN- 615	Plasminogen activator receptor	Human	Mouse	mAb	No	Cancer

2125	New Antigen Receptor	Lysozyme	Chicken	Nurse shark	IgNAR	No	None (antibody-antigen recognition)
21/16	PBLA8	Flu virus	Viena	Manaa	m A h	No	Viral infection
2115	Fab HC19 Murine	hemagglutinin	virus	Mouse	mAb	NO	(Influenza)
2VXT	reference antibody 125-2H FAB	Interleukin-18	Human	Mouse	mAb	No	Autoimmune disorders
2W9E	ICSM 18 FAB fragment	Prion protein fragment	Human	Mouse	mAb	No	Prion diseases
3EOA	Efalizumab FAB fragment	Integrin alpha-L I domain	Human	Human (humanized)	mAb	Yes (Raptiva)	Autoimmune disorders
3HMX	Ustekinumab FAB	Interleukin-12	Human	Human	mAb	Yes (Stelara)	Autoimmune disorders
3MXW	Anti-Shh 5E1 chimera FAB fragment	Sonic Hedgehog N- terminal domain	Human	Human/Mou se	mAb	No	Cancer
3RVW	4C1 FAB	DER P 1 allergen	Dust mite	Mouse	mAb	No	Asthma
4DN4	CNTO888 FAB	MCP-1	Human	Human	mAb	Yes (Carlumab)	Cancer
4FQI	CR9114 FAB	H5N1 Influenza virus hemagglutinin	Virus	Human	mAb	No	Viral infection (Influenza)
4G6J	Canakinuma b antibody fragment	Interleukin-1 beta	Human	Human	mAb	Yes (Ilaris)	Arthritis/Autoimmunity
4G6M	Gevokizuma b antibody fragment	Interleukin-1 beta	Human	Human (humanized)	mAb	Yes (no brand name)	Arthritis/Autoimmunity
4GXU	1F1 antibody	1918 H1 Hemagglutinin	Virus	Human	mAb	No	Viral infection (Influenza)
3EO1	GC-1008 FAB fragment	Transforming Growth Factor- Beta 3	Human	Human	mAb	Yes (fresolimuma b)	Cancer
3G6D	CNTO607 FAB	Interleukin-13	Human	Human	mAb	No	Asthma
3HI6	AL-57 FAB fragment	Integrin alpha-L I domain	Human	Human	mAb	No	Autoimmune disorders
3L5W	C836 FAB	Interleukin-13	Human	Human (humanized)	mAb	No	Asthma
3V6Z	FAB E6	Capsid protein assembly	Virus	Human	mAb	No	Viral infection (Hepatitis B)

<sup>1</sup>Trade names for therapeutic antibodies in benchmark cases are listed in parentheses.

Table 5.3 Additional details and references for new antibody-antigen affinity cases
---

Complex	K <sub>D</sub> , nM	$\Delta G$ , kcal/mol	Temperature,°C	Method <sup>1</sup>	Reference
1 <b>S</b> 78	500	-8.45	20	BLI	(412)
2DD8	20	-10.50	25	SPR	(413)
2FJG	20	-10.92	37	SPR	(414)
3MJ9	8	-11.05	25	SPR	(415)
3U7Y	160	-9.27	25	SPR	(416)
3WD5	0.115	-13.56	25	SPR	(417)
4ETQ	0.18	-13.29	25	SPR	(418)
4FP8	430	-8.83	30	BLI	(419)
4M5Z	10	-11.10	30	BLI	(420)

5C7X	0.007	-15.22	25	SPR	(421)
5CBA	0.3715	-13.38	37	SPR	(422)
5GRJ	0.0421	-14.15	25	SPR	(423)
5HGG	0.054	-14.01	25	SPR	(424)
5HYS	10	-10.91	25	SPR	(425)
5KOV	1.87	-11.91	25	SPR	(426)
5014	0.019	-14.62	25	SPR	(427)
5SV3	0.627	-12.55	25	SPR	(428)
5WHK	1.81	-11.93	25	SPR	(429)
5WUX	0.0274	-14.41	25	SPR	(430)
5Y9J	0.995	-12.28	25	SPR	(431)
6A0Z	53	-9.92	25	SPR	(432)
6B0S	178	-9.11	22	BLI	(433)
4POU	157	-9.28	25	ITC	(434)
4Y7M	1.61	-11.79	20	SPR	(435)
5VNW	430	-8.68	25	SPR	(436)
6EY6	8	-11.05	25	BLI	(437)
6DBG	0.15	-13.40	25	SPR	(438)
6CWG	0.1	-13.64	25	SPR	(439)

<sup>1</sup>Experimental affinity measurement method used. SPR = surface plasmon resonance, BLI = bio-layer interferometry, ITC = isothermal titration calorimetry.

Case	CDRH1 <sup>1</sup>	CDRH2 <sup>1</sup>	CDRH3 <sup>1</sup>	CDRL1	CDRL2	CDRL3
1AHW	KASGFNIKDYYMH	LIDPENGNTI	ARDNSYYFDY	KASQDIRKYLN	YYATSLAD	LQHGESPYT
1DQJ	SVTGDSVTSDYWS	YISYSGSTY	ASWGGDV	RASQSISNNLH	KYASQSIS	QQSNSWPYT
1E6J	KASGYTFTSYTMH	YINPSSGYSN	SRPVVRLGYNFDY	SASSSVSYMH	YEISKLAS	QQWNYPFT
1JPS	AASGFNIKEYYMH	LIDPEQGNTI	ARDTAAYFDY	RASRDIKSYLN	YYATSLAE	LQHGESPWT
1MLC	KATGYTFSTYWIE	ILPGSGST	ARGDGNYGY	RASQSISNNLH	KYVSQSSS	QQSNSWPRT
1S78	AASGFTFTDYTMD	DVNPNSGGSI	ARNLGPSFYFDY	KASQDVSIGVA	YSASYRYT	QQYYIYPYT
1VFB	TVSGFSLTGYGVN	MIWGDGNTD	ARERDYRLDY	RASGNIHNYLA	YYTTTLAD	QHFWSTPRT
1WEJ	TASGFNIKDTYMH	RIDPASGNTK	AGYDYGNFDY	RASGNIHNYLA	YNAKTLAD	QHFWSTPWT
2DD8	KASGGTFSSYTIS	GITPILGIAN	ARDTVMGGMDV	GGNNIGSKSVH	YDDSDRPS	QVWDSSSDYV
2FD6	KASGYSFTNFYIH	WIFHGSDNTE	ARWGPHWYFDV	SASSSVSYMH	FEISKLAS	QQWNYPFT
2FJG	AASGFTISDYWIH	ITPAGGYT	ARFVFFLPYAMDY	RASQDVSTAVA	YSASFLYS	QQSYTTPPT
2125	VVRDSRCVLSTG	-	KPESRYGSYDAVCAALNDQ			
2VIS	TVSGFLLISNGVH	VIWAGGNTN	ARDFYDYDVFYYAMDY	RSSTGAVTTSNYAN	GGTNNRAP	ALWYSNHWV
2VXT	KASGYSFTDYFIY	DIDPYNGDTS	ARGLRF	RASQDIGSKLY	YATSSLDS	LQYASSPYT
2W9E	KASRNTFTDYNLD	NVYPNNGVTG	ALYYYDVSY	SASSSVSYMH	YDTSKLAS	HQWRSNPYT
3EO1	KASGYTFSSNVIS	GVIPIVDIAN	ASTLGLVLDAMDY	RASQSLGSSYLA	YGASSRAP	QQYADSPIT
3EOA	AASGYSFTGHWMN	MIHPSDSETR	ARGIYFYGTTYFDY	RASKTISKYLA	YSGSTLQS	QQHNEYPLT
3G6D	AASGFTFNSYWIN	GIAYDSSNTL	ARGLGAFHWDMQPDY	SGDNIGGTFVS	YDDNDRPS	GTWDMVTNNV
3HI6	AASGFTFSRYVMW	YIWPSGGNTY	ASSYDFWSNAFDI	RASQSIGSYLN	YAASSLQS	QQSYSTPS
3HMX	KGSGYSFTTYWLG	IMSPVDSDIR	ARRRPGQGYFDF	RASQGISSWLA	YAASSLQS	QQYNIYPYT
3L5W	SFSGFSLSTYGMGVG	HIWWDDVKR	ARMGSDYDVWFDY	RASKSISKYLA	YSGSTLQS	QQHNEYPYT
3MJ9	TVSGISLSDYGVH	IIGHAGGTD	ARHFYTYFDV	SGDKLSDVYVH	YEDNRRPS	QSWDGTNSAWV
3MXW	KGSGYTFIDEALH	VIRPYSGETN	ARDWERGDFFDY	KASQSVSNDLT	YYASNRYT	QQDYGSPPT
3RJQ	TASGRISSSYDMG	AISWSGGTTD	AAKWRPLRYSDYPSNSDYYD			

Table 5.4 Antibody CDR loop sequences of benchmark cases.

3RVW	TVTGYSITSDYAWN	YISYSGTTS	GRTGVYRYPERAPY	KASQDIYSYLS	YRANRLIT	LQYDEFPYT
3SE8	RASGYNFRDYSIH	WIKPLWGAVS	VRRGSCDYCGDFPWQY	KASQGGNAMT	YDTSRRAS	QQFEF
3U7Y	RASGYEFLNCPIN	WLKPRGGAVN	TRGKYCTARDYYNWDFEH	RTSQSGSLA	YSGSTRAA	QQYEF
3V6Z	AASGFTFSSYGMS	TISSGGNYIY	TREGAYSGSSSYPMDY	KSSQSVLYSSNQKNYLA	YWASTRES	HQYLSSYMYT
3WD5	AASGFTFDDYAMH	AITWNSGHID	AKVSYLSTASSLDY	RASQGIRNYLA	YAASTLQS	QRYNRAPYT
4DN4	KASGGTFSSYGIS	GIIPIFGTAN	ARYDGIYGELDF	RASQSVSDAYLA	YDASSRAT	HQYIQLHSFT
4DW2	SASGFTFSRYAMS	SITNGGSTY	ERGELTYAMDY	RASSTVSFHYLH	YATSNLAS	QHYSAYPRT
4ETQ	KASGYSFNFYWMH	MIDPSESESR	TRSNYRYDYFDV	SASSSVSYMY	YDTSNLAS	QQWTSYPLT
4FP8	VGSGSSFGESTLSYYAVS	IINAGGGDID	AKHMSMQQVVSAGWERADLVG DAFDV	QASQDIRKFLN	YDASNLQR	QQYDGLPFT
4FQI	KSSGGTSNNYAIS	GISPIFGSTA	ARHGNYYYYSGMDV	SGSDSNIGRRSVN	YSNDQRPS	AAWDDSLKGAV
4G6J	AASGFTFSVYGMN	IIWYDGDNQY	ARDLRTGPFDY	RASQSIGSSLH	KYASQSFS	HQSSSLPFT
4G6M	SFSGFSLSTSGMGVG	HIWWDGDES	ARNRYDPPWFVD	RASQDISNYLS	YYTSKLHS	LQGKMLPWT
4GXU	AASGFTFSSYAMH	VISYDGRNKY	ARELLMDYYDHIGYSPGPT	SGSSSNIGSYTVN	YSLNQRPS	AAWDDSLSAHVV
4M3K	AASGSISSITTMG	LINSVGDTT	NAFMSTNSGRTGSF			
4M5Z	AVSGYSISSNYYWG	SIYHSGSTY	ARHVRSGYPDTAYYFDK	GGNNIGTKVLH	YDDSDRPS	QVWDISTDQAV
4POU	AASGYPHPYLHMG	AMDSGGGGTL	AAGGYQLRDRTYGH			
4Y7M	AASGFTFEDYAIG	CISNLDGSTY	AAVNAQGIYCTDYIIGPYGMDY			
5C7X	AASGFTFSSYWMN	GIENKYAGGATY	ARGFGTDF	SGDSIGKKYAY	YKKRPS	SAWGDKGMV
5CBA	KASGGTFSSYAIS	GIIPIFGTAN	AREPDYYDSSGYYPIDAFDI	TGTSSDVGAYDWVS	FDVNNRPS	SSYTRRDTYV
5E5M	AASGSTISSVAVG	TSSTSSTTAT	KTGLTN			
5GRJ	AASGFTFSSYIMM	SIYPSGGITF	ARIKLGTVTTVDY	TGTSSDVGGYNYVS	YDVSNRPS	SSYTSSSTRV
5HGG	AASGFTLDSYAIG	CISASGGSTN	AADHPGLCTSESGRRRYLEV			
5HYS	AVSGYSITSGYSWN	SITYDGSTN	ARGSHYFGHWHFAV	RASQSVDYDGDSYMN	YAASYLES	QQSHEDPYT
5JMO	AASGFTFSSYSMY	SINRVGSNTD	AVGMYAAPPW			
5KOV	TVSGFSLIDYGVH	VIWTGGSTD	GRPYYGNVMDY	RASQDISNYLN	YYTSRLHS	QQGNTFPPT
5014	KASGYTFTNYWVV	SIHPRDSDAR	ARLSQVSGWSPWVGP	RASQSISVSLN	YAASRLQS	QETYSDLMYT
501R	TVSGGSVSSGSSYWT	YTSYSGSTK	ARDRFDVASGSSFDF	RASQSISNYLN	YAASSLGS	QQSYGSPT
5SV3	TASGRTLGDYGVA	VISRSTIITD	AVIANPVYATSRNSDDYGH			
5VNW	AASGYISDAYYMG	TITHGTNTY	AVLETRSYSFRY			
5WHK	AASGFTFSEYAMG	SIGSSGGQTK	ARLAIGDSY	TGTGSDVGSYNLVS	YGDSQRPS	ASYAGSGIYV
5WK3	KGSGYSFTSYWIG	IIDPSDSDTR	ARVGPADVWDSFDY	KSSQSVLLSPWNSNQLA	YGASTRES	QQYYLIPST
5WUX	AASGYVFTDYGMN	WINTYIGEPI	ARGYRSYAMDY	KASQNVGTNVA	YSASFLYS	QQYNIYPLT
5X0T	VASGFTFSNFWMN	EIRLKSNNYATH	TSYDYEY	KASENVGTYVS	YGASNRYT	GQSYSYPFT
5Y9J	KASGGTFNNNAIN	GIIPMFGTAK	ARSRDLLLFPHHALSP	QGDSLRSYYAS	YGKNNRPS	SSRDSSGNHWV
6A0Z	KATGYTFSGHWIE	EILPGSGNIH	ARLGTTAVERDWYFDV	KASQNVGTHLA	YSASYRYS	QQYNNFPLT
6A77	AASGFTFSTYDMS	TINSNGGSTY	AREALLRPPYYALDY	GASENIYGALT	YGAINLAD	QNVLSTPFT
6AL0	AASGFTFSNYGMA	SISAGGDKTY	AKTSRVYFDY	KRSTGNIGSNYVN	YRDDKRPD	HSYSSGIV
6B0S	AASGFTFSSYSMN	SITSSSSYIY	ARDPGIAAADNHWFDP	SGDKLGDKYAC	YQDTKRPS	QAWDSSTVV
6BPC	TASGFTFSDYYMA	NINYDGSTPD	ARETVVGSFDY	KASQNVGTNVA	YSASYRYS	QQYNSYPYT
6CWG	AASGRDFSMYMLA	AIMCSGGGGGTY	AASTTYCSATTYSSDRLYDF			
6DBG	AASGHTYSTYCMG	RINVGGSSTW	TLHRFCNTWSLGTLNV			
6EY6	AASGRTFSSYVMG	AISWSGGSIH	VAGFAGYGSFTSRSARDSDKYDY			
60C3	SVSGVSVTSDIYYWT	YIFYNGDTN	ARGTEDLGYCSSGSCPNH	RPSQNIRSFLN	YAASNLQS	QQSYNTPPT

<sup>1</sup>For camelid/single-domain antibodies, CDR loop sequences are in the CDRH1/H2/H3 columns. IgNAR case 2I25 does not have a CDRH2-equivalent loop.



**Figure 5.1 Docking and affinity benchmark composition.** (A) Composition of antibody-antigen benchmark depicted here through classification by antibody type, antigen type, and docking difficulty. Separation of categories by stage of benchmark inclusion (BM5, New) highlights differing proportions for some categories, increasing overall diversity. (B) Averages of CDR3 length for sdAb, heavy, and light chains. Standard deviation for each group is shown in error bars. (C) Interface sizes for antibody-antigen complexes in (Å<sup>2</sup>), with mAb cases shown as black points, sdAb cases as gray points, and mean size as black bar. (D) Experimentally determined binding affinities, with mean shown as black bar. K<sub>D</sub> scale shown on right for reference, corresponding to steady-state affinities for 25°C, which was the most frequently used temperature in reported affinity measurements.

#### 5.3.2 Binding conformational changes

We calculated conformational changes of benchmark complexes at multiple levels, comparing unbound and bound antibody structures by variable domains, antibody CDRs, and individual antibody residues. Binding I-RMSDs, a key metric for predicting docking difficulty, were also calculated for antibody and antigen structures. Binding RMSDs of antibody CDRs and variable domains highlight significant differences in, and extraordinary cases of, structural change of antibodies upon binding within the benchmark (Figure 5.2A). Although most antibody CDRs remain relatively static upon antigen binding (RMSD < 1 Å; 311 out of 360 CDR loops), several Medium and Difficult cases exhibit notable CDR-specific conformational changes (3-7 Å) that pose unusual challenges for docking prediction. Most of these dramatic shifts occur in CDR3 loops, which show the highest RMSDs on average for each antibody chain type. However, unexpectedly large conformational changes were also found in CDR1 and CDR2 loops for both light and heavy chains. Several CDR1 loops in sdAb chains exhibited these shifts, with the largest in 4POU (434) at 3.9 Å (Figure 5.2A). When comparing mAb and sdAb CDR RMSDs, conformational changes trended higher in sdAb chains for all CDRs on average, with sdAb CDR1 chain RMSDs significantly higher than mAb heavy chains (p = 0.006). Figure 5.2B revealed higher antibody I-RMSDs for sdAb cases, showing that this observed difference is not limited to individual CDR loops. Conversely, antigen I-RMSDs of sdAb chains trended lower than mAb chains, though these changes were not statistically significant.

To further investigate patterns in antibody and antigen binding conformational changes, and whether they co-occur or are mutually exclusive, we compared antigen versus antibody I-RMSDs for all test cases (**Figure 5.2C**). While antibody and antigen I-RMSDs were distributed broadly, the cases with larger antibody or antigen I-RMSD values (>2 Å) generally had lower I-RMSD values for the other side of the interface.

Antibody residues in each variable domain were re-numbered, allowing for calculations of average I-RMSD and interface contacts for residues present in mAb heavy, mAb light, or sdAb chains (Figure 5.3). sdAb chains showed the highest average I-RMSDs for many segments of the variable domain, including CDR1, portions of CDR3, and several framework regions (AHo renumbered residues 43-56, 80-90). An example of an sdAb CDR1 with pronounced binding conformational change, from test case 4POU, is shown in Figure 5.4A. Surprisingly, sdAb CDR3 RMSDs show a decrease for residues near the apex of the loop before increasing again, a noticeable deviation from mAb CDRH3 conformational changes. To determine whether certain antibody amino acids are associated with higher or lower conformational changes upon binding, RMSDs of antibody residues near the antigen interface were compared by amino acid type (Figure 5.5). While most residues did not display significant differences in binding RMSDs, glycine and proline exhibited significantly larger conformational changes, whereas tyrosine and tryptophan were associated with smaller conformational changes. These intriguing trends warrant further investigation and could possibly be incorporated into predictive antibody modeling and docking algorithms.

**Figure 5.3** highlights sdAb and mAb residues that form an interface contact in 50% or more of benchmark cases that included the residue. Indicated with black bars in each plot, interface contacts were only found in CDR loops, but the amount and residue positions of these contacts differed by antibody chain type. mAb heavy chains contained 17 interface contacts, including seven in CDRH3. Only 9 interface contacts were found in mAb light chains, largely due to fewer interface contacts in CDRL2 and CDRL3. sdAb chains showed 22 prevalent interface contact residues, with 16 in sdAb CDR3 alone. Significantly longer sdAb CDR3s in the benchmark may help to account for increased sdAb CDR3 contacts, but structural comparisons of mAb and sdAb CDR3s suggest key differences (**Figure 5.4B**). sdAb case 5HGG (424) contains a relatively long CDR3 of 20 aa that folds over the immunoglobulin domain, offering a stark contrast to mAb case 5CBA (440), in which the CDRH3 of identical length extends from the heavy chain. sdAb CDR3 folding in 5HGG is influenced by a disulfide bond between CDR2 and CDR3. Interloop disulfide bonds are relatively common in sdAbs (441), and were found in four sdAb benchmark cases (4Y7M, 5HGG, 6CWG, 6DBG). sdAbs without an interloop disulfide bond showed larger average conformational changes than sdAbs with an interloop disulfide bond in the first eight CDR3 residues (**Table 5.5**). These unexpected CDR3 orientations in sdAb chains help to explain the lower RMSDs shown at the apex of sdAb CDR3 loops when compared to mAb CDRH3 loops.

Newly introduced cases also increased structural diversity among benchmark cases that bind to highly similar or identical antigens, as shown in a superposition of benchmark mAbs on Influenza hemagglutinin (**Figure 5.4C**). This antibody-antigen benchmark update adds four hemagglutinin antigen cases (4FP8, 4M5Z, 6A0Z, 6OC3) to the three cases present in BM5 (2VIS, 4FQI, 4GXU).



**Figure 5.2 Binding conformational changes of antibody-antigen benchmark cases.** (A) RMSDs of CDR loops and variable domains of sdAb, heavy, and light chains. sdAb chains (N=12) were plotted independently of mAb heavy and light chains (N=54). For clarity, RMSD values are capped at 4 Å in this plot, and four CDR RMSDs > 4 Å are not shown. These values occur in sdAb CDR3 (5.0 Å; 4M3K), mAb CDRH2 (6.7 Å; 2FJG), CDRL1 (5.3 Å; 5WHK), and CDRL2 (5.9 Å; 5C7X) loops. Median RMSD for each group is shown as a black bar. (B) I-RMSDs of sdAb and mAb cases by antibody alone, antigen alone, and the entire interface. The plot includes all sdAb (N=13) and mAb (N=54) benchmark cases. Median I-RMSD for each group is shown as a black bar. (C) Antigen versus antibody I-RMSD for all test cases, categorized by benchmark version (BM5, New) and docking difficulty, with sdAbs shown separately (12 of 13 in New).



Figure 5.3 Comparison of residue-level conformational changes by antibody chain type. Average residue-level conformational changes (RMSDs) for (A) light, (B) heavy, and (C) sdAb chains. Red lines represent average RMSD at each residue, re-numbered according to the AHo scheme (408). CDR regions are highlighted with gray shading. Residues were included in the analysis only when present in ~40% or more cases for a given chain type (N  $\geq$  5 for sdAb, N  $\geq$  20 for heavy/light). As a result, residue gaps vary in size between chain types, particularly in CDR3 loops. Plots are aligned numerically for easier comparisons, maintaining the same x-axis except for the last AHo re-numbered residue before the gap in CDR3 numbering (114 for heavy, 111 for light, and 117 for sdAb). Residue number 60 in the light chain is marked with an asterisk because it was not present in any mAb benchmark cases, but was included for consistency with other panels. Black rectangles above the x-axis indicate which AHo re-numbered residue.

Α 4POU interface **Unbound CDR1** Bound CDR1 Antigen 4FP8 С В 2VIS 6A0Z 4GXU **5CBA CDRH3** 5HGG CDR3 6OC3

**Figure 5.4 Structural diversity of benchmark cases.** (A) Visualization of superposed bound and unbound interface of case 4POU (434), with flexibility of CDR1 highlighted. Unbound CDR1 is colored orange, bound CDR1 is colored red, and the rest of the antibody chains are in grey. Unbound and bound antigen chains are colored light blue. (B) Superposed antibody chains for test cases 5CBA (440) and 5HGG (424), with CDR3 loops highlighted. 5CBA CDRH3 is colored light green, while 5HGG CDR3 is colored magenta. The disulfide bond between CDR2 and CDR3 in 5HGG is visualized by yellow sticks. (C) Structures of seven benchmark cases with Influenza hemagglutinin as antigen (419, 420, 432, 442-445). All cases were aligned into the same reference frame using two chains from the 4FQI antigen, which is depicted in green and cyan with transparent surfaces. Heavy and light chains of each antihemagglutinin mAb were given the same color for better visualization of individual antibody structures.



Figure 5.5 Binding conformational changes of antibody residues near the antigen interface by amino acid. (A) Binding conformational changes, calculated for each residue as  $C\alpha$  atom I-RMSDs after superposing unbound and bound Fv domains, are shown with amino acids identified by their three letter codes. Asterisks indicate significant differences between I-RMSDs of that amino acid and the combined I-RMSDs of all other amino acids (\*, p < 0.05; \*\*\*, p < 0.001), calculated by Wilcoxon rank sum test. 81

RMSD values greater than 2.5 Å are not shown, allowing for better visualization of differences in average I-RMSD by residue. (B) The same boxplot as (A), with y-axis scaled to show the outlier I-RMSDs. Due to residue preferences of antibody interface residues, the numbers of I-RMSD values vary substantially among the amino acids. Amino acids with N < 10: CYS, MET. Amino acids with N between 25 and 50: ALA, GLN, GLU, HIS, LYS, PRO, VAL. Amino acids with N between 51 and 100: ARG, ILE, LEU, PHE, TRP. Amino acids with N between 101 and 200: ASN, ASP, GLY, THR. Amino acids with N > 200: SER, TYR.

CDR3	sdAb Average	No disulfide	Disulfide average	No disulfide	Disulfide
residue <sup>1</sup>	RMSD (Å)	average RMSD (Å) <sup>2</sup>	RMSD $(Å)^3$	count <sup>4</sup>	count <sup>5</sup>
107	0.36	0.41	0.26	8	4
108	0.46	0.58	0.23	8	4
109	0.47	0.61	0.18	8	4
110	1.1	1.45	0.41	8	4
111	1.08	1.40	0.51	7	4
112	1.69	2.28	0.65	7	4
113	2.15	3.17	0.62	6	4
114	1.62	2.45	0.58	5	4
115	0.65	0.78	0.56	3	4
116	0.81	0.69	0.93	3	3
117	0.74	0.60	0.83	2	3
130	0.76	0.42	1.10	3	3
131	0.76	0.44	1.08	3	3
132	1.04	0.57	1.40	3	4
133	1.55	2.08	0.90	5	4
134	1.32	1.76	0.67	6	4
135	1.13	1.35	0.75	7	4
136	0.91	1.09	0.58	7	4
137	1.22	1.63	0.39	8	4
138	0.75	0.93	0.40	8	4

#### Table 5.5 sdAb CDR3 average RMSDs for subsets with or without interloop disulfide.

<sup>1</sup>Residue number, by AHo numbering scheme (408).

<sup>2</sup>Average RMSD by CDR3 residue for sdAb structures with no interloop disulfide (N=8).

<sup>3</sup>Average RMSD by CDR3 residue for sdAb structures with interloop disulfide (N=4).

<sup>4</sup>Number of sdAb structures without an interloop disulfide that include the CDR3 residue.

<sup>5</sup>Number of sdAb structures with an interloop disulfide that include the CDR3 residue.

# 5.3.3 Global docking prediction

For an initial assessment of the test cases and their docking difficulty, we performed global protein docking simulations with the unbound structures as input. To assess shared or divergent patterns in docking performance across the benchmark from different methods, we used two global

docking programs: ZDOCK (165, 393) and ClusPro (396). For ClusPro docking, the antibodyantigen potential was selected, as it was reported to improve performance on this class of complexes (172). For both ZDOCK and ClusPro, framework regions of antibodies were blocked or masked during docking to favor models with antibody CDR loops in the interface. The ClusPro server returned between 20 and 30 models for each case, with a median of 30 models, while ZDOCK generated 3600 models per case.

