# Survey of Job Shop Scheduling Techniques

Albert Jones, Ph.D. National Institute of Standards and Technology Building 220, Room A127 Gaithersburg, MD 20899-0001 Phone: (301) 975-3554, Fax: (301) 258-9749 jonesa@cme.nist.gov

Luis C. Rabelo, Ph.D. SDM Program Massachusetts Institute of Technology Cambridge, MA 02139-0254 Phone: (216) 447-5203, Fax: (216) 447-5196 email: lrabelo@mit.edu

## Introduction

In the United States today, there are approximately 40,000 factories producing metal-fabricated parts. These parts end up in a wide variety of products sold here and abroad. These factories employ roughly 2 million people and ship close to \$3 billion worth of products every year. The vast majority of these factories are what we call "job shops", meaning that the flow of raw and unfinished goods through them is completely random. Over the years, the behavior and performance of these job shops have been the focus of considerable attention in the Operations Research (OR) literature. Research papers on topics such as factory layout, inventory control, process control, production scheduling, and resource utilization can be found in almost every issue of every OR journal on the market today. The most popular of these topics is production (often referred to as job shop) scheduling. Job shop scheduling can be thought of as the allocation of resources over a specified time to perform a predetermined collection of tasks. Job shop scheduling has received this large amount of attention, because it has the potential to dramatically decrease costs and increase throughput, thereby, profits.

A large number of approaches to the modeling and solution of these job shop scheduling problems have been reported in the OR literature, with varying degrees of success. These approaches revolve around a series of technological advances that have occurred over that last 35 years. These include mathematical programming, dispatching rules, expert systems, neural networks, genetic algorithms, and inductive learning. In this article, we take an evolutionary view in describing how these technologies have been applied to job shop scheduling problems. To do this, we discuss a few of the most important contributions in each of these technology areas and the most recent trends.

## Mathematical techniques

Mathematical programming has been applied extensively to job shop scheduling problems. Problems have been formulated using integer programming (Balas 1965, 1967), mixed-integer programming (Balas 1969, 1970), and dynamic programming (Srinivasan 1971). Until recently, the use of these approaches has been limited because scheduling problems belong to the class of NP-complete problems. To overcome these deficiencies, a group of researchers began to decompose the scheduling problem into a number of subproblems, proposing a number of techniques to solve them. In addition, new solution techniques, more powerful heuristics, and the computational power of modern

computers have enabled these approaches to be used on larger problems. Still, difficulties in the formulation of material flow constraints as mathematical inequalities and the development of generalized software solutions have limited the use of these approaches.

#### **Decomposition strategies**

Davis and Jones (1988) proposed a methodology based on the decomposition of mathematical programming problems that used both Benders-type (Benders 1960) and Dantzig/Wolfe-type (Dantzig and Wolfe, 1960) decompositions. The methodology was part of closed-loop, real-time, two-level hierarchical shop floor control system. The top-level scheduler (i.e., the supremal) specified the earliest start time and the latest finish time for each job. The lower level scheduling modules (i.e., the infimals) would refine these limit times for each job by detailed sequencing of all operations. A multicriteria objective function was specified that included tardiness, throughput, and process utilization costs. The decomposition was achieved by first reordering the constraints of the original problem to generate a block angular form, then transforming that block angular form into a hierarchical tree structure. In general, N subproblems would result plus a constraint set that contained partial members of each of the subproblems. The latter was termed the "coupling " constraints, and included precedence relations and material handling. The supremal unit explicitly considered the coupling constraints, while the infimal units considered their individual decoupled constraint sets. The authors pointed out that the inherent stochastic nature of job shops and the presence of multiple, but often conflicting, objectives made it difficult to express the coupling constraints using exact mathematical relationships. This made it almost impossible to develop a general solution methodology. To overcome this, a new real-time simulation methodology was proposed in (Davis and Jones, 1988) to solve the supremal and infimal problems.

Gershwin (1989) used the notion of temporal decomposition to propose a mathematical programming framework for analysis of production planning and scheduling. This framework can be characterized as hierarchical and multi-layer. The problem formulations to control events at higher layers ignored the details of the variations of events occurring at lower layers. The problem formulations at the lower layers view the events at the higher layers as static, discrete events. Scheduling is actually carried out in bottom three layers so that the production requirements imposed by the planning layers can be met. First, a hedging point is found by solving a dynamic programming problem. This hedging point is the number of excess goods that should be produced to compensate for future equipment failures. This hedging point is used to formulate a linear programming problem to determine instantaneous production rates. These rates are then used to determine the actual schedule (which parts to make and when). A variety of approaches are under investigation for generating schedules.

#### Enumerative techniques and Lagrangian relaxation

Two popular solution techniques for integer-programming problems are branch-and-bound and Lagrangian relaxation. Branch-and-bound is an enumerative technique (Agin 1966, Lawler and Wood 1966). Summarizing Morton and Pentico (1993), "The basic idea of branching is to conceptualize the problem as a decision tree. Each decision choice point - a node - corresponds to a partial solution. From each node, there grow a number of new branches, one for each possible decision. This branching process continues until leaf nodes, that cannot branch any further, are reached. These leaf nodes are solutions to the scheduling problem". Although efficient bounding and pruning procedures have been developed to speed up the search, this is still a very computational intensive procedure for solving large scheduling problems. If the integer constraint is the main problem, then why not remove that (Shapiro 1979). Lagrangian relaxation solves integer-programming problems by omitting specific integer-valued constraints and adding the corresponding costs (due to these omissions and/or relaxations) to the objective function. As with branch and bound, Lagrangian relaxation is computationally expensive for large scheduling problems.

