

TECHNICAL RESEARCH REPORT

Hierarchical Production Planning for Complex Manufacturing Systems

by A. Mehra, I. Minis, and J.M. Proth

T.R. 94-36



*Sponsored by
the National Science Foundation
Engineering Research Center Program,
the University of Maryland,
Harvard University,
and Industry*

Hierarchical Production Planning for Complex Manufacturing Systems

Anshu Mehra and Ioannis Minis

Mechanical Engr. Dept. and
Institute for Systems Research,
University of Maryland

Jean M. Proth

INRIA-Lorraine,
4 rue Marconi, Metz 2000
57070 Metz, France

ABSTRACT

A hierarchical approach to production planning for complex manufacturing systems is presented. A single facility comprising of a number of work-centers that produce multiple part types is considered. The planning horizon includes a sequence of time periods, and the demand for all part types is assumed to be known. The production planning problem consists of minimizing the holding costs for all part types as well as the work-in-process, and backlogging cost for the end items. We present a two-level hierarchy that is based on aggregating parts to part families, work-centers to manufacturing cells and time periods to aggregate time periods. The solution at the aggregate level is imposed as a constraint to the detailed level problem which employs a decomposition based on manufacturing cells. This architecture uses a rolling horizon strategy to perform the production management function. We have employed perturbation analysis techniques to adjust certain parameters of the optimization problems at the detailed level to reach a near-optimal detailed production plan.

INTRODUCTION

Production planning consists of determining the quantities of products to manufacture in a sequence of time periods in order to optimize a certain criterion while satisfying constraints such as capacity of resources. This problem is further complicated by endogenous (e.g. resource failures) as well as exogenous (e.g. delayed receipts of raw material and unexpected changes in demand) random events. The resulting optimization problem is extremely large and complex. Maxwell et al. (1983) stress the need to develop and implement a generic production planning system by stating: *“Billions of dollars are wasted in US each year by manufacturers of discrete parts because of inadequate procedures for controlling inventory and production.”*

Two distinct approaches to production planning have been adopted in the literature. The first is a *monolithic approach*, wherein the entire problem is formulated as a large mixed-integer linear programming-type problem. The second is a *hierarchical approach* which partitions the global problem into a series of sub-problems that correspond to different hierarchical levels and are solved sequentially, such that the solution at each level imposes constraints on the solution of the subsequent lower level. The major advantages of the hierarchical approach are: (i) reduction of complexity, and (ii) gradual absorption of random events; for other advantages see (Dempster et al., 1981).

Hax and Meal (1975) conducted some pioneering work in the area of hierarchical production planning (HPP). They considered a multi-plant firm with four decision levels. The highest level performs the distribution of products to individual plants, and is solved once. The next three decision levels address the management of a single plant. The concept of a rolling horizon is employed to perform repetitive open-loop optimizations. Bitran and Hax (1977, 1981) further refined Hax and Meal’s methodology for a single stage production system and Bitran et al. (1982) extended this methodology to a two-stage process. These models considered three levels of aggregation for products: (i) items, (ii) families, and (iii) types. Graves (1982) adopts a different approach to the problem and introduces feedback between decision levels.

Following the product aggregation scheme of Hax and Meal, the problem is formulated as a monolithic mixed integer program, which is decomposed by Lagrangian relaxation. The best values of the Lagrangian multipliers are determined by an iterative procedure which may be interpreted as a feedback mechanism in the Hax-Meal framework. Aardal et al. (1990) present another hybrid approach to the mixed integer problem suggested by Graves. It is emphasized that the Hax and Meal model, and its extensions, are based on a special cost structure, and the proposed aggregation schemes are relevant to only a particular class of production systems. Furthermore, no spatial and time decompositions of the production system are proposed.

Axsater (1981), Hillion et al. (1988), Meier (1989), Nagi (1991), Thompson and Davis (1990) and Thompson et al. (1993) address a double aggregation scheme over products and machines, where at the top layer the entities of relevance are machine cells/departments and part families. This class of literature does not address aggregation of time periods. Mehra et al. (1993) have addressed aggregation of time periods, but for single work-center case. Gabby (1975), Erschler (1986), Saad (1990), Pienkosz and Toczowski (1993) derive necessary and sufficient conditions for the consistent disaggregation of the Hax and Meal model. The sub-optimality issues have been addressed by Dempster et al. (1981).

Designing an HPP system which is *optimal* for a given manufacturing system remains an open problem. To address this issue we propose a two-level hierarchy which is based on aggregating parts to part families, work-centers to manufacturing cells and time periods to aggregate time periods. The optimization problems at the aggregate and detailed level of the proposed hierarchy have been formulated. The solution at the aggregate level is imposed as a constraint to the detailed level problem which employs a decomposition based on manufacturing cells. This architecture uses a rolling horizon strategy to perform the production planning. We adjust some parameters at the detailed level to obtain the best possible performance of the HPP system. These parameters are typically related to local decisions at each level, which cannot be de-

rived from the solution of the previous upper level. We propose an iterative approach that is closely linked to *Perturbation Analysis* (PA) to determine near optimal values of these parameters.

We determine cost gradient estimates with respect to parameters of interest and use them to obtain near-optimal values of the parameters. We have presented a numerical example to demonstrate the effectiveness of the proposed hierarchical approach. To the best of our knowledge, it is the first time that simultaneous aggregation of parts, work-centers, and time periods is considered, and PA is employed in HPP systems.

HIERARCHICAL PRODUCTION PLANNING

In this section, we present a two-level problem to illustrate the proposed hierarchical production planning approach, which is based on a triple aggregation scheme: (i) parts are grouped into part families, (ii) work-centers are grouped into manufacturing cells, and (iii) elementary periods are aggregated into sub-periods.

The two levels in the proposed hierarchical structure are (i) the *aggregate* (or high) level that plans production for part families on cells within sub-periods and (ii) the *detailed* (or low) level that computes the disaggregation of the high level solution to determine the production plan for part types on work-centers during elementary periods. At the high level, the criterion is both to minimize inventory costs related to inter-cell work-in-process (WIP) as well as to families of finished parts, and backloging costs for families of finished parts. At the low level, the criterion is to minimize both the inventory and WIP costs of finished part types as well as backloging costs of finished parts. Thus, we introduce a hierarchy among inventories; priority is given to inter-cell and end-product inventories. The disaggregation is performed over a short term horizon, z ($z < H$ where H is the planning horizon). Production planning is usually performed on a rolling horizon basis. That is, although the high level solution is computed over H , only a part of it (over z) is implemented,

followed by a recomputation of the plan after z based on the actual state of the system. This is done in order to incorporate the future demands/forecasts progressively.