Performance of these docking algorithms on the benchmark is shown in **Figure 5.6**, with detailed results given in **Table 5.6**. As expected, performance on the Rigid cases was higher than for the Medium and Difficult classes of cases. Both algorithms exhibited comparable performance on the benchmark overall, with ClusPro showing moderately higher success rates on the benchmark (34% success for top 10, 45% success for top 30), although ZDOCK produced more Medium accuracy or higher models (22% success for top 30, versus 16% for ClusPro). While some cases, particularly in the Rigid subset, had successful predictions from both methods (e.g. 6DBG, 3MXW), there are also many cases of divergent performance between docking methods. Cases of Medium difficulty, particularly below 5CBA in **Figure 5.6** (corresponding to binding I-RMSD > 1.5 Å), were more challenging, though both ZDOCK and ClusPro produced models of Acceptable accuracy for several cases above this I-RMSD threshold.

To determine docking performance when larger sets of models are considered, success rates were computed for ZDOCK for up to 2000 models per case (Figure 5.7A). In ZDOCK, models of both Acceptable and Medium accuracy increased with the total number of models generated, with success rates of approximately 66% for Medium accuracy models and 90% for Acceptable accuracy models, given that 2000 models were allowed per case. This result suggests the need for model selection and refinement methods to improve the overall rankings of near-

native predictions, which are evidently identified within larger sets of models, but ranked below incorrect predictions. To assess possible differential performance on new cases versus the previous set of BM5 cases, comparison of ZDOCK performance between these subsets of cases is shown in **Figure 5.7B**. Only cases with Rigid docking difficulty were evaluated to avoid bias due to a higher proportion of Medium and Difficult category cases in the new benchmark (**Figure 5.1A**). No major overall difference in docking success was observed between the BM5 and new Rigid cases, though minor reductions of success were seen for docking model ranks < 100.

As antigen glycosylation can be an important factor in antibody-antigen recognition, including antibody recognition of viral glycoproteins, we tested ZDOCK using a subset of benchmark cases with glycans present in the experimentally determined structures of the unbound antigens. ZDOCK results were compared with the same structures with glycans removed (**Table 5.7**), corresponding to the ZDOCK results reported above. The presence of glycans did not markedly affect ZDOCK's performance, though future studies could explore this in more detail on a larger set of cases, for instance using modeled N-glycans (446, 447), and possible improvement of N-glycan parameters in ZDOCK.



**Figure 5.6 Docking performance on the antibody-antigen benchmark.** Global docking methods ZDOCK and ClusPro were assessed for predictive modeling of all 67 antibody-antigen cases, using unbound proteins as input. SnugDock was used to perform local docking perturbations on the same set of

cases (except IgNAR case 2I25) using unbound proteins aligned to the bound antibody-antigen complex as input, to test for the prediction of binding funnels at the native binding site. To compare with SnugDock score, ZRANK2 was also tested for binding funnel prediction through rescoring of the set of all locally perturbed models from SnugDock for each case. The top 10 (T10) and top 30 (T30) ranked models from global docking methods and the top 1 (T1) and top 10 (T10) ranked models from local perturbations were assessed for near-native predictions using CAPRI criteria of High, Medium, or Acceptable accuracy. Cases are classified by docking difficulty and sorted by I-RMSD within each classification. Newly added cases and sdAbs are also indicated. Success rates, calculated as the percent of test cases with near-native predictions from the corresponding method for a given range of top-ranked models, are shown at the bottom, with bars colored according to CAPRI accuracy.

	ClusPro A	ntibody	ZDOCK 3.0.2			
Case	Тор	Тор	Тор	Тор		
Case	Acceptable	Medium	Acceptable	Medium		
1AHW	3	-	210	443		
1DQJ	-	-	460	-		
1E6J	1	-	23	56		
1JPS	3	-	853	972		
1MLC	-	-	30	30		
1VFB	2	2	27	562		
1WEJ	1	1	28	28		
2FD6	19	19	5	18		
2I25	-	-	5	5		
2VIS	-	-	413	732		
4G6J	-	-	329	329		
4G6M	2	-	6	9		
2VXT	4	-	3	3		
3EOA	-	-	164	164		
3RVW	1	-	127	497		
2W9E	19	2	11	108		
3MXW	1	7	1	1		
4DN4	3	-	32	164		
3HMX	2	-	5	5		
4FQI	-	1	123	1569		
4GXU	-	-	3377	-		
3EO1	-	1	349	1255		
3G6D	18	-	54	-		
3HI6	-	-	173	-		
3L5W	8	2	107	175		
3V6Z	-	-	449	-		
1 <b>S</b> 78	-	-	205	-		
3MJ9	-	-	766	2592		
3SE8	-	-	8	8		
3U7Y	-	-	67	67		
3WD5	6	-	228	1062		

Table 5.6 Global docking	ranks of top	Acceptable and	Medium models.
--------------------------	--------------	----------------	----------------

4FP8	21	-	-	-		
4M5Z	1	-	4	78		
4Y7M	-	-	5	32		
5GRJ	-	-	211	-		
5JMO	12	12	13	13		
5014	-	-	29	106		
501R	15	-	308	308		
5SV3	-	-	37	37		
5WK3	17	17	23	28		
5WUX	-	-	248	459		
5X0T	-	-	28	28		
5Y9J	-	-	231	1631		
6A77	10	-	5	5		
6B0S	5	-	2513	-		
6BPC	-	-	1391	-		
6CWG	-	-	23	81		
6DBG	1	4	3	3		
6OC3	-	-	205	424		
2DD8	11	-	199	-		
3RJQ	-	-	416	-		
4M3K	-	-	512	512		
4ETQ	6	-	6	6		
4POU	-	-	12	65		
5CBA	9	-	6	412		
5E5M	12	-	92	-		
5HGG	-	-	5	-		
5HYS	-	-	2923	-		
5KOV	-	-	38	54		
5VNW	-	-	855	-		
5WHK	-	-	186	-		
6A0Z	-	-	3419	-		
6AL0	-	-	-	-		
6EY6	_	-	2568	-		
2FJG	4	-	57	431		
4DW2	-	-	648	-		
5C7X	-	-	1100	-		

Docking model assessments of Acceptable and Medium are based on CAPRI docking criteria (197), and rank given for Acceptable or Medium denotes the highest-ranked docking model with Acceptable or higher accuracy, or Medium or higher accuracy, respectively. Top Acceptable or Medium hits in the top 30 ranked docking models for any case are highlighted in green. "-" indicates no predictions of that accuracy within the full sets of models from the docking method.



**Figure 5.7 Comparison of docking success rates in ZDOCK models.** Success rates from top 1 to top 2000 ranked models from ZDOCK are shown for (A) All cases and (B) Rigid-body cases. (A) Success rates for CAPRI criteria of Acceptable accuracy or higher, and Medium accuracy or higher, are compared. (B) Success rates for Acceptable accuracy or higher models for new rigid-body cases (N=23) in comparison with rigid-body cases from BM5 (N=21).

			No Glycans			With Glycans				
Casa	Catagomy	# Clyconal	#	Rank	#	Rank	#	Rank	#	Rank
Case Ca	Category	# Orycans	$Acc^2$	Acc <sup>3</sup>	Med <sup>2</sup>	Med <sup>3</sup>	$Acc^2$	Acc <sup>3</sup>	Med <sup>2</sup>	Med <sup>3</sup>
4FP8	Rigid	16	0	-	0	-	0	-	0	-
3U7Y	Rigid	10	10	67	3	67	11	34	3	103
1S78	Rigid	3	4	205	0	-	4	251	0	-
3SE8	Rigid	10	9	8	1	8	12	25	0	-
3MJ9	Rigid	3	7	766	1	2592	2	1424	0	-
3RJQ	Medium	8	4	416	0	-	1	368	0	-
5HYS	Medium	2	1	2923	0	-	0	-	0	-
6A0Z	Medium	2	1	3419	0	-	0	-	0	-
2FJG	Difficult	2	13	57	3	431	15	253	2	376

Table 5.7 Comparison of ZDOCK results with or without glycans removed in unbound antigen.

<sup>1</sup>Number of N-glycans present in the unbound antigen PDB file.

<sup>2</sup>Number of Acceptable or higher accuracy (# Acc) or Medium or higher accuracy (# Med) models in the 3600 ZDOCK models for that test case.

<sup>3</sup>Rank of first Acceptable or higher accuracy (Rank Acc) or Medium or higher accuracy (Rank Med) in ZDOCK models for the test case. "-" denotes no models with that accuracy in the ZDOCK models.

#### 5.3.4 Local docking perturbations

In addition to global rigid-body docking simulations, we performed local docking with the SnugDock algorithm (178), to test for the presence of binding energy funnels near the native complexes with the unbound antibody and antigen structures as input (**Figure 5.6**; additional details in **Table 5.8**). This algorithm samples side chains and CDR backbone conformations during the docking search, thus providing a modeling method distinct from global rigid-body docking algorithms we tested. For comparison of binding funnel detection versus the SnugDock Rosetta interface score, all SnugDock models were scored and re-ranked using the ZRANK2 algorithm (180). Unbound structures were superposed onto bound complex structures prior to SnugDock simulations, giving the low initial I-RMSDs for most cases (**Table 5.8**). The goal of these simulations was to detect binding energy funnels and near-native top-ranked models in the context

of local flexible perturbations, and not to improve I-RMSD or accuracy over the input (i.e. superposed unbound structures).

Numerous antibody-antigen cases had highly ranked near-native models in **Figure 5.6**, indicative of binding energy funnels in the SnugDock simulations, while others were more challenging for local docking. SnugDock and ZRANK2 scoring resulted in different performance for several cases, including case 6CWG, where a High accuracy model was ranked in the top 10 for SnugDock, and case 3MXW, where a High accuracy model was ranked #1 by ZRANK2 score. Two representative SnugDock score versus I-RMSD plots, both from the Rigid docking category and with relatively low binding conformational changes, are shown in **Figure 5.8**. Plots of SnugDock score versus I-RMSD for rigid cases (**Figure 5.9**) and medium or difficult cases (**Figure 5.10**) further show the range of success in reaching binding energy funnels. As can be seen in **Figure 5.8**, even Rigid cases can exhibit vastly different binding energy funnels in SnugDock, bearing further investigation; the affinity for the complex in test case 4FP8, with a less distinctive funnel in **Figure 5.8B**, is relatively low, and previous analysis showed that lower affinities may be associated with lower docking success (186).
Case	Input I-	Snugdock	Top Rank	Top Rank	Top Rank	ZRANK2
	$RMSD^1$	I-RMSD <sup>2</sup>	Acceptable <sup>3</sup>	Medium <sup>3</sup>	High <sup>3</sup>	I-RMSD <sup>2</sup>
5JMO	0.39	0.98	1	1	1	1.08
1WEJ	0.58	8.01	4	11	1001	5.52
4FP8	0.57	8.08	3	1001	1001	8.08
3EOA	3.01	5.86	16	1001	1001	11.43
6DBG	0.51	1.11	1	1	1001	2.01
6BPC	1.96	1.23	1	1	1001	13.32
3MXW	0.82	2.59	1	3	3	0.98
3RVW	0.63	7.7	4	5	1001	2.56
1JPS	0.63	2.02	1	1	1001	2.12
4G6M	0.58	1.77	1	1	1001	2.3
5SV3	0.58	2.8	1	2	1001	11.57
5WK3	0.97	2.47	1	1	1001	3.54
6OC3	0.85	13.68	3	13	1001	13.68
1MLC	0.85	2.4	1	1	1001	1.92
4G6J	2	14.37	2	4	1001	2.4
5014	0.78	5.07	4	51	1001	3.81
1AHW	0.94	2.03	1	3	1001	5.35
6B0S	6.33	6.3	1001	1001	1001	10.91
3HMX	2.23	2.13	1	2	1001	3.25
4M5Z	1.27	2.03	1	1	1001	9.5
1DQJ	1.12	14.1	4	36	86	9.07
6CWG	0.79	1.28	1	1	4	5.57
4GXU	1.93	4.6	11	1001	1001	9.14
2VIS	3.75	5.03	1001	1001	1001	11.44
5WUX	0.95	14.84	11	1001	1001	5.15
4DN4	2.52	2.62	1	16	1001	2.41
3U7Y	0.93	6.96	5	35	1001	2.37
4Y7M	0.85	7.78	6	30	30	11.94
501R	1.06	5.38	2	2	1001	1.82
3WD5	1.02	4.66	2	5	1001	4.43
6A77	2.91	8.81	12	251	1001	9.54
5Y9J	1.08	2.54	1	1	1001	14.73
1VFB	1.03	12.4	4	5	1001	12.51
1E6J	1.15	2.35	1	1	1001	2.35
2FD6	1.99	12.69	44	51	1001	13.15
4FQI	2.04	9.95	3	42	1001	2.67
2W9E	1.25	1.88	1	1	1001	7.07
1S78	1.78	3.01	1	4	1001	2.03
5GRJ	1.98	3.35	1	1001	1001	5.86
3SE8	1.36	11.67	16	16	1001	5.53
5X0T	1.37	3.98	1	5	1001	3.33
2VXT	1.48	10.53	2	2	1001	2.41

 Table 5.8 SnugDock local perturbation performance by test case.

4.99	15.9	127	1001	1001	14.69
0.56	1.09	1	1	19	4.15
0.79	2.71	1	1	1001	1.49
0.91	8.4	1001	1001	1001	12.13
0.85	8.23	12	1001	1001	10.15
1.14	15.1	2	1001	1001	9.83
1.36	10.11	5	1001	1001	15.88
1.98	3.22	1	29	1001	3.37
1.53	11.64	7	1001	1001	11.78
2.67	2.89	1	23	1001	9.56
1.54	10.5	4	1001	1001	5.63
1.95	13.17	8	12	1001	8.73
1.82	5.49	22	249	1001	9.12
1.87	8.29	5	20	1001	2.23
1.92	13.38	4	1001	1001	4.79
1.86	8.26	2	59	1001	6.9
4.55	14	8	1001	1001	3.51
1.93	4.35	3	1001	1001	3.61
3.41	8.79	23	1001	1001	4.17
2.22	3.55	1	1001	1001	4.22
4.04	8.38	75	1001	1001	15.49
2.52	10.06	2	1001	1001	8.17
6.83	11.06	1001	1001	1001	6.08
2.67	4.73	2	31	1001	5.17
	$\begin{array}{c} 4.99\\ 0.56\\ 0.79\\ 0.91\\ 0.85\\ 1.14\\ 1.36\\ 1.98\\ 1.53\\ 2.67\\ 1.54\\ 1.95\\ 1.82\\ 1.87\\ 1.92\\ 1.86\\ 4.55\\ 1.93\\ 3.41\\ 2.22\\ 4.04\\ 2.52\\ 6.83\\ 2.67\end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

<sup>1</sup>I-RMSD from bound complex for pose input to SnugDock, prior to minimization and ensemble generation. <sup>2</sup>I-RMSD from bound complex for top-ranked model from SnugDock local perturbation based on

SnugDock's interface score or ZRANK2 score.

<sup>3</sup>Top-ranked SnugDock model at indicated level of CAPRI accuracy, based on SnugDock interface score. As 1000 models were generated, a rank of 1001 denotes no models within the SnugDock set with that CAPRI accuracy. Figure provided by Jing Zhou.



Α

**Figure 5.8 SnugDock binding funnels for two benchmark cases.** Rosetta interface score versus I-RMSD is shown for rigid-body docking cases (A) 5JMO (Nb14 camelid/Furin) and (B) 4FP8 (C05 Fab/Influenza H3 hemagglutinin). Each point represents one of the 1000 SnugDock models generated for each case. Point types and colors represent model accuracy according to CAPRI criteria (197), as detailed in the legend of panel A.



**Figure 5.9 SnugDock binding funnels for Rigid benchmark cases.** Docking perturbations were performed on all cases in SnugDock, generating 1000 models for each case. In each plot, Rosetta interface score versus interface RMSD between model and bound structure is shown. Point types represent model accuracies based on CAPRI criteria (197): Incorrect (gray circles), Acceptable (pink circles), Medium (red squares), High (dark red triangles).



**Figure 5.10 SnugDock binding funnels for Medium and Difficult benchmark cases.** Docking perturbations were performed on all cases in SnugDock, generating 1000 models for each case. For each case, Rosetta interface score versus interface RMSD between model and bound structure is shown. Point types represent model accuracies based on CAPRI criteria (197): Incorrect (gray circles), Acceptable (pink circles), Medium (red squares), High (dark red triangles).

## 5.3.5 Binding affinity prediction

Given the implications for protein interface design, several algorithms have been reported that predict binding affinities based on protein structures, though comprehensive comparative assessments of these methods in the context of varied antibody-antigen interfaces have been limited (210). We tested a variety of affinity prediction functions and interface descriptors to determine performance, as well as similarities between functions, on the set of 51 experimentally determined antibody-antigen affinities (**Figure 5.11**). Bound structures of the cases were used as input to affinity calculations, and all structures were pre-processed using the Rosetta "score" application (399) to ensure consistent atom naming, remove double occupancy atoms, and add missing side chain atoms, without performing side chain packing of intact residues or backbone minimization. Additionally, complexes scored by Rosetta Energy Function 2015 ("REF15") (204) or the "beta16" Rosetta scoring function were pre-processed through constrained minimization using the FastRelax protocol in Rosetta (397).

Resultant correlations between scores and experimental  $\Delta$ Gs varied widely, up to approximately r = 0.46 for REF15. Some relatively simple descriptors, such as  $\Delta$ ASA (r = 0.17) and hydrogen bonding energetics computed by PyRosetta (403) (HBOND2; r = 0.29) suggest that interface size and specific energetic components are key determinants of antibody affinity, and composite functions with weighted combinations of terms, such as ZRANK (179) (r = 0.32), the beta\_nov16 score function (an update of REF15) (r = 0.40), and the antibody-antigen potential of FireDock (448) (FIREDOCK\_AB; r = 0.37), performed comparatively well. Interestingly, two statistical contact potential functions alone also exhibited relatively high predictive performance on this set: TB (405) (r = 0.33) and T2 (404) (r = 0.42). With the exception of  $\Delta$ ASA, all of the correlations with  $\Delta$ G noted above were significant (p < 0.05), ranging from p = 0.04 (HBOND2) to p = 0.0007 (REF15) (**Table 5.9**).

Given the previously noted relationship between interface size and affinity, which was specifically exhibited for complexes with limited binding conformational change (217), we examined whether this relationship is observed for the antibody-antigen complexes in this set, and also evaluated any changes in correlations for the top-performing scoring function (REF15) when stratifying the cases according to I-RMSD (**Figure 5.12**). Unexpectedly, we found that higher correlations with experimentally determined  $\Delta$ Gs were observed for the subset of cases with high binding conformational change (I-RMSD  $\geq 1$  Å), for both  $\Delta$ ASA and REF15. Though several outlier points were observed, it does not appear that these were responsible for this effect, nor were the ranges of  $\Delta G$  values, or number of data points, markedly different between subsets of cases. To examine this effect in more detail, we calculated  $\Delta G$  correlations with  $\Delta ASA$  and REF15 based on stratification using antibody or antigen I-RMSD alone (**Table 5.10**). Based on this analysis, greater antibody conformational change is associated with the higher correlations observed for the cases with high binding conformational change (antibody I-RMSD  $\geq$  1.0), with correlations with experimental  $\Delta Gs$  of r = 0.63 for  $\Delta ASA$  and r = 0.74 for REF15. Despite the lower number of cases (N=15), p-values for both correlations were significant (p = 0.01, 0.002 for  $\Delta ASA$  and REF15, respectively; **Table 5.10**), and more significant than p-values for the corresponding subset identified by overall I-RMSD (p = 0.16, 0.004 for  $\Delta ASA$  and REF15).



Figure 5.11 Affinity predictions on benchmark cases. A set of 20 scoring functions representing protein modeling functions, statistical potentials, and interface descriptors were used to assess the 51 affinity prediction cases, and are compared with experimentally determined  $\Delta$ Gs (kcal/mol). For heatmap generation, scores for each term were scaled and clustered hierarchically to assess similarities between functions; the dendrogram is shown above the heatmap.  $\Delta$ ASA values were negated to facilitate comparison with energetic functions and  $\Delta$ G. Pearson correlation coefficients with experimental  $\Delta$ Gs for all scoring methods are shown at the bottom.

Function/Term <sup>1</sup>	Correlation	p-value
ΔASA	0.18	0.205
REF15	0.46	0.001
beta_nov16	0.41	0.003
ZRANK	0.32	0.021
ZRANK2	0.28	0.047
LISA	0.27	0.059
ZAPP	0.24	0.086
PYDOCK_TOT	0.17	0.236
dDFIRE	0.19	0.193
PISA	0.36	0.010
FIREDOCK_AB	0.37	0.008
ROSETTADOCK	0.24	0.091
T2	0.40	0.003
CP_TB	0.34	0.014
LK_SOLV	0.21	0.136
ELE	0.24	0.087
DDG_W	0.18	0.212
HBOND2	0.29	0.039
VDW	0.12	0.393
DCOMPLEX	0.16	0.263
Ab I-RMSD	0.16	0.256
Ag I-RMSD	-0.08	0.564
I-RMSD	0.08	0.566

Table 5.9 Pearson correlation, and correlation p-value, of functions/terms with experimentally determined  $\Delta$ Gs.

<sup>1</sup>Functions and terms are from structural analysis and scoring algorithms applied to the bound complex structures, as detailed in the Methods. Ab I-RMSD, Ag I-RMSD and I-RMSD are root-mean-squared distances between unbound and bound interface residue Cα atoms, corresponding to binding conformational change of antibody (Ab I-RMSD), antigen (Ag I-RMSD), or both (I-RMSD).



Figure 5.12 I-RMSD,  $\Delta$ ASA, and Rosetta REF15 scores versus experimentally determined  $\Delta$ Gs. (A) I-RMSD, (B) negative  $\Delta$ ASA, and (C) Rosetta REF15 score is compared with experimentally determined  $\Delta$ G for 51 cases. For (B) and (C), cases are stratified by I-RMSD, with open circles representing cases with I-RMSD < 1 Å (N=28), while filled circles represent cases with I-RMSD ≥ 1 Å (N=23). Dashed lines show linear fits, excluding labeled outliers in (A) and (B), and including only filled circles in (B) and (C). Pearson correlation coefficients are shown on the right of each plot.

Subset	$\mathbf{N}^1$	$\Delta ASA^2$	p-value	REF15 <sup>2</sup>	p-value
I-RMSD < 1 Å	28	0.12	0.55	0.38	0.047
$I-RMSD \ge 1 \text{ Å}$	23	0.31	0.16	0.57	0.004
Antibody I-RMSD < 1 Å	36	-0.05	0.77	0.27	0.1
Antibody I-RMSD≥1 Å	15	0.63	0.01	0.74	0.002
Antigen I-RMSD < 1 Å	31	0.036	0.84	0.38	0.035
Antigen I-RMSD ≥ 1 Å	20	0.31	0.18	0.56	0.01

Table 5.10 Correlations with experimental  $\Delta G$  values for  $\Delta ASA$  and Rosetta REF15 stratified by I-RMSD.

<sup>1</sup>Number of test cases in subset.

<sup>2</sup>Pearson correlation coefficients between negated  $\Delta$ ASA or Rosetta REF15 score with experimental  $\Delta$ G values for subset of affinity benchmark. Correlations greater than 0.5, and corresponding p-values, are shown in bold.

#### 5.4 Discussion

This antibody-antigen benchmark represents an expanded set of non-redundant and structurally diverse complexes, which can provide robust and challenging tests for docking and affinity prediction algorithms. Initial assessments of the benchmark not only reflect its diversity, but also raise important questions regarding conformational changes and predictive challenge that can be addressed by further research.

As antibody-antigen docking algorithms are evaluated and compared using benchmarks (172, 388, 389, 449), it is critical to understand how to improve predictive success and areas needing improvement. Initial benchmarking of rigid-body and flexible docking reported here highlights the challenge posed by this class of interfaces. Even for local docking perturbations, several cases did not have top-ranked hits or prominent binding funnels, including Rigid cases with minimal binding conformational changes. Possible factors underlying lack of local docking successes include low binding affinity, which as noted before can negatively impact success (186),

as well as suboptimal initial positioning for local docking perturbation searches, CDR loop sampling away from key bound-like conformations, or lack of explicit sampling of antigen backbone flexibility. Additionally, accurate modeling of glycosylation of antigens, including viral glycoproteins, is another possible avenue to reduce false positive docking predictions, or detect favorable docking poses near the native binding site, by representing glycan surface masking and glycan-antibody interactions.

Correlations of predictive methods with experimental affinities show that the antibodyantigen benchmark provides challenging tests of affinity prediction. Unexpectedly, success of affinity predictions was found to be higher for antibody-antigen complexes with larger conformational changes, both for REF15, the affinity predictor with highest overall correlation with  $\Delta G$ , and  $\Delta ASA$ . This contradicts previous results based on a more general set of proteinprotein interaction affinities (217), which were later used to generate a linear model incorporating I-RMSD and  $\Delta ASA$  to predict affinities (220). The reasons for this intriguing shift with respect to previous studies are unclear, but the focus on antibody-antigen complexes in this benchmark may have played a substantial role. Future studies, possibly with explicit dynamics simulations, can be used to understand how the predictability of antibody-antigen affinities may be influenced by binding mobility and conformational changes.

Past successes in structure-based antibody design have demonstrated the utility of rational design methods (159, 416, 450-452), yet previous benchmarking of several computational design methods showed relatively low correlations with experimental data (210). This new and updated set of affinities provides a wide range of interfaces and  $\Delta G$  values, representing a complementary dataset to an existing benchmark with binding energy changes in antibody-antigen interfaces ( $\Delta\Delta Gs$ ) (210); these can collectively be used to assess predictive performance of antibody affinity

prediction methods, or an element of assessments that incorporate distinct datasets of affinities for other classes of interactions with experimentally determined structures (207, 453, 454).

Future efforts to develop advanced predictive antibody docking methods can utilize this expanded antibody-antigen benchmark to train and test algorithms. Recent studies have indicated that machine learning approaches trained on protein interface data can lead to improved selection of protein-protein docking models (183, 455, 456), and this set can provide test cases to assess, and potentially optimize, predictive success for antibody-antigen complexes with these methods. Algorithm development to address challenging areas including "cross-docking" to predict antibody targets using docking simulations (457, 458) and integrative epitope prediction (459-462) likewise can be facilitated by this benchmark.

Prospective benchmark updates can incorporate antibody-antigen structural information that does not adhere to the limits that were set when building this benchmark, offering more options to assess docking and affinity predictions as well as modeling accuracy. For antibody-antigen complexes without a corresponding experimentally determined unbound antibody or antigen structure, the missing component could be modeled and docked to the structure of its binding partner. Indeed, use of a modeled unbound antibody structure for docking represents a more common predictive modeling scenario, as an experimentally determined unbound antibody structure is not likely to be available for a given antibody sequence. We plan to identify new antibody-antigen cases from the PDB on a regular basis, and we will include these on the benchmark site as a pre-release prior to the next major benchmark version, to aid algorithm development and testing by the community.

# Chapter 6: A curated dataset of antibody-antigen affinities and structures to facilitate development of affinity prediction algorithms

# Abstract

Antibody recognition of antigens represents a unique class of immune system interactions, with binding affinities ranging from micromolar to picomolar. Current methods for in silico prediction of antibody-antigen affinity leave room for improvement, and more robust affinity predictors would aid antibody design and therapeutic development. Several databases report antibody-antigen affinities, but the lack of a dedicated and updated affinity dataset has limited the resources needed for training and developing new algorithms. Here, we present a curated dataset of antibody-antigen binding affinities paired with experimentally determined complex structures, allowing us to demonstrate the utility of this data for training new or optimized functions for affinity prediction. 401 binding free energy ( $\Delta G$ ) values in kcal/mol were included in this dataset following extensive searches connecting affinity data in the literature to available antibody-antigen complex structures in the Protein Data Bank (PDB). Scoring functions from Rosetta and other programs showed promising correlations to antibody-antigen affinity values even though they were not optimized or designed solely for affinity prediction. Using the affinity dataset for training, existing scoring functions with weights optimized through multilinear regression showed improved correlations to antibody-antigen affinities. These retrained models were evaluated with cross-validation and an independent test set. This dataset represents an improved and centralized resource for antibody-antigen affinity prediction, and it can be used to test new approaches or optimize current methods. Information about the full dataset can be found at https://piercelab.ibbr.umd.edu/Draft antibody antigen affinity dataset.xlsx.