#### **Recent trends**

Model-Based Optimization (MBO) is an optimization approach that uses mathematical expressions (e.g., constraints and inequalities) to model scheduling problems as mixed integer (non) linear programs (MINLP's) (Zentner et al., 1994). A set of methods such as linear programming, branch-and-bound, and decomposition techniques are used to search the scenario space of solutions. Due to the advances in computer technologies, the computation times are becoming very practical. According to Subrahmanyam et al. (1996) "For problems of moderate size, solutions of type

D are given." Solutions of type D are optimal solutions of the maximum desirability possible within the constraints of operation. These approaches are being enhanced by the development of English-like "scheduling languages" and high-level graphical interfaces. The scheduling languages support the developing of the mathematical formulations with minimum intervention from the user.

## **Dispatching rules**

Dispatching rules have been applied consistently to scheduling problems. They are procedures designed to provide good solutions to complex problems in real-time. The term dispatching rule, scheduling rule, sequencing rule, or heuristic are often used synonymously (Panwalker and Islander 1977, Blackstone et al., 1982, Baker 1974). Dispatching rules have been classified mainly according to the performance criteria for which they have been developed. Wu (1987) categorized dispatching rules into several classes. Class 1 contains simple priority rules, which are based on information related to the jobs. Sub-classes are based on the particular piece of information used. Example classes include those based on processing times (such as shortest processing time (SPT)), due dates (such as earliest due date (EDD)), slack (such as minimum slack (MINSLACK)), and arrival times (such as first-in firstout (FIFO)). Class 2 consists of combinations of rules from class one. The particular rule that is implemented can now depend on the situation that exists on the shop floor. A typical example of a rule in this class is, for example, SPT until the queue length exceeds 5, then switch to FIFO. This prohibits jobs with large processing times from staying in the queue for long periods. Class 3 contains rules that are commonly referred to as Weight Priority Indexes. The idea here is to use more than one piece of information about the jobs to determine the schedule. Pieces of information are assigned weights to reflect their relative importance. Usually, an objective function f(x) is defined. For example,  $f(x) = weight_1 * Processing Time of Job(x) + weight_2 * (Current Time - Due Date of Job(x))$ . Then, any time new sequence is needed, the function f(x) is evaluated for each job x in the queue. The jobs are ranked based on this evaluation.

During the last 30 years, the performance of a large number of these rules has been studied extensively using simulation techniques (Montazer and Van Wassenhove, 1990). These studies have been aimed at answering the question: If you want to optimize a particular performance criterion, which rule should you choose? Most of the early work concentrated on the shortest processing time rule (SPT). Conway and Maxwell (1967) were the first to study the SPT rule and its variations. They found that, although some individual jobs could experience prohibitively long flow times, the SPT rule minimized the mean flow time for all jobs. They also showed that SPT was the best choice for optimizing the mean value of other basic measures such as waiting time and system utilization. Many similar investigations have been carried out to determine the dispatching rule which optimizes a wide range of job-related (such as due date and tardiness) and shop-related (such as throughput and utilization) performance measures. This problem of selecting the best dispatching rule for a given performance measure continues to be a very active area of research. However, the research has been expanded to include the possibility of switching rules to address an important problem: error recovery. Two early efforts to address error recovery were conducted by Bean and Birge (1986) and Saleh (1988). Both developed heuristic rules to smooth-out disruptions to the original schedule, thereby creating a match-up with that schedule. Bean and Birge (1986) based their heuristic on Turnpike Theory (McKenzie 1976) to optimize a generalized cost function. Saleh showed that he could minimize duration of the disruption by switching the objective function from mean flow time to makespan based on disjunctive graphs (Adams et al., 1988).

## Artificial intelligence (AI) techniques

Starting in the early 80s, a series of new technologies were applied to job shop scheduling problems. They fall under the general title of artificial intelligence (AI) techniques and include expert systems, knowledge-based systems, and several search techniques. Expert and knowledge-based systems were quite prevalent in the early and mid 1980s. They have four main advantages. First, and perhaps most important, they use both quantitative and qualitative knowledge in the decision-making process. Second, they are capable of generating heuristics that are significantly more complex than the simple dispatching rules described above. Third, the selection of the best heuristic can be based on information about the entire job shop including the current jobs, expected new jobs, and the current status of resources, material transporters, inventory, and personnel. Fourth, they capture complex relationships in elegant new data structures and contain special techniques for powerful manipulation of the information in these data structures There are, however, serious disadvantages. They can be time consuming to build and verify, as well as difficult to maintain and change. Moreover, since they generate only feasible solutions, it is rarely possible to tell how close that solution is to the optimal solution. Finally, since they are tied directly to the system they were built to manage, there is no such thing as a generic AI system.