The aggregation approach and the hierarchical production planning problem are discussed below.

Problem formulation

We consider a single manufacturing facility consisting of a number of manufacturing resources or work-centers. A number of part types are manufactured in the facility. Each part type is produced following a certain routing or a sequence of operations, each operation is performed on a particular work-center. The processing time of each operation is given; set-up times are ignored in this treatment.

We limit ourselves to a two-level hierarchy in order to present the proposed approach design method and to evaluate its effectiveness. In this paper we assume that the temporal aggregation has been defined *a priori*. Furthermore, the aggregation of parts to families and machines to cells satisfy the following: (i) Parts belonging to the same family follow a common sequence of cells during their manufacture. Note that it is not necessary for parts to follow the same machine sequence. (ii) Parts belonging to the same family have similar processing times within each cell. This requirement will result in homogeneity of flows. We consider a planning horizon H consisting of Z *sub-periods* (aggregate time periods). Each sub-period is divided into z *elementary periods* (detailed time periods) of duration T each. Parts are due and new orders arrive at the beginning of these elementary periods. T is supposed to be much larger than the maximal amount of time needed to perform one operation. Furthermore, $z \times T$, which is the duration of a sub-period, is much larger than the time needed to complete the set of operations performed on any type of part in any cell. We also assume that: (i) at most, one operation related to a given part is performed during one elementary period, and (ii) a part visits at most one cell during each sub-period.

The production planning problem consists of determining the number of opera-

tions of each part type to be performed on the work-centers during each elementary period of the planning horizon in order to minimize the total inventory holding and backloging costs. The notation employed in this problem is as follows:

(i) $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$ is the set of N part types, (ii) $\mathcal{M} = \{m_1, m_2, \dots, m_M\}$ is the set of M work-centers, (iii) \mathcal{R}_j is the manufacturing routing, i.e. the sequence of operations that part type p_j must undergo. The total number of operations required for part type p_j is denoted by n_j , (iv) $I_{j,w}$ is the holding cost of one unit of part type p_j for one elementary period after its w -th operation is completed, (v) B_j is the backloging cost of one unit of finished part type p_j for one elementary period, (vi) $t_{j,w}$ is the processing time of w -th operation of part type p_j , (vii) $\delta_{j,w}(i)$ is the indicator function; assumes a value equal to 1 if machine m_i is required for w -th operation of part type p_j , and 0 otherwise, (viii) d_j^k is the external demand of part type p_j at the beginning of the k -th elementary period, and which must be satisfied by the end of that period, (ix) $s_{j,0}^k$ is the raw-material inventory for part type p_j at the end of k -th elementary period, (x) $s_{j,w}^0$ is the initial inventory of part type p_j at the end of w -th operation (at the beginning of the first elementary period), (xi) $s_{j,w}^k$ is the inventory of part type p_j at the end of k -th elementary period and at the end of w -th operation after satisfying the demand; negative values indicate backlog, (xii) $u_{j,w}^k$ is the production volume of part type p_j related to the w -th operation on p_j , during the k -th elementary period.

The problem can be formally modeled by the following linear program:

(Problem \mathcal{P}_1) *minimize*

$$\sum_{k=1}^{Z \times z} \sum_{j=1}^N \left\{ \left(\sum_{w=1}^{n_j} I_{j,w} \times [s_{j,w}^k]^+ \right) + B_j \times [-s_{j,n_j}^k]^+ \right\} \quad (1)$$

subject to:

$$s_{j,w}^k = s_{j,w}^{k-1} + u_{j,w}^k - u_{j,w+1}^k; \quad (2)$$

$$s_{j,n_j}^k = s_{j,n_j}^{k-1} + u_{j,n_j}^k - d_j^k; \quad (3)$$

$$\sum_{j=1}^N \sum_{w=1}^{n_j} \delta_{j,w}(i) \times u_{j,w}^k \times t_{j,w} \leq T; \quad (4)$$

$$u_{j,w}^k \leq s_{j,w-1}^{k-1}; \quad (5)$$

$$u_{j,w}^k \geq 0; \quad (6)$$

for $j = 1, 2, \dots, N, k = 1, 2, \dots, Z \times z$ in constraints (2) – (6);

for $w = 1, 2, \dots, n_j$ in constraints (5) – (6);

for $w = 1, 2, \dots, n_j - 1$ in constraint (2);

for $i = 1, 2, \dots, M$ in constraint (4);

The symbol $[\bullet]^+$ implies $\max\{0, \bullet\}$. The objective function (1) consists of both (i) the inventory holding costs of work-in-process (for $w = 1, 2, \dots, n_j - 1$) and finished parts ($w = n_j$) and (ii) the backlogging costs of finished parts over the entire planning horizon. Constraints (2) and (3) represent the state equations of the inventory levels. Note that the inventory level at any operation of a part is only updated at the end of the elementary period. Also, no more than one operation can take place on a specific part in a single elementary period. The capacity constraints on the work-centers for all the elementary periods of the horizon are represented by (4). Constraint (5) denotes that the total instances of an operation on a certain part within an elementary period cannot exceed the number of parts contained in the upstream buffer at the end of the previous elementary period. Finally, (6) represents the non-negativity of production.

Note that although the objective function in (1) is non-linear, it can be transformed into a linear one by adding constraints as shown below.

$$\text{minimize} \quad \sum_{k=1}^{Z \times z} \sum_{j=1}^N \sum_{w=1}^{n_j} P_{j,w}^k$$

subject to :

$$P_{j,w}^k \geq I_{j,w} \times s_{j,w}^k; \quad w = 1, 2, \dots, n_j \quad (7)$$

$$P_{j,n_j}^k \geq -B_j \times s_{j,n_j}^k; \quad (8)$$

$$P_{j,w}^k \geq 0; \quad w = 1, 2, \dots, n_j \quad (9)$$

for $j = 1, 2, \dots, N, k = 1, 2, \dots, Z \times z$ in constraints (7) – (9).

where $P_{j,w}^k$ is the inventory/backlogging cost related to the w -th operation and to the k -th elementary period.