## **6.1 Introduction**

Antibody-antigen recognition is a unique class of protein-protein interactions, and an improved understanding of what dictates these interactions can lead to advances in protein design and therapeutics development. Powerful and promising in silico tools have helped alleviate experimental resource constraints, including computational prediction of antibody-antigen interfaces through protein docking and of antibody-antigen interaction strength through affinity prediction. As discussed in Chapter 5, a variety of methods have been developed to predict proteinprotein binding affinities (201, 202, 218, 220, 380, 463, 464). Some of these methods utilized a benchmark of 144 affinities for training (217), where each protein-protein interface of the corresponding structure was scored by a predictive physics-based or knowledge-based potential and then compared to experimentally determined  $\Delta G$  values (often in kcal/mol). Despite the progress made in addressing affinity prediction, the applicability of potentials trained on general protein-protein interactions may be limited when specifically predicting antibody-antigen affinity. Less than 200 protein-protein binding affinities have been used for training from the affinity benchmark and several docking benchmarks (186, 195), and a small fraction of those cases represent antibody-antigen interactions, forming a small potential training set for predicting antibody-antigen  $\Delta G$  values. The unique properties of antibody-antigen interfaces have been explicitly addressed by protein docking algorithms such as ClusPro (172) and FireDock (181), showing how antibody-antigen recognition can be treated as separate from all protein-protein interactions. Though predictions of affinities specific to antibody-antigen recognition have been studied recently, the suitability of current datasets for benchmarking remains unclear.

Another study noted the dearth of antibody-antigen complexes and affinities currently available for benchmarking and presented Docking Benchmark 5.5 (BM5.5), an expanded set of

antibody-antigen cases for docking and affinity prediction, with full details in Chapter 5 (465). The expanded set of antibody-antigen affinities consisted of 51 cases, which were evaluated by a variety of functions that score protein-protein interface structures, reporting correlations between resulting scores and experimentally determined  $\Delta G$  values. The correlations found were modest at best, with Pearson correlation coefficients of less than 0.5, suggesting that there is substantial room for improvement in predicting antibody-antigen affinities. In response, recent work has presented CSM-AB, an antibody-antigen affinity prediction algorithm using graph-based machine learning methods (214) that improved on correlations with  $\Delta G$  values when training on a dataset of over 400 affinities obtained from PDBbind (216). Yet correlations were still modest despite improvements, showing that more algorithmic development is needed to address this difficult problem, possibly through optimization of existing functions tested on BM5.5 (465). At the same time, the affinity data used for training had to be assembled from a yearly-updated database of protein-protein binding affinities (216); other datasets dedicated to antibody-antigen complexes are available and frequently updated, but at most include affinity data occasionally (153). This situation sits in stark contrast with antibody-antigen  $\Delta\Delta G$  values, where several curated datasets (209, 210) have spurred more advancements in  $\Delta\Delta G$  prediction through machine learning (211, 213, 466). Without similar dedication to a dataset of antibody-antigen  $\Delta G$  values, future advances in algorithm development of affinity predictors may be limited by the suboptimal size and accessibility of validated  $\Delta G$  values that can be used for benchmarking. Furthermore, this deficit undercuts previous presentations of  $\Delta G$  and  $\Delta \Delta G$  prediction as equally valuable efforts that may work in tandem to increase understanding of protein-protein interactions, which could aid the design of therapeutics (467).

Here, we present a curated and annotated dataset of 401 antibody-antigen  $\Delta G$  values, each validated in the literature with a documented temperature and a corresponding complex structure. This dataset represents a wide range of  $\Delta G$  values from micromolar to picomolar, along with broad diversity in antigen origin and type of antibody. Correlations between experimental affinities and scoring functions REF15, beta16, IRAD, and ZRANK were calculated, showing modest predictive performance on the expanded dataset. The dataset was also broken into subsets based on method of affinity measurement and structure resolution, but correlations between  $\Delta G$  values and each scoring function did not differ substantially by subset. To demonstrate the utility of this curated dataset for improving antibody-antigen affinity prediction, we assessed the correlations of individual terms from tested scoring functions, then attempted to optimize the weights of these terms for antibody-antigen affinity prediction using stepwise regression and cross-validation. Several retrained models showed modest improvements in  $\Delta G$  correlations with existing scoring functions, and additional methods of regression trained on the affinity dataset did not yield the same improvements. Though these models showed higher correlations with the affinity dataset, they showed lower correlations than existing functions to an independent test set of antibodyantigen affinities after scoring the corresponding modeled complexes. However, two of the top models did show improved correlations with a smaller set of neutralization data during this validation stage. This dataset represents the largest and most detailed standalone database for antibody-antigen  $\Delta G$  values, providing a useful resource that can stimulate development of affinity prediction algorithms.

## 6.2 Methods

## 6.2.1 Collection of cases in the antibody-antigen affinity dataset

Each case in this affinity dataset consists of an experimentally determined antibody-antigen affinity value (calculated as  $\Delta G$ , in kcal/mol) and a corresponding antibody-antigen complex structure deposited in the PDB; (336). Counterintuitively, collection of these cases started with antibody-antigen structures in PDB, followed by manual inspection of the reported methods, results, and supplemental files of any publications reporting the affinity of a particular complex. In most cases, a structure was linked to a publication through its PDB web page; if affinity of this complex was measured, it was often reported in this publication. Antibody-antigen complex structures were also visualized in PyMOL (Schrödinger), in part to confirm correspondence between a given structure and affinity value that would form a case. For PDB structures released prior to 2020, SAbDab (153) and PDBbind (216) datasets were compared to connect any structures listed in SAbDab to any affinity values for the same structure listed in PDBbind. The list of PDB codes found in both databases formed an initial list of affinity values and were investigated individually. Antibody-antigen complex structures listed in SAbDab with search criteria of protein antigen and <3.3 Å resolution, and that did not have a corresponding affinity value in PDBbind, were manually inspected for possible affinity values. Each affinity value was confirmed through literature searches; papers that reported an affinity corresponding to a PDB structure either reported the affinity value and structure or reported an affinity value exclusively, with the structure published in a subsequent paper.

Affinity values for corresponding structures were only included in the dataset if there was a documented temperature for the affinity measurement in the supporting literature for the structure or in literature that the paper directly cited in a description of its methods. If the temperature was listed as RT or room temperature, it was assumed to be 25°C. The experimentally determined dissociation constant ( $K_D$ ) and listed temperature were used to calculate  $\Delta G$  (in kcal/mol) for each case. More specifically,  $\Delta G$  was calculated using the equation  $\Delta G = RTlnK_D$ , where R is the gas constant and T is the temperature in degrees Kelvin. In four cases, an experimentally determined association constant ( $K_A$ ) that corresponded to a particular structure was used to calculate  $K_D$  through the equation  $K_D = 1/K_A$ , which was then used to calculate  $\Delta G$  in the same fashion. Structures with a resolution of <3.25 Å, a documented affinity value and temperature in the literature, and a match between the antibody-antigen complex used in structure and the antibody-antigen complex used to measure affinity were placed on a preliminary list of new affinity cases. Possible redundancy amongst new cases and between new cases and BM5.5 cases (465) was evaluated using BLAST (362). The sequence identity set as a threshold for non-redundant cases, except where antigens are already non-redundant (<80% sequence identity) even though heavy chain variable domain sequence identity is >90%.

## 6.2.2 Curation of antibody-antigen affinity dataset

Several characteristics of the affinity measurement for each new case were collected to add background information and enable comparisons between affinity measurement methods. Recorded information from supporting literature included the method of affinity measurement, often surface plasmon resonance (SPR), bio-layer interferometry (BLI), or isothermal titration calorimetry (ITC), organization of the affinity measurement (whether the antibody or antigen was immobilized), and additional notes about comparisons between the complex structure and the complex affinity. The resolution of the corresponding structure was collected from a batch download of PDB structure information.

The number of non-water hetero atoms (HETATMs) in the interface of each new affinity case was found using an in-house Perl script. All atoms within 6 Å of an atom of its binding partner were defined as interface contacts, and included contacts between non-water HETATMs, between non-water HETATMs and amino acid atoms, and between amino acid atoms. This list of interface contacts was used to define the level of HETATM involvement, as measured by the percentage of interface contacts that include a HETATM. This method was primarily utilized as a mechanism to initially screen affinity dataset structures for substantial HETATM interface contacts. The permissive detection of HETATM contacts (i.e. not restricted to previously defined antibodyantigen interface residues) likely inflated the percentage of HETATM contacts in some cases. Manual inspection of structures in PyMOL was ultimately used to determine the degree and impact of non-water HETATM contacts in antibody-antigen interfaces. Any affinity case with a substantial percentage of contacts (>30%) that involved non-water HETATMs with a suspected effect on antibody recognition in the corresponding literature was defined as a "high HETATM" case. Five cases were classified as "high HETATM", with substantial involvement originating from either N-linked glycan (4JM2, 4TVP, 5CEZ, 6J11) or ATP (7DC8) HETATMs in the interface. These cases were not utilized in the analysis and optimization of affinity predictors, but have been included in the affinity dataset for future affinity prediction efforts that may integrate HETATMs.

Two cases, corresponding to structures 7A29 and 7D2Z, satisfied the criteria outlined here, but were excluded due to uncertainty in correspondence between the resolved structure of the complex and the experimentally determined  $\Delta G$  value. We confirmed an affinity measurement for 7A29, a SARS-CoV-2 spike trimer in complex with a nanobody, but it was unclear if this measurement would correspond to binding the receptor binding domain (RBD) in an up or down conformation (468). Since the structure contained both conformations of the RBD and these conformations showed different numbers of interface residues contacting the nanobody, we decided to remove the case. We also confirmed an affinity measurement for 7D2Z, a SARS-CoV-2 RBD in complex with a nanobody, but found a portion of an N-terminal signal peptide and a C-terminal cleaved 3C protease site on the RBD that were resolved in the structure and formed a substantial part of the antibody-antigen interface (469). Since additional C-terminal tags appeared to be included in RBD protein used for the affinity measurement, this discrepancy led to uncertainty regarding the correspondence between structure and affinity, and ultimately the case was removed.

## 6.2.3 Analysis of affinity predictors

Prior to refinement and scoring of affinity cases, any hydrogens, double occupancy, or HETATMs in all chains of structures were removed. The "high HETATM" cases were not included in this analysis. Affinity cases were refined using Rosetta FastRelax (weekly release 2020.28) (397) with two variations: constrained sidechain atoms and unconstrained sidechain atoms. In both variations of the refinement protocol, backbone atoms were constrained. A typical command in Rosetta for the refinement stage is as follows, with the "coord\_constrain\_sidechains" flag removed when refining with unconstrained sidechains:

-ignore\_zero\_occupancy false

```
-ignore_unrecognized_res
```

-relax:constrain\_relax\_to\_start\_coords

-relax:coord\_constrain\_sidechains

-relax:ramp\_constraints false

-relax:fast

-ex1
-ex2
-use\_input\_sc
-no\_his\_his\_pairE
-no\_optH false
-flip\_HNQ
-renumber\_pdb F
-overwrite
-nstruct 1

Following refinement, the binding score of each case was determined by REF15 and beta16 scoring functions (204) in Rosetta using the score application (weekly release 2020.28). Each case was scored as a complex and separated by components. The binding score for each case was calculated as the difference between the complex score and the sum of the scores of each component when separated. Pearson correlation coefficients between experimental affinity values and binding scores were calculated in R (366), both for the total score and for individual terms of each scoring function. Integrated Residue- and Atom-based potentials for Docking (IRAD; (182)) and ZRANK (179) functions were also used to separately score affinity cases by running each application on the command line in a Unix environment. Pearson correlation coefficients between experimental affinity values and IRAD scores, ZRANK scores, IRAD individual terms, or ZRANK individual terms were also calculated in R.

6.2.4 Comparison of correlations in affinity subsets

Affinity cases were separated by resolution into four subsets:  $\leq 2 \text{ Å}$ ,  $2 \text{ Å} < x \leq 2.5 \text{ Å}$ , 2.5 Å, 2.5 Å, and > 3 Å. To separate affinity cases by method of affinity measurement, four subsets

were determined: Measurement by SPR, measurement by BLI, measurement by ITC, and measurement by a method that was not solely defined by one of the other three techniques and classified as "Other". Techniques in the "Other" subset for measuring affinity include enzyme-linked immunosorbent assay (ELISA), competition inhibition assays, radioligand binding assays, kinetic exclusion assay (KinExA), spectroscopy, radioactive iodination and ultracentrifugation, fluorescence polarization, stopped-flow fluorescence kinetics, and microscale thermophoresis. Correlations between experimental affinity values and REF15 binding scores in each subset were calculated in R.

#### 6.2.5 Additional case information in affinity dataset

The full table for the affinity dataset contains additional annotations for each case. The reference for each affinity found in the literature is shown in the "Reference" column, often showing the PubMed ID number (pubmed.ncbi.nlm.nih.gov) for the publication. In several cases, the documented temperature for a given affinity measurement was gathered from a publication referenced by the publication that reported the affinity. When this situation occurred, the reference for the temperature was placed in a separate column as a PubMed ID to signify that a second reference was needed to calculate  $\Delta G$ . Affinity measurements were found to be oriented in a variety of ways in the literature, sometimes with the antigen immobilized and the antibody as the analyte, and sometimes vice versa. To provide more description of affinity measurements beyond naming the technique, we added a "Method notes" column. This column was primarily used to indicate what component of the measurement was immobilized if applicable, with the common note "Immobilized IgG" showing that the antibody was immobilized. Immobilized antigen and antibody as analyte was treated as the default orientation for affinity measurements in this table, so any case with "Method notes" left blank had affinity measured in this format. PDB validation

metrics are also included in the table to annotate cases with measurements of structure quality beyond resolution. These metrics, which are placed in separate columns, include Rfree, Clashscore, Ramachandran outliers, Sidechain outliers, and RSRZ outliers. PDB validation reports for each case were downloaded from rcsb.org in xml format and parsed to retrieve the above metrics using in-house Bash scripts. If "N/A" is listed as a value in any of these five columns, the validation metric with this value was not calculated for this case.

#### 6.2.5 Individual terms in REF15 and beta16 scoring functions

In addition to the binding score, calculated values for each individual term in these Rosetta scoring functions were obtained for all antibody-antigen affinity cases except for "high HETATM" cases (N = 396). REF15 includes a total of 19 individual weighted terms, while beta16 includes a total of 26 individual weighted terms. However, only eight terms in REF15 and eleven terms in beta16 contributed to binding score calculated from each scoring function, and only those terms are listed and described here. The values for all other individual terms were zero or approximately zero when the sum of scores from complex components were subtracted from the score of the entire complex structure. Most terms are shared between REF15 and beta16; if a particular term is only present in one scoring function, it is indicated below. The " $\Delta$ " placed before each term in the results was included to signify that these scoring terms were calculated from the difference between scores of the complex structure and the sum of its components. This symbol is not present in the original description of these scoring terms.

Brief descriptions of REF15 and beta16 terms (204, 470):

fa\_atr – Lennard-Jones attractive energy between two atoms on different residues fa\_rep – Lennard-Jones repulsive energy between two atoms on different residues fa\_sol – Gaussian exclusion implicit solvation energy lk\_ball (beta16 only) – Anisotropic contribution to the solvation

lk\_ball\_iso (beta16 only) – Same as fa\_sol

lk ball wtd (REF15 only) – weighted sum of lk ball and lk ball iso

lk\_ball\_bridge (beta16 only) – Solvation bonus from bridging waters, as measured from interactions of polar atoms

lk\_ball\_bridge\_uncpl (beta16 only) – Same as lk\_ball\_bridge, but the value is uncoupled with dGfree

fa\_elec - Energy of interaction between two nonbonded charged atoms

hbond\_lr\_bb - Energy of long-range backbone-backbone hydrogen bonds

hbond bb sc – Energy of backbone-sidechain hydrogen bonds

hbond\_sc – Energy of sidechain-sidechain hydrogen bonds

6.2.6 Individual terms in IRAD and ZRANK scoring functions

Calculated values of individual terms for IRAD and ZRANK scoring functions were output by IRAD while scoring complex structures on the command line in a Unix environment. These terms include atom-based and residue-based contact potentials, numbers of atoms or residues used to calculate potentials, and individual ZRANK terms that include van der Waals and electrostatic components. IRAD also output electrostatic terms from the Zlab Affinity for Protein-Protein interaction (ZAPP; (380)) free energy function. These terms were considered as input for subsequent regression analyses even though correlations between ZAPP scores and  $\Delta G$  values were not assessed.

Brief descriptions of IRAD and ZRANK terms (179, 182):

potlatot - total of potential 1 from Zhang et al. (471), actual center of mass

pot2atot - total of potential 2 from Zhang et al., actual center of mass

pot1ptot - total of potential 1 from Zhang et al., parameterized center of mass

pot2ptot - total of potential 2 from Zhang et al., parameterized center of mass

npot1atot - number of atom pairs identified in calculation of pot1atot

npot2atot - number of atom pairs identified in calculation of pot2atot

npot1ptot - number of atom pairs identified in calculation of pot1ptot

npot2ptot - number of atom pairs identified in calculation of pot2ptot

presa – Chakrabarti and Janin's potentials based on interface propensities (472) with Yang's contacting definition, actual center of mass

presp – Chakrabarti and Janin's potentials based on interface propensities with Yang's contacting definition, parameterized center of mass

nresa - number of interface residues identified in calculation of presa

nresp - number of interface residues identified in calculation of presp

potctot – Yang's potential (471) with 4.5 Å cutoff used for defining contacts

npotctot – number of atom pairs identified in calculation of potctot

presc – Chakrabarti and Janin's potentials with 4.5 Å cutoff used for defining contacts

nresc – number of interface residues identified in calculation of presc

pottbtot – Tobi and Bahar potential for protein-protein docking (405)

npottbtot - number of interactions identified in calculation of pottbtot

prestb – Chakrabarti and Janin's potentials with contacting definition from Tobi and Bahar potential

nrestb - number of interface residues identified in calculation of nrestb

potllstot – Lu, Lu, and Skolnick potential for protein-protein interactions (473)

npotllstot - number of interactions identified in calculation of potllstot

potg1tot – Potential from Glaser et al. for interface residue contact preferences (474), with a cutoff of 6 Å for contacting definition

potg2tot – Normalized potential from Glaser et al. for interface residue contact preferences, with a cutoff of 6 Å for contacting definition

npotgtot - number of interactions identified in calculation of potential from Glaser et al.

- presg Chakrabarti and Janin's potentials with contacting definition from Glaser et al. potential
- nresg number of interface residues identified in calculation of presg
- acont4 number of atom contacts within 4 Å cutoff
- acont5 number of atom contacts within 5 Å cutoff
- acont6 number of atom contacts within 6 Å cutoff

solvlk – Lazaridis-Karplus implicit solvation model, with cutoff distance optimized to match Rosetta (475)

- vdw\_atr van der Waals attractive force from ZRANK
- vdw\_rep van der Waals repulsive force from ZRANK
- elec\_sra short-range electrostatics attractive force from ZRANK
- elec srr short-range electrostatics repulsive force from ZRANK
- elec lra long-range electrostatics attractive force from ZRANK
- elec lrr long-range electrostatics repulsive force from ZRANK

ace – statistical contact potential of atomic contact energies derived from monomeric protein structures (476)

iface – statistical contact potential of interface atomic contact energies derived from proteinprotein complex structures (477)

- elec\_sra\_x short-range electrostatics attractive force from ZAPP (380)
- elec\_srr\_x short-range electrostatics repulsive force from ZAPP
- elec\_lra\_x long-range electrostatics attractive force from ZAPP
- elec\_lrr\_x long-range electrostatics repulsive force from ZAPP

## 6.2.7 Regression analysis and cross-validation

Stepwise regression was conducted in R on all sets of terms using the StepAIC() command of the MASS package (478), with selection of terms set in both directions. Optimization was performed on eight different sets of input terms, either from one scoring function or terms combined from several scoring functions, to predict  $\Delta G$  values of antibody-antigen affinity cases. Following selection of input terms for each set, linear regression with stepwise selected terms was conducted in R using the train() command of the caret package (479). Ridge, Least Absolute Shrinkage and Selection Operator (LASSO), and Elastic net regressions were conducted in R on all sets of terms using the train() command of the caret package, with glmnet as a dependency (480). All methods of cross-validation during regression (leave-one-out, 5-fold, or 10-fold) were specified and automated using the trainControl() command of the caret package. To ensure reproducibility, a constant seed of 123 was set with set.seed() immediately prior to any regression analysis.

Sets of input terms from REF15/beta16, IRAD, and ZRANK were combined in several different ways for regression analysis. Descriptions of these sets were broken into categories based on the source of input terms.

#### 6.2.8 REF15/beta16-based sets

The eight REF15 terms that contributed to binding score calculations were used as a set of input terms. The eleven beta16 terms that contributed to binding score calculations formed a separate set of input terms.

## 6.2.9 IRAD/ZRANK-based sets

Individual terms output by IRAD were combined into several different sets. One set only included eight terms from IRAD and ZRANK that previously formed an optimized scoring function for ranking protein docking models (182). These terms include vdw\_atr, vdw\_rep, elec\_s, elec\_l, ace, iface, potctot, and npotctot. The terms elec\_s and elec\_l are separate sums of short-range and long-range electrostatic force calculations from IRAD. elec\_s values were determined by calculating the sum of terms elec\_sra and elec\_srr output by IRAD. elec\_l values were determined by calculating the sum of terms elec\_lra and elec\_lrr output by IRAD. Another set was a combination of terms output by IRAD that calculated potentials and all ZRANK terms, totaling 23 terms in all ("IRAD + pot"). Finally, a larger set included all values output by IRAD as input terms, including calculated potentials, counts of residues or interactions determined by potentials, ZRANK individual terms, and ZAPP electrostatic terms, totaling 43 terms in all ("IRAD + pot/count").

# 6.2.9 Composite REF15/IRAD sets

Additional sets were formed with combinations of individual terms from different sources, namely REF15 and IRAD. One set of terms combined the eight REF15 terms contributing to binding score with all calculated potentials output by IRAD for a broader set of 22 individual terms ("REF15 + pot"). The other composite set contained the REF15 and IRAD potential terms, but also added the terms for counts output by IRAD, leading to a set of 38 terms ("REF15 + pot/count").

## 6.2.10 Data for independent test set

60 VRC01-class germline antibodies were isolated from transgenic mice, namely VRC01<sup>gHL</sup> mice expressing both heavy and light chains of broadly neutralizing antibody VRC01 germline version, which were immunized with HIV-1 Env by the lab of Yuxing Li (University of Maryland IBBR) and characterized for binding affinity to and neutralization of six HIV-1 isolates (481). Antibody heavy and light chain sequences, isolate gp120/gp41 sequences, binding affinity data,

and neutralization data were kindly shared by Lin Lei and Andrey Galkin in Yuxing Li's lab. Binding affinities were quantified in ELISA assays using area under the curve (auc) calculations. These auc values were negated and  $log_{10}$ -transformed to better match the range and direction of predicted  $\Delta G$  values and binding scores. Neutralization values were quantified as IC50s in ELISA and reported in  $\mu g/ml$  concentration. These values were log10-transformed prior to correlations with predicted  $\Delta G$  values and binding scores.

#### 6.2.11 Modeling and scoring of antibody-antigen complexes from independent test set

Antibody-antigen complexes were modeled using a structural template of SOSIP from the 45 01dG5 isolate (482, 483) bound to a VRC01-class intermediate antibody (481), a structure that was determined through cryo-EM and kindly shared by Andrey Galkin. This structure was also modified to complete the antibody-antigen interface by grafting a gp120 loop (residues 59-67) from a VRC01-bound SOSIP structure (PDB code: 5FYK (45)) to the cryo-EM structure, which did not have this loop resolved. To graft this loop, 5FYK and the cryo-EM structure were first superposed and visualized in PyMOL (Schrödinger) to ensure that the resolved loop would fit into the disordered region of the cryo-EM structure. This loop was then extracted from 5FYK and included with the coordinates of the cryo-EM structure, forming a new template for modeling. Provided sequences for intermediate antibodies and isolates were modeled as a trimeric SOSIP construct bound to three intermediate antibodies through a custom homology modeling pipeline using the MODELLER program (484). For stages of refinement and scoring, two gp120 chains and one antibody (heavy and light chains) were extracted from the full model to reduce computation time while keeping the antibody-antigen interface intact. Each model was then refined by Rosetta FastRelax, with both backbone and sidechain atoms constrained. REF15 binding score was calculated for modeled gp120-antibody complexes using the same procedure outlined for scoring antibody-antigen affinity case structures. IRAD scores of modeled structures were also calculated as previously outlined.

#### 6.2.12 Predictions of optimized scoring functions on independent test set

Performance of retrained scoring functions on an independent test set was assessed by calculating correlations between negated and  $log_{10}$ -transformed affinity or  $log_{10}$ -transformed neutralization values and the calculated scores of modeled structures of gp120-antibody complexes. To calculate predicted  $\Delta G$  values for each modeled complex, scores of individual terms output by REF15 and IRAD were combined according to the components included in each retrained function, multiplied by the corresponding weight for each term derived during model optimization, then added together with the intercept value for a given function. Pearson correlation coefficients and p-values between predicted  $\Delta G$  values and experimental affinity or neutralization values were calculated in R. Correlations of REF15 and IRAD scores with affinity or neutralization values were calculated in the same fashion, and these correlations were used as a baseline for comparisons to performance with retrained functions.

#### **6.3 Results**

#### 6.3.1 Dataset assembly and diversity

This curated dataset of antibody-antigen affinities includes 401 cases that have both an experimental  $\Delta G$  value for antibody-antigen affinity and a corresponding structure of the antibodyantigen complex in the PDB. Five of these cases form a small but unique set of antibody-antigen complexes with a substantial number of hetero atom (HETATM) interface contacts, often due to the presence of N-glycans in the interface. These cases were not utilized in antibody-antigen affinity prediction for this research, but they may be useful to include in any future prediction algorithms that explicitly address HETATMs, and especially N-glycans, as part of their analysis. A summary of cases in the antibody-antigen affinity dataset is shown in **Figure 6.1**, which illustrates the diversity of this dataset. Affinity values were obtained from structures that were released in the PDB between 1998 and July 2021, covering a variety of antigens and methods of affinity measurement. In 177 cases, or about 44% of this affinity dataset, an affinity value for a corresponding antibody-antigen structure had not been reported in any commonly-used database of  $\Delta G$  values, including PDBbind (216). A sizeable portion of cases that do not overlap with PDBbind were released in 2020 or 2021, making their absence from PDBbind expected since the database currently contains affinity values through 2019 (pdbbind.org.cn, accessed February 2022). However, approximately 100  $\Delta G$  values in this affinity dataset were released before 2020 and have not been reported by PDBbind.

In the full dataset, 96 structures contain a nanobody, comprising ~24% of affinity cases in this database, which is a higher proportion of nanobody affinities than contained in Docking Benchmark 5.5 (BM5.5) (~18%; (465)). This expanded affinity dataset also includes a higher proportion of non-human antigens (~60%) than BM5.5 (~53%), another indication of increased benchmark diversity. A noticeable increase in viral antigens was observed as well, with a higher proportion of viral antigens in new cases (~35% of dataset) when compared with BM5.5 (~24%). A large majority of experimental affinity values were measured with surface plasmon resonance (SPR), bio-layer interferometry (BLI), or isothermal titration calorimetry (ITC) methods, with the largest number of affinity values collected using SPR (~62% of dataset). Any trends in the changing composition of evaluated antibody-antigen affinities may reflect broader trends in favored techniques for affinity measurement and in the biological focus of structural determination, underscoring the need for an up-to-date set of antibody-antigen affinity data.

This dataset contains 51 experimental affinity values previously reported in BM5.5 (465), and we wanted to examine any shifts in the distribution of affinity values or structural resolution following the expansion of cases, which may affect the suitability of this dataset for antibodyantigen affinity prediction. Additional cases exhibited a broader range of  $\Delta G$  values than BM5.5 cases, but median values were largely unchanged, showing that this expanded dataset largely resembles the dataset of affinity values established previously (**Figure 6.2A**). The median structural resolution of cases from BM5.5 was also unchanged when compared to newly added cases (**Figure 6.2B**). Though the median resolution of additional cases was slightly lower, these cases also expanded the range of structure resolutions observed in the dataset, notably through eight cases between 1.1 and 1.5 Å. The resolution of antibody-antigen complex structures showed a large range of values from ~1.2 Å to ~3.3 Å, with the median resolution at 2.5 Å. Affinity values from BM5.5 were previously correlated with scoring functions or potentials analyzing the structures of antibody-antigen interfaces (465), and little change in median affinity or resolution suggests that this larger affinity dataset can also be used to assess methods of affinity prediction.