#### Expert/knowledge-based systems

Expert and knowledge-based systems consist of two parts: a knowledge base, and inference engine to operate on that knowledge base. Formalizations of the "knowledge" that human experts use -- rules, procedures, heuristics, and other types of abstractions -- are captured in the knowledge base. Three types of knowledge are usually included: procedural, declarative, and meta. Procedural knowledge is domain-specific problem solving knowledge. Declarative knowledge provides the input data defining the problem domain. Meta knowledge is knowledge about how to use the procedural and declarative knowledge to actually solve the problem. Several data structures have been utilized to represent the knowledge in the knowledge base including semantic nets, frames, scripts, predicate calculus, and production rules. The inference engine selects a strategy to apply to the knowledge bases to solve the problem at hand. It can be forward chaining (data driven) or backward chaining (goal driven).

ISIS (Fox 1983) was the first major expert system aimed specifically at job shop scheduling problems. ISIS used a constraint-directed reasoning approach with three constraint categories: organizational goals, physical limitations and causal restrictions. Organizational goals considered objective functions based on due-date and work-in-progress. Physical limitations referred to situations where a resource had limited processing capability. Procedural constraints and resource requirements were typical examples of the third category. Several issues with respect to constraints were considered such as constraints in conflict, importance of a constraint, interactions of constraints, constraint generation and constraint obligation. ISIS used a three level, hierarchical, constraint-directed search. Orders were selected at level 1. A capacity analysis was performed at level 2 to determine the availability of the resources required by the order. Detailed scheduling was performed at level 3. ISIS also provided for the capability to interactively construct and alter schedules. In this capacity, ISIS utilized its constraint knowledge to maintain the consistency of the schedule and to identify scheduling decisions that would result in poorly satisfied constraints.

Wysk et al. (1986) developed an integrated expert system/simulation scheduler called MPECS. The expert system used both forward and backward chaining to select a small set of potentially good rules from predefined set of dispatching rules and other heuristics in the knowledge base. These rules optimized a single performance measure, although that measure could change from one scheduling period to the next. The selected rules were then evaluated one at a time using a deterministic simulation of a laboratory manufacturing system. After all of the rules were evaluated, the best rule was implemented on the laboratory system. Data could be gathered about how the rule actually performed and used to update the knowledge base off-line. They were able to show that periodic rescheduling makes the system more responsive and adaptive to a changing environment. MPECS was important for several reasons. It was the first hybrid system to make decisions based on the actual feedback from the shop floor. It incorporated some learning into its knowledge base to improve future decisions. The same systems could be used to optimize several different performance measures. Finally, it utilized a new multi-step approach to shop floor scheduling.

Other examples of expert/knowledge-based scheduling systems developed OPIS (Opportunistic Intelligent Scheduler) (Smith 1995), and SONIA (Le Pape 1995).

### **Distributed AI: agents**

Due to the limited knowledge and the problem solving ability of a single expert or knowledge based system, these AI approaches have difficulty solving large scheduling problems as well. To address this, AI researchers have also begun to develop distributed scheduling system approaches (Parunak et al., 1985). They have done this by an application of their well-known "divide and conquer" approach. This requires a problem decomposition technique, such as those described above, and the development of different expert/knowledge-based systems that can cooperate to solve the overall problem (Zhang and Zhang, 1995). The AI community's answer is the "agent" paradigm. An agent is a unique software process operating asynchronously with other agents. Agents are complete knowledge-based systems by themselves. The set of agents in a system may be heterogeneous with respect to long-

term knowledge, solution-evaluation criteria, or goals, as well as languages, algorithms, hardware requirements. Integrating agents selected from a "library" creates a multi-agent system.

For example, one such multi-agent system could involve two types of agents: tasks and resources. Each task agent might be responsible for scheduling a certain class of tasks such as material handling, machining, or inspection, on those resources capable of performing those tasks. This can be done using any performance measure related to tasks, such as minimize tardiness, and any solution technique. Each resource agent might be responsible for a single resource or a class of resources. Task agents must send their resource requests to the appropriate resource agent, along with the set of operations to be performed by that resource (Daouas et al., 1995). Upon receipt of such a request, the resource agent must generate a new schedule using its own performance measures, such as maximize utilization, which includes this request. The resource agent will use the results to decide whether to accept this new request or not. To avoid the situation where no resource will accept a request, coordination mechanisms must be developed. There are, now, no general guidelines for the design and implementation of this coordination. Therefore, the debates about centralized vs. decentralized approaches to job shop scheduling go on. The agents' formalism may provide an answer to these debates.

### Artificial neural networks

Neural networks, also called connectionist or distributed parallel processing models, have been studied for many years in an attempt to mirror the learning and prediction abilities of human beings. Neural network models are distinguished by network topology, node characteristics, and training or learning rules. An example of a three-layer, feed-forward neural network is shown in Figure 1.



Figure 1. An example of a three-layer, feed-forward neural network

Supervised learning neural networks

Through exposure to historical data, supervised learning neural networks attempt to capture the desired relationships between inputs and the outputs. Back-propagation is the most popular and widely used supervised training procedure. Back-propagation (Rumelhart et al., 1986, Werbos 1995) applies the gradient-descent technique in the feed-forward network to change a collection of weights so that some cost function can be minimized. The cost function, which is only dependent on weights (**W**) and training patterns, is defined by:

$$\mathbf{C}(\mathbf{W}) = \frac{1}{2} \mathbf{S} \left( \mathbf{T}_{ij} \cdot \mathbf{O}_{ij} \right) \tag{1}$$

where the **T** is the target value, **O** is the output of the network, **i** represents the output nodes, and **j** represents the training patterns.