Although this LP problem is simple to formulate, there are several reasons why it cannot be solved easily or implemented in practice: (i) the LP is of a very large dimension for a typical manufacturing system and planning horizon, (ii) detailed information about the demand of part types is not known for the entire horizon, (iii) the demand is subject to change due to order cancellations and acceptance of new orders, and thus the monolithic formulation requires frequent recomputation and excessive computation time, (iv) this formulation does not allow random events to be absorbed with a computation effort proportional to the impact of the random event, and (v) this formulation does not allow different criteria to be used at different levels of hierarchy.

These problems necessitate a hierarchical approach to the production planning problem. In addition, the hierarchical structure is parallel to a corporate management hierarchy and thus provides significant assistance to the overall management function. The following paragraphs outline the hierarchical approach.

Aggregate level problem

The problem at the aggregate level considers the production planning of part families on manufacturing cells during sub-periods of the planning horizon (H). Notation related to the aggregate level is as follows:

(i) $\mathcal{F} = \{\bar{p}_1, \bar{p}_2, \dots, \bar{p}_{\bar{N}}\}$ is the set of \bar{N} families; note that $\bar{N} \leq N$, (ii) $\mathcal{C} = \{c_1, c_2, \dots, c_{\bar{M}}\}$ is the set of \bar{M} manufacturing cells; note that $\bar{M} \leq M$, (iii) $T \times z$ is the length of a sub-period, (iv) $\widehat{\mathcal{R}}_f$ is the macro-manufacturing process, i.e. the sequence of macro-operations a part family \bar{p}_f must undergo, where, \bar{n}_f represents the number of macro-operations, (v) $\bar{I}_{f,q}$ is the inventory cost of one part unit in representative family \bar{p}_f after the q -th macro-operation of \bar{p}_f is completed, (vi) \bar{B}_f is the backlogging cost of one part unit in representative family \bar{p}_f at the final production stage, (vii) $\bar{u}_{f,q}^\kappa$ is the production volume of parts in family \bar{p}_f related to q -th macro-operation during the κ -th sub-period, (viii) $\tau_{f,q}^\kappa$ is the processing time related to q -th macro-operation of one unit of part type in representative family \bar{p}_f during the κ -th sub-period, (ix) D_f^κ is the demand of parts in family \bar{p}_f at the beginning of the κ -th sub-period and is given from

$$D_f^\kappa = \sum_{k=(\kappa-1) \times z + 1}^{\kappa \times z} \sum_{j \in \bar{p}_f} d_j^k. \quad (10)$$

The macro-manufacturing processes (macro-routing) are developed from aggregating work-centers into cells according to the aggregation rules presented above. In order to obtain this aggregation, we consider the routing of each part $j \in \bar{p}_f$, \mathcal{R}_j , and we define a set of \bar{n}_f subsets (sub-routings) $\{\mathcal{R}_j^q; q = 1, 2, \dots, \bar{n}_f\}$ where q is the rank of the macro-operation in the macro-manufacturing process of family \bar{p}_f . Note that several value of q may correspond to the same cell. Each sub-routing \mathcal{R}_j^q corresponds to one of the cells visited by part p_j , and alternatively, each sub-routing corresponds to one macro-operation of its family \bar{p}_f . That is:

$$\mathcal{R}_j = \bigcup_{q=1}^{\bar{n}_f} \mathcal{R}_j^q, \text{ and } \mathcal{R}_j^a \cap \mathcal{R}_j^b = \emptyset, \text{ for } a \neq b.$$

Note that \mathcal{R}_j^q corresponds to the q -th macro-operation of family \bar{p}_f , if $j \in \bar{p}_f$. Let $\widehat{\mathcal{I}}(j, q)$ be the last operation of the sub-routing \mathcal{R}_j^q and $\mathcal{I}(j, q)$ be the set of operations of \mathcal{R}_j^q , except the last one.

The computation of inventory/backlogging costs and processing times for families can be obtained from Eqs. (11) - (13). In fact, these parameters depend on the production volume of the part type $u_{j,w}^k$. However, since the production volumes are known only after solving the detailed level problem, we use the weighted average of the costs and processing times with respect to the parts demand, i.e.,

$$\tau_{f,q}^\kappa = \left(D_f^\kappa \times \left\| \mathcal{R}_j^q \right\| \right)^{-1} \sum_{k=(\kappa-1) \times z + 1}^{\kappa \times z} \sum_{j \in \bar{p}_f} \left(d_j^k \times \sum_{w \in \mathcal{R}_j^q} t_{j,w} \right) \quad (11)$$

$$\bar{B}_f = \left(\sum_{k=1}^{Z \times z} \sum_{j \in \bar{p}_f} d_j^k \right)^{-1} \sum_{k=1}^{Z \times z} \sum_{j \in \bar{p}_f} \left(d_j^k \times B_j \right) \quad (12)$$

$$\bar{I}_{f,q} = \left(\sum_{k=1}^{Z \times z} \sum_{j \in \bar{p}_f} d_j^k \right)^{-1} \sum_{k=1}^{Z \times z} \sum_{j \in \bar{p}_f} \left(d_j^k \times I_{j,w} \right); \quad (13)$$

for $w \in \widehat{\mathcal{I}}(j, q)$ in Eqn. (13).

The symbol $\|\cdot\|$ represents cardinality. The aggregate level problem can be formally stated as a linear program similar to \mathcal{P}_1 . The problem consists of determining the production $\bar{u}_{f,q}^\kappa$ over the entire planning horizon for all families at each macro-operation. The production parameters at this level (inventory/backlogging cost, processing time) correspond to macro-operations of families over sub-periods. The criterion at the aggregate level is minimization of: (i) the inventory holding costs of work-in-process between cells and finished families, and (ii) the backlogging costs of finished families, over the entire planning horizon. The inventory is estimated at the end of each sub-period (instead of elementary periods in \mathcal{P}_1) for all part families $f \in \mathcal{F}$ at the end of each macro-operation. The capacity constraints of the aggregate level problem are applied to cells for all sub-periods of the planning horizon. The aggregate level problem is presented in Chu et al. (1993).