Figure 6.1 Summary of diversity in the antibody-antigen affinity dataset. Broad categories of antigen origin, antibody type, and affinity measurement method are shown, with cases from Docking Benchmark 5.5 (BM5.5; N = 51) and new additions (N = 350) separated into different columns. Types of antibodies include Fab, single-chain variable fragment (scFv), and single domain antibody (sdAb).



Figure 6.2 Ranges of  $\Delta G$  values and structural resolution in affinity dataset. (A) Experimental affinity values plotted for cases originating from BM5.5 (N = 51) and cases added to form the affinity dataset (N = 350). (B) Resolution of structures for antibody-antigen affinity cases originating in BM5.5 and added to affinity dataset. Median values for  $\Delta G$  and resolution in Å are shown as black bars for all subsets.

## 6.3.2 Performance of existing scoring functions as affinity predictors

To assess existing scoring functions and potentials as possible affinity predictors, antibodyantigen structures were refined, scored, and evaluated for correlations between predicted score of the complex and the corresponding  $\Delta G$  value. Rosetta scoring functions REF15 and beta16 (204), integrated residue-based and atom-based potentials for docking (IRAD) (182), and ZRANK (179) algorithms were used to score all complexes in the dataset following refinement with Rosetta FastRelax (397), with the exception of the five affinity cases that had substantial HETATM interface contacts. First, methods of FastRelax refinement were tested: complexes with backbone and sidechain atoms constrained, and complexes with only backbone atoms constrained. These

В

variations on refinement were named "constrained sidechains" and "unconstrained sidechains", respectively. Following refinement, both sets of complexes were scored with REF15 or beta16 and correlations with experimental affinity values were determined (**Table 6.1**). Pearson correlation coefficients between binding scores derived from Rosetta and affinity values were between R = 0.27 and R = 0.29, a notable drop from the correlation found when  $\Delta G$  values from BM5.5 were correlated with REF15 scores (R = 0.45; (465)). However, the p-values for Pearson correlation coefficients with the complete affinity dataset showed higher significance (p << 0.001) than the p-value for the Pearson correlation coefficient between REF15 binding score and affinity values from BM5.5 (p = 0.001), suggesting that both scoring functions show some promise in predicting antibody-antigen affinity values. These correlations were also clearly higher with refinement than without (R = 0.09-0.13), showing that a refinement protocol can mitigate possible structural errors and better reflect experimental  $\Delta G$ .

Though Pearson correlation coefficients were similar overall, scoring and refinement with constrained sidechains was the combination that produced the highest Pearson correlation coefficients for both scoring functions ( $R \sim 0.29$ ). Based on this result, complexes were refined with constrained sidechains prior to scoring with other functions or potentials included in this analysis. Correlations of IRAD and ZRANK scores were checked for correlations with the affinity dataset, finding Pearson correlation coefficients near REF15 correlations ( $R \sim 0.3$ ; **Table 6.1**). Both scoring functions were originally designed to rank or re-rank models from rigid-body docking algorithms, likely making their scoring term weights suboptimal for antibody-antigen affinity prediction, along with the weights of individual terms in REF15 and beta16. In addition, some contact potentials and counts were calculated by IRAD and tested for predictive performance but excluded from the published scoring function, providing more terms that can be evaluated
specifically for antibody-antigen affinity prediction The tested scoring functions showed some promise as affinity predictors, leading to a more detailed analysis of scoring terms and how they correlate with  $\Delta G$  values individually.

Scoring	Sidechain	Pearson
function	refinement	correlation^
REF15	No refinement	0.09
REF15	Constrained	0.29***
REF15	Not constrained	0.27***
beta16	No refinement	0.13*
beta16	Constrained	0.29***
beta16	Not constrained	0.27***
IRAD	Constrained	0.30***
ZRANK	Constrained	0.29***

Table 6.1 Correlations of REF15, beta16, IRAD, and ZRANK scores with ΔG values.

^Scores for all cases in the affinity dataset except for five high-HETATM cases were correlated with  $\Delta G$  values (N = 396). P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\*\* < 0.001

### 6.3.3 Correlations of predictors by annotation subset

To better understand the correlations between scores and affinity values, results for REF15, beta16, IRAD, and ZRANK were individually plotted (**Figure 6.3A-D**). As suggested by Pearson correlation coefficients around R = 0.3, the highest affinity values, or lowest  $\Delta G$  values, tended to have the lowest calculated scores, or an estimate of higher free energy. Each scatterplot shows noticeable and consistent variation around this trend, including some lower affinity cases with higher calculated scores and some higher affinity cases with lower calculated scores. One case (7D2Z) showed substantially higher scores than suggested by its affinity value, and was later removed from the dataset after the match between affinity and structure was found to be uncertain, as detailed in the methods. Since cases in the affinity dataset represented various methods of affinity measurement and a range of structure resolutions, we separated cases into subsets based on this information and determined correlations with  $\Delta G$  values for each subset. REF15 score

correlations to  $\Delta G$  values showed some differences in subsets of affinity measurement method, which included SPR, BLI, ITC, and Other (Table 6.2). Correlations for the SPR and BLI subsets were significant for all affinity predictors and exhibited a range of Pearson correlation coefficients (R = 0.255-0.367) that were largely consistent with correlations to the entire dataset. While all correlations for the ITC subset were also relatively consistent with full dataset correlations (R =0.255-0.39), these correlations were found to be significant only for IRAD and ZRANK. Cases with a method of measurement classified as Other showed no significance between  $\Delta G$  values and scores of affinity predictors when tested as a subset, and Pearson correlation coefficients for REF15 and beta16 scores were notably lower (R = 0.06-0.07) than correlations with the entire dataset. Though it suggests that measurement method may influence affinity prediction by REF15 and beta16, this observation comes with caveats that include the nature of this subset (combination of several unrelated measurement methods) and its small size (N = 29). Decreased correlations in this small subset are also likely to be influenced by cases that deviate from the trendline in the entire dataset, if these cases happen to use a measurement method classified as Other. The case 6C9U may reduce the Pearson correlation coefficient substantially on its own in this subset, as it has a low  $\Delta G$  value but a higher REF15 score. IRAD and ZRANK correlations for this subset (R = 0.239 - 0.295) were not significant, but much closer to the correlation for the entire dataset.

The affinity dataset was also separated into subsets by structural resolution, specifically groups of <=2 Å, 2 Å < x <= 2.5 Å, 2.5 Å < x <= 3 Å, and >3 Å (**Table 6.3**). These subsets did not lead to an even distribution of cases, largely because the >3 Å subset was less than half the size of any other subset (N = 40), but the resulting subsets were more evenly distributed than cases separated by methods of affinity measurement. All resolution subsets showed significant p-values for Pearson correlation coefficients between  $\Delta G$  values and all affinity predictors, suggesting that

one resolution subset does not have a dramatic impact on overall correlations. The observed correlations by subset also suggest that lower resolution structures do not impede antibody-antigen affinity prediction by these functions, as the subset with the lowest resolution (>3 Å) showed some of the highest Pearson correlation coefficients (R = 0.327-0.389), which were occasionally higher than correlations with the highest resolution subset (<=2 Å; R = 0.331-0.352). The lowest correlations were found for cases between 2.5 and 3 Å (R = 0.218-0.257), which may be due to its size (N = 145) and the greater likelihood of including cases that show higher deviations from the overall trend. Overall, these results show that correlations between  $\Delta G$  values and affinity predictors were largely unaffected by crucial characteristics of structure and affinity measurement, demonstrating the consistency and quality of cases in this affinity dataset.



**Figure 6.3 Predictive performance of existing scoring functions.** Correlation plots for prediction of  $\Delta G$  values by scoring functions (A) REF15 (R = 0.29), (B) beta16 (R = 0.29), (C) IRAD (R = 0.30), and (D) ZRANK (R = 0.29). The trend for each plot is shown as a red line.

Measurement	REF15	beta16	IRAD	ZRANK	# of
method	correlation^	correlation^	correlation^	correlation^	cases
SPR	0.35***	0.338***	0.284***	0.287***	245
BLI	0.255*	0.282*	0.367***	0.317**	82
ITC	0.271	0.25	0.342*	0.39*	40
Other	0.069	0.064	0.295	0.239	29

Table 6.2 Scoring function correlations with affinity values by measurement method.

^P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01, \*\*\* < 0.001

Structure	REF15	beta16	IRAD	ZRANK	# of
resolution (Å)	correlation^	correlation^	correlation^	correlation^	cases
<= 2 Å	0.463***	0.44***	0.352***	0.331**	89
> 2 and $<= 2.5$ Å	0.27**	0.284**	0.30***	0.329***	122
> 2.5 and <= 3 Å	0.257**	0.252**	0.25**	0.218**	145
> 3 Å	0.327*	0.329*	0.359*	0.382*	40

Table 6.3 Scoring function correlations with affinity values by structure resolution.

^P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01, \*\*\* < 0.0016.3.4 Correlations of individual terms with  $\Delta G$  values

Individual scoring terms for tested affinity predictors were also calculated and assessed independently for correlations with  $\Delta G$  values. Since the IRAD scoring function contains individual terms from ZRANK (182), only REF15, beta16 and IRAD scoring terms were investigated directly to minimize redundancy. The REF15 scoring function contains eight individual terms that contributed to calculated binding score, less than half the total number of terms (204). Beta16 contains eleven terms that contributed to binding score, including several terms that overlap with REF15, highlighting the similarity of these scoring functions. Correlations between  $\Delta G$  values and assessed REF15 and beta16 terms are shown in **Table 6.4**. Unsurprisingly, Pearson correlation coefficients for individual terms were lower than the complete REF15 scoring function, with the change in van der Waals attractive force ( $\Delta$ fa atr) showing the highest correlation among individual terms (R = 0.23). Other terms exhibited lower but significant correlations with  $\Delta G$  values, including the change in sidechain-sidechain hydrogen bond potential ( $\Delta$ hbond sc), the change in electrostatics ( $\Delta$ fa elec), and the change in van der Waals repulsive force ( $\Delta$ fa rep). More terms in beta16 showed significant correlations, due to the addition of terms that measure the change in solvation ( $\Delta$ lk ball,  $\Delta$ lk ball iso). Interestingly, a weighted combination of these terms did not show significant correlations in REF15 scores ( $\Delta$ lk ball wtd),

possibly because this combination also included solvation components that did not show significant correlations with  $\Delta G$  values on their own ( $\Delta lk$  ball bridge,  $\Delta lk$  ball bridge uncpl).

Though the published IRAD scoring function contains eight terms (182), a total of 43 terms were output during scoring of antibody-antigen structures. The score for each term output by IRAD was also correlated with  $\Delta G$  values, finding significant correlations with 24 terms (Table 6.5). A ZRANK term measuring van der Waals attractive force (vdw atr) exhibited one of the highest Pearson correlation coefficients (R = 0.21), corroborating the findings of the correlation for  $\Delta fa$  atr in REF15. Eight other terms output by IRAD showed significant positive correlations and were generally atom-based or residue-based potentials, including IFACE (iface; (477)), potential 2 from Zhang et al. (pot2atot; (471)), and Tobi-Bahar (pottbtot; (405)). All other terms showed significant negative correlations with  $\Delta G$  values, and these terms tended to be counts output during the calculation of potentials. Of the terms included in the existing IRAD scoring function, only vdw atr, elec s (a combination of elec sra and elec srr), iface, and potctot showed significant correlations with  $\Delta G$  values, suggesting that a reexamined and optimized set of terms from IRAD could improve prediction of antibody-antigen affinities in this dataset. Since some terms output by IRAD had similar characteristics and showed very similar correlations, we assessed how the scores of individual terms correlated with each other through hierarchical clustering and visualization (Figure 6.4). A heatmap of this correlation data highlights the similarity in scores for several large groups of terms output by IRAD, including one with 12 terms of potentials and another set of 18 terms primarily consisting of counts. Terms with significant individual correlations with  $\Delta G$  values also tended to cluster together. Surprisingly, the vdw atr term from ZRANK showed high correlations with the counts of interface atoms in specified distance cutoffs (acont4-6). Short-range electrostatics terms clustered together, but also with the solvlk term, which measures solvation and

not electrostatics in interfaces. These clustering results show both intuitive and unexpected similarities in the scores of terms output by IRAD, demonstrating patterns in these output terms and potentially informing optimization of affinity prediction using these terms in combination.

REF15 scoring	Pearson	beta16 scoring term	Pearson
term	correlation^		correlation^
$\Delta fa_atr$	0.23***	$\Delta fa_atr$	0.23***
$\Delta$ hbond_sc	0.18***	$\Delta$ hbond_sc	0.18***
$\Delta fa_elec$	0.12*	$\Delta lk_ball_iso$	0.14**
Δfa_rep	-0.11*	$\Delta lk_ball$	-0.13**
$\Delta fa_{sol}$	-0.08	$\Delta$ fa_elec	0.11*
$\Delta$ hbond_bb_sc	0.05	∆fa_rep	-0.11*
$\Delta lk_ball_wtd$	-0.04	$\Delta fa_{sol}$	-0.08
$\Delta$ hbond_lr_bb	-0.03	$\Delta lk_ball_bridge_uncpl$	0.06
		$\Delta$ hbond_bb_sc	0.05
		$\Delta$ lk_ball_bridge	0.04
		∆hbond_lr_bb	-0.03

 Table 6.4 Correlations of REF15 and beta16 scoring terms with affinity values.

^P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01, \*\*\* < 0.001

<b>Fable 6.5 Correlations of IRAD and ZRAN</b>	K scoring terms	with affinity values.
--	-----------------	-----------------------

Scoring term	Pearson correlation^
npot2atot	-0.24***
ZRANK vdw_atr	0.22***
acont4	-0.21***
nresa	-0.21***
acont6	-0.20***
acont5	-0.20***
nresc	-0.18***
nresp	-0.17***
nprestb	-0.17***
npotctot/npotllstot&	-0.16**
npot1ptot	-0.16**
pot2atot	0.15**
pottbtot	0.15**
iface	0.14**
npot2ptot	-0.13**
potg2tot	-0.13**
solvlk	0.13**
ZAPP elec_sra_x	0.12*
ZRANK elec_sra	0.12*

ZAPP elec_srr_x	-0.11*
potllstot	0.11*
potctot	0.11*
ZRANK elec_srr	-0.10*

 $^{P}$ -value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01, \*\*\* < 0.001 <sup>&</sup>Numbers of interactions counted by npotctot and npotllstot were identical for all complex structures, and resulted in identical Pearson correlation coefficients



Figure 6.4 Heatmap of correlations between terms output by IRAD. Listed terms were ordered through hierarchical clustering of squared Pearson correlation coefficients between scores. Black lines on the y-axis of the heatmap highlight terms that have a significant correlation with  $\Delta G$  values. A black line shows that all terms within the range of the black line share a significant correlation (p < 0.05) with  $\Delta G$  values.

### 6.3.5 Selection and retraining of input terms for antibody-antigen affinity prediction

Analysis of correlations between antibody-antigen affinity values and existing scoring functions revealed the potential to improve predictions of antibody-antigen affinities following optimization of the weights for individual scoring terms. As a proof of concept for utilizing this dataset to improve antibody-antigen  $\Delta G$  prediction, we utilized several regression methods in R to investigate optimized sets of scoring terms from REF15, beta16, IRAD, and ZRANK, either from one function or a combination of these functions. Eight sets of terms were implemented, and the scores of these terms on affinity cases not classified as high HETATM (N = 396) were directly compared with  $\Delta G$  values through multilinear regression. Four of these sets consisted of input terms from a single scoring function, with one retrained model for each scoring function. Two other sets of terms built on the integration of ZRANK terms into the published IRAD scoring function by adding other potentials output by IRAD ("IRAD + pot") or all other terms output by IRAD ("IRAD + pot/count"). Lastly, two more sets directly combined REF15 terms either with potentials output by IRAD ("REF15 + pot") or with potentials and counts output by IRAD ("REF15 + pot/count"). It should be noted that terms from Rosetta scoring functions and ZRANK terms were not combined with each other in a set, as the characteristics of included terms (e.g. van der Waals, electrostatics) substantially overlapped.

On each set of terms, we performed stepwise linear regression analysis to select a reduced set of terms with individual weights (**Table 6.6**). The number of input terms varied widely by model, from seven or eight terms in retrained existing functions to several dozen terms in composite sets. In all stepwise regression analyses, at least two terms in the input set were not selected in the retrained model, and the number of selected terms (not including the intercept) ranged from four (REF15, beta16) to 13 (IRAD + pot/count, REF15 + pot/count). Since the scores of input terms were not centered or scaled prior to analysis, the weights determined by stepwise regression are not indicative of the importance of individual terms. Interestingly, some individual terms with significant correlations with  $\Delta G$  values were not included in retrained models, while other terms with lower correlations such as  $\Delta fa_s$  and prestb were selected in multiple models.

A few models also contain terms output by IRAD that are correlated as shown in **Figure 6.4**, suggesting that multicollinearity between terms may be present following regression, especially in retrained models with higher numbers of terms.

Each set of terms selected by stepwise regression was then tested for antibody-antigen affinity prediction through several variations of cross-validation (leave-one-out, 5-fold, 10-fold). The root-mean-square error (RMSE) and Pearson correlation coefficients were computed and compared with the same metrics of the existing scoring functions, showing that several retrained models yielded improvements in  $\Delta G$  prediction (**Table 6.7**). Models with terms from several sources tended to show the most improvement, though an increase in number of terms may contribute to this trend. However, retrained REF15 and beta16 models with just four selected terms also showed lower RMSE and higher Pearson correlation coefficients following cross-validation. These models contain two terms ( $\Delta$ fa atr,  $\Delta$ hbond sc) with highly significant correlations with  $\Delta G$  values, suggesting that even a minimal set of impactful terms can increase predictive performance when retrained. Two composite models also included the selected terms from REF15/beta16, along with four potentials output by IRAD (REF15 + pot) or with  $\Delta$ fa rep, five potentials output by IRAD, and three counts output by IRAD (REF15 + pot/count). Another composite model with greater predictive performance (IRAD + pot/count) did not include REF15 terms as input, but stepwise regression still selected several electrostatics and solvation terms, suggesting that these types of terms may be important components for antibody-antigen affinity prediction.

Following these results, predictions from leave-one-out cross-validation were plotted for the top performing retrained models (**Figure 6.5**). The RMSE and correlation coefficients produced by these models were comparable to graph-based antibody-antigen prediction models trained on a dataset of similar size (214), suggesting that training scoring functions with this antibody-antigen affinity dataset can meaningfully improve predictive performance. However, these retrained models still leave ample room for improvement of antibody-antigen affinity prediction, with substantial differences between predicted and actual  $\Delta G$  values shown for several cases in each correlation plot. Though correlations improved with retraining, the modest correlations of these models to  $\Delta G$  are exemplified by stark differences in the ranges of predicted and actual  $\Delta G$  values. Even the top retrained models contain just a four-log  $\Delta G$  range in predictions (-10 to -14 kcal/mol), far narrower than the nine-log  $\Delta G$  range in this dataset (-16 to -7 kcal/mol). The worst predictions of  $\Delta G$  were concentrated in cases with the highest ( $\Delta G < -14$  kcal/mol) and lowest ( $\Delta G > -10$  kcal/mol) affinities, suggesting that the retrained models are limited in their ability to recognize properties of antibody-antigen interfaces that lead to very strong or very weak affinities within the context of this dataset.

This result could be attributed to suboptimal selection of input terms by stepwise regression, and additional methods of regression were tested on the same sets of input terms to help address this possibility. Ridge, Least Absolute Shrinkage and Selection Operator (LASSO), and Elastic net regression were also used to develop retrained models. These regression methods can either penalize or remove multicollinear terms, which could be especially beneficial for reweighing the largest sets of input terms. Though models with removed terms and different weights were generated, training and cross-validation showed RMSE and correlation coefficients that did not match the top models from stepwise regression, and in some cases showed decreased correlations when compared to existing functions (**Table 6.8**). Models trained with ridge regression, where all input terms were kept but penalized for multicollinearity, showed the lowest correlations overall, especially when training the largest sets of input terms. LASSO and Elastic

net regression often converged following training and cross-validation, finding the same set of selected terms and similar correlations. The performance of these regression methods suggests that selection of terms by stepwise regression did not result in suboptimal retrained functions. Instead, the narrower range of  $\Delta G$  predictions may be due to limitations inherent to the scores output by these terms, which could be missing key characteristics in the energetics of antibody-antigen interfaces, despite measurements of van der Waals, solvation, and electrostatic forces. Overall, the top performing models were generating using stepwise regression, and the improvements in affinity prediction highlight the utility of this antibody-antigen affinity dataset in future algorithmic developments.

Individual	REF15	beta16	IRAD	ZRANK	IRAD	IRAD +	REF15	<b>REF15</b> +
term*					+ pot	pot/count	+ pot	pot/count
(Intercept)	-9.04	-9.03	-9.51	-9.61	-9.55	-9.51	-9.04	-8.99
ace	Х	Х	0.022	-	-	-	Х	Х
acont4	Х	Х	Х	Х	Х	-0.024	Х	-
elec_l	Х	Х	0.015	Х	Х	Х	Х	Х
elec_lra	Х	Х	Х	0.013	0.011	0.009	Х	Х
elec_lrr	Х	Х	Х	0.023	0.019	0.019	Х	Х
elec_s	Х	Х	0.005	Х	Х	Х	Х	Х
elec_sra	Х	Х	Х	0.005	0.005	0.006	Х	Х
elec_srr	Х	Х	Х	0.006	0.005	0.006	Х	Х
fa_atr	0.051	0.052	Х	Х	Х	Х	0.062	0.078
fa_elec	0.042	0.039	Х	Х	Х	Х	0.032	0.031
fa_rep	-	-	Х	Х	Х	Х	-	0.059
fa_sol	0.057	0.057	Х	Х	Х	Х	0.056	0.068
hbond_sc	0.136	0.145	Х	Х	Х	Х	0.136	0.146
iface	Х	Х	-	Х	-	-0.001	Х	Х
npot1atot	Х	Х	Х	Х	Х	-	Х	0.08
npot2atot	Х	Х	Х	Х	Х	-0.064	Х	-
npotgtot	Х	Х	Х	Х	Х	0.074	Х	0.062
nresa	Х	Х	Х	Х	Х	-	Х	-0.036
pot2atot	Х	Х	Х	Х	0.064	0.047	0.09	0.082
potctot	Х	Х	0.025	Х	-	-	-	-
potg2tot	Х	Х	Х	Х	-	-0.016	-	-0.013
potllstot	Х	Х	Х	Х	-	-	-0.027	-0.035
presa	Х	Х	Х	Х	-0.09	-0.104	-	-
presc	Х	Х	Х	Х	-	-	-0.105	-0.103
prestb	Х	Х	Х	Х	0.119	0.12	0.128	0.125
solvlk	X	X	X	Х	-	-0.063	Х	X
vdw_atr	Х	Х	0.016	0.025	0.021	-	Х	X

Table 6.6 Terms and weights of retrained models selected by stepwise regression.

\*For each term, the weight given by the model following stepwise regression is shown. Terms that were not present in one or more models are indicated as "X" (this term was not included as input for this model) or "-" (this term was included as input for this model, but was not selected during stepwise regression).

Linear model	Validation	RMSE^	Pearson
			correlation^
REF15 bind score	N/A	1.62	0.29
beta16 bind score	N/A	1.62	0.29
IRAD score	N/A	1.61	0.3
ZRANK score	N/A	1.62	0.29
REF15	Leave-one-out	1.61	0.31
REF15	5-fold cross	1.6	0.34
REF15	10-fold cross	1.6	0.35
beta16	Leave-one-out	1.61	0.31
beta16	5-fold cross	1.6	0.34
beta16	10-fold cross	1.6	0.35
IRAD	Leave-one-out	1.63	0.26
IRAD	5-fold cross	1.63	0.3
IRAD	10-fold cross	1.63	0.29
ZRANK	Leave-one-out	1.63	0.26
ZRANK	5-fold cross	1.64	0.3
ZRANK	10-fold cross	1.63	0.3
IRAD + pot	Leave-one-out	1.62	0.29
IRAD + pot	5-fold cross	1.62	0.32
IRAD + pot	10-fold cross	1.62	0.32
IRAD + pot/count	Leave-one-out	1.6	0.34
IRAD + pot/count	5-fold cross	1.59	0.37
IRAD + pot/count	10-fold cross	1.6	0.35
REF15 + pot	Leave-one-out	1.59	0.34
REF15 + pot	5-fold cross	1.59	0.36
REF15 + pot	10-fold cross	1.58	0.37
REF15 + pot/count	Leave-one-out	1.58	0.36
REF15 + pot/count	5-fold cross	1.57	0.38
REF15 + pot/count	10-fold cross	1.57	0.39

Table 6.7  $\Delta G$  prediction of models retrained through stepwise regression.

^RMSE and Pearson correlation coefficients are highlighted in bold if both metrics were improved for a retrained model when compared to existing scoring functions



**Figure 6.5 Predictive performance of top retrained models.** Correlation plots of  $\Delta G$  predictions for retrained models (A) IRAD + pot/count (R = 0.34), (B) REF15 + pot (R = 0.34), and (C) REF15 + pot/count (R = 0.36) from leave-one-out cross-validation. The trend for each plot is shown as a red line.

	RE	REF15 + pot		IRAD + pot/count			REF15 + pot/count		
Regression	RMSE	R^	# of	RMSE	R^	# of	RMSE	R^	# of
			terms			terms			terms
Ridge	1.63	0.27	22	1.64	0.28	43	1.63	0.29	38
LASSO	1.62	0.3	5	1.63	0.31	3	1.62	0.31	4
Elastic net	1.62	0.3	5	1.63	0.3	11	1.62	0.31	4

Table 6.8 Affinity prediction of models retrained with Ridge, LASSO, Elastic net regression following 5-fold cross-validation.

^Pearson correlation coefficient

#### 6.3.6 Performance of top retrained models on independent test set

Although models retrained through multilinear regression improved antibody-antigen affinity predictions within the presented affinity dataset, its ability to predict affinity of antibody-antigen complexes not used for training was uncertain. We tested the best performing retrained models (IRAD + pot/count, REF15 + pot, and REF15 + pot/count) on an independent set of experimental affinity values and compared the resultant correlations with predictions made by existing scoring functions REF15 and IRAD. This independent test set contains the affinity values of VRC01-class intermediate antibodies to HIV-1 gp120 from six different isolates (481). Each affinity value was paired with a model of the gp120-antibody complex, which was then scored by each retrained model, REF15, and IRAD. Neutralization data of modeled complexes were also available, but correlations to retrained models were examined on a smaller scale due to a reduced dataset and concerns about distance between model and template sequences. Correlation coefficients between retrained model scores and affinity values (R = 0.3-0.41) were lower than REF15 (R = 0.56) and IRAD (R = 0.54) for the entire dataset (N = 360), with all correlations shown in **Table 6.9**.

The use of this independent set comes with several caveats, as models trained on highresolution antibody-antigen structures may be less effective in predicting the affinity modeled complexes, which also may have modeling inaccuracies due to varying degrees of sequence identity with the structural template. Interestingly, correlations between neutralization data for isolate 45 01dG5 and scores of retrained models REF15 + pot and REF15 + pot/count were higher than correlations established by REF15 and IRAD (Table 6.10). The predictions of these models may be more robust or generalizable than IRAD + pot/count, which showed little to no correlation between  $\Delta G$  predictions and antibody neutralization. Though this dataset is small (N = 48), these improvements in correlations suggest that retrained models may better distinguish between a range of antibody-antigen interfaces, which in this case differed widely by neutralization of 45 01dG5. The plots for REF15 + pot and REF15 + pot/count showed their significant correlations to neutralization data, but also revealed individual complex scores that deviate from these trends (Figure 6.6). In both plots, complexes for antibodies that best neutralize 45 01dG5 (log IC50 < -1) have a substantial range of predicted  $\Delta G$  values around the trendline, showing that these retrained models could not consistently distinguish between complexes with the highest neutralization values and complexes with the lowest neutralization values (log IC50 > 0) in this dataset. Despite yielding higher correlations with neutralization data, REF15 + pot and REF15 + pot/count plots demonstrate that these improvements on existing scoring functions are quite limited. Overall, predictive performance of retrained models on an independent test set showed little to no improvements over existing scoring functions, showing the importance of testing affinity predictors carefully and extensively with data withheld during training.