After the network propagates the input values to the output layer, the error between the desired output and actual output will be "back-propagated" to the previous layer. In the hidden layers, the error for each node is computed by the weighted sum of errors in the next layer's nodes. In a three-layered network, the next layer means the output layer. The activation function is usually a sigmoid function with the weights modified according to

$$\mathbf{D}\mathbf{W}_{ij} = \mathbf{h} \mathbf{X}_{j} (\mathbf{1} \cdot \mathbf{X}_{j}) (\mathbf{T}_{j} \cdot \mathbf{X}_{j}) \mathbf{X}_{i}$$
<sup>(2)</sup>

or

$$\mathbf{D}\mathbf{W}_{ij} = \mathbf{h} \mathbf{X}_{j} \left(\mathbf{1} - \mathbf{X}_{j}\right) \left(\mathbf{S} \ \mathbf{d}_{k} \mathbf{W}_{jk}\right) \mathbf{X}_{i}$$
(3)

where  $W_{jk}$  is weight from node i to node (e.g., neuron) j, **h** is the learning rate,  $X_j$  is the output of node j,  $T_j$  is the target value of node j, and **d**<sub>k</sub> is the error function of node k.

If  $\mathbf{j}$  is in the output layer, Eq. (2) is used. If  $\mathbf{j}$  is the hidden layers, Eq. (3) is used. The weights are updated to reduce the cost function at each step. The process continues until the error between the predicted and the actual outputs is smaller than some predetermined tolerance.

Rabelo (1990) was the first to use back-propagation neural nets to solve job shop scheduling problems with several job types, exhibiting different arrival patterns, process plans, precedence sequences and batch sizes. Training examples were generated to train the neural network to select the correct characterization of the manufacturing environments suitable for various scheduling policies and the chosen performance criteria. In order to generate training samples, a performance simulation of the dispatching rules available for the manufacturing system was carried out. The neural networks were trained for problems involving 3, 4, 5, 8, 10, and 20 machines. To carry out this training, a special, input-feature space was developed. This space contained both job characteristics (such as types, number of jobs in each type, routings, due dates, and processing times) and shop characteristics (such as number of machines and their capacities). The output of the neural network represented the relative ranking of the available dispatching rules for that specific scheduling problem and the selected performance criteria. The neural networks were tested in numerous problems and their performance (in terms of minimizing Mean Tardiness) was always better than each single dispatching rule (25% to 50%).

#### **Relaxation models**

Neural networks based on relaxation models are defined by energy functions. They are pre-assembled systems that relax from input to output along a predefined energy contour. Hopfield neural networks (Hopfield and Tank 1985) are a classical example of a relaxation model that has been used to solve some classic, textbook scheduling problems (Foo and Takefuji, 1988). Two-dimensional Hopfield networks were used to solve 4-job, 3-machine problems and 10-job, 10-machine problems (Zhou et al., 1990). They were extended in (Lo and Bavarian, 1991) to 3 dimensions to represent jobs (i=1,...,I), machines j=1,...,J), and time (m=1,...,M). In each case, the objective was to minimize the makespan, total time to complete all jobs, which is defined as

$$E = S S S (v_{ijm}) (m + T_{ij} - 1)$$

(4)

where  $v_{ijm}$  is the output (1 or 0) of neuron ijm, and  $T_{ij}$  is the time required by  $j^{th}$  resource (e.g., machine) to complete the  $i^{th}$  job.

Due to a large number of variables involved in generating a feasible schedule, these approaches tend to be computationally inefficient and frequently generate infeasible solutions. Consequently, they have not been used to solve realistic scheduling problems.

### **Temporal reinforcement learning**

It was noted above that supervised learning neural networks attempt to capture the desired relationships between inputs and the outputs through exposure to training patterns. However, for some problems, the desired response may not always be available during the time of learning. When, the desired response is obtained, changes to the neural network are performed by assessing penalties for the scheduling actions previously decided by the neural network. As summarized by Tesauro (1992), "In the simplest form of this paradigm, the learning system passively observes a temporal sequence of input states that eventually leads to a final reinforcement or reward signal (usually a scalar). The learning system's task in this case is to predict expected reward given an observation of an input state or sequence of input states. The system may also be set up so that it can generate control signals that influence the sequence of states." For scheduling, the learning task is to produce an scheduling action that will lead to minimizing (or maximizing) the performance measure (e.g., makespan, tardiness) based on the state of the system (e.g., inventories, machine status, routings, due dates, layouts). Several procedures have been developed to train neural networks when the desired response is not available during the time of learning. Rabelo et al. (1994) utilized a procedure developed by Watkins (1989), denominated Q-learning, to implement a scheduling system to solve dynamic job shop scheduling problems. The scheduling system was able to follow trends in the shop floor and select a dispatching rule that provided the maximum reward according to performance measures based on tardiness and flow time. On the other hand, Zhang and Dietterich (1996) utilized a procedure developed by Sutton (1988) called  $TD(\lambda)$  to schedule payload processing of NASA's space shuttle program.

### Neighborhood search methods

Neighborhood search methods are very popular. Neighborhood search methods provide good solutions and offer possibilities to be enhanced when combined with other heuristics. Wilkerson and Irwin (1971) developed one of the first neighborhood procedures. This method iteratively added small changes ("perturbations") to an initial schedule, which is obtained by any heuristic. Conceptually similar to hill climbing, these techniques continue to perturb and evaluate schedules until there is no improvement in the objective function. When this happens, the procedure is ended. Popular techniques that belong to this family include Tabu search, simulated annealing, and genetic algorithms. Each of these has its own perturbation methods, stopping rules, and methods for avoiding local optimum.