Detailed level problem

The detailed level problem consists of determining the production plan for part types on work-centers during elementary periods of the first sub-period. Recall that we solve only the first sub-period in our rolling horizon approach to determine the production volume $u_{j,w}^k$ (for the w -th operation of part type p_j , during the k -th elementary period; $k = 1, 2, \dots, z$), respecting the aggregate high level solution $\bar{u}_{f,q}^1$. The criterion is to minimize the holding cost of all part types (end items and work-in-process), and backlogging costs for end items over the first sub-period. In addition, the low level problem adopts a decomposition-based approach in which the overall planning problem is solved by considering the manufacturing cells **independently**. This is also referred to as spatial decomposition and is consistent with CIM industrial practice. In order to make the problem definition concise we have defined some sets presented below.

$\mathcal{G}(j, v)$ represents the set of operations for part p_j that are performed in cell c_v . Let $J(j, v)$ be an indicator function which takes the value of 1 if the last operation in the routing of part p_j is performed in cell c_v and 0 otherwise.

Since, the demand for part types is known only for the cell that is in charge of the final stage of the manufacturing process, we define the parameters $\alpha_{j,q}^k (j \in \bar{p}_f)$ as the ratios of the expected production $u_{j,w}^k$ of part type p_j to the expected production of parts in family \bar{p}_f during the k -th elementary period, where $k = (\kappa - 1) \times z + 1, \dots, \kappa \times z$.

$$\alpha_{j,q}^k = \frac{u_{j,w}^k}{\bar{u}_{f,q}^k}; \quad w \in \widehat{\mathcal{I}}(j, q) \quad (14)$$

Note that:

$$\sum_{j \in \bar{p}_f} \alpha_{j,q}^k = 1; \quad (15)$$

For **each** cell c_v , where $c_v \in \mathcal{C}$, the detailed level problem can be formally stated by the following linear program:

(Problem $\mathcal{P}_2(v)$)

$$\begin{aligned} \text{minimize } & \sum_{k=1}^z \sum_{j=1}^N \left(\left\{ \sum_{w \in \mathcal{G}(j, v)} I_{j,w} \times [s_{j,w}^k]^+ \right\} + \right. \\ & \left. B_j \times J(j, v) \times [-s_{j,n_j}^k]^+ \right) \end{aligned} \quad (16)$$

subject to :

$$s_{j,w}^k = s_{j,w}^{k-1} + u_{j,w}^k - u_{j,w+1}^k; \quad (17)$$

for $w \in \mathcal{G}(j, v)$ such that $(w+1)$ -th operation is not performed in c_v

$$s_{j,w}^k = s_{j,w}^{k-1} + u_{j,w}^k - d_j^k; \quad (18)$$

if the last operation of \mathcal{R}_j is performed in c_v

$$s_{j,w}^k = s_{j,w}^{k-1} + u_{j,w}^k - \alpha_{j,q+1}^k \times \bar{u}_{f,q+1}^1; \quad (19)$$

for $w \neq n_j$ and $w \in \widehat{\mathcal{I}}(j, q)$; q being the index of macro-operation for part p_j performed in c_v

$$\sum_{j=1}^N \sum_{w \in \mathcal{G}(j, v)} \left\{ \partial_{j,w}(i) \times u_{j,w}^k \times t_{j,w} \right\} \leq T; \quad i \in c_v \quad (20)$$

$$u_{j,w}^k \leq s_{j,w-1}^{k-1}; \quad w \in \mathcal{G}(j, v) \quad (21)$$

$$\sum_{k=1}^z \sum_{j \in \bar{p}_f} u_{j,n_j}^k = \bar{u}_{f,\bar{n}_f}^1; \quad f \in \mathcal{F} \quad (22)$$

if the last macro-operation of $\widehat{\mathcal{R}}_f$ is performed in c_v ;

$$\sum_{k=1}^z u_{j,w}^k = \bar{u}_{f,q}^1 \times \sum_{k=1}^z \alpha_{j,q}^k; \quad (23)$$

for all macro-operations of $\widehat{\mathcal{R}}_f$ performed in c_v ; $w \in \widehat{\mathcal{I}}(j, q)$

$$u_{j,w}^k \geq 0; \quad w \in \mathcal{G}(j, v) \quad (24)$$

for $j = 1, 2, \dots, N$; $k = 1, 2, \dots, z$ and $j \in \bar{p}_f$ in constraints (17)-(19), (21), (23), (24);
for $k = 1, 2, \dots, z$ in constraints (20),(22).

The objective function (16) consists of: (i) the inventory holding costs of work-in-process, and (ii) the holding and backlogging costs of finished parts, over the first sub-period. Constraints (17) to (19) represent the state equations of the inventory levels. The capacity constraints on the work-centers m_i belonging to cell c_v for all elementary periods of the horizon are represented by (20). Constraint (21) denotes that the total instances of a certain operation in an elementary period cannot exceed the number of parts present in the upstream buffer at the end of the previous elementary period. Constraints (22) and (23) denote the production volume constraints imposed by the aggregate level. Constraint (22) represents the cumulative production of family \bar{p}_f related to the n_j -th operation for the entire first sub-period. This value is imposed by the solution of the aggregate level problem. Constraint (22) is applied only if macro-operation \bar{n}_f is performed in cell c_v . Constraint (23) represents the production of part p_j related to the last operation in sub-routing \mathcal{R}_j^q for the entire first sub-period. This value is determined from the solution of the aggregate level problem and the parameter alpha for that part type summed over the elementary periods of the first sub-period. Constraint (24) represents the non-negativity of production.

In this formulation, we observe that the demand is known only for the cell that is in charge of the final stage of the manufacturing process. Thus, we use the parameters $\alpha_{j,q}^k$ (see constraint 19) to determine the demand requirement of parts during the intermediate macro-operations. A perturbation analysis-based algorithm is used to determine the near-optimal value of these parameters. This algorithm is presented below.

GRADIENT ESTIMATION USING PA

As seen in the previous section, the optimization problems \mathcal{P}_2 at the detailed level include the control parameters $\alpha_{j,q}^k$ the value of which are not known *a priori*. We seek to obtain the near-optimal value of these parameters by adjusting them iteratively. To obtain the new set of $\alpha_{j,q}^k$, which reduces the cost function, we evaluate the gradient of the cost function with respect to these parameters at each iteration.