Isolate	REF15	IRAD +	REF15 +	REF15	IRAD	% identity&
	+ pot	pot/count	pot/count	bind score	score	
dG5	0.39**	0.43***	0.46***	0.62***	0.56***	95
dG5_K278T	0.36**	0.36**	0.37**	0.64***	0.6***	95
dH5	0.33**	0.27*	0.34**	0.5***	0.58***	88
BG505	0.33**	0.23	0.31*	0.51***	0.46***	78
BG505_T278A	0.55***	0.48***	0.55***	0.64***	0.6***	78
426c	0.53***	0.49***	0.46***	0.64***	0.56***	73
All	0.36***	0.3***	0.4***	0.56***	0.54***	N/A

Table 6.9 Correlations of top retrained models with affinities in independent test set.

^P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01, \*\*\* < 0.001. All other tested correlations were not significant.

<sup>&</sup>Percent identity of gp120 structure template sequence to the isolate gp120 sequence that was modeled in complex with antibody structures.

Table 6.10 Correlations of top retrained models with neutralization data of 45\_01dG5 isolate in independent test set.

Affinity predictor	Pearson correlation^
REF15 + pot	0.45**
IRAD + pot/count	-0.1
REF15 + pot/count	0.29*
REF15 (existing function)	0.21
IRAD (existing function)	0.13

^P-value significance of all correlation coefficients is shown in asterisks. \* < 0.05, \*\* < 0.01. All other tested correlations were not significant.



Figure 6.6 Predictions of retrained models with significant correlations to 45\_01dG5 neutralization data. Correlation plots between 45\_01dG5 neutralization data and  $\Delta G$  predictions of gp120-antibody complex models from (A) REF15 + pot (R = 0.45) and (B) REF15 + pot/count (R = 0.29). The trend for each plot is shown as a red line.

### 6.4 Discussion

We have collected and annotated a large dataset of 401 antibody-antigen affinities and corresponding structures. This dataset represents a diverse group of antibody types, antigen targets, and calculated  $\Delta G$  values. We found modest correlations between the  $\Delta G$  values and scores from four scoring functions, suggesting that optimization of terms in these functions could improve antibody-antigen affinity prediction. To examine the potential to improve predictions, we trained and tested linear regression models of input terms from Rosetta scoring functions REF15 and beta16, docking model scoring functions IRAD and ZRANK, or a combination of both. Retraining and cross-validation of input terms generated several linear models with improved correlations with  $\Delta G$  values in the affinity dataset, and the largest improvements tended to occur when retraining models with input terms from multiple sources. Finally, top retrained models showed lower correlations than existing scoring functions with affinity values of modeled HIV gp120antibody complexes in an independent test set, but two models did show higher correlations with a limited set of neutralization values. This in-depth regression analysis provided an example of how this curated set of affinities can be utilized for future development of algorithms for antibodyantigen  $\Delta G$  prediction.

The antibody-antigen affinity dataset presented here represents a dedicated effort to compile and annotate high-quality affinities and structures, but the tremendous diversity of this dataset poses clear challenges to reliable affinity prediction. Given the modest correlations observed in this study and elsewhere, it is possible that these challenges may stem in part from limitations inherent to including data with variable affinity measurement and structure determination, both in method and relative quality. Apart from different methods of affinity measurement, variability in documentation of affinity measurements such as proper controls in the literature may make reported affinities in this dataset less reliable (485). The potential reliability of measurements was not directly evaluated in this study, but the corresponding literature used to confirm antibody-antigen affinities could be assessed in future work to recommend best or common practices for reporting experimental affinities. At the same time, inconsistent documentation in the literature prevented this curated dataset from growing even larger, potentially offering a more robust dataset. We found numerous cases with an antibody-antigen structure and affinity value, but that could not be included because there was no documented temperature for the affinity measurement in the literature. Based on this experience, it would be beneficial for future studies to include a recorded temperature when presenting affinity values, as adopting this standard would allow for more  $\Delta G$  values to be calculated accurately and for more cases to be included in subsequent updates of this affinity dataset.

In addition, the nature of structure determination and interface properties for a given complex may have led to aberrant predicted scores, which could affect the strength of correlations with experimental  $\Delta G$  values. Metrics of structure quality such as Rfree and Ramachandran outliers were included to annotate cases in the affinity dataset, but these metrics were not explicitly tested for their possible contributions to changes in predictive performance. Though refinement with FastRelax alone may minimize or resolve aberrant determination of antibody-antigen interfaces that could affect scoring by affinity predictors, any applications of this affinity dataset may filter cases based on quality metrics prior to training, offering flexibility to various demands of research efforts in algorithmic development. Future research on the affinity dataset could also explore utilizing PDB\_REDO, a repository and server that aims to optimize PDB structures by correcting crystallographic errors and improving fits to electron density (486). Analyzing optimized antibody-antigen interfaces could reduce the risk of spurious scores due to variable

structure quality, which may better capture the utility of tested scoring functions for affinity prediction. In addition, future research should test affinity predictors in more practical applications useful for therapeutic design, such as predicting affinities of high-quality antibody-antigen complex models from *in silico* screening.

Despite retraining of scoring terms with this affinity dataset, cross-validation of these multilinear regression models yielded minor improvements to the modest Pearson correlation coefficients of existing scoring functions. These results highlight potential limitations to this method of antibody-antigen affinity prediction. First, structural analysis and scoring of antibodyantigen interfaces by affinity predictors only involves the bound complex. However,  $\Delta G$  values stem from thermodynamic free energy shifts from an unbound-unbound complex state to the bound interface, and the strength of  $\Delta G$  is determined by the free energy demands in forming a particular interaction, which can be influenced by associated conformational changes (217, 220) and by entropy loss through desolvation and the formation of hydrogen bonds (487). Considering just the bound state in affinity prediction may fail to capture important thermodynamic properties, adding substantial noise to  $\Delta G$  predictions and affecting the proper training of multilinear regression models. Second, complex scoring functions that can contain colinear terms may be at a higher risk of overtraining on antibody-antigen affinity prediction, as suggested by the lower correlations of top retrained models to experimental  $\Delta G$  values in the independent test set. This risk of overtraining has been previously discussed when combining separate potentials into a function for selection of near-native docking models (488). To mitigate overtraining, future developments of predictive scoring functions could implement additional assessments of trained models, such as testing correlations on a portion of the affinity dataset withheld from training and cross-validation. Regression methods could also be used only on a set of input terms more restricted by evidence of collinearity, where a clustering analysis of correlations between terms (see **Figure 6.4**) may reduce the number of terms utilized for training. Despite these suggestions for assessing future scoring functions with this dataset, there are caveats regarding the prediction of affinities in the independent test set that should be mentioned. Linear models trained on high-resolution and nonredundant antibody-antigen structures may naturally be less adept when scoring modeled complexes, which have various levels of sequence identity between model and template that could result in modelling errors. At the same time, distinguishing the affinities of highly similar antibodies to the same epitope could require a different set of scoring terms that may not have been selected in top retrained models. Overall, considering the applicability of trained models to specific situations of antibody-antigen affinity prediction is crucial for any future developments.

Though optimization of a scoring function composed of linearly weighted terms may improve correlations with antibody-antigen affinities, non-linear models generated from various machine learning methods may be better suited to affinity predictions. Several algorithms have already shown promise in  $\Delta G$  and  $\Delta \Delta G$  prediction (206, 214, 215), providing additional templates for affinity prediction algorithms that can be trained on this affinity dataset. The predictive performance of multilinear models could also improve when training scoring terms that can better distinguish antibody-antigen interface properties that lead to the wide range of observed  $\Delta G$ values. Future efforts in modeling with this affinity dataset could incorporate an even broader set of scoring potentials, such as those listed in CCharPPI (400, 488), as well as antibody-specific interface features that can be output from Rosetta (159). A weighted scoring function that improves predictions of  $\Delta G$  values would become a valuable tool for understanding and characterizing antibody-antigen interactions *in silico*. However, the presented models were not tested on  $\Delta\Delta G$ prediction, which could have demonstrated additional utility for therapeutic antibody design efforts to increase affinity to a desired antigen using a complex structure as a template. These retrained models, or any that utilize the affinity dataset, should also be evaluated for  $\Delta\Delta G$  prediction through validation with one or more available datasets (209, 210). As examined with CSM-AB (214), top performing affinity prediction models could also be tested in ranking of docking predictions to see if improvements in affinity prediction translated to protein docking. In testing affinity predictors on docking success, this research may help determine if docking and affinity prediction can be addressed with one algorithm or may require different foundations for algorithm development. In any case, this affinity dataset provides an expanded training set that can help advance antibody-antigen affinity prediction.

## Chapter 7: Summary and future directions

The research in this thesis examined approaches of vaccine design and computational analysis for HCV and expanded the resources available for computational modeling and prediction of all antibody-antigen interactions. We were able to design and validate multiple iterations of a soluble and secreted HCV vaccine candidate, providing new insights into how HCV E1E2 heterodimer can be scaffolded to maintain native-like antigenicity and immunogenicity while avoiding the more difficult process of extracting E1E2 from the viral membrane. Computational analysis of antibody interactions with HCV led to novel predictions of sequence contributions to neutralization and heterodimerization, which could have broad implications for future directions of HCV vaccine design. This interest in antibody-antigen interactions was also applied to the development of an antibody-antigen benchmark, expanding resources for the entire research community to facilitate improved antibody-antigen docking and affinity predictions.

The vaccine design work presented in this thesis (chapters 2 and 3) represents an intriguing proof of concept that is now being pursued as an HCV vaccine candidate. The success of certain scaffolded E1E2 assemblies can provide insights into the determinants of proper glycoprotein presentation in a heterodimeric and hetero-hexameric state, despite the recent report of the E1E2 heterodimer structure. In collaboration with Alexander Ploss and others, we have begun to quantify the immunogenicity of top sE1E2 vaccine candidates, with promising results showing an increase in the neutralization of heterologous isolates compared with mbE1E2 (489). Multimerization or multivalent scaffolding of these sE1E2 designs is also being explored, and the current success of some designs provides a template for targeted HCV scaffold designs, either through fixed backbone or *de novo* techniques. As discussed in chapter 2, we expect to assess the viability of sE1E2 as a platform for HCV isolates from multiple genotypes, with incorporation of additional

isolates and consensus sequences shown in chapter 3. Since sE1E2 is easier to produce while still recognized by HCV antibodies, this antigen may become a useful resource for the HCV community in further characterization of HCV antigenicity and immunogenicity.

Research involving predictions of important HCV E1E2 residues based on computational analysis (chapter 4) can be expanded in several directions. We can engage with collaborators to have our long list of predicted polymorphisms and E1E2 interface residues be tested and validated experimentally. Verifying the neutralization change for predicted polymorphisms and a loss of E1E2 assembly for predicted heterodimerization residues would help elucidate the sequence changes most impactful to vaccine design and determine possible mechanisms of immune evasion. We can also utilize additional mutagenesis or neutralization data to make more predictions using the same techniques, with any overlap in predicted contributors potentially providing further support for those residues. Validated polymorphisms can be tested against a larger set of antibodies to elucidate the true breadth of neutralization change, which could also offer insights into the actual mechanism, be it a change in the glycan shield, a shift in glycoprotein dynamics, or an alteration of receptor binding. In addition, the impact of predicted polymorphisms on E1E2 flexibility and dynamics can be addressed with molecular dynamics simulations on a small set of polymorphisms found to induce the largest changes in antibody neutralization.

An updated antibody-antigen benchmark and expanded antibody-antigen affinity dataset (chapters 5 and 6) can lead to some exciting developments in docking and affinity prediction while providing useful resources for the community. We plan to update both datasets frequently, keeping pace with rapidly accumulating structure and affinity data. With more cases available, development of protein-protein docking or affinity prediction algorithms will utilize a more robust and diverse set of prediction challenges, which will likely result in improved predictors. We can work to better understand the differential docking and affinity prediction success observed, possibly revealing trends in performance on the benchmark and defining the limitations of current algorithms. The lab is also taking benchmark structures in new directions of prediction, modeling some complexes in the benchmark with AlphaFold-Multimer (490, 491) to assess and potentially improve its ability to correctly predict antibody-antigen interfaces. This research could result in improved docking algorithms, either by optimizing current algorithms such as ZDOCK for antibody-antigen docking, or by identifying native antibody-antigen interfaces through deep learning. We can continue to improve predictions of  $\Delta G$  using the affinity dataset through integration of more scoring functions or potentials, as well as the application of non-linear predictive methods with machine learning. This research can inspire further developments in assembling and presenting benchmarks, incorporating key computational techniques such as antibody modeling and design or  $\Delta\Delta G$  prediction that allow for testing of current algorithms with a curated and updated set of cases.

# **Publication Information**

## Peer-reviewed publications related to the dissertation:

- 1. **Guest JD**, Wang R, Elkholy KH, Chagas A, Chao KL, Cleveland TE, Kim YC, Keck Z-Y, Marin A, Yunus AS, Mariuzza RA, Andrianov AK, Toth EA, Foung SKH, Pierce BG, Fuerst TR. "Design of a native-like secreted form of the hepatitis C virus E1E2 heterodimer." *PNAS* (2021) 118(3): e2015149118.
- Guest JD\*, Vreven T\*, Zhou J, Moal I, Jeliazkov JR, Gray JJ, Weng Z, Pierce BG. "An expanded benchmark for antibody-antigen docking and affinity prediction reveals insights into antibody recognition determinants." *Structure* (2021) 29(6): 606-621.e5. *\*Joint first authors*

### Other peer-reviewed publications:

- 1. Wang R, Suzuki S, **Guest JD**, Heller B, Andrianov AK, Marin A, Mariuzza RA, Keck Z-Y, Foung SKH, Yunus AS, Pierce BG, Toth EA, Ploss A, Fuerst TR. "Induction of broadly neutralizing antibodies using a secreted form of the hepatitis C virus E1E2 heterodimer as a vaccine candidate." *PNAS* (2022) 119(11): e2112008119.
- Wu D, Kolesnikov A, Yin R, Guest JD, Gowthaman R, Shmelev A, Serdyuk Y, Efimov GA, Pierce BG, Mariuzza RA. "Structural assessment of two HLA-A2-restricted SARS-CoV-2 spike epitopes by public and private T cell receptors." *Nat Commun.* (2022) 13(1): 19.
- 3. **Guest JD**, Pierce BG. "Structure-Based and Rational Design of a Hepatitis C Virus Vaccine." *Viruses* (2021) 13(5): 837.
- 4. Yin R, **Guest JD**, Taherzadeh G, Gowthaman R, Mittra I, Quackenbush J, Pierce BG. "Structural and energetic profiling of SARS-CoV-2 receptor binding domain antibody recognition and the impact of circulating variants." *PLoS Comput Biol.* (2021) 17(9): e1009380.
- Gowthaman R, Guest JD, Yin R, Adolf-Bryfogle J, Schief WR, Pierce BG. "CoV3D: a database of high resolution coronavirus protein structures." *Nucleic Acids Res.* (2021) 49(D1): D282-287.
- Salas JH, Urbanowicz RA, Guest JD, Frumento N, Figueroa A, Clark KE, Keck Z-Y, Cowton VM, Cole SJ, Patel AH, Fuerst TR, Drummer HE, Major M, Tarr AW, Ball JK, Law M, Pierce BG, Foung SKH, Bailey JR. "An antigenically diverse, representative panel of envelope glycoproteins for HCV vaccine development". *Gastroenterology* (2021) S0016-5085(21)03624-6.
- Lensink MF, Brysbaert G, Mauri T, et al., Guest JD, et al., Zhang S, Zhu X, Wodak SJ. "Prediction of protein assemblies, the next frontier: The CASP14-CAPRI experiment." *Proteins* (2021) 89(12): 1800-1823.
- 8. Pierce BG, Keck Z-Y, Wang R, Lau P, Garagusi KJ, Elkholy KH, Toth EA, Urbanowicz RA, **Guest JD**, Agnihotri P, Kerzic MC, Marin A, Andrianov AK, Ball JK, Mariuzza RA, Fuerst TR, Foung SKH. "Structure-Based Design of Hepatitis C Virus E2

Glycoprotein Improves Serum Binding and Cross-Neutralization." *J Virol.* (2020) 94(22): e00704-20.

- Lin S, Yang S, He J, Guest JD, Ma Z, Yang L, Pierce BG, Tang Q, Zhang Y-J. "Zika virus NS5 protein antagonizes type I interferon production via blocking TBK1 activation." *Virology* (2019) 527: 180-87.
- Urbanowicz RA, Wang R, Schiel JE, Keck Z-Y, Kerzic MC, Lau P, Rangarajan S, Garagusi KJ, Tan L, Guest JD, Ball JK, Pierce BG, Mariuzza RA, Foung SKH, Fuerst TR. "Antigenicity and Immunogenicity of Differentially Glycosylated HCV E2 Envelope Proteins Expressed in Mammalian and Insect Cells." *J Virol.* (2019) 93: e01403-18.
- 11. Keck Z-Y, Pierce BG, Lau P, Lu J, Wang Y, Underwood A, Bull RA, Prentoe J, Velazquez-Moctezuma R, Walker MR, Luciani F, Guest JD, Fauvelle C, Baumert TF, Bukh J, Lloyd AR, Foung SKH. "Broadly neutralizing antibodies from an individual that naturally cleared multiple hepatitis C virus infections uncover molecular determinants for E2 targeting and vaccine design." *PLoS Pathog.* (2019) 15(5): e1007772.
- Lensink MF, Brysbaert G, Nadzirin N, et al., Guest JD, et al., Honorato RV, Bonvin AMJJ, Wodak SJ. "Blind prediction of homo- and hetero-protein complexes: The CASP13-CAPRI experiment." *Proteins* (2019) 87(12): 1200-1221.
- 13. Guest JD, Pierce BG. "Computational Modeling of Hepatitis C Virus Envelope Glycoprotein Structure and Recognition." *Front. Immunol.* (2018) 9: 1117.

# Bibliography

- 1. C. Trepo, A brief history of hepatitis milestones. *Liver Int* **34 Suppl 1**, 29-37 (2014).
- 2. H. Laugi, Discovery of Hepatitis C Virus: 2020 Nobel Prize in Medicine. *Euroasian J Hepatogastroenterol* **10**, 105-108 (2020).
- 3. Q. Ding, M. von Schaewen, A. Ploss, The impact of hepatitis C virus entry on viral tropism. *Cell Host Microbe* **16**, 562-568 (2014).
- 4. S. H. Moosavy *et al.*, Epidemiology, transmission, diagnosis, and outcome of Hepatitis C virus infection. *Electron Physician* **9**, 5646-5656 (2017).
- 5. S. Zaltron, A. Spinetti, L. Biasi, C. Baiguera, F. Castelli, Chronic HCV infection: epidemiological and clinical relevance. *BMC Infect Dis* **12 Suppl 2**, S2 (2012).
- 6. F. Ansaldi, A. Orsi, L. Sticchi, B. Bruzzone, G. Icardi, Hepatitis C virus in the new era: perspectives in epidemiology, prevention, diagnostics and predictors of response to therapy. *World J Gastroenterol* **20**, 9633-9652 (2014).
- 7. R. K. Dhiman, G. S. Grover, M. Premkumar, Hepatitis C elimination: a Public Health Perspective. *Curr Treat Options Gastroenterol* **17**, 367-377 (2019).
- 8. T. Wilkins, M. Akhtar, E. Gititu, C. Jalluri, J. Ramirez, Diagnosis and Management of Hepatitis C. *Am Fam Physician* **91**, 835-842 (2015).
- 9. G. J. Dore, S. Bajis, Hepatitis C virus elimination: laying the foundation for achieving 2030 targets. *Nat Rev Gastroenterol Hepatol* **18**, 91-92 (2021).
- 10. WHO (2017) Global Hepatitis Report 2017. (World Health Organization, Geneva).
- 11. T. Kish, A. Aziz, M. Sorio, Hepatitis C in a New Era: A Review of Current Therapies. *P T* **42**, 316-329 (2017).
- 12. M. P. Manns et al., Hepatitis C virus infection. Nat Rev Dis Primers 3, 17006 (2017).
- 13. R. Schinazi, P. Halfon, P. Marcellin, T. Asselah, HCV direct-acting antiviral agents: the best interferon-free combinations. *Liver Int* **34 Suppl 1**, 69-78 (2014).
- A. Geddawy, Y. F. Ibrahim, N. M. Elbahie, M. A. Ibrahim, Direct Acting Anti-hepatitis C Virus Drugs: Clinical Pharmacology and Future Direction. *J Transl Int Med* 5, 8-17 (2017).
- 15. A. Al-Khazraji *et al.*, Identifying Barriers to the Treatment of Chronic Hepatitis C Infection. *Dig Dis* **38**, 46-52 (2020).

- 16. J. M. Pawlotsky, Hepatitis C Virus Resistance to Direct-Acting Antiviral Drugs in Interferon-Free Regimens. *Gastroenterology* **151**, 70-86 (2016).
- O. Falade-Nwulia, M. S. Sulkowski, A. Merkow, C. Latkin, S. H. Mehta, Understanding and addressing hepatitis C reinfection in the oral direct-acting antiviral era. *J Viral Hepat* 25, 220-227 (2018).
- 18. A. Yeung *et al.*, Population-level estimates of hepatitis C reinfection post scale-up of direct-acting antivirals among people who inject drugs. *J Hepatol* **76**, 549-557 (2022).
- 19. B. Roche, A. Coilly, J. C. Duclos-Vallee, D. Samuel, The impact of treatment of hepatitis C with DAAs on the occurrence of HCC. *Liver Int* **38 Suppl 1**, 139-145 (2018).
- 20. N. Scott *et al.*, The case for a universal hepatitis C vaccine to achieve hepatitis C elimination. *BMC Med* **17**, 175 (2019).
- 21. C. M. Walker, A. Grakoui, Hepatitis C virus: why do we need a vaccine to prevent a curable persistent infection? *Curr Opin Immunol* **35**, 137-143 (2015).
- 22. A. L. Cox, MEDICINE. Global control of hepatitis C virus. *Science* **349**, 790-791 (2015).
- 23. S. Plotkin, History of vaccination. Proc Natl Acad Sci USA 111, 12283-12287 (2014).
- 24. M. Ghattas, G. Dwivedi, M. Lavertu, M. G. Alameh, Vaccine Technologies and Platforms for Infectious Diseases: Current Progress, Challenges, and Opportunities. *Vaccines (Basel)* **9** (2021).
- 25. J. L. Hsu, A brief history of vaccines: smallpox to the present. *S D Med* **Spec no**, 33-37 (2013).
- 26. L. R. Baden *et al.*, Efficacy and Safety of the mRNA-1273 SARS-CoV-2 Vaccine. *N* Engl J Med **384**, 403-416 (2021).
- 27. C. B. Creech, S. C. Walker, R. J. Samuels, SARS-CoV-2 Vaccines. *JAMA* **325**, 1318-1320 (2021).
- 28. Y. Dong *et al.*, The way of SARS-CoV-2 vaccine development: success and challenges. *Signal Transduct Target Ther* **6**, 387 (2021).
- 29. D. R. Burton, L. M. Walker, Rational Vaccine Design in the Time of COVID-19. *Cell Host Microbe* **27**, 695-698 (2020).
- 30. C. Rueckert, C. A. Guzman, Vaccines: from empirical development to rational design. *PLoS Pathog* **8**, e1003001 (2012).
- 31. M. Pizza *et al.*, Identification of vaccine candidates against serogroup B meningococcus by whole-genome sequencing. *Science* **287**, 1816-1820 (2000).

- 32. V. Masignani, R. Rappuoli, M. Pizza, Reverse vaccinology: a genome-based approach for vaccine development. *Expert Opin Biol Ther* **2**, 895-905 (2002).
- R. Rappuoli, M. J. Bottomley, U. D'Oro, O. Finco, E. De Gregorio, Reverse vaccinology 2.0: Human immunology instructs vaccine antigen design. *J Exp Med* 213, 469-481 (2016).
- 34. B. S. Graham, M. S. A. Gilman, J. S. McLellan, Structure-Based Vaccine Antigen Design. *Annu Rev Med* **70**, 91-104 (2019).
- D. W. Kulp, W. R. Schief, Advances in structure-based vaccine design. *Curr Opin Virol* 3, 322-331 (2013).
- 36. S. B. Sable, M. Kalra, I. Verma, G. K. Khuller, Tuberculosis subunit vaccine design: the conflict of antigenicity and immunogenicity. *Clin Immunol* **122**, 239-251 (2007).
- 37. S. Mahanty, A. Prigent, O. Garraud, Immunogenicity of infectious pathogens and vaccine antigens. *BMC Immunol* **16**, 31 (2015).
- 38. C. L. Hsieh *et al.*, Structure-based design of prefusion-stabilized SARS-CoV-2 spikes. *Science* **369**, 1501-1505 (2020).
- 39. K. S. Corbett *et al.*, SARS-CoV-2 mRNA vaccine design enabled by prototype pathogen preparedness. *Nature* **586**, 567-571 (2020).
- 40. Q. D. Nguyen, K. Kikuchi, B. Maity, T. Ueno, The Versatile Manipulations of Self-Assembled Proteins in Vaccine Design. *Int J Mol Sci* **22** (2021).
- 41. J. S. McLellan *et al.*, Structure-based design of a fusion glycoprotein vaccine for respiratory syncytial virus. *Science* **342**, 592-598 (2013).
- 42. R. Derking, R. W. Sanders, Structure-guided envelope trimer design in HIV-1 vaccine development: a narrative review. *J Int AIDS Soc* **24 Suppl 7**, e25797 (2021).
- 43. X. Wu *et al.*, Rational design of envelope identifies broadly neutralizing human monoclonal antibodies to HIV-1. *Science* **329**, 856-861 (2010).
- 44. Y. Li *et al.*, Mechanism of neutralization by the broadly neutralizing HIV-1 monoclonal antibody VRC01. *J Virol* **85**, 8954-8967 (2011).
- 45. G. B. Stewart-Jones *et al.*, Trimeric HIV-1-Env Structures Define Glycan Shields from Clades A, B, and G. *Cell* **165**, 813-826 (2016).
- 46. R. W. Sanders *et al.*, A next-generation cleaved, soluble HIV-1 Env trimer, BG505 SOSIP.664 gp140, expresses multiple epitopes for broadly neutralizing but not non-neutralizing antibodies. *PLoS Pathog* **9**, e1003618 (2013).

- 47. H. The Lancet Gastroenterology, The hunt for a vaccine for hepatitis C virus continues. *Lancet Gastroenterol Hepatol* **6**, 253 (2021).
- 48. M. Naderi *et al.*, Hepatitis C virus and vaccine development. *Int J Mol Cell Med* **3**, 207-215 (2014).
- 49. K. Page *et al.*, Randomized Trial of a Vaccine Regimen to Prevent Chronic HCV Infection. *N Engl J Med* **384**, 541-549 (2021).
- 50. M. Lavie, A. Goffard, J. Dubuisson, Assembly of a functional HCV glycoprotein heterodimer. *Curr Issues Mol Biol* **9**, 71-86 (2007).
- 51. P. Falson *et al.*, Hepatitis C Virus Envelope Glycoprotein E1 Forms Trimers at the Surface of the Virion. *J Virol* **89**, 10333-10346 (2015).
- 52. B. D. Lindenbach, C. M. Rice, The ins and outs of hepatitis C virus entry and assembly. *Nat Rev Microbiol* **11**, 688-700 (2013).
- 53. A. Kumar *et al.*, Structural insights into hepatitis C virus receptor binding and entry. *Nature* **598**, 521-525 (2021).
- 54. K. G. Hadlock *et al.*, Human monoclonal antibodies that inhibit binding of hepatitis C virus E2 protein to CD81 and recognize conserved conformational epitopes. *J Virol* **74**, 10407-10416 (2000).
- 55. M. Law *et al.*, Broadly neutralizing antibodies protect against hepatitis C virus quasispecies challenge. *Nat Med* **14**, 25-27 (2008).
- 56. J. D. Guest, B. G. Pierce, Structure-Based and Rational Design of a Hepatitis C Virus Vaccine. *Viruses* **13** (2021).
- 57. J. D. Guest, B. G. Pierce, Computational Modeling of Hepatitis C Virus Envelope Glycoprotein Structure and Recognition. *Front Immunol* **9**, 1117 (2018).
- 58. T. B. Ruwona, E. Giang, T. Nieusma, M. Law, Fine mapping of murine antibody responses to immunization with a novel soluble form of hepatitis C virus envelope glycoprotein complex. *J Virol* **88**, 10459-10471 (2014).
- 59. L. Cao *et al.*, Functional expression and characterization of the envelope glycoprotein E1E2 heterodimer of hepatitis C virus. *PLoS Pathog* **15**, e1007759 (2019).
- 60. J. Prentoe *et al.*, Antigenic and immunogenic evaluation of permutations of soluble hepatitis C virus envelope protein E2 and E1 antigens. *PLoS One* **16**, e0255336 (2021).
- 61. A. T. de la Peña *et al.*, Structure of the hepatitis C virus E1E2 glycoprotein complex. *bioRxiv* 10.1101/2021.12.16.472992, 2021.2012.2016.472992 (2021).