### Tabu search

The basic idea of Tabu search (Glover 1989, 1990) is to explore the search space of all feasible scheduling solutions by a sequence of moves. A move from one schedule to another schedule is made by evaluating all candidates and choosing the best available, just like gradient-based techniques. Some moves are classified as tabu (i.e., they are forbidden) because they either trap the search at a local optimum, or they lead to cycling (repeating part of the search). These moves are put onto something called the Tabu List, which is built up from the history of moves used during the search. These tabu moves force exploration of the search space until the old solution area (e.g., local optimum) is left behind. Another key element is that of freeing the search by a short term memory function that provides "strategic forgetting". Tabu search methods have been evolving to more advanced frameworks that includes longer term memory mechanisms. These advanced frameworks are sometimes referred as Adaptive Memory Programming (AMP, Glover 1996).

Tabu search methods have been applied successfully to scheduling problems and as solvers of mixed integer programming problems. Nowicki and Smutnicki (Glover 1996) implemented tabu search methods for job shop and flow shop scheduling problems. Vaessens (Glover 1996) showed that tabu search methods (in specific job shop

scheduling cases) are superior over other approaches such as simulated annealing, genetic algorithms, and neural networks.

#### Simulated annealing

Simulated annealing is based on the analogy to the physical process of cooling and recrystalization of metals. The current state of the thermodynamic system is analogous to the current scheduling solution, the energy equation for the thermodynamic system is analogous to the objective function, and the ground state is analogous to the global optimum. In addition to the global energy J, there is a global temperature T, which is lowered as the iterations progress. Using this analogy, the technique randomly generates new schedules by sampling the probability distribution of the system (Kirkpatrick et al., 1983):

$$\mathbf{p}_{j} \mathbf{\mu} \exp(-\mathbf{T}(\mathbf{D}\mathbf{J}_{best} - \mathbf{D}\mathbf{J}_{j})/\mathbf{K})$$

where  $P_j$  represents the probability of making move j from among the neighborhood choices.  $DJ_{best}$  represents the improvement of the objective function for the best choice, and  $DJ_j$  represents the improvement for choice j. K is a normalization factor. Since increases of energy can be accepted, the algorithm is able to escape local minima.

(5)

Simulated annealing has been applied effectively to job shop scheduling problems. Vakharia and Chang (1990) developed a scheduling system based on simulated annealing for manufacturing cells. Jeffcoat and Bulfin (1993) applied simulated annealing to a resource-constrained scheduling problem. Their computational results indicated that the simulated annealing procedure provided the best results in comparison with other neighborhood search procedures.

### Genetic algorithms

Genetic algorithms (GA) are an optimization methodology based on a direct analogy to Darwinian natural selection and mutations in biological reproduction. In principle, genetic algorithms encode a parallel search through concept space, with each process attempting coarse-grain hill climbing (Goldberg 1988). Instances of a concept correspond to individuals of a species. Induced changes and recombinations of these concepts are tested against an evaluation function to see which ones will survive to the next generation. The use of genetic algorithms requires five components:

- **1.** A way of encoding solutions to the problem -- fixed length string of symbols.
- **2.** An evaluation function that returns a rating for each solution.
- **3.** A way of initializing the population of solutions.
- **4.** Operators that may be applied to parents when they reproduce to alter their genetic composition such as crossover (i.e., exchanging a randomly selected segment between parents), mutation (i.e., gene modification), and other domain specific operators.
- 5. Parameter setting for the algorithm, the operators, and so forth.

A number of approaches have been utilized in the application of genetic algorithms (GA) to job shop scheduling problems (Davis 1985, Goldberg and Lingle 1985, Starkweather et al., 1992):

- 1. Genetic algorithms with blind recombination operators have been utilized in job shop scheduling. Their emphasis on relative ordering schema, absolute ordering schema, cycles, and edges in the offsprings will arise differences in such blind recombination operators.
- **2.** Sequencing problems have been addressed by mapping their constraints to a Boolean satisfiability problem using partial payoff schemes. This scheme has produced good results for very simple problems.
- **3.** Heuristic genetic algorithms have been applied to job shop scheduling. In these genetic schemes, problem specific heuristics are incorporated in the recombination operators (such as optimization operators based).

Starkweather et al. (1993) were the first to use genetic algorithms to solve a dual -criteria job shop scheduling problem in a real production facility. Those criteria were the minimization of average inventory in the plant and the minimization of the average waiting time for an order to be selected. These criteria are negatively correlated (The larger the inventory, the shorter the wait; the smaller the inventory, the longer the wait.). To represent the production/shipping optimization problem, a symbolic coding was used for each member (chromosome) of the population. In this scheme, customer orders are represented by discrete integers. Therefore, each member of the population is a permutation of customer orders. The Genetic Algorithm used to solve this problem was based on blind recombinant operators. This recombination operator emphasizes information about the relative order of the elements in the permutation, because this impacts both inventory and waiting time. A single evaluation function (a weighted sum of the two criteria) was utilized to rank each member of the population. That ranking was based on an on-line simulation of the plant operations. This approach generated schedules that produced inventory levels and waiting times that were acceptable to the plant manager. In addition, the integration of the genetic algorithm with the on-line simulation made it possible to react to system dynamics.