To estimate these derivatives efficiently we use Perturbation Analysis (PA). Consider the optimal value of cost function (16), say Y , which is a function of vector of random parameters α . $Y = Y(\alpha)$ where $\alpha = [\alpha_{j,q}^k]; q = 1, 2, \dots, \bar{n}_f, j \in f \in \mathcal{F}, k = 1, 2, \dots, Z \times z$. While using PA the optimization is done by recursive reassignment of the control parameter vector α to the series of random vectors $\{\alpha_m\}_{m=0}^{\infty}$ as follows:

$$\alpha^{m+1} = \alpha^m + A_m \widehat{\nabla}_{\alpha} Y(\alpha^m).$$

Here $\widehat{\nabla}_{\alpha} Y(\alpha^m)$ is the estimate of the gradient and $A_m = \text{diag}[a_m^1, \dots, a_m^k]$.

Blum (1954a, 1954b) has shown that if we define the sequences $\{a_m^i\}_{m=0}^{\infty}$, $i = 1 \dots k$ such that $\lim_{m \rightarrow \infty} a_m^i = 0, \forall i; \sum_{m=1}^{\infty} a_m^i = \infty, \forall i; \sum_{m=1}^{\infty} (a_m^i)^2 < \infty, \forall i$ and $\widehat{\nabla}_{\alpha} Y(\alpha)$ is an unbiased estimate of $\nabla_{\alpha} Y(\alpha)$, then the sequence $\{\alpha_m\}_{m=0}^{\infty}$ converges asymptotically to at least a local optimum of Y . For our simulation we chose $a_m^i = a_0^i/m, \forall i$ which satisfies all the conditions discussed above.

The gradient estimation proposed in this paper is carried out in two steps. We first perturbed the parameters $\beta_{j,q}^{\kappa}$

$$\beta_{j,q}^{\kappa} = \sum_{k=(\kappa-1) \times z+1}^{k=\kappa \times z} \alpha_{j,q}^k$$

of each sub-period and estimate the value of the gradient $\widehat{\nabla}_{\beta} Y(\beta^m)$ such that the inventories of parts at the end of a macro-operation remain same at the end of all successive sub-periods. Secondly, we perturb parameters α and estimate the value of the gradient $\widehat{\nabla}_{\alpha} Y(\alpha^m)$ such that the part inventories at the end of a macro-operation

remain invariant at the end of that sub-period, and all successive sub-periods. The small perturbations are represented by $\Delta\beta_{j,q}^\kappa$ and $\Delta\alpha_{j,q}^k$ for the first and second step, respectively. In this paper, we present the procedure for the gradient estimation of the second step. The gradient estimation of the first step can be performed in a similar manner.

Gradient estimation w.r.t parameters α

Consider the parameters $\Delta\alpha_{j,q+1}^k$, such that Eq. (25) is satisfied.

$$\sum_{k=(\kappa-1)\times z+1}^{k=\kappa\times z} \Delta\alpha_{j,q}^k = 0; \quad \forall j, \kappa, q \quad (25)$$

Note that $\Delta\alpha_{j,q+1}^k$ may have both positive and negative values. Let $\Delta s_{j,w}^k (w \in \widehat{\mathcal{I}}(j, q))$ be the change in inventory $s_{j,w}^k$ due to perturbations $\Delta\alpha_{j,q+1}^k$. It will be seen below that due to Eq. (25), the effect of $\Delta\alpha_{j,q+1}^k; k = (\kappa - 1) \times z + 1, \dots, \kappa \times z$ is limited to inventories of part types related to operations $w \in \mathcal{R}_j^q$ which are performed during the κ -th sub-period. The change in inventory $s_{j,w}^k (w \in \widehat{\mathcal{I}}(j, q))$ due to perturbations $\Delta\alpha_{j,q+1}^k (k = (\kappa - 1) \times z + 1, \dots, \kappa \times z)$ can be obtained from Eq. (19) as follows:

$$\begin{aligned} \Delta s_{j,w}^{(\kappa-1)\times z+1} &= -\Delta\alpha_{j,q+1}^{(\kappa-1)\times z+1} \times \bar{u}_{f,q+1}^\kappa \\ \Delta s_{j,w}^{(\kappa-1)\times z+2} &= -\left(\Delta\alpha_{j,q+1}^{(\kappa-1)\times z+1} + \right. \\ &\quad \left. \Delta\alpha_{j,q+1}^{(\kappa-1)\times z+2}\right) \times \bar{u}_{f,q+1}^\kappa \\ &\dots \quad \dots \end{aligned}$$

Furthermore from Eq. (25)

$$\begin{aligned} \Delta s_{j,w}^{\kappa\times z} &= -\left(\Delta\alpha_{j,q+1}^{(\kappa-1)\times z+1} + \Delta\alpha_{j,q+1}^{(\kappa-1)\times z+2} + \right. \\ &\quad \left. \Delta\alpha_{j,q+1}^{(\kappa-1)\times z+3} + \dots + \Delta\alpha_{j,q+1}^{\kappa\times z}\right) \times \bar{u}_{f,q+1}^\kappa \\ &= 0 \end{aligned}$$

Substituting the new value of inventories in criterion (16) and taking derivatives, we get

$$\partial Y / \partial \alpha_{j,q+1}^{(\kappa-1)\times z+1} \approx$$

$$\begin{aligned}
& - \left(\left[s_{j,w}^{(\kappa-1) \times z+1} + \Delta s_{j,w}^{(\kappa-1) \times z+1} \right]^\diamond + \left[s_{j,w}^{(\kappa-1) \times z+2} + \Delta s_{j,w}^{(\kappa-1) \times z+2} \right]^\diamond + \dots + \right. \\
& \quad \left. \left[s_{j,w}^{\kappa \times z-1} + \Delta s_{j,w}^{\kappa \times z-1} \right]^\diamond \right) \times \bar{u}_{f,q+1}^\kappa \times I_{j,w} \\
\partial Y / \partial \alpha_{j,q+1}^{(\kappa-1) \times z+2} & \approx \\
& - \left(\left[s_{j,w}^{(\kappa-1) \times z+2} + \Delta s_{j,w}^{(\kappa-1) \times z+2} \right]^\diamond + \dots + \right. \\
& \quad \left. \left[s_{j,w}^{\kappa \times z-1} + \Delta s_{j,w}^{\kappa \times z-1} \right]^\diamond \right) \times \bar{u}_{f,q+1}^\kappa \times I_{j,w} \\
& \dots \\
& \dots \\
\partial Y / \partial \alpha_{j,q+1}^{\kappa \times z} & \approx 0
\end{aligned}$$

where

$$\left[s_{j,w}^k + \Delta s_{j,w}^k \right]^\diamond = \begin{cases} 1 & \text{if } s_{j,w}^k + \Delta s_{j,w}^k \geq 0 \text{ and } \Delta s_{j,w}^k \neq 0 \\ 0 & \text{if } \Delta s_{j,w}^k = 0 \\ -\infty & \text{if } s_{j,w}^k + \Delta s_{j,w}^k \leq 0 \end{cases}$$