- 62. B. G. Pierce, Z. Y. Keck, S. K. Foung, Viral evasion and challenges of hepatitis C virus vaccine development. *Curr Opin Virol* **20**, 55-63 (2016).
- 63. B. G. Pierce *et al.*, Global mapping of antibody recognition of the hepatitis C virus E2 glycoprotein: Implications for vaccine design. *Proc Natl Acad Sci U S A* **113**, E6946-E6954 (2016).
- 64. M. M. Santoro, C. F. Perno, HIV-1 Genetic Variability and Clinical Implications. *ISRN Microbiol* **2013**, 481314 (2013).
- 65. R. Pejchal *et al.*, A potent and broad neutralizing antibody recognizes and penetrates the HIV glycan shield. *Science* **334**, 1097-1103 (2011).
- 66. D. B. Smith (2022) (Flaviviridae Study Group, <u>https://talk.ictvonline.org/ictv\_wikis/flaviviridae/w/sg\_flavi/634/table-1---confirmed-hcv-genotypes-subtypes-march-2022</u>).
- 67. D. B. Smith *et al.*, Expanded classification of hepatitis C virus into 7 genotypes and 67 subtypes: updated criteria and genotype assignment web resource. *Hepatology* **59**, 318-327 (2014).
- 68. D. Moradpour, F. Penin, C. M. Rice, Replication of hepatitis C virus. *Nat Rev Microbiol* 5, 453-463 (2007).
- 69. M. H. Powdrill *et al.*, Contribution of a mutational bias in hepatitis C virus replication to the genetic barrier in the development of drug resistance. *Proc Natl Acad Sci U S A* **108**, 20509-20513 (2011).
- 70. W. Shi *et al.*, Recombination in hepatitis C virus: identification of four novel naturally occurring inter-subtype recombinants. *PLoS One* **7**, e41997 (2012).
- A. W. Tarr *et al.*, Genetic Diversity Underlying the Envelope Glycoproteins of Hepatitis C Virus: Structural and Functional Consequences and the Implications for Vaccine Design. *Viruses* 7, 3995-4046 (2015).
- 72. K. Tsukiyama-Kohara, M. Kohara, Hepatitis C Virus: Viral Quasispecies and Genotypes. *Int J Mol Sci* **19** (2017).
- 73. J. P. Messina *et al.*, Global distribution and prevalence of hepatitis C virus genotypes. *Hepatology* **61**, 77-87 (2015).
- 74. M. Lavie, X. Hanoulle, J. Dubuisson, Glycan Shielding and Modulation of Hepatitis C Virus Neutralizing Antibodies. *Front Immunol* **9**, 910 (2018).
- 75. J. Prentoe *et al.*, Hypervariable region 1 and N-linked glycans of hepatitis C regulate virion neutralization by modulating envelope conformations. *Proc Natl Acad Sci U S A* **116**, 10039-10047 (2019).

- 76. F. Helle *et al.*, The neutralizing activity of anti-hepatitis C virus antibodies is modulated by specific glycans on the E2 envelope protein. *J Virol* **81**, 8101-8111 (2007).
- 77. R. A. Urbanowicz *et al.*, Antigenicity and Immunogenicity of Differentially Glycosylated Hepatitis C Virus E2 Envelope Proteins Expressed in Mammalian and Insect Cells. *J Virol* **93** (2019).
- E. V. LeBlanc, Y. Kim, C. J. Capicciotti, C. C. Colpitts, Hepatitis C Virus Glycan-Dependent Interactions and the Potential for Novel Preventative Strategies. *Pathogens* 10 (2021).
- 79. F. Helle *et al.*, Role of N-linked glycans in the functions of hepatitis C virus envelope proteins incorporated into infectious virions. *J Virol* **84**, 11905-11915 (2010).
- 80. A. Goffard *et al.*, Role of N-linked glycans in the functions of hepatitis C virus envelope glycoproteins. *J Virol* **79**, 8400-8409 (2005).
- 81. T. R. Fuerst, B. G. Pierce, Z. Y. Keck, S. K. H. Foung, Designing a B Cell-Based Vaccine against a Highly Variable Hepatitis C Virus. *Front Microbiol* **8**, 2692 (2017).
- J. Prentoe, J. Bukh, Hypervariable Region 1 in Envelope Protein 2 of Hepatitis C Virus: A Linchpin in Neutralizing Antibody Evasion and Viral Entry. *Front Immunol* 9, 2146 (2018).
- 83. J. Prentoe, R. Velazquez-Moctezuma, S. K. Foung, M. Law, J. Bukh, Hypervariable region 1 shielding of hepatitis C virus is a main contributor to genotypic differences in neutralization sensitivity. *Hepatology* **64**, 1881-1892 (2016).
- 84. Z. Y. Keck *et al.*, Antibody Response to Hypervariable Region 1 Interferes with Broadly Neutralizing Antibodies to Hepatitis C Virus. *J Virol* **90**, 3112-3122 (2016).
- 85. N. A. Brasher *et al.*, B cell immunodominance in primary hepatitis C virus infection. *J Hepatol* **72**, 670-679 (2020).
- 86. Z. Y. Keck *et al.*, Analysis of a highly flexible conformational immunogenic domain a in hepatitis C virus E2. *J Virol* **79**, 13199-13208 (2005).
- 87. R. Gopal *et al.*, Probing the antigenicity of hepatitis C virus envelope glycoprotein complex by high-throughput mutagenesis. *PLoS Pathog* **13**, e1006735 (2017).
- 88. Z. Y. Keck *et al.*, Broadly neutralizing antibodies from an individual that naturally cleared multiple hepatitis C virus infections uncover molecular determinants for E2 targeting and vaccine design. *PLoS Pathog* **15**, e1007772 (2019).
- 89. E. H. Augestad *et al.*, Global and local envelope protein dynamics of hepatitis C virus determine broad antibody sensitivity. *Sci Adv* **6**, eabb5938 (2020).

- 90. E. H. Augestad, J. Bukh, J. Prentoe, Hepatitis C virus envelope protein dynamics and the link to hypervariable region 1. *Curr Opin Virol* **50**, 69-75 (2021).
- 91. J. L. M. Law *et al.*, Role of the E2 Hypervariable Region (HVR1) in the Immunogenicity of a Recombinant Hepatitis C Virus Vaccine. *J Virol* **92** (2018).
- 92. B. G. Pierce *et al.*, Structure-based design of hepatitis C virus E2 glycoprotein improves serum binding and cross-neutralization. *J Virol* 10.1128/JVI.00704-20 (2020).
- 93. Z. Y. Keck *et al.*, Non-random escape pathways from a broadly neutralizing human monoclonal antibody map to a highly conserved region on the hepatitis C virus E2 glycoprotein encompassing amino acids 412-423. *PLoS Pathog* **10**, e1004297 (2014).
- 94. A. W. Tarr *et al.*, Determination of the human antibody response to the epitope defined by the hepatitis C virus-neutralizing monoclonal antibody AP33. *J Gen Virol* **88**, 2991-3001 (2007).
- 95. H. Pantua *et al.*, Glycan shifting on hepatitis C virus (HCV) E2 glycoprotein is a mechanism for escape from broadly neutralizing antibodies. *J Mol Biol* **425**, 1899-1914 (2013).
- 96. Y. Li *et al.*, Structural basis for penetration of the glycan shield of hepatitis C virus E2 glycoprotein by a broadly neutralizing human antibody. *J Biol Chem* **290**, 10117-10125 (2015).
- 97. R. T. Chung *et al.*, Human monoclonal antibody MBL-HCV1 delays HCV viral rebound following liver transplantation: a randomized controlled study. *Am J Transplant* **13**, 1047-1054 (2013).
- 98. R. Velazquez-Moctezuma *et al.*, Mechanisms of Hepatitis C Virus Escape from Vaccine-Relevant Neutralizing Antibodies. *Vaccines (Basel)* **9** (2021).
- 99. R. Velazquez-Moctezuma, A. Galli, M. Law, J. Bukh, J. Prentoe, Hepatitis C Virus-Escape Studies for Human Monoclonal Antibody AR4A Reveal Isolate-Specific Resistance and a High Barrier to Resistance. *The Journal of infectious diseases* 219, 68-79 (2019).
- 100. Z. Y. Keck *et al.*, Mutations in hepatitis C virus E2 located outside the CD81 binding sites lead to escape from broadly neutralizing antibodies but compromise virus infectivity. *J Virol* **83**, 6149-6160 (2009).
- 101. G. J. Babcock *et al.*, High-throughput sequencing analysis of post-liver transplantation HCV E2 glycoprotein evolution in the presence and absence of neutralizing monoclonal antibody. *PLoS One* **9**, e100325 (2014).
- 102. J. R. Bailey *et al.*, Naturally selected hepatitis C virus polymorphisms confer broad neutralizing antibody resistance. *J Clin Invest* **125**, 437-447 (2015).

- R. El-Diwany *et al.*, Extra-epitopic hepatitis C virus polymorphisms confer resistance to broadly neutralizing antibodies by modulating binding to scavenger receptor B1. *PLoS Pathog* 13, e1006235 (2017).
- L. N. Wasilewski *et al.*, A Hepatitis C Virus Envelope Polymorphism Confers Resistance to Neutralization by Polyclonal Sera and Broadly Neutralizing Monoclonal Antibodies. J Virol 90, 3773-3782 (2016).
- 105. J. Y. Lee *et al.*, Apolipoprotein E likely contributes to a maturation step of infectious hepatitis C virus particles and interacts with viral envelope glycoproteins. *J Virol* **88**, 12422-12437 (2014).
- 106. M. T. Catanese *et al.*, Ultrastructural analysis of hepatitis C virus particles. *Proc Natl Acad Sci U S A* **110**, 9505-9510 (2013).
- M. Sidorkiewicz, Hepatitis C Virus Uses Host Lipids to Its Own Advantage. *Metabolites* 11 (2021).
- 108. C. Fauvelle *et al.*, Apolipoprotein E Mediates Evasion From Hepatitis C Virus Neutralizing Antibodies. *Gastroenterology* **150**, 206-217 e204 (2016).
- 109. F. Wrensch *et al.*, Hepatitis C Virus (HCV)-Apolipoprotein Interactions and Immune Evasion and Their Impact on HCV Vaccine Design. *Front Immunol* **9**, 1436 (2018).
- L. Riva, J. Dubuisson, Similarities and Differences Between HCV Pseudoparticle (HCVpp) and Cell Culture HCV (HCVcc) in the Study of HCV. *Methods Mol Biol* 1911, 33-45 (2019).
- 111. J. D. Duncan, R. A. Urbanowicz, A. W. Tarr, J. K. Ball, Hepatitis C Virus Vaccine: Challenges and Prospects. *Vaccines (Basel)* 8 (2020).
- 112. H. M. Shepard, G. L. Phillips, D. T. C, M. Feldmann, Developments in therapy with monoclonal antibodies and related proteins. *Clin Med (Lond)* **17**, 220-232 (2017).
- 113. R. L. Stanfield, I. A. Wilson, Antibody Structure. *Microbiol Spectr* 2 (2014).
- 114. I. Sela-Culang, V. Kunik, Y. Ofran, The structural basis of antibody-antigen recognition. *Front Immunol* **4**, 302 (2013).
- 115. A. K. Mishra, R. A. Mariuzza, Insights into the Structural Basis of Antibody Affinity Maturation from Next-Generation Sequencing. *Front Immunol* 9, 117 (2018).
- 116. I. Sela-Culang, S. Alon, Y. Ofran, A systematic comparison of free and bound antibodies reveals binding-related conformational changes. *J Immunol* **189**, 4890-4899 (2012).
- 117. Y. Barrios, P. Jirholt, M. Ohlin, Length of the antibody heavy chain complementarity determining region 3 as a specificity-determining factor. *J Mol Recognit* **17**, 332-338 (2004).
- 118. S. Muyldermans, Nanobodies: natural single-domain antibodies. *Annu Rev Biochem* **82**, 775-797 (2013).
- 119. I. Jovcevska, S. Muyldermans, The Therapeutic Potential of Nanobodies. *BioDrugs* **34**, 11-26 (2020).
- 120. L. S. Mitchell, L. J. Colwell, Analysis of nanobody paratopes reveals greater diversity than classical antibodies. *Protein Eng Des Sel* **31**, 267-275 (2018).
- 121. P. Deszynski *et al.*, INDI-integrated nanobody database for immunoinformatics. *Nucleic Acids Res* **50**, D1273-D1281 (2022).
- 122. X. Li *et al.*, Comparative Analysis of Immune Repertoires between Bactrian Camel's Conventional and Heavy-Chain Antibodies. *PLoS One* **11**, e0161801 (2016).
- 123. J. W. Stave, K. Lindpaintner, Antibody and antigen contact residues define epitope and paratope size and structure. *J Immunol* **191**, 1428-1435 (2013).
- 124. M. Lustrek *et al.*, Epitope predictions indicate the presence of two distinct types of epitope-antibody-reactivities determined by epitope profiling of intravenous immunoglobulins. *PLoS One* **8**, e78605 (2013).
- 125. M. H. Van Regenmortel, What is a B-cell epitope? *Methods Mol Biol* **524**, 3-20 (2009).
- 126. J. L. Sanchez-Trincado, M. Gomez-Perosanz, P. A. Reche, Fundamentals and Methods for T- and B-Cell Epitope Prediction. *J Immunol Res* **2017**, 2680160 (2017).
- 127. L. Potocnakova, M. Bhide, L. B. Pulzova, An Introduction to B-Cell Epitope Mapping and In Silico Epitope Prediction. *J Immunol Res* **2016**, 6760830 (2016).
- F. Chen, N. Tzarum, I. A. Wilson, M. Law, VH1-69 antiviral broadly neutralizing antibodies: genetics, structures, and relevance to rational vaccine design. *Curr Opin Virol* 34, 149-159 (2019).
- 129. Q. Wang, L. Zhang, Broadly neutralizing antibodies and vaccine design against HIV-1 infection. *Front Med* 14, 30-42 (2020).
- M. G. Pauthner, L. Hangartner, Broadly Neutralizing Antibodies to Highly Antigenically Variable Viruses as Templates for Vaccine Design. *Curr Top Microbiol Immunol* 428, 31-87 (2020).
- 131. T. Zhou, K. Xu, Structural Features of Broadly Neutralizing Antibodies and Rational Design of Vaccine. *Adv Exp Med Biol* **1075**, 73-95 (2018).
- 132. S. A. Griffith, L. E. McCoy, To bnAb or Not to bnAb: Defining Broadly Neutralising Antibodies Against HIV-1. *Front Immunol* **12**, 708227 (2021).

- 133. L. B. Shrestha, N. Tedla, R. A. Bull, Broadly-Neutralizing Antibodies Against Emerging SARS-CoV-2 Variants. *Front Immunol* **12**, 752003 (2021).
- 134. P. D. Kwong, J. R. Mascola, Human antibodies that neutralize HIV-1: identification, structures, and B cell ontogenies. *Immunity* **37**, 412-425 (2012).
- 135. N. M. C. De Jong, A. Aartse, M. J. Van Gils, D. Eggink, Development of broadly reactive influenza vaccines by targeting the conserved regions of the hemagglutinin stem and head domains. *Expert Rev Vaccines* **19**, 563-577 (2020).
- 136. D. Pinto *et al.*, Broad betacoronavirus neutralization by a stem helix-specific human antibody. *Science* **373**, 1109-1116 (2021).
- 137. P. Zhou *et al.*, A human antibody reveals a conserved site on beta-coronavirus spike proteins and confers protection against SARS-CoV-2 infection. *Sci Transl Med* **14**, eabi9215 (2022).
- 138. M. F. Good, S. K. Yanow, Cryptic epitope for antibodies should not be forgotten in vaccine design. *Expert Rev Vaccines* **15**, 675-676 (2016).
- 139. M. Yuan *et al.*, A highly conserved cryptic epitope in the receptor binding domains of SARS-CoV-2 and SARS-CoV. *Science* **368**, 630-633 (2020).
- 140. B. R. West *et al.*, Structural Basis of Pan-Ebolavirus Neutralization by a Human Antibody against a Conserved, yet Cryptic Epitope. *mBio* **9** (2018).
- J. Li *et al.*, Structural and Functional Characterization of a Cross-Reactive Dengue Virus Neutralizing Antibody that Recognizes a Cryptic Epitope. *Structure* 26, 51-59 e54 (2018).
- 142. A. S. Dingens *et al.*, High-resolution mapping of the neutralizing and binding specificities of polyclonal sera post-HIV Env trimer vaccination. *Elife* **10** (2021).
- 143. C. A. Cottrell *et al.*, Mapping the immunogenic landscape of near-native HIV-1 envelope trimers in non-human primates. *PLoS Pathog* **16**, e1008753 (2020).
- 144. A. R. Rees, Understanding the human antibody repertoire. *MAbs* 12, 1729683 (2020).
- 145. B. North, A. Lehmann, R. L. Dunbrack, Jr., A new clustering of antibody CDR loop conformations. *J Mol Biol* **406**, 228-256 (2011).
- J. Adolf-Bryfogle, Q. Xu, B. North, A. Lehmann, R. L. Dunbrack, Jr., PyIgClassify: a database of antibody CDR structural classifications. *Nucleic Acids Res* 43, D432-438 (2015).
- M. Dondelinger *et al.*, Understanding the Significance and Implications of Antibody Numbering and Antigen-Binding Surface/Residue Definition. *Front Immunol* 9, 2278 (2018).

- 148. B. D. Weitzner, R. L. Dunbrack, Jr., J. J. Gray, The origin of CDR H3 structural diversity. *Structure* 23, 302-311 (2015).
- M. L. Fernandez-Quintero *et al.*, Characterizing the Diversity of the CDR-H3 Loop Conformational Ensembles in Relationship to Antibody Binding Properties. *Front Immunol* 9, 3065 (2018).
- 150. R. J. Blackler *et al.*, Antigen binding by conformational selection in near-germline antibodies. *J Biol Chem* 10.1016/j.jbc.2022.101901, 101901 (2022).
- 151. W. Wang *et al.*, Conformational selection and induced fit in specific antibody and antigen recognition: SPE7 as a case study. *J Phys Chem B* **117**, 4912-4923 (2013).
- 152. B. D. Weitzner, J. J. Gray, Accurate Structure Prediction of CDR H3 Loops Enabled by a Novel Structure-Based C-Terminal Constraint. *J Immunol* **198**, 505-515 (2017).
- 153. J. Dunbar *et al.*, SAbDab: the structural antibody database. *Nucleic Acids Res* **42**, D1140-1146 (2014).
- 154. S. Ferdous, A. C. R. Martin, AbDb: antibody structure database-a database of PDBderived antibody structures. *Database (Oxford)* **2018** (2018).
- 155. C. T. Schoeder *et al.*, Modeling Immunity with Rosetta: Methods for Antibody and Antigen Design. *Biochemistry* **60**, 825-846 (2021).
- 156. J. Graves *et al.*, A Review of Deep Learning Methods for Antibodies. *Antibodies (Basel)* 9 (2020).
- 157. B. Abanades, G. Georges, A. Bujotzek, C. M. Deane, ABlooper: Fast accurate antibody CDR loop structure prediction with accuracy estimation. *Bioinformatics* 10.1093/bioinformatics/btac016 (2022).
- 158. J. A. Ruffolo, J. Sulam, J. J. Gray, Antibody structure prediction using interpretable deep learning. *Patterns (N Y)* **3**, 100406 (2022).
- 159. J. Adolf-Bryfogle *et al.*, RosettaAntibodyDesign (RAbD): A general framework for computational antibody design. *PLoS Comput Biol* **14**, e1006112 (2018).
- 160. J. C. Almagro et al., Antibody modeling assessment. Proteins 79, 3050-3066 (2011).
- A. Teplyakov *et al.*, Antibody modeling assessment II. Structures and models. *Proteins* 82, 1563-1582 (2014).
- 162. I. A. Vakser, Protein-protein docking: from interaction to interactome. *Biophys J* **107**, 1785-1793 (2014).
- 163. S. Y. Huang, Search strategies and evaluation in protein-protein docking: principles, advances and challenges. *Drug Discov Today* **19**, 1081-1096 (2014).

- 164. D. W. Ritchie, G. J. Kemp, Protein docking using spherical polar Fourier correlations. *Proteins* **39**, 178-194 (2000).
- R. Chen, L. Li, Z. Weng, ZDOCK: an initial-stage protein-docking algorithm. *Proteins* 52, 80-87 (2003).
- 166. D. Kozakov, R. Brenke, S. R. Comeau, S. Vajda, PIPER: an FFT-based protein docking program with pairwise potentials. *Proteins* **65**, 392-406 (2006).
- 167. D. W. Ritchie, V. Venkatraman, Ultra-fast FFT protein docking on graphics processors. *Bioinformatics* **26**, 2398-2405 (2010).
- H. A. Gabb, R. M. Jackson, M. J. Sternberg, Modelling protein docking using shape complementarity, electrostatics and biochemical information. *J Mol Biol* 272, 106-120 (1997).
- 169. V. Venkatraman, Y. D. Yang, L. Sael, D. Kihara, Protein-protein docking using regionbased 3D Zernike descriptors. *BMC Bioinformatics* **10**, 407 (2009).
- 170. M. Estrin, H. J. Wolfson, SnapDock-template-based docking by Geometric Hashing. *Bioinformatics* **33**, i30-i36 (2017).
- 171. J. J. Gray *et al.*, Protein-protein docking with simultaneous optimization of rigid-body displacement and side-chain conformations. *J Mol Biol* **331**, 281-299 (2003).
- 172. R. Brenke *et al.*, Application of asymmetric statistical potentials to antibody-protein docking. *Bioinformatics* **28**, 2608-2614 (2012).
- 173. N. A. Marze, S. S. Roy Burman, W. Sheffler, J. J. Gray, Efficient flexible backbone protein-protein docking for challenging targets. *Bioinformatics* **34**, 3461-3469 (2018).
- 174. A. May, M. Zacharias, Protein-protein docking in CAPRI using ATTRACT to account for global and local flexibility. *Proteins* **69**, 774-780 (2007).
- 175. M. Krol, R. A. Chaleil, A. L. Tournier, P. A. Bates, Implicit flexibility in protein docking: cross-docking and local refinement. *Proteins* **69**, 750-757 (2007).
- 176. I. H. Moal, P. A. Bates, SwarmDock and the use of normal modes in protein-protein docking. *Int J Mol Sci* **11**, 3623-3648 (2010).
- C. Dominguez, R. Boelens, A. M. Bonvin, HADDOCK: a protein-protein docking approach based on biochemical or biophysical information. *J Am Chem Soc* 125, 1731-1737 (2003).
- 178. A. Sircar, J. J. Gray, SnugDock: paratope structural optimization during antibody-antigen docking compensates for errors in antibody homology models. *PLoS computational biology* **6**, e1000644 (2010).

- 179. B. Pierce, Z. Weng, ZRANK: reranking protein docking predictions with an optimized energy function. *Proteins* **67**, 1078-1086 (2007).
- 180. B. Pierce, Z. Weng, A combination of rescoring and refinement significantly improves protein docking performance. *Proteins* **72**, 270-279 (2008).
- 181. N. Andrusier, R. Nussinov, H. J. Wolfson, FireDock: fast interaction refinement in molecular docking. *Proteins* **69**, 139-159 (2007).
- 182. T. Vreven, H. Hwang, Z. Weng, Integrating atom-based and residue-based scoring functions for protein-protein docking. *Protein Sci* **20**, 1576-1586 (2011).
- 183. I. H. Moal *et al.*, IRaPPA: information retrieval based integration of biophysical models for protein assembly selection. *Bioinformatics* 10.1093/bioinformatics/btx068 (2017).
- 184. M. Zacharias, Accounting for conformational changes during protein-protein docking. *Current opinion in structural biology* **20**, 180-186 (2010).
- C. Pons, S. Grosdidier, A. Solernou, L. Perez-Cano, J. Fernandez-Recio, Present and future challenges and limitations in protein-protein docking. *Proteins* 78, 95-108 (2010).
- T. Vreven *et al.*, Updates to the Integrated Protein-Protein Interaction Benchmarks: Docking Benchmark Version 5 and Affinity Benchmark Version 2. *J Mol Biol* 427, 3031-3041 (2015).
- 187. I. A. Vakser, Challenges in protein docking. Curr Opin Struct Biol 64, 160-165 (2020).
- 188. H. Madaoui, R. Guerois, Coevolution at protein complex interfaces can be detected by the complementarity trace with important impact for predictive docking. *Proc Natl Acad Sci U S A* **105**, 7708-7713 (2008).
- 189. J. Janin *et al.*, CAPRI: a Critical Assessment of PRedicted Interactions. *Proteins* **52**, 2-9 (2003).
- 190. J. Janin, Assessing predictions of protein-protein interaction: the CAPRI experiment. *Protein Sci* 14, 278-283 (2005).
- Y. Gao, D. Douguet, A. Tovchigrechko, I. A. Vakser, DOCKGROUND system of databases for protein recognition studies: unbound structures for docking. *Proteins* 69, 845-851 (2007).
- 192. P. J. Kundrotas *et al.*, Dockground: A comprehensive data resource for modeling of protein complexes. *Protein Sci* 27, 172-181 (2018).
- 193. H. Hwang, T. Vreven, J. Janin, Z. Weng, Protein-protein docking benchmark version 4.0. *Proteins* **78**, 3111-3114 (2010).

- 194. H. Hwang, B. Pierce, J. Mintseris, J. Janin, Z. Weng, Protein-protein docking benchmark version 3.0. *Proteins* **73**, 705-709 (2008).
- 195. J. Mintseris *et al.*, Protein-Protein Docking Benchmark 2.0: an update. *Proteins* **60**, 214-216 (2005).
- R. Chen, J. Mintseris, J. Janin, Z. Weng, A protein-protein docking benchmark. *Proteins* 52, 88-91 (2003).
- R. Mendez, R. Leplae, L. De Maria, S. J. Wodak, Assessment of blind predictions of protein-protein interactions: current status of docking methods. *Proteins* 52, 51-67 (2003).
- 198. A. Azimzadeh, M. H. Van Regenmortel, Antibody affinity measurements. *J Mol Recognit* **3**, 108-116 (1990).
- 199. M. Malmqvist, Surface plasmon resonance for detection and measurement of antibodyantigen affinity and kinetics. *Curr Opin Immunol* **5**, 282-286 (1993).
- 200. P. L. Kastritis, A. M. Bonvin, On the binding affinity of macromolecular interactions: daring to ask why proteins interact. *J R Soc Interface* **10**, 20120835 (2013).
- 201. I. H. Moal, J. Fernandez-Recio, Intermolecular Contact Potentials for Protein-Protein Interactions Extracted from Binding Free Energy Changes upon Mutation. *J Chem Theory Comput* 9, 3715-3727 (2013).
- 202. R. Raucci, E. Laine, A. Carbone, Local Interaction Signal Analysis Predicts Protein-Protein Binding Affinity. *Structure* **26**, 905-915 e904 (2018).
- 203. T. Siebenmorgen, M. Zacharias, Evaluation of Predicted Protein-Protein Complexes by Binding Free Energy Simulations. *J Chem Theory Comput* **15**, 2071-2086 (2019).
- 204. R. F. Alford *et al.*, The Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *J Chem Theory Comput* **13**, 3031-3048 (2017).
- 205. R. Akbar *et al.*, Progress and challenges for the machine learning-based design of fit-forpurpose monoclonal antibodies. *MAbs* **14**, 2008790 (2022).
- 206. M. Wang, Z. Cang, G. W. Wei, A topology-based network tree for the prediction of protein-protein binding affinity changes following mutation. *Nat Mach Intell* **2**, 116-123 (2020).
- I. H. Moal, J. Fernandez-Recio, SKEMPI: a Structural Kinetic and Energetic database of Mutant Protein Interactions and its use in empirical models. *Bioinformatics* 28, 2600-2607 (2012).