These applications have emphasized the utilization of genetic algorithms as a "solo" technique. This has limited the level of complexity of the problems solved and their success. Recent research publications have demonstrated the sensitivity of genetic algorithms to the initial population. When the initial population is generated randomly, genetic algorithms are shown to be less efficient that the annealing-type algorithms, but better than the heuristic methods alone. However, if the initial population is generated by a heuristic, the genetic algorithms become as good as, or better than the annealing-type algorithms. In addition, integration with other search procedures (e.g., tabu search) has enhanced the capabilities of both. This result is not surprising, as it is consistent with results from non-linear optimization. Simply stated, if you begin the search close to the optimal solution you are much more likely to get the optimum than if you begin the search far away.

## **Fuzzy logic**

Fuzzy set theory has been utilized to develop hybrid scheduling approaches. Fuzzy set theory can be useful in modeling and solving job shop scheduling problems with uncertain processing times, constraints, and set-up times. These uncertainties can be represented by fuzzy numbers that are described by using the concept of an interval of confidence. These approaches usually are integrated with other methodologies (e.g., search procedures, constraint relaxation). For example, Slany (1994) stresses the imprecision of straight-forward methods presented in the mathematical approaches and introduces a method known as fuzzy constraint relaxation, which is integrated with a knowledge-based scheduling system. His system was applied to a steel manufacturing plant. Grabot and Geneste (1994) use fuzzy logic principles to combine dispatching rules for multi-criteria problems. On the other hand, Krucky (1994) addresses the problem of minimizing setup times of a medium-to-high product mix production line using fuzzy logic. The heuristic, fuzzy logic based algorithm described helps determine how to minimize setup time by clustering assemblies into families of products that share the same setup by balancing a product's placement time between multiple-high-speed placement process steps. Tsujimura et al. (1993) presented a hybrid system, which uses fuzzy set theory to model the processing times of a flow shop scheduling facility. Triangular Fuzzy Numbers (TFNs) are used to represent these processing times. Each job is defined by two TFNs, a lower bound and an upper bound. A branch and bound procedure is utilized to minimize makespan.

## **Reactive Scheduling**

Reactive scheduling is generally defined as the ability to revise or repair a complete schedule that has been "overtaken" by events on the shop floor (Zweben et al., 1995). Such events include rush orders, excessive delays, and broken resources. There are two approaches: reactive repair and the proactive adjustment. In reactive repair, the scheduling system waits until an event has occurred before it attempts to recover from that event. The match-up techniques described in section 3 fall into this category. Proactive adjustment requires a capability to monitor the system continuously, predict the future evolution of the system, do contingency planning for likely events, and generate new schedules, all during the execution time of the current schedule. The work of Wysk et al. (1986) and Davis and Jones (1988) fall into this category. Approaches that are more recent utilize artificial intelligence and knowledge-based methodologies (Smith 1995). Still most of the AI approaches propose a quasi-deterministic view of

the system, i.e., a stochastic system featuring implicit and/or explicit causal rules. The problem formulation used does not recognize the physical environment of the shop floor domain where interference not only leads to readjustment of schedules but also imposes physical actions to minimize them.

## Learning in Scheduling

The first step in developing a knowledge base is knowledge acquisition. This in itself is a two step process: get the knowledge from knowledge sources and store that knowledge in digital form. Much work has been done in the area of knowledge acquisition, such as protocol analysis and interactive editing (Shaw et al., 1992). Knowledge sources may be human experts, simulation data, experimental data, databases, and text. In scheduling problems, the knowledge sources are likely to be human experts or simulation data. To extract knowledge from these two sources, the machine learning technique that learns from examples (data) becomes a promising tool. Inductive learning is a state classification process. If we view the state space as a hyperplane, the training data (consisting of conditions and decisions) can be represented as points on the hyperplane. The inductive learning algorithm seeks to draw lines on the hyperplane based on the training data to divide the plane into several areas within which the same decision (conclusion) will be made.

One algorithm that has been implemented in inductive aids and expert system shells is that developed by Quinlan (1986), called Iterative Dichotomister 3 or ID3. ID3 uses examples to induce production rules (e.g. IF ... THEN ...), which form a simple decision tree. Decision trees are one way to represent knowledge for the purpose of classification. The nodes in a decision tree correspond to attributes of the objects to be classified, and the arcs are alternative values for these attributes. The end nodes of the tree (leaves) indicate classes to which groups of objects belong. Each example is described by attributes and a resulting decision. To determine a good attribute to partition the objects into classes, entropy is employed to measure the information content of each attribute, and then rules are derived through a repetitive decomposition process that minimizes the overall entropy. The entropy value of attribute  $A_k$  can be defined as

$$\mathbf{H}(\mathbf{A}_{k}) = \sum_{j=1}^{M_{k}} \mathbf{P}(\mathbf{a}_{kj}) \left\{ -\sum_{i=1}^{N} \mathbf{P}(\mathbf{c}_{i}|\mathbf{a}_{kj}) \log_{2} \mathbf{P}(\mathbf{c}_{i}|\mathbf{a}_{kj}) \right\}$$

(6)

where  $H(A_k)$  is the entropy value of attribute  $A_{ks} P(a_{kj})$  is the probability of attribute k being at its j<sup>th</sup> value,  $P(c_i|a_{kj})$  is the probability that the class value is  $c_i$  when attribute k is at its j<sup>th</sup> value,  $M_k$  is the total number of values for attribute  $A_k$  and N is the total number of different classes (outcomes).