Note that the inventory at the end of the $\kappa \times z$ -th elementary period does not change by $\Delta \alpha_{j,q+1}^{(\kappa-1) \times z+1}, \dots, \Delta \alpha_{j,q+1}^{(\kappa) \times z}$ i.e. $\Delta s_{j,w}^{\kappa \times z} = 0$. Hence, the impact of $\Delta \alpha_{j,q+1}^k$ which satisfies Eq. (25) is local and it effects only the production volume of the part-types corresponding to the macro-operation and sub-period under consideration.

To illustrate the procedure of gradient estimation with respect to variable $\Delta \alpha_{j,q}^k$, we consider the following example. Let for part type $j \in \bar{p}_f$, $s_{j,w}^1 = 80$, $s_{j,w}^2 = 50$, $s_{j,w}^3 = 20$, $z = 3$, $\bar{u}_{f,q+1}^1 = 100$. If $\Delta \alpha_{j,q+1}^1 = 0.2$, $\Delta \alpha_{j,q+1}^2 = -0.15$ and $\Delta \alpha_{j,q+1}^3 = -0.05$, then from the above equations, the values of change in inventories are - $\Delta s_{j,w}^1 = -(0.2) \times 100 = -20$; $\Delta s_{j,w}^2 = -(0.2 - 0.15) \times 100 = -5$; $\Delta s_{j,w}^3 = 0$, and the value of gradients are - $\partial Y / \partial \alpha_{j,q+1}^1 \approx -200$; $\partial Y / \partial \alpha_{j,q+1}^2 \approx -100$; $\partial Y / \partial \alpha_{j,q+1}^3 \approx 0$.

Algorithm to update α iteratively

We present an algorithm to update the value of the parameters α . We first evaluate perturbations $\Delta \beta_{j,q}^\kappa$ and then perturbations $\Delta \alpha_{j,q}^k$.

Algorithm APA

A. Update parameters β

B. Update parameters α

Step 1: Estimate Gradient with respect to parameters $\alpha_{j,q}^k$

for all $\bar{p}_f \in \mathcal{F}$, $q = 1, 2, \dots, \bar{n}_f - 1$ and $\kappa = 1, 2, \dots, Z$

- choose a value of variable ss ($ss = \text{stepsize}$)
- estimate the value of $\Delta s_{j,w}^{(\kappa-1) \times z + 1}, \dots, \Delta s_{j,w}^{(\kappa-1) \times z + z - 1}$ ($w \in \widehat{\mathcal{I}}(j, q)$)
- estimate the value of $\partial Y / \partial \alpha_{j,q+1}^{(\kappa-1) \times z + 1}, \dots, \partial Y / \partial \alpha_{j,q+1}^{(\kappa-1) \times z + z - 1}$

Step 2: Sort Positive and Negative Gradients

Let

$$\overline{\overline{W}}_q^\kappa(j) \leftarrow \left\{ k : j \in \bar{p}_f, \partial Y / \partial \alpha_{j,q}^k \geq 0, \alpha_{j,q}^k \neq 0; \right. \\ \left. \forall k \in \{(\kappa - 1) \times z + 1, \dots, \kappa \times z\} \right\}$$

$$\overline{W}_q^\kappa(j) \leftarrow \left\{ k : j \in \bar{p}_f, \partial Y / \partial \alpha_{j,q}^k < 0; \right. \\ \left. \forall k \in \{(\kappa - 1) \times z + 1, \dots, \kappa \times z\} \right\}$$

$$W_q^\kappa(j) \leftarrow \min \left\{ \min_{k \in \overline{\overline{W}}_q^\kappa(j)} (1 - \alpha_{j,q}^k), \right. \\ \left. \min_{k \in \overline{W}_q^\kappa(j)} \alpha_{j,q}^k \right\}$$

Step 3: Estimation of $\Delta \alpha_{j,q}^k$

$$\Delta \alpha_{j,q}^k \leftarrow \begin{cases} \frac{\partial Y / \partial \alpha_{j,q}^k}{\sum_{k \in \overline{\overline{W}}_q^\kappa(j)} \partial Y / \partial \alpha_{j,q}^k} \min(ss, W_q^\kappa(j)) & \text{if } k \in \overline{\overline{W}}_q^\kappa(j) \\ \frac{\partial Y / \partial \alpha_{j,q}^k}{\sum_{k \in \overline{W}_q^\kappa(j)} \partial Y / \partial \alpha_{j,q}^k} \min(ss, W_q^\kappa(j)) & \text{if } k \in \overline{W}_q^\kappa(j) \\ 0 & \text{if } \partial Y / \partial \alpha_{j,q}^k \geq 0 \\ & \text{and } \alpha_{j,q}^k = 0 \end{cases}$$

- $\alpha_{j,q}^k = (\alpha_{j,q}^k + \Delta \alpha_{j,q}^k) \times \left(1 + \frac{\Delta \beta_{j,q}^\kappa}{\beta_{j,q}^\kappa}\right)$

Step 4: Analyze

if $\alpha_{j,q}^k > 1$ then

- $ss \leftarrow ss/2;$
- $W_q^\kappa(j) \leftarrow W_q^\kappa(j)/2;$

- go back to Step 2

ALGORITHM FOR HPP

In this section we present the algorithm used to solve the production planning problem. In this approach, the aggregate level optimization problem is solved first. The detailed level problems are then solved for all manufacturing cells at the beginning of each sub-period. The gradient estimates of the cost function $Y(\boldsymbol{\alpha})$ are carried out by the algorithm APA. A rolling horizon technique is employed to adjust the production plan according to the most recent forecast of the demand.

The algorithm for the hierarchical production planning based on the notations given in previous sections is presented below.