- 208. J. Jankauskaite, B. Jimenez-Garcia, J. Dapkunas, J. Fernandez-Recio, I. H. Moal, SKEMPI 2.0: an updated benchmark of changes in protein-protein binding energy, kinetics and thermodynamics upon mutation. *Bioinformatics* **35**, 462-469 (2019).
- T. Sulea, V. Vivcharuk, C. R. Corbeil, C. Deprez, E. O. Purisima, Assessment of Solvated Interaction Energy Function for Ranking Antibody-Antigen Binding Affinities. *J Chem Inf Model* 56, 1292-1303 (2016).
- S. Sirin, J. R. Apgar, E. M. Bennett, A. E. Keating, AB-Bind: Antibody binding mutational database for computational affinity predictions. *Protein Sci* 25, 393-409 (2016).
- 211. Y. Kurumida, Y. Saito, T. Kameda, Predicting antibody affinity changes upon mutations by combining multiple predictors. *Sci Rep* **10**, 19533 (2020).
- X. Liu, Y. Luo, P. Li, S. Song, J. Peng, Deep geometric representations for modeling effects of mutations on protein-protein binding affinity. *PLoS Comput Biol* 17, e1009284 (2021).
- 213. D. E. Pires, D. B. Ascher, mCSM-AB: a web server for predicting antibody-antigen affinity changes upon mutation with graph-based signatures. *Nucleic Acids Res* 44, W469-473 (2016).
- 214. Y. Myung, D. E. V. Pires, D. B. Ascher, CSM-AB: graph-based antibody-antigen binding affinity prediction and docking scoring function. *Bioinformatics* 10.1093/bioinformatics/btab762 (2021).
- 215. Y. Myung, C. H. M. Rodrigues, D. B. Ascher, D. E. V. Pires, mCSM-AB2: guiding rational antibody design using graph-based signatures. *Bioinformatics* **36**, 1453-1459 (2020).
- 216. R. Wang, X. Fang, Y. Lu, S. Wang, The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *J Med Chem* **47**, 2977-2980 (2004).
- 217. P. L. Kastritis *et al.*, A structure-based benchmark for protein-protein binding affinity. *Protein Sci* **20**, 482-491 (2011).
- 218. L. C. Xue, J. P. Rodrigues, P. L. Kastritis, A. M. Bonvin, A. Vangone, PRODIGY: a web server for predicting the binding affinity of protein-protein complexes. *Bioinformatics* **32**, 3676-3678 (2016).
- 219. A. Vangone, A. M. Bonvin, Contacts-based prediction of binding affinity in proteinprotein complexes. *Elife* **4**, e07454 (2015).
- 220. J. Janin, A minimal model of protein-protein binding affinities. *Protein Sci* 23, 1813-1817 (2014).

- 221. Y. Waheed, M. Siddiq, Z. Jamil, M. H. Najmi, Hepatitis elimination by 2030: Progress and challenges. *World J Gastroenterol* **24**, 4959-4961 (2018).
- 222. S. Buhler, R. Bartenschlager, Promotion of hepatocellular carcinoma by hepatitis C virus. *Dig Dis* **30**, 445-452 (2012).
- 223. R. Bartenschlager *et al.*, Critical challenges and emerging opportunities in hepatitis C virus research in an era of potent antiviral therapy: Considerations for scientists and funding agencies. *Virus Res* **248**, 53-62 (2018).
- 224. H. Midgard *et al.*, Hepatitis C reinfection after sustained virological response. *J Hepatol* **64**, 1020-1026 (2016).
- 225. J. R. Bailey, E. Barnes, A. L. Cox, Approaches, Progress, and Challenges to Hepatitis C Vaccine Development. *Gastroenterology* **156**, 418-430 (2019).
- 226. S. B. Cashman, B. D. Marsden, L. B. Dustin, The Humoral Immune Response to HCV: Understanding is Key to Vaccine Development. *Front Immunol* **5**, 550 (2014).
- 227. D. Sepulveda-Crespo, S. Resino, I. Martinez, Hepatitis C virus vaccine design: focus on the humoral immune response. *J Biomed Sci* **27**, 78 (2020).
- L. Kong, K. N. Jackson, I. A. Wilson, M. Law, Capitalizing on knowledge of hepatitis C virus neutralizing epitopes for rational vaccine design. *Current opinion in virology* 11, 148-157 (2015).
- 229. F. Penin, J. Dubuisson, F. A. Rey, D. Moradpour, J. M. Pawlotsky, Structural biology of hepatitis C virus. *Hepatology* **39**, 5-19 (2004).
- 230. D. Lapa, A. R. Garbuglia, M. R. Capobianchi, P. Del Porto, Hepatitis C Virus Genetic Variability, Human Immune Response, and Genome Polymorphisms: Which Is the Interplay? *Cells* **8** (2019).
- 231. L. Cocquerel, C. Wychowski, F. Minner, F. Penin, J. Dubuisson, Charged residues in the transmembrane domains of hepatitis C virus glycoproteins play a major role in the processing, subcellular localization, and assembly of these envelope proteins. *J Virol* 74, 3623-3633 (2000).
- 232. A. Op De Beeck *et al.*, The transmembrane domains of hepatitis C virus envelope glycoproteins E1 and E2 play a major role in heterodimerization. *J Biol Chem* **275**, 31428-31437 (2000).
- A. Bianchi, S. Crotta, M. Brazzoli, S. K. Foung, M. Merola, Hepatitis C virus e2 protein ectodomain is essential for assembly of infectious virions. *Int J Hepatol* 2011, 968161 (2011).
- 234. J. G. Haddad *et al.*, Identification of Novel Functions for Hepatitis C Virus Envelope Glycoprotein E1 in Virus Entry and Assembly. *J Virol* **91** (2017).

- 235. G. Vieyres, J. Dubuisson, T. Pietschmann, Incorporation of hepatitis C virus E1 and E2 glycoproteins: the keystones on a peculiar virion. *Viruses* **6**, 1149-1187 (2014).
- 236. C. C. Colpitts, P. L. Tsai, M. B. Zeisel, Hepatitis C Virus Entry: An Intriguingly Complex and Highly Regulated Process. *Int J Mol Sci* **21** (2020).
- 237. M. B. Zeisel, D. J. Felmlee, T. F. Baumert, Hepatitis C virus entry. *Curr Top Microbiol Immunol* **369**, 87-112 (2013).
- 238. P. Pileri et al., Binding of hepatitis C virus to CD81. Science 282, 938-941 (1998).
- 239. E. Scarselli *et al.*, The human scavenger receptor class B type I is a novel candidate receptor for the hepatitis C virus. *EMBO J* **21**, 5017-5025 (2002).
- 240. V. J. Kinchen *et al.*, Plasma deconvolution identifies broadly neutralizing antibodies associated with hepatitis C virus clearance. *J Clin Invest* **130**, 4786-4796 (2019).
- 241. N. Tzarum, I. A. Wilson, M. Law, The Neutralizing Face of Hepatitis C Virus E2 Envelope Glycoprotein. *Front Immunol* **9**, 1315 (2018).
- 242. Y. Wang, Z. Y. Keck, S. K. Foung, Neutralizing antibody response to hepatitis C virus. *Viruses* **3**, 2127-2145 (2011).
- 243. M. D. Colbert *et al.*, Broadly Neutralizing Antibodies Targeting New Sites of Vulnerability in Hepatitis C Virus E1E2. *J Virol* **93** (2019).
- 244. A. I. Flyak *et al.*, HCV Broadly Neutralizing Antibodies Use a CDRH3 Disulfide Motif to Recognize an E2 Glycoprotein Site that Can Be Targeted for Vaccine Design. *Cell Host Microbe* **24**, 703-716 e703 (2018).
- 245. I. Vasiliauskaite *et al.*, Conformational Flexibility in the Immunoglobulin-Like Domain of the Hepatitis C Virus Glycoprotein E2. *mBio* **8** (2017).
- 246. L. J. Stroh, K. Nagarathinam, T. Krey, Conformational Flexibility in the CD81-Binding Site of the Hepatitis C Virus Glycoprotein E2. *Front Immunol* **9**, 1396 (2018).
- 247. E. Giang *et al.*, Human broadly neutralizing antibodies to the envelope glycoprotein complex of hepatitis C virus. *Proc Natl Acad Sci U S A* **109**, 6205-6210 (2012).
- 248. F. Chen *et al.*, Antibody Responses to Immunization With HCV Envelope Glycoproteins as a Baseline for B-Cell-Based Vaccine Development. *Gastroenterology* **158**, 1058-1071 e1056 (2020).
- 249. C. Fauvelle *et al.*, Hepatitis C virus vaccine candidates inducing protective neutralizing antibodies. *Expert Rev Vaccines* **15**, 1535-1544 (2016).
- 250. H. Zazrin, H. Shaked, J. H. Chill, Architecture of the hepatitis C virus E1 glycoprotein transmembrane domain studied by NMR. *Biochim Biophys Acta* **1838**, 784-792 (2014).

- 251. L. Kong *et al.*, Hepatitis C virus E2 envelope glycoprotein core structure. *Science* **342**, 1090-1094 (2013).
- 252. L. Kong *et al.*, Structure of Hepatitis C Virus Envelope Glycoprotein E1 Antigenic Site 314-324 in Complex with Antibody IGH526. *J Mol Biol* **427**, 2617-2628 (2015).
- 253. R. Spadaccini *et al.*, Structural characterization of the transmembrane proximal region of the hepatitis C virus E1 glycoprotein. *Biochim Biophys Acta* **1798**, 344-353 (2010).
- 254. N. Tzarum *et al.*, Genetic and structural insights into broad neutralization of hepatitis C virus by human VH1-69 antibodies. *Sci Adv* **5**, eaav1882 (2019).
- 255. H. M. Yassine *et al.*, Hemagglutinin-stem nanoparticles generate heterosubtypic influenza protection. *Nat Med* **21**, 1065-1070 (2015).
- 256. L. Rutten *et al.*, Structure-Based Design of Prefusion-Stabilized Filovirus Glycoprotein Trimers. *Cell Rep* **30**, 4540-4550 e4543 (2020).
- 257. J. L. Slon-Campos *et al.*, A protective Zika virus E-dimer-based subunit vaccine engineered to abrogate antibody-dependent enhancement of dengue infection. *Nat Immunol* **20**, 1291-1298 (2019).
- 258. M. Kanekiyo, B. S. Graham, Next-Generation Influenza Vaccines. *Cold Spring Harb Perspect Med* 10.1101/cshperspect.a038448 (2020).
- 259. L. D. Jones, M. A. Moody, A. B. Thompson, Innovations in HIV-1 Vaccine Design. *Clin Ther* 10.1016/j.clinthera.2020.01.009 (2020).
- 260. Y. Lu, J. P. Welsh, J. R. Swartz, Production and stabilization of the trimeric influenza hemagglutinin stem domain for potentially broadly protective influenza vaccines. *Proc Natl Acad Sci U S A* **111**, 125-130 (2014).
- 261. E. Kim *et al.*, Microneedle array delivered recombinant coronavirus vaccines: Immunogenicity and rapid translational development. *EBioMedicine* 10.1016/j.ebiom.2020.102743, 102743 (2020).
- 262. W. Tai *et al.*, A recombinant receptor-binding domain of MERS-CoV in trimeric form protects human dipeptidyl peptidase 4 (hDPP4) transgenic mice from MERS-CoV infection. *Virology* **499**, 375-382 (2016).
- 263. Y. C. Chang *et al.*, Efficacy of heat-labile enterotoxin B subunit-adjuvanted parenteral porcine epidemic diarrhea virus trimeric spike subunit vaccine in piglets. *Appl Microbiol Biotechnol* **102**, 7499-7507 (2018).
- 264. P. Leblanc *et al.*, VaxCelerate II: rapid development of a self-assembling vaccine for Lassa fever. *Hum Vaccin Immunother* **10**, 3022-3038 (2014).

- 265. W. Cai *et al.*, Expression, purification and immunogenic characterization of hepatitis C virus recombinant E1E2 protein expressed by Pichia pastoris yeast. *Antiviral research* **88**, 80-85 (2010).
- 266. M. Logan *et al.*, Native Folding of a Recombinant gpE1/gpE2 Heterodimer Vaccine Antigen from a Precursor Protein Fused with Fc IgG. *J Virol* **91** (2017).
- 267. B. Wen *et al.*, Signal peptide replacements enhance expression and secretion of hepatitis C virus envelope glycoproteins. *Acta Biochim Biophys Sin (Shanghai)* **43**, 96-102 (2011).
- Z. Keck *et al.*, Cooperativity in virus neutralization by human monoclonal antibodies to two adjacent regions located at the amino terminus of hepatitis C virus E2 glycoprotein. J Virol 87, 37-51 (2013).
- 269. Z. Y. Keck *et al.*, Hepatitis C virus E2 has three immunogenic domains containing conformational epitopes with distinct properties and biological functions. *J Virol* **78**, 9224-9232 (2004).
- 270. Z. Y. Keck *et al.*, Human monoclonal antibodies to a novel cluster of conformational epitopes on HCV E2 with resistance to neutralization escape in a genotype 2a isolate. *PLoS Pathog* **8**, e1002653 (2012).
- 271. Z. Y. Keck *et al.*, Affinity maturation of a broadly neutralizing human monoclonal antibody that prevents acute hepatitis C virus infection in mice. *Hepatology* **64**, 1922-1933 (2016).
- 272. S. E. Boyken *et al.*, De novo design of protein homo-oligomers with modular hydrogenbond network-mediated specificity. *Science* **352**, 680-687 (2016).
- 273. L. Gonzalez, Jr., D. N. Woolfson, T. Alber, Buried polar residues and structural specificity in the GCN4 leucine zipper. *Nat Struct Biol* **3**, 1011-1018 (1996).
- 274. P. S. Huang *et al.*, High thermodynamic stability of parametrically designed helical bundles. *Science* **346**, 481-485 (2014).
- 275. T. J. Broering *et al.*, Identification and characterization of broadly neutralizing human monoclonal antibodies directed against the E2 envelope glycoprotein of hepatitis C virus. *J Virol* **83**, 12473-12482 (2009).
- 276. Z. Y. Keck *et al.*, Human monoclonal antibody to hepatitis C virus E1 glycoprotein that blocks virus attachment and viral infectivity. *J Virol* **78**, 7257-7263 (2004).
- 277. J. Lebowitz, M. S. Lewis, P. Schuck, Modern analytical ultracentrifugation in protein science: a tutorial review. *Protein Sci* **11**, 2067-2079 (2002).
- 278. T. M. Laue, B. D. Shah, T. M. Ridgeway, S. L. Pelletier, *Analytical ultracentrifugation in biochemistry and polymer science* (Royal Society of Chemistry, 1992).

- 279. A. K. Andrianov *et al.*, Supramolecular assembly of Toll-like receptor 7/8 agonist into multimeric water-soluble constructs enables superior immune stimulation in vitro and in vivo. *ACS Appl. Bio Mater.* **3**, 3187-3195 (2020).
- 280. A. K. Andrianov, Y. Y. Svirkin, M. P. LeGolvan, Synthesis and biologically relevant properties of polyphosphazene polyacids. *Biomacromolecules* **5**, 1999-2006 (2004).
- A. K. Andrianov, A. Marin, T. R. Fuerst, Molecular-Level Interactions of Polyphosphazene Immunoadjuvants and Their Potential Role in Antigen Presentation and Cell Stimulation. *Biomacromolecules* 17, 3732-3742 (2016).
- 282. J. N. Glover, S. C. Harrison, Crystal structure of the heterodimeric bZIP transcription factor c-Fos-c-Jun bound to DNA. *Nature* **373**, 257-261 (1995).
- 283. B. E. Willcox *et al.*, Production of soluble alphabeta T-cell receptor heterodimers suitable for biophysical analysis of ligand binding. *Protein Sci* **8**, 2418-2423 (1999).
- 284. A. Berthelmann, J. Lach, M. A. Grawert, M. Groll, J. Eichler, Versatile C(3)-symmetric scaffolds and their use for covalent stabilization of the foldon trimer. *Org Biomol Chem* **12**, 2606-2614 (2014).
- 285. D. Wrapp *et al.*, Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* **367**, 1260-1263 (2020).
- 286. G. Wang, R. N. de Jong, E. T. J. van den Bremer, P. Parren, A. J. R. Heck, Enhancing Accuracy in Molecular Weight Determination of Highly Heterogeneously Glycosylated Proteins by Native Tandem Mass Spectrometry. *Anal Chem* **89**, 4793-4797 (2017).
- 287. R. E. Iacob, I. Perdivara, M. Przybylski, K. B. Tomer, Mass spectrometric characterization of glycosylation of hepatitis C virus E2 envelope glycoprotein reveals extended microheterogeneity of N-glycans. *J Am Soc Mass Spectrom* **19**, 428-444 (2008).
- 288. A. V. Gandhi, M. R. Pothecary, D. L. Bain, J. F. Carpenter, Some Lessons Learned From a Comparison Between Sedimentation Velocity Analytical Ultracentrifugation and Size Exclusion Chromatography to Characterize and Quantify Protein Aggregates. *J Pharm Sci* 106, 2178-2186 (2017).
- 289. Z. Y. Keck *et al.*, Mapping a region of hepatitis C virus E2 that is responsible for escape from neutralizing antibodies and a core CD81-binding region that does not tolerate neutralization escape mutations. *J Virol* **85**, 10451-10463 (2011).
- Y. Wang *et al.*, Affinity maturation to improve human monoclonal antibody neutralization potency and breadth against hepatitis C virus. *J Biol Chem* 286, 44218-44233 (2011).
- 291. A. G. Khan *et al.*, Structure of the core ectodomain of the hepatitis C virus envelope glycoprotein 2. *Nature* **509**, 381-384 (2014).

- 292. A. Wahid *et al.*, Disulfide bonds in hepatitis C virus glycoprotein E1 control the assembly and entry functions of E2 glycoprotein. *J Virol* **87**, 1605-1617 (2013).
- 293. J. Dubuisson, C. M. Rice, Hepatitis C virus glycoprotein folding: disulfide bond formation and association with calnexin. *J Virol* **70**, 778-786 (1996).
- 294. P. Angel, M. Karin, The role of Jun, Fos and the AP-1 complex in cell-proliferation and transformation. *Biochim Biophys Acta* **1072**, 129-157 (1991).
- 295. M. Karin, Z. Liu, E. Zandi, AP-1 function and regulation. *Curr Opin Cell Biol* **9**, 240-246 (1997).
- 296. K. A. Swanson *et al.*, A respiratory syncytial virus (RSV) F protein nanoparticle vaccine focuses antibody responses to a conserved neutralization domain. *Sci Immunol* **5** (2020).
- 297. J. Marcandalli *et al.*, Induction of Potent Neutralizing Antibody Responses by a Designed Protein Nanoparticle Vaccine for Respiratory Syncytial Virus. *Cell* **176**, 1420-1431 e1417 (2019).
- 298. A. M. Owsianka *et al.*, Identification of conserved residues in the E2 envelope glycoprotein of the hepatitis C virus that are critical for CD81 binding. *J Virol* **80**, 8695-8704 (2006).
- 299. K. McCaffrey *et al.*, An Optimized Hepatitis C Virus E2 Glycoprotein Core Adopts a Functional Homodimer That Efficiently Blocks Virus Entry. *J Virol* **91** (2017).
- 300. M. Q. Marin *et al.*, Optimized Hepatitis C Virus (HCV) E2 Glycoproteins and their Immunogenicity in Combination with MVA-HCV. *Vaccines (Basel)* **8** (2020).
- B. G. Pierce *et al.*, Structure-Based Design of Hepatitis C Virus Vaccines That Elicit Neutralizing Antibody Responses to a Conserved Epitope. *J Virol* 91, e01032-01017 (2017).
- 302. L. He *et al.*, Proof of concept for rational design of hepatitis C virus E2 core nanoparticle vaccines. *Sci Adv* **6**, eaaz6225 (2020).
- 303. H. Freedman, M. R. Logan, J. L. Law, M. Houghton, Structure and Function of the Hepatitis C Virus Envelope Glycoproteins E1 and E2: Antiviral and Vaccine Targets. *ACS Infect Dis* **2**, 749-762 (2016).
- 304. J. P. Julien *et al.*, Crystal structure of a soluble cleaved HIV-1 envelope trimer. *Science* **342**, 1477-1483 (2013).
- 305. T. J. Ruckwardt *et al.*, Safety, tolerability, and immunogenicity of the respiratory syncytial virus prefusion F subunit vaccine DS-Cav1: a phase 1, randomised, open-label, dose-escalation clinical trial. *Lancet Respir Med* 10.1016/S2213-2600(21)00098-9 (2021).

- 306. A. W. Reinke, R. A. Grant, A. E. Keating, A synthetic coiled-coil interactome provides heterospecific modules for molecular engineering. *J Am Chem Soc* **132**, 6025-6031 (2010).
- 307. N. R. Zaccai *et al.*, A de novo peptide hexamer with a mutable channel. *Nat Chem Biol* 7, 935-941 (2011).
- 308. D. A. Lindhout, J. R. Litowski, P. Mercier, R. S. Hodges, B. D. Sykes, NMR solution structure of a highly stable de novo heterodimeric coiled-coil. *Biopolymers* **75**, 367-375 (2004).
- 309. L. He *et al.*, Approaching rational epitope vaccine design for hepatitis C virus with metaserver and multivalent scaffolding. *Sci Rep* **5**, 12501 (2015).
- 310. M. Kanekiyo *et al.*, Self-assembling influenza nanoparticle vaccines elicit broadly neutralizing H1N1 antibodies. *Nature* **499**, 102-106 (2013).
- 311. J. Jardine *et al.*, Rational HIV immunogen design to target specific germline B cell receptors. *Science* **340**, 711-716 (2013).
- 312. L. He *et al.*, Presenting native-like trimeric HIV-1 antigens with self-assembling nanoparticles. *Nat Commun* **7**, 12041 (2016).
- 313. E. Lamazares *et al.*, A Heterologous Viral Protein Scaffold for Chimeric Antigen Design: An Example PCV2 Virus Vaccine Candidate. *Viruses* **12** (2020).
- 314. A. C. Walls *et al.*, Elicitation of Potent Neutralizing Antibody Responses by Designed Protein Nanoparticle Vaccines for SARS-CoV-2. *Cell* **183**, 1367-1382 e1317 (2020).
- 315. B. Zakeri *et al.*, Peptide tag forming a rapid covalent bond to a protein, through engineering a bacterial adhesin. *Proc Natl Acad Sci U S A* **109**, E690-697 (2012).
- 316. L. Li, J. O. Fierer, T. A. Rapoport, M. Howarth, Structural analysis and optimization of the covalent association between SpyCatcher and a peptide Tag. *J Mol Biol* **426**, 309-317 (2014).
- K. D. Brune *et al.*, Dual Plug-and-Display Synthetic Assembly Using Orthogonal Reactive Proteins for Twin Antigen Immunization. *Bioconjug Chem* 28, 1544-1551 (2017).
- Y. F. Kang *et al.*, Rapid Development of SARS-CoV-2 Spike Protein Receptor-Binding Domain Self-Assembled Nanoparticle Vaccine Candidates. *ACS Nano* 15, 2738-2752 (2021).
- 319. A. A. Cohen *et al.*, Mosaic nanoparticles elicit cross-reactive immune responses to zoonotic coronaviruses in mice. *Science* **371**, 735-741 (2021).

- 320. A. A. Cohen *et al.*, Construction, characterization, and immunization of nanoparticles that display a diverse array of influenza HA trimers. *PLoS One* **16**, e0247963 (2021).
- 321. J. H. Salas *et al.*, An antigenically diverse, representative panel of envelope glycoproteins for HCV vaccine development. *Gastroenterology* 10.1053/j.gastro.2021.10.005 (2021).
- 322. T. H. Carlsen *et al.*, Breadth of neutralization and synergy of clinically relevant human monoclonal antibodies against HCV genotypes 1a, 1b, 2a, 2b, 2c, and 3a. *Hepatology* **60**, 1551-1562 (2014).
- 323. R. R. Meyerhoff *et al.*, HIV-1 Consensus Envelope-Induced Broadly Binding Antibodies. *AIDS Res Hum Retroviruses* **33**, 859-868 (2017).
- 324. K. Sliepen *et al.*, Structure and immunogenicity of a stabilized HIV-1 envelope trimer based on a group-M consensus sequence. *Nat Commun* **10**, 2355 (2019).
- 325. X. Ping *et al.*, Generation of a broadly reactive influenza H1 antigen using a consensus HA sequence. *Vaccine* **36**, 4837-4845 (2018).
- 326. H. Sun, J. H. Sur, S. Sillman, D. Steffen, H. L. X. Vu, Design and characterization of a consensus hemagglutinin vaccine immunogen against H3 influenza A viruses of swine. *Vet Microbiol* **239**, 108451 (2019).
- 327. C. Kuiken, P. Hraber, J. Thurmond, K. Yusim, The hepatitis C sequence database in Los Alamos. *Nucleic Acids Res* **36**, D512-516 (2008).
- 328. C. Kuiken, K. Yusim, L. Boykin, R. Richardson, The Los Alamos hepatitis C sequence database. *Bioinformatics* **21**, 379-384 (2005).
- 329. C. Combet *et al.*, euHCVdb: the European hepatitis C virus database. *Nucleic Acids Res* **35**, D363-366 (2007).
- 330. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* **30**, 772-780 (2013).
- 331. E. Paradis, J. Claude, K. Strimmer, APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics* **20**, 289-290 (2004).
- 332. J. D. Thompson, T. J. Gibson, F. Plewniak, F. Jeanmougin, D. G. Higgins, The CLUSTAL\_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res* **25**, 4876-4882 (1997).
- M. Gouy, S. Guindon, O. Gascuel, SeaView version 4: A multiplatform graphical user interface for sequence alignment and phylogenetic tree building. *Mol Biol Evol* 27, 221-224 (2010).
- 334. N. Ogata, H. J. Alter, R. H. Miller, R. H. Purcell, Nucleotide sequence and mutation rate of the H strain of hepatitis C virus. *Proc Natl Acad Sci U S A* **88**, 3392-3396 (1991).

- 335. M. Yanagi, R. H. Purcell, S. U. Emerson, J. Bukh, Transcripts from a single full-length cDNA clone of hepatitis C virus are infectious when directly transfected into the liver of a chimpanzee. *Proc Natl Acad Sci U S A* **94**, 8738-8743 (1997).
- 336. H. M. Berman et al., The Protein Data Bank. Nucleic Acids Research 28, 235-242 (2000).
- 337. J. D. Guest *et al.*, Design of a native-like secreted form of the hepatitis C virus E1E2 heterodimer. *Proc Natl Acad Sci U S A* **118** (2021).
- 338. A. Y. Lam, E. Pardon, K. V. Korotkov, W. G. J. Hol, J. Steyaert, Nanobody-aided structure determination of the EpsI:EpsJ pseudopilin heterodimer from Vibrio vulnificus. *J Struct Biol* **166**, 8-15 (2009).
- 339. D. Elnatan *et al.*, Symmetry broken and rebroken during the ATP hydrolysis cycle of the mitochondrial Hsp90 TRAP1. *Elife* **6** (2017).
- 340. M. Janiak *et al.*, Hepatitis C virus (HCV) genotype 1b displays higher genetic variability of hypervariable region 1 (HVR1) than genotype 3. *Sci Rep* **9**, 12846 (2019).
- 341. H. Zhang, A. A. Quadeer, M. R. McKay, Evolutionary modeling reveals enhanced mutational flexibility of HCV subtype 1b compared with 1a. *iScience* **25**, 103569 (2022).
- 342. W. S. Baker, S. Negi, W. Braun, C. H. Schein, Producing physicochemical property consensus alphavirus protein antigens for broad spectrum vaccine design. *Antiviral Res* **182**, 104905 (2020).
- 343. X. Qiu, V. R. Duvvuri, J. Bahl, Computational Approaches and Challenges to Developing Universal Influenza Vaccines. *Vaccines (Basel)* **7** (2019).
- 344. L. Feneant, S. Levy, L. Cocquerel, CD81 and hepatitis C virus (HCV) infection. *Viruses* 6, 535-572 (2014).
- 345. J. Zhang *et al.*, CD81 is required for hepatitis C virus glycoprotein-mediated viral infection. *J Virol* **78**, 1448-1455 (2004).
- 346. A. W. Tarr *et al.*, Immunization with a synthetic consensus hepatitis C virus E2 glycoprotein ectodomain elicits virus-neutralizing antibodies. *Antiviral Res* **160**, 25-37 (2018).
- 347. S. Munshaw *et al.*, Computational reconstruction of Bole1a, a representative synthetic hepatitis C virus subtype 1a genome. *J Virol* **86**, 5915-5921 (2012).
- 348. H. E. Drummer, Challenges to the development of vaccines to hepatitis C virus that elicit neutralizing antibodies. *Front Microbiol* **5**, 329 (2014).
- 349. J. C. Meunier *et al.*, Isolation and characterization of broadly neutralizing human monoclonal antibodies to the e1 glycoprotein of hepatitis C virus. *J Virol* **82**, 966-973 (2008).