The attribute with the minimum entropy value will be selected as a node in the decision tree to partition the objects. The arcs out of this node represent different values of this attribute. If all the objects in an arc belong to one class, the partition process stops. Otherwise, another attribute will be identified using entropy values to further partition the objects that belong to this arc. This partition process continues until all the objects in an arc are in the same class. Before applying this algorithm, all attributes that have continuous values need to be transformed to discrete values.

In the context of job shop scheduling, the attributes represent system status and the classes represent the dispatching rules. Very often, the attribute values are continuous. Yih (1988) proposed a trace-driven knowledge acquisition (TDKA) methodology to deal with continuous data and to avoid the problems occurring in verbally interviewing human experts. TDKA learns scheduling knowledge from expert schedulers without a dialogue with them. There are three steps in this approach. In Step 1, an interactive simulator is developed to mimic the system of interest. The expert will interact with this simulator and make decisions. The entire decision making process will be recorded in the simulator and can be repeated for later analysis. The series of system information and the corresponding decision collected is called a "trace." Step 2 analyzes the "trace" and forms classification rules to partition the trace into groups. The partition process stops when most of the cases in each group use the same dispatching rule (error rate is below the threshold defined by the knowledge engineer). Then, the decision rules are formed. The last step is to verify the generated rules. The resulting rule base is used to schedule jobs in the

simulator. If it performs as well as or better than the expert, the process stops. Otherwise, the threshold value is increased, and the process returns to Step 2.

As the job shop operates over time, it is important to be able to modify the knowledge contained in these rule bases. Chiu (1994) looks at knowledge modification for job shop scheduling problems by a framework of dynamic scheduling schemes that explores routing flexibility and handles uncertainties. The proposed methodology includes three modules: discrete-event simulation, instance generation, and incremental induction. First, a simulation module is developed to implement the dynamic scheduling scheme, to generate training examples, and to evaluate the methodology. Second, in an instance-generation module, the searching of good training examples is successfully fulfilled by a genetic algorithm. Finally, in an incremental-induction module, a tolerance-based incremental learning algorithm is proposed to allow continuous learning and facilitate knowledge modification. This algorithm uses entropy values to select attributes to partition the examples where the attribute values are continuous. The tolerance is used to maintain the stability of the existing knowledge while the new example is introduced. The decision tree will not be reconstructed unless there is enough momentum from the new data, that is, the change of the entropy value becomes significant. The experimental results showed that the tolerance-based incremental learning algorithm cannot only reduce the frequency of modifications, but also enhances the generalization ability of the resulting decision tree in a distributed job shop environment.

### **Theory of Constraints**

The Theory of Constraints (TOC) developed by Eliyahu Goldratt (1990, 1992) is the underlying philosophy for synchronized manufacturing. Goldratt (1990) defined synchronized manufacturing as any systematic method that attempts to move material quickly and smoothly through the production process in concert with market demand. A core concept to TOC is the idea that a few critical constraints exist. Goldratt contends that there is only one constraint in a system at any given time. As defined by Dettmer (1997), a constraint is "any element of a system or its environment that limits the output of the system". A constraint will prevent increases in throughput regardless of improvements made to the system. The best schedule is obtained by focusing on the planning and scheduling of these constraint operations. In essence, the constraint operations become the basis from which the entire schedule is derived. TOC has several important concepts and principles. Among them (Goldratt 1990,1992):

- **1.** Systems function like chains.
- **2.** The system optimum is not the sum of the local optima.
- 3. The effect-cause-effect method identifies constraints.
- **4.** System constraints can be either physically or policy.
- 5. Inertia is the worst enemy of a process of ongoing improvement.
- **6.** Throughput is the rate at which the entire system generates money through sales.
- 7. Inventory is all the money the system invests in things it intends to sell.
- 8. Operating expense is all the money the system spends turning inventory into throughput.

The general process of TOC is as follows (Goldratt 1990):

- **1.** Identify the systems' constraints.
- 2. Decide how to exploit the system's constraints.
- **3.** Subordinate everything else to the above decision.
- 4. Elevate the system's constraints.
- **5.** If in the previous steps a constraint have been broken, go back to Step1, but do not allow inertia to cause a system constraint.

TOC has been successfully applied to scheduling problems (Academic and Industrial) (Advanced Manufacturing Research, Inc. 1996). Its tools that comprised five distinct logic trees (explained extensively in (Dettmer1997)) are the Current Reality Tree, the Evaporating Cloud Diagram, the Future Reality Tree, the Prerequisite Tree, and the Transition Tree. These trees are tied to the Categories of Legitimate Reservation (that provide the logic to guide the

construction of the trees). These tools have not only been used in production scheduling but also in other enterprise functions such as marketing and sales.

### Summary and conclusions

Since job shop scheduling problems fall into the class of NP-complete problems, they are among the most difficult to formulate and solve. Operations Research analysts and engineers have been pursuing solutions to these problems for more than 35 years, with varying degrees of success.