Step 1: Estimate part-families attributes

- $\bar{I}_{f,q}, \bar{B}_f, \tau_{f,q}^\kappa$ and D_f^κ for all $\bar{p}_f \in \mathcal{F}, q \in \widehat{\mathcal{R}}_f$ and $\kappa = 1, 2, \dots, Z$

Step 2: Initialize

- parameters α_{j,q_f}^k
- $a \leftarrow 1$ (first sub-period)
- $iter \leftarrow 1$
- $sum(iter) \leftarrow 0$
- $stepsize \leftarrow initialstepsize$

Step 3: Solve Aggregate level problem

- for the production values $\bar{u}_{f,q}^\kappa$ of part families in sub-periods $\kappa = a, a+1, \dots, a+Z-1$

Step 4: Solve detailed level problem for a -th sub-period

- for each cell $c_v \in \mathcal{C}$ individually
- for production values $u_{j,w}^k$ of part-types in elementary periods
 $k = (a - 1) \times z + 1, (a - 1) \times z + 2, \dots, a \times z$
- $sum(iter) \leftarrow sum(iter) +$ criterion value at detail level for all cell $c_v \in \mathcal{C}$

Step 5: Roll the planning horizon by one sub-period

- **If** $a = Z$ **Then** go to Step 6 **Else** $a \leftarrow a + 1$, estimate $s_{j,w}^k$ for all $p_j \in \mathcal{P}, w \in \mathcal{R}_j$ and $k = (a - 1) \times z + 1, (a - 1) \times z + 2, \dots, a \times z$, go to Step 3.

Step 6: Stopping Criterion

- **if** $\left\| \frac{sum(iter) - sum(iter-1)}{sum(iter-1)} \right\| \leq \epsilon$ where $\epsilon \leq 0.01$ **then** stop. Production value of part types obtained at the last iteration is the solution. **Else** go to Step 7

Step 7: Iterate

- $iter \leftarrow iter + 1$
- $a \leftarrow 1$
- $sum(iter) \leftarrow 0$
- $stepsize \leftarrow initialstepsize / iter$
- estimate new value for parameters α_{j,q_f}^k from algorithm APA
- go to Step 4

A NUMERICAL EXAMPLE

To demonstrate performance of the model and the algorithm a numerical example is considered in this section. The sample production planning problem comprises six part types and six work-centers, which are to be managed over a horizon

13 sub-periods. Each sub-period consists of 4 elementary periods. The routing details of part types 1-6 are: $\langle m_2 m_3 m_6 \rangle$, $\langle m_3 m_4 m_6 \rangle$, $\langle m_1 m_3 m_5 \rangle$, $\langle m_1 m_4 m_5 \rangle$, $\langle m_1 m_3 m_6 m_4 \rangle$ and $\langle m_1 m_2 m_5 m_3 \rangle$ respectively. Part-types p_1 and p_2 ; p_3 and p_4 ; p_5 and p_6 are grouped together to form families f_1 , f_2 and f_3 , respectively. Similarly, work-centers m_1 ; m_2, m_3 and m_4 ; m_5 and m_6 are grouped together to form cell c_1, c_2 and c_3 respectively. For this grouping the macro-routing of parts in family 1-3 are $\langle c_2 c_3 \rangle$, $\langle c_1 c_2 c_3 \rangle$ and $\langle c_1 c_2 c_3 c_2 \rangle$ respectively.

The coefficients (constants) $I_{j,w}$, B_j , $t_{j,w}$ and d_j^k of the LP were obtained using a random number generator. $s_{j,0}^k = 9999$ for $j = 1, \dots, 6$ and $k = 1, \dots, 13$; $s_{j,w}^0 = 100$ for $j = 1, \dots, 6$ and $w = 1, 2, \dots, n_j - 1$. Note that a high value of $s_{j,0}^k$ is selected assuming that a large quantity of raw-material is available for all part types. The cost coefficients and processing times were also obtained from the random number generator. The capacity of all work-centers (T) was selected to be 750 units. The cost coefficients, the processing times ($\bar{I}_{f,q}$, \bar{B}_f , $\tau_{f,q}^k$) and the initial inventory states at the aggregate level are estimated from the low level parameters $I_{j,w}$, B_j , $t_{j,w}$, d_j^k , $s_{j,0}^k$ and $s_{j,w}^0$ using relations (11 - 13). For the first iteration, the value of α_j^q is assumed to be 0.125 for $\bar{p}_f \ni j, \forall j, \forall q$.

The numerical results for this example were obtained by executing the two-level HPP model of section 4 on the SUN/SPARC station. The programs were coded in C and the LPs were solved by XMP (1981). The HPP problem is solved on a rolling horizon basis, i.e., the aggregate level LP is first solved for the entire planning horizon for all aggregate time-periods. Its solution is passed through the corresponding parameters to the detailed level LPs (see section). These LPs are solved for each cell separately. The planning horizon at the detailed level is one sub-period. The monolithic LP for problem $\mathcal{P}1$ was also solved. For comparison purposes it has 2080 decision variables and 3120 structural constraints.

The number of structural constraints and decision variable given by the software are 39 and 52 for cell c_1 , 104 and 80 for cell c_2 , 59 and 48 for cell c_3 , and 351 and 234

for the aggregate LP.

The criterion value and the computational times obtained by solving production planning problem by the monolithic approach (Chu et al., 1993) and the HPP approach are presented in Table 1. It is emphasized that the optimality of the solution depends on the selection of the initial step size (Fu, 1994), which is critical for the PA portion of the algorithm. The step size should go to zero at a rate that is slow enough to avoid premature convergence to the wrong value and simultaneously fast enough to eventually converge. The results show that for this example the converged solution was 3.17% away from the "best possible criterion value" that can be obtained if monolithic problem is solved.

However, these results should be considered as preliminary. A thorough evaluation of the approach is currently underway.

Table 1

Comparison of HPP approach vs. monolithic approach

	optimal value of criterion ($\times 10^3$)	computational time (minutes)
monolithic approach	193.85	32.0
hierarchical approach (after first iteration)	277.02	2.4
hierarchical approach (converged)	200.05	28.2

CONCLUSION

This paper presents a hierarchical approach for solving the production planning

problems. Aggregation of part types, machines and time-periods is considered. A two-level hierarchy is presented. At the aggregate level, the production planning of part families on manufacturing cells over aggregate elementary periods of the entire planning horizon is performed. At the detailed level, the production planning of parts on work-centers over the time periods of the first aggregate time period is performed. The rolling horizon strategy is utilized. Perturbation analysis is employed to adjust some control parameters at the detailed level. An algorithm is presented that provides a “near-optimal” solution of the HPP problem.