- 350. J. Dubuisson *et al.*, Formation and intracellular localization of hepatitis C virus envelope glycoprotein complexes expressed by recombinant vaccinia and Sindbis viruses. *J Virol* **68**, 6147-6160 (1994).
- 351. C. Yi *et al.*, Junctional and somatic hypermutation-induced CX4C motif is critical for the recognition of a highly conserved epitope on HCV E2 by a human broadly neutralizing antibody. *Cell Mol Immunol* **18**, 675-685 (2021).
- 352. L. J. Stroh, T. Krey, HCV Glycoprotein Structure and Implications for B-Cell Vaccine Development. *Int J Mol Sci* **21** (2020).
- 353. H. E. Drummer, I. Boo, A. L. Maerz, P. Poumbourios, A conserved Gly436-Trp-Leu-Ala-Gly-Leu-Phe-Tyr motif in hepatitis C virus glycoprotein E2 is a determinant of CD81 binding and viral entry. *J Virol* **80**, 7844-7853 (2006).
- 354. F. Aleman *et al.*, Immunogenetic and structural analysis of a class of HCV broadly neutralizing antibodies and their precursors. *Proc Natl Acad Sci U S A* **115**, 7569-7574 (2018).
- 355. J. R. Bailey *et al.*, Broadly neutralizing antibodies with few somatic mutations and hepatitis C virus clearance. *JCI Insight* **2** (2017).
- 356. R. A. Urbanowicz *et al.*, A Diverse Panel of Hepatitis C Virus Glycoproteins for Use in Vaccine Research Reveals Extremes of Monoclonal Antibody Neutralization Resistance. *J Virol* **90**, 3288-3301 (2015).
- 357. H. E. Drummer, P. Poumbourios, Hepatitis C virus glycoprotein E2 contains a membrane-proximal heptad repeat sequence that is essential for E1E2 glycoprotein heterodimerization and viral entry. *J Biol Chem* **279**, 30066-30072 (2004).
- 358. A. W. Tarr *et al.*, Characterization of the hepatitis C virus E2 epitope defined by the broadly neutralizing monoclonal antibody AP33. *Hepatology* **43**, 592-601 (2006).
- 359. D. X. Johansson *et al.*, Human combinatorial libraries yield rare antibodies that broadly neutralize hepatitis C virus. *Proc Natl Acad Sci U S A* **104**, 16269-16274 (2007).
- 360. A. W. Tarr *et al.*, An alpaca nanobody inhibits hepatitis C virus entry and cell-to-cell transmission. *Hepatology* **58**, 932-939 (2013).
- T. D. Schiano *et al.*, Monoclonal antibody HCV-AbXTL68 in patients undergoing liver transplantation for HCV: results of a phase 2 randomized study. *Liver Transpl* 12, 1381-1389 (2006).
- 362. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic Local Alignment Search Tool. *J Mol Biol* **215**, 403-410 (1990).
- 363. K. El Omari *et al.*, Unexpected structure for the N-terminal domain of hepatitis C virus envelope glycoprotein E1. *Nat Commun* **5**, 4874 (2014).

- 364. T. Kortemme, D. E. Kim, D. Baker, Computational alanine scanning of protein-protein interfaces. *Sci STKE* **2004**, pl2 (2004).
- 365. T. Kortemme, D. Baker, A simple physical model for binding energy hot spots in proteinprotein complexes. *Proc Natl Acad Sci U S A* **99**, 14116-14121 (2002).
- 366. R Core Team (2018) R: A Language and Environment for Statistical Computing. (R Foundation for Statistical Computing, Vienna, Austria).
- 367. H. Wickham, ggplot2: Elegant Graphics for Data Analysis (Springer-Verlag New York, 2016).
- 368. P. Murrell, *R Graphics* (CRC Press, ed. 2nd, 2005).
- 369. Z. Zhang, Reshaping and aggregating data: an introduction to reshape package. *Ann Transl Med* **4**, 78 (2016).
- 370. A. Owsianka *et al.*, Monoclonal antibody AP33 defines a broadly neutralizing epitope on the hepatitis C virus E2 envelope glycoprotein. *J Virol* **79**, 11095-11104 (2005).
- 371. C. Harman *et al.*, A view of the E2-CD81 interface at the binding site of a neutralizing antibody against hepatitis C virus. *J Virol* **89**, 492-501 (2015).
- 372. T. J. Morin *et al.*, Human monoclonal antibody HCV1 effectively prevents and treats HCV infection in chimpanzees. *PLoS Pathog* **8**, e1002895 (2012).
- 373. M. Castelli *et al.*, A Biologically-validated HCV E1E2 Heterodimer Structural Model. *Sci Rep* **7**, 214 (2017).
- H. Freedman *et al.*, Computational Prediction of the Heterodimeric and Higher-Order Structure of gpE1/gpE2 Envelope Glycoproteins Encoded by Hepatitis C Virus. *J Virol* 91 (2017).
- 375. Y. Ciczora, N. Callens, F. Penin, E. I. Pecheur, J. Dubuisson, Transmembrane domains of hepatitis C virus envelope glycoproteins: residues involved in E1E2 heterodimerization and involvement of these domains in virus entry. *J Virol* **81**, 2372-2381 (2007).
- E. Mashiach, R. Nussinov, H. J. Wolfson, FiberDock: a web server for flexible inducedfit backbone refinement in molecular docking. *Nucleic acids research* 38, W457-461 (2010).
- 377. T. M. Cheng, T. L. Blundell, J. Fernandez-Recio, pyDock: electrostatics and desolvation for effective scoring of rigid-body protein-protein docking. *Proteins* **68**, 503-515 (2007).
- 378. D. Kozakov, K. H. Clodfelter, S. Vajda, C. J. Camacho, Optimal clustering for detecting near-native conformations in protein docking. *Biophys J* **89**, 867-875 (2005).

- 379. S. Lorenzen, Y. Zhang, Identification of near-native structures by clustering protein docking conformations. *Proteins* **68**, 187-194 (2007).
- 380. T. Vreven, H. Hwang, B. G. Pierce, Z. Weng, Prediction of protein-protein binding free energies. *Protein science : a publication of the Protein Society* **21**, 396-404 (2012).
- 381. I. Anishchenko, P. J. Kundrotas, A. V. Tuzikov, I. A. Vakser, Protein models: the Grand Challenge of protein docking. *Proteins* **82**, 278-287 (2014).
- 382. I. Anishchenko, P. J. Kundrotas, A. V. Tuzikov, I. A. Vakser, Protein models docking benchmark 2. *Proteins* 83, 891-897 (2015).
- 383. T. Bohnuud *et al.*, A benchmark testing ground for integrating homology modeling and protein docking. *Proteins* **85**, 10-16 (2017).
- 384. V. M. Chauhan, S. Islam, A. Vroom, R. Pantazes, Development and Analyses of a Database of Antibody – Antigen Complexes. *Computer Aided Chemical Engineering* 44, 2113-2118 (2018).
- 385. S. Mahajan *et al.*, Benchmark datasets of immune receptor-epitope structural complexes. *BMC Bioinformatics* **20**, 490 (2019).
- 386. U. Kulkarni-Kale, S. Raskar-Renuse, G. Natekar-Kalantre, S. A. Saxena, Antigen-Antibody Interaction Database (AgAbDb): a compendium of antigen-antibody interactions. *Methods Mol Biol* **1184**, 149-164 (2014).
- 387. M. N. Nguyen, M. R. Pradhan, C. Verma, P. Zhong, The interfacial character of antibody paratopes: analysis of antibody-antigen structures. *Bioinformatics* **33**, 2971-2976 (2017).
- 388. K. Krawczyk, T. Baker, J. Shi, C. M. Deane, Antibody i-Patch prediction of the antibody binding site improves rigid local antibody-antigen docking. *Protein Eng Des Sel* **26**, 621-629 (2013).
- 389. F. Ambrosetti, B. Jimenez-Garcia, J. Roel-Touris, A. Bonvin, Modeling Antibody-Antigen Complexes by Information-Driven Docking. *Structure* **28**, 119-129 e112 (2020).
- 390. J. M. Rini, U. Schulze-Gahmen, I. A. Wilson, Structural evidence for induced fit as a mechanism for antibody-antigen recognition. *Science* **255**, 959-965 (1992).
- 391. I. A. Wilson, R. L. Stanfield, Antibody-antigen interactions: new structures and new conformational changes. *Curr Opin Struct Biol* **4**, 857-867 (1994).
- 392. M. Torchala, I. H. Moal, R. A. Chaleil, J. Fernandez-Recio, P. A. Bates, SwarmDock: a server for flexible protein-protein docking. *Bioinformatics* **29**, 807-809 (2013).
- 393. B. G. Pierce, Y. Hourai, Z. Weng, Accelerating protein docking in ZDOCK using an advanced 3D convolution library. *PLoS One* **6**, e24657 (2011).

- 394. K. Wiehe *et al.*, ZDOCK and RDOCK performance in CAPRI rounds 3, 4, and 5. *Proteins* **60**, 207-213 (2005).
- 395. B. Al-Lazikani, A. M. Lesk, C. Chothia, Standard conformations for the canonical structures of immunoglobulins. *J Mol Biol* **273**, 927-948 (1997).
- 396. D. Kozakov *et al.*, The ClusPro web server for protein-protein docking. *Nat Protoc* **12**, 255-278 (2017).
- 397. M. D. Tyka *et al.*, Alternate states of proteins revealed by detailed energy landscape mapping. *J Mol Biol* **405**, 607-618 (2011).
- 398. S. J. Hubbard, J. M. Thornton (1993) NACCESS. (University College London, Department of Biochemistry and Molecular Biology).
- 399. J. K. Leman *et al.*, Macromolecular modeling and design in Rosetta: recent methods and frameworks. *Nat Methods* **17**, 665-680 (2020).
- 400. I. H. Moal, B. Jimenez-Garcia, J. Fernandez-Recio, CCharPPI web server: computational characterization of protein-protein interactions from structure. *Bioinformatics* **31**, 123-125 (2015).
- 401. Y. Yang, Y. Zhou, Specific interactions for ab initio folding of protein terminal regions with secondary structures. *Proteins* **72**, 793-803 (2008).
- 402. S. Viswanath, D. V. Ravikant, R. Elber, Improving ranking of models for protein complexes with side chain modeling and atomic potentials. *Proteins* **81**, 592-606 (2013).
- 403. S. Chaudhury, S. Lyskov, J. J. Gray, PyRosetta: a script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* **26**, 689-691 (2010).
- 404. D. Tobi, Designing coarse grained-and atom based-potentials for protein-protein docking. *BMC Struct Biol* **10**, 40 (2010).
- 405. D. Tobi, I. Bahar, Optimal design of protein docking potentials: efficiency and limitations. *Proteins* **62**, 970-981 (2006).
- 406. T. Lazaridis, M. Karplus, Effective energy function for proteins in solution. *Proteins-Structure Function and Genetics* **35**, 133-152 (1999).
- 407. S. Liu, C. Zhang, H. Zhou, Y. Zhou, A physical reference state unifies the structurederived potential of mean force for protein folding and binding. *Proteins* **56**, 93-101 (2004).
- 408. A. Honegger, A. Pluckthun, Yet another numbering scheme for immunoglobulin variable domains: an automatic modeling and analysis tool. *J Mol Biol* **309**, 657-670 (2001).

- 409. J. Dunbar, C. M. Deane, ANARCI: antigen receptor numbering and receptor classification. *Bioinformatics* **32**, 298-300 (2016).
- 410. R. L. Stanfield, H. Dooley, P. Verdino, M. F. Flajnik, I. A. Wilson, Maturation of shark single-domain (IgNAR) antibodies: evidence for induced-fit binding. *J Mol Biol* **367**, 358-372 (2007).
- 411. A. C. R. Martin, C. T. Porter (2009) ProFit. (University College London).
- 412. W. H. Lua, S. K. Gan, D. P. Lane, C. S. Verma, A search for synergy in the binding kinetics of Trastuzumab and Pertuzumab whole and F(ab) to Her2. *NPJ Breast Cancer* **1**, 15012 (2015).
- 413. P. Prabakaran *et al.*, Structure of severe acute respiratory syndrome coronavirus receptorbinding domain complexed with neutralizing antibody. *J Biol Chem* **281**, 15829-15836 (2006).
- 414. W. C. Liang *et al.*, Cross-species vascular endothelial growth factor (VEGF)-blocking antibodies completely inhibit the growth of human tumor xenografts and measure the contribution of stromal VEGF. *J Biol Chem* **281**, 951-961 (2006).
- 415. P. Verdino *et al.*, Molecular insights into gammadelta T cell costimulation by an anti-JAML antibody. *Structure* **19**, 80-89 (2011).
- 416. R. Diskin *et al.*, Increasing the potency and breadth of an HIV antibody by using structure-based rational design. *Science* **334**, 1289-1293 (2011).
- 417. S. Hu *et al.*, Comparison of the inhibition mechanisms of adalimumab and infliximab in treating tumor necrosis factor alpha-associated diseases from a molecular view. *J Biol Chem* **288**, 27059-27067 (2013).
- 418. M. H. Matho *et al.*, Structural and biochemical characterization of the vaccinia virus envelope protein D8 and its recognition by the antibody LA5. *J Virol* **86**, 8050-8058 (2012).
- 419. D. C. Ekiert *et al.*, Cross-neutralization of influenza A viruses mediated by a single antibody loop. *Nature* **489**, 526-532 (2012).
- 420. M. Hong *et al.*, Antibody recognition of the pandemic H1N1 Influenza virus hemagglutinin receptor binding site. *J Virol* **87**, 12471-12480 (2013).
- 421. S. Steidl, O. Ratsch, B. Brocks, M. Durr, E. Thomassen-Wolf, In vitro affinity maturation of human GM-CSF antibodies by targeted CDR-diversification. *Mol Immunol* **46**, 135-144 (2008).
- 422. B. J. Fennell *et al.*, CDR-restricted engineering of native human scFvs creates highly stable and soluble bifunctional antibodies for subcutaneous delivery. *mAbs* **5**, 882-895 (2013).

- 423. K. Liu *et al.*, Structural basis of anti-PD-L1 monoclonal antibody avelumab for tumor therapy. *Cell Res* 27, 151-153 (2017).
- 424. T. Kromann-Hansen *et al.*, A Camelid-derived Antibody Fragment Targeting the Active Site of a Serine Protease Balances between Inhibitor and Substrate Behavior. *J Biol Chem* **291**, 15156-15168 (2016).
- 425. R. K. Jensen *et al.*, Structure of the omalizumab Fab. *Acta Crystallogr F Struct Biol Commun* **71**, 419-426 (2015).
- 426. W. A. Bogdanoff *et al.*, Structure of a Human Astrovirus Capsid-Antibody Complex and Mechanistic Insights into Virus Neutralization. *J Virol* **91** (2017).
- 427. P. T. Beernink, S. Giuntini, I. Costa, A. H. Lucas, D. M. Granoff, Functional Analysis of the Human Antibody Response to Meningococcal Factor H Binding Protein. *mBio* **6**, e00842 (2015).
- 428. P. M. Legler *et al.*, Stability of isolated antibody-antigen complexes as a predictive tool for selecting toxin neutralizing antibodies. *mAbs* **9**, 43-57 (2017).
- 429. J. A. Kenniston *et al.*, Structural basis for pH-insensitive inhibition of immunoglobulin G recycling by an anti-neonatal Fc receptor antibody. *J Biol Chem* **292**, 17449-17460 (2017).
- 430. J. U. Lee, W. Shin, J. Y. Son, K. Y. Yoo, Y. S. Heo, Molecular Basis for the Neutralization of Tumor Necrosis Factor alpha by Certolizumab Pegol in the Treatment of Inflammatory Autoimmune Diseases. *Int J Mol Sci* **18** (2017).
- 431. W. Shin *et al.*, BAFF-neutralizing interaction of belimumab related to its therapeutic efficacy for treating systemic lupus erythematosus. *Nat Commun* **9**, 1200 (2018).
- 432. Q. Lin *et al.*, Structural Basis for the Broad, Antibody-Mediated Neutralization of H5N1 Influenza Virus. *J Virol* **92** (2018).
- 433. S. W. Scally *et al.*, Rare PfCSP C-terminal antibodies induced by live sporozoite vaccination are ineffective against malaria infection. *J Exp Med* **215**, 63-75 (2018).
- 434. S. W. Fanning, R. Walter, J. R. Horn, Structural basis of an engineered dual-specific antibody: conformational diversity leads to a hypervariable loop metal-binding site. *Protein Eng Des Sel* **27**, 391-397 (2014).
- 435. V. S. Nguyen *et al.*, Inhibition of type VI secretion by an anti-TssM llama nanobody. *PLoS One* **10**, e0122187 (2015).
- 436. C. McMahon *et al.*, Yeast surface display platform for rapid discovery of conformationally selective nanobodies. *Nat Struct Mol Biol* **25**, 289-296 (2018).

- 437. Y. Duhoo *et al.*, Camelid nanobodies used as crystallization chaperones for different constructs of PorM, a component of the type IX secretion system from Porphyromonas gingivalis. *Acta Crystallogr F Struct Biol Commun* **73**, 286-293 (2017).
- 438. R. W. Gene *et al.*, High affinity anti-Internalin B VHH antibody fragments isolated from naturally and artificially immunized repertoires. *J Immunol Methods* **416**, 29-39 (2015).
- 439. M. J. Rudolph *et al.*, Contribution of an unusual CDR2 element of a single domain antibody in ricin toxin binding affinity and neutralizing activity. *Protein Eng Des Sel* **31**, 277-287 (2018).
- 440. C. Tu *et al.*, A Combination of Structural and Empirical Analyses Delineates the Key Contacts Mediating Stability and Affinity Increases in an Optimized Biotherapeutic Single-chain Fv (scFv). *J Biol Chem* **291**, 1267-1276 (2016).
- 441. K. E. Conrath, U. Wernery, S. Muyldermans, V. K. Nguyen, Emergence and evolution of functional heavy-chain antibodies in Camelidae. *Dev Comp Immunol* 27, 87-103 (2003).
- 442. S. Bangaru *et al.*, A Site of Vulnerability on the Influenza Virus Hemagglutinin Head Domain Trimer Interface. *Cell* **177**, 1136-1152 e1118 (2019).
- 443. D. Fleury, S. A. Wharton, J. J. Skehel, M. Knossow, T. Bizebard, Antigen distortion allows influenza virus to escape neutralization. *Nat Struct Biol* **5**, 119-123 (1998).
- 444. C. Dreyfus *et al.*, Highly conserved protective epitopes on influenza B viruses. *Science* **337**, 1343-1348 (2012).
- 445. T. Tsibane *et al.*, Influenza human monoclonal antibody 1F1 interacts with three major antigenic sites and residues mediating human receptor specificity in H1N1 viruses. *PLoS Pathog* **8**, e1003067 (2012).
- 446. A. Bohne-Lang, C. W. von der Lieth, GlyProt: in silico glycosylation of proteins. *Nucleic acids research* **33**, W214-219 (2005).
- 447. J. W. Labonte, J. Adolf-Bryfogle, W. R. Schief, J. J. Gray, Residue-centric modeling and design of saccharide and glycoconjugate structures. *J Comput Chem* **38**, 276-287 (2017).
- 448. N. Andrusier, R. Nussinov, H. J. Wolfson, FireDock: Fast interaction refinement in molecular docking. *Proteins* (2007).
- 449. B. D. Weitzner *et al.*, Modeling and docking of antibody structures with Rosetta. *Nat Protoc* **12**, 401-416 (2017).
- 450. S. M. Lippow, K. D. Wittrup, B. Tidor, Computational design of antibody-affinity improvement beyond in vivo maturation. *Nat Biotechnol* **25**, 1171-1176 (2007).
- 451. J. R. Willis *et al.*, Redesigned HIV antibodies exhibit enhanced neutralizing potency and breadth. *J Clin Invest* **125**, 2523-2531 (2015).

- 452. T. Sulea, G. Hussack, S. Ryan, J. Tanha, E. O. Purisima, Application of Assisted Design of Antibody and Protein Therapeutics (ADAPT) improves efficacy of a Clostridium difficile toxin A single-domain antibody. *Scientific reports* **8**, 2260 (2018).
- 453. T. Borrman *et al.*, ATLAS: A database linking binding affinities with structures for wild-type and mutant TCR-pMHC complexes. *Proteins* **85**, 908-916 (2017).
- 454. Z. Liu *et al.*, PDB-wide collection of binding data: current status of the PDBbind database. *Bioinformatics* **31**, 405-412 (2015).
- 455. P. Gainza *et al.*, Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. *Nat Methods* 10.1038/s41592-019-0666-6 (2019).
- 456. X. Wang, G. Terashi, C. W. Christoffer, M. Zhu, D. Kihara, Protein Docking Model Evaluation by 3D Deep Convolutional Neural Networks. *Bioinformatics* 10.1093/bioinformatics/btz870 (2019).
- 457. K. P. Kilambi, J. J. Gray, Structure-based cross-docking analysis of antibody-antigen interactions. *Scientific reports* **7**, 8145 (2017).
- 458. A. Lopes *et al.*, Protein-protein interactions in a crowded environment: an analysis via cross-docking simulations and evolutionary information. *PLoS Comput Biol* **9**, e1003369 (2013).
- 459. C. K. Hua *et al.*, Computationally-driven identification of antibody epitopes. *eLife* **6** (2017).
- 460. K. Krawczyk, X. Liu, T. Baker, J. Shi, C. M. Deane, Improving B-cell epitope prediction and its application to global antibody-antigen docking. *Bioinformatics* **30**, 2288-2294 (2014).
- M. C. Jespersen, S. Mahajan, B. Peters, M. Nielsen, P. Marcatili, Antibody Specific B-Cell Epitope Predictions: Leveraging Information From Antibody-Antigen Protein Complexes. *Front Immunol* 10, 298 (2019).
- 462. S. Pittala, C. Bailey-Kellogg, Learning context-aware structural representations to predict antigen and antibody binding interfaces. *Bioinformatics* **36**, 3996-4003 (2020).
- 463. Y. X. Yang, P. Wang, B. T. Zhu, Importance of interface and surface areas in proteinprotein binding affinity prediction: A machine learning analysis based on linear regression and artificial neural network. *Biophys Chem* **283**, 106762 (2022).
- 464. M. Agostino, S. O. Pohl, A. Dharmarajan, Structure-based prediction of Wnt binding affinities for Frizzled-type cysteine-rich domains. *J Biol Chem* **292**, 11218-11229 (2017).
- 465. J. D. Guest *et al.*, An expanded benchmark for antibody-antigen docking and affinity prediction reveals insights into antibody recognition determinants. *Structure* **29**, 606-621 e605 (2021).

- 466. J. Wee, K. Xia, Persistent spectral based ensemble learning (PerSpect-EL) for proteinprotein binding affinity prediction. *Brief Bioinform* 10.1093/bib/bbac024 (2022).
- 467. M. M. Gromiha, K. Yugandhar, S. Jemimah, Protein-protein interactions: scoring schemes and binding affinity. *Curr Opin Struct Biol* **44**, 31-38 (2017).
- 468. T. F. Custodio *et al.*, Selection, biophysical and structural analysis of synthetic nanobodies that effectively neutralize SARS-CoV-2. *Nat Commun* **11**, 5588 (2020).
- 469. H. Yao *et al.*, A high-affinity RBD-targeting nanobody improves fusion partner's potency against SARS-CoV-2. *PLoS Pathog* **17**, e1009328 (2021).
- 470. H. Park *et al.*, Simultaneous Optimization of Biomolecular Energy Functions on Features from Small Molecules and Macromolecules. *J Chem Theory Comput* 12, 6201-6212 (2016).
- 471. Y. Zhang, A. Kolinski, J. Skolnick, TOUCHSTONE II: a new approach to ab initio protein structure prediction. *Biophys J* **85**, 1145-1164 (2003).
- 472. P. Chakrabarti, J. Janin, Dissecting protein-protein recognition sites. *Proteins* **47**, 334-343 (2002).
- 473. H. Lu, L. Lu, J. Skolnick, Development of unified statistical potentials describing protein-protein interactions. *Biophys J* 84, 1895-1901 (2003).
- 474. F. Glaser, D. M. Steinberg, I. A. Vakser, N. Ben-Tal, Residue frequencies and pairing preferences at protein-protein interfaces. *Proteins* **43**, 89-102 (2001).
- 475. B. Kuhlman *et al.*, Design of a novel globular protein fold with atomic-level accuracy. *Science* **302**, 1364-1368 (2003).
- 476. C. Zhang, G. Vasmatzis, J. L. Cornette, C. DeLisi, Determination of atomic desolvation energies from the structures of crystallized proteins. *J Mol Biol* **267**, 707-726 (1997).
- 477. J. Mintseris *et al.*, Integrating statistical pair potentials into protein complex prediction. *Proteins* **69**, 511-520 (2007).
- 478. W. N. Venables, B. D. Ripley, *Modern Applied Statistics with S* (Springer, New York, ed. Fourth, 2002).
- 479. M. Kuhn, Building Predictive Models in R Using the caret Package. *Journal of Statistical Software* 28, 1 26 (2008).
- 480. J. Friedman, T. Hastie, R. Tibshirani, Regularization Paths for Generalized Linear Models via Coordinate Descent. *J Stat Softw* **33**, 1-22 (2010).
- 481. L. Lei (2020) HIGH-RESOLUTION ANALYSIS OF HIV ENVELOPE-SPECIFIC ANTIBODY RESPONSES TO ACCELERATE RATIONAL IMMUNOGEN DESIGN.

in *Cell Biology & Molecular Genetics* (University of Maryland, Digital Repository at the University of Maryland).

- 482. X. Wu *et al.*, Selection pressure on HIV-1 envelope by broadly neutralizing antibodies to the conserved CD4-binding site. *J Virol* **86**, 5844-5856 (2012).
- 483. R. M. Lynch *et al.*, HIV-1 fitness cost associated with escape from the VRC01 class of CD4 binding site neutralizing antibodies. *J Virol* **89**, 4201-4213 (2015).
- 484. A. Sali, T. L. Blundell, Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol* **234**, 779-815 (1993).
- 485. I. Jarmoskaite, I. AlSadhan, P. P. Vaidyanathan, D. Herschlag, How to measure and evaluate binding affinities. *Elife* **9** (2020).
- 486. R. P. Joosten, F. Long, G. N. Murshudov, A. Perrakis, The PDB\_REDO server for macromolecular structure model optimization. *IUCrJ* **1**, 213-220 (2014).
- 487. H. Akiba, K. Tsumoto, Thermodynamics of antibody-antigen interaction revealed by mutation analysis of antibody variable regions. *J Biochem* **158**, 1-13 (2015).
- 488. D. Barradas-Bautista, I. H. Moal, J. Fernandez-Recio, A systematic analysis of scoring functions in rigid-body protein docking: The delicate balance between the predictive rate improvement and the risk of overtraining. *Proteins* **85**, 1287-1297 (2017).
- 489. R. Wang *et al.*, Induction of broadly neutralizing antibodies using a secreted form of the hepatitis C virus E1E2 heterodimer as a vaccine candidate. *Proc Natl Acad Sci U S A* 119, e2112008119 (2022).
- 490. R. Evans *et al.*, Protein complex prediction with AlphaFold-Multimer. *bioRxiv* 10.1101/2021.10.04.463034, 2021.2010.2004.463034 (2022).
- 491. R. Yin, B. Y. Feng, A. Varshney, B. G. Pierce, Benchmarking AlphaFold for protein complex modeling reveals accuracy determinants. *bioRxiv* 10.1101/2021.10.23.465575, 2021.2010.2023.465575 (2021).