The proposed hierarchical scheme permits the computation of aggregate as well as detailed production plans when detailed demand/forecast information is not known (or considered necessary) at high management levels with long planning horizons. It also allows absorption of random events without frequent needs of recomputation. For faster computation, parallel processing can be also employed to solve the set of independent detailed level problems concurrently.

ACKNOWLEDGEMENTS

This work is supported by the Institute for Systems Research of the University of Maryland under grant # NSFD CD 8803012.

REFERENCES

K. Aardal, T. Larsson, “A Benders decomposition based heuristic for the hierarchical production planning problem,” *European Journal of Operational Research*, 45, pp. 1990.

Axsater, S., “Aggregation of Product Data for Hierarchical Production Planning,” *Operation Research*, 29, pp. 744-755, 1981.

G.R. Bitran, E.A. Haas, A.C. Hax, “Hierarchical Production Planning: A Two Stage System,” *Operations Research*, 30, pp. 232-251, 1982.

G.R. Bitran, A.C. Hax, “On the Design of Hierarchical Planning Systems,” *Decision Sciences*, 8, pp. 28-55, 1977.

G.R. Bitran, A.C. Hax, "Disaggregation and Resource Allocation using Convex Knapsack Problems with Bounded Variables," *Management Science*, 27, pp. 717-743, 1981.

J. R. Blum, "Approximation methods which converge with probability one," *Annals of Math. Stat.*, 25, pp. 382-386, 1954a.

J. R. Blum, "Multidimensional stochastic approximation methods," *Annals of Math. Stat.*, 25, pp. 737-744, 1954b.

X. R. Cao and Y. C. Ho, "Sensitivity analysis and optimization of throughput in a production line with blocking," *IEEE Transactions on Automatic Control*, 32, pp. 959-967, 1987.

C. Chu, M. Fu, J. Herrmann, A. Mehra, I. Minis, J.M. Proth, "Perturbation Analysis and Parallel Computation in Hierarchical Production Planning," NSF proposal, ISR, University of Maryland, College Park, Maryland, 1993.

M.A.H. Dempster, M.L. Fisher, B. Lageweg, L. Jansen, J.K. Lenstra, A.H.G. Rinnoy Kan, "Analytical evaluation of Hierarchical Planning Systems," *Operations Research*, 29, pp. 707-716, 1981.

J. Erscheler, G. Fontan, C. Merce, "Consistency of the Disaggregation Process in Hierarchical Planning," *Operations Research*, 34, pp. 464-469, 1986.

M. C. Fu, "Optimization via simulation: a review," to appear in *Annals of Operations Research*, 1994.

Gabby, H., "A Hierarchical Approach to Production Planning," TR-120, Operations Research Center, M.I.T, Cambridge, MA, 1975.

S.C. Graves, "Using Lagrangian Techniques to solve Hierarchical Production Planning Problems," *Management Science*, 28, pp. 260-275, 1982.

A. Mehra, R. Nagi, J.M. Proth, "Temporal Aggregation in Hierarchical Production Planning," Technical Report, University of Maryland, College Park, Maryland, also submitted to *European Journal of Operational Research*, 1993.

A.C. Hax, H.C. Meal, "Hierarchical Integration of Production Planning and

Scheduling,” *Studies in the Management Sciences*, M.A. Geisler, ed., 1, Logistics, North Holland - American Elsevier, 1975.

H. Hillion, K. Meier, J.M. Proth, “Production Subsystems and Part-families: The Top Level Model,” *Hierarchical Production Planning Systems*, Operational Research’87, G.K. Rand, ed., Elsevier Science Publishers B.V. (North-Holland), IFORS, 1988.

Y. C. Ho and X. R. Cao, *Discrete Event Dynamic Systems and Perturbation Analysis*, Kluwer Academic, 1991.

Y. C. Ho, “Performance evaluation and perturbation analysis of discrete event dynamic systems,” *IEEE Transactions on Automatic Control*, 32, pp. 563-572, 1987.

Y. C. Ho, “Extensions of infinitesimal perturbation analysis,” *IEEE Transactions on Automatic Control*, 33, pp. 427-438, 1988.

W. B. Gong and Y. C. Ho, “Smoothed (conditional) perturbation analysis of discrete event dynamical systems,” *IEEE Transactions on Automatic Control*, 32, pp. 858-866, 1987.

W. B. Gong, “Smoothed perturbation analysis of Markovian queuing networks,” *Proc. American Control Conference*, pp. 456-461, 1988.

W. Maxwell, J.A. Muckstadt, J. Thomas, J. VanderEecken, “A modeling framework for planning and control production in discrete parts manufacturing systems and assembly systems,” *Interface*, 3, 1983.

K. Meier, *Commande Hierarchisee d’un Systeme de Production*, Ph.D. Thesis, University of Metz, France, 1989.

R. Nagi, *Design and Operation of Hierarchical Production Management Systems*, Ph.D. Thesis, University of Maryland, College Park, 1989.

K. Pienkosz and E. Toczyłowski, “On aggregation of items in single-stage production systems with limited inventory levels,” *Operations Research*, 41, pp. 419-426, 1993.

G. H. Saad, “Hierarchical production-planning systems: extensions and modifi-

cations," *OR; The Journal of the Operational Research Society*, 41, pp. 609-624, 1990.

R. Suri, "Perturbation analysis: the state of the art and research issues explained via the G/G/1 queue," *Proceedings of the IEEE*, 77, pp. 114-137, 1989.

R. Suri, "Infinitesimal perturbation analysis for general discrete event systems," *Journal of the ACM*, 34, pp. 686-717, 1987.

A. D. Thompson, D. T. Watanabe and W. J. Davis, "A comparative study of aggregate production planning strategies under conditions of uncertainty and cyclic product demands," *International journal of production research*, 31, pp. 1957-1979, 1993.

A. D. Thompson and W. J. Davis, "An integrated approach for modeling uncertainty in aggregate production planning," *IEEE transactions on systems, man and cybernetic*, 20, pp. 1000-1012, 1990.

XMP or R. E. Marsten, "The design of the XMPLP Library," *Transactions on Mathematical Software*, 7, pp. 481-497, 1981.