While they are difficult to solve, job shop scheduling problems are among the most important because they impact the ability of manufacturers to meet customer demands and make a profit. They also impact the ability of autonomous systems to optimize their operations, the deployment of intelligent systems, and the optimizations of communications systems. For this reason, operations research analysts and engineers will continue this pursuit well into the next century.

## References

[1] Adams, J., E. Balas and D. Zawack (1988), "The shifting bottleneck procedure for job shop scheduling," *Management Science*: 34 (3), 391-401.

[2] Advanced Manufacturing Research, Inc. (1996), "Advanced planning and scheduling systems: just a fad or a breakthrough in manufacturing and supply chain management?," *The Report on Manufacturing*, December 1996.

[3] Agin, N. (1966), "Optimum seeking with branch and bound," *Management Science*, 13, 176-185.

[4] Baker, K. (1974), *Introduction to Sequencing and Scheduling*, New York: John Wiley & Sons.

[5] Balas, E. (1965), "An additive algorithm for solving linear programs with zero-one variables," *Operations Research*, 13: 517-546.

[6] Balas, E. (1967), "Discrete programming by the filter method," *Operations Research*, 15, 915-957.

[7] Balas, E. (1969), "Machine sequencing via disjunctive graphs: An implicit enumeration algorithm," *Operations Research*, 17: 1-10.

[8] Balas, E. (1970), "Machine sequencing: disjunctive graphs and degree-constrained subgraphs," *Naval Research Logistics Quarterly*, 17, 941-957.

[9] Bean, J. and J. Birge (1986), "Match-up real-time scheduling," *NBS Special Publication*, 724: 197-212.

[10] Benders, J. (1960), "Partitioning procedures for solving mixed-variables mathematical programming problems," *Numersche Mathematik*, 4 (3): 238-252.

[11] Blackstone, J., D. Phillips and G. Hogg (1982), "A state-of-the-art survey of dispatching rules for manufacturing job shop operations," *International Journal of Production Research*, 20 (1): 27-45.

[12] Chiu, C. (1994), "A Learning-Based Methodology for Dynamic Scheduling in Distributed Manufacturing Systems," Ph.D. Dissertation, Purdue University.

[13] Conway, R. and W. Maxwell (1967), *Theory of Scheduling*, Reading, Massachusetts: Addison-Wesley.

[14] Dantzig, G. and P. Wolfe (1960), "Decomposition principles for linear programs," *Naval Research Logistics Quarterly*, 8 (1): 101-111.

[15] Daouas, T., K. Ghedira and J. Muller (1995), "Distributed flow shop scheduling problem versus local optimization," *Proceedings of the First International Conference on Multi-Agent Systems*, Cambridge, Massachusetts: MIT Press, 1995

[16] Davis, L. (1985), "Job shop scheduling with genetic algorithms," *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, Carnegie Mellon University, 136-140.

[17] Davis, W. and A. Jones (1988), "A real-time production scheduler for a stochastic manufacturing

environment. International Journal of Computer Integrated Manufacturing," 1 (2): 101-112.

[18] Dettmer, W. (1997), *Goldratt's Theory of Constraints: A Systems Approach to Continuous Improvement*, Milwaukee, Wisconsin: Quality Press.

[19] Foo, Y. and Y. Takefuji (1988), "Stochastic neural networks for solving job-shop scheduling: Part 2 Architecture and simulations," *Proceedings of the IEEE International Conference on Neural Networks*, published by IEEE TAB: II283-II290. [20] Fox, M. (1983), "Constraint-Directed Search: A case study of Job Shop Scheduling," Ph.D. Dissertation, Carnegie-Mellon University.

[21] Gershwin, S. (1989), "Hierarchical flow control: a framework for scheduling and planning discrete events in manufacturing systems," *Proceedings of IEEE Special Issue on Discrete Event Systems*, 77: 195-209.

[22] Glover, F. (1989), "Tabu search - Part I," ORSA Journal on Computing, 1 (3): 190-206.

[23] Glover, F. (1990), "Tabu search - Part II," ORSA Journal on Computing, 2 (1): 4-32.

[24] Glover, F. (1996), "Tabu search and adaptive memory programming - advances, applications and

challenges," To appear in: *Interfaces in Computer Science and Operations Research*, The Netherlands: Kluwer Academic Publishers.

[25] Goldberg, D. and R. Lingle (1985), "Alleles, loci, and the traveling salesman problem," *Proceedings of the of the International Conference on Genetic Algorithms and Their Applications*, Carnegie Mellon University, 162-164.
[26] Goldberg, D. (1988), *Genetic Algorithms in Search Optimization and Machine Learning*, Menlo Park: California: Addison-Wesley.

[27] Goldratt, E. (1990), *Theory of Constraints*, Great Barrington, Massachusetts: North River Press.

[28] Goldratt, E. (1992), *The Goal*, Great Barrington, Massachusetts: North River Press.

[29] Grabot, B. and L. Geneste (1994), "Dispatching rules in scheduling: a fuzzy approach," *International Journal of Production Research*, 32 (4): 903-915.

[30] Hopfield, J. and D. Tank (1985), "Neural computation of decisions in optimization problems," *Biological Cybernetics*, 52: 141-152.

[31] Jeffcoat, D. and R. Bulfin (1993), "Simulated annealing for resource-constrained scheduling," *European Journal of Operational Research*, 70: 43-51.

[32] Kirkpatrick, S., C. Gelatt and M. Vecchi (1983), "Optimization by simulated annealing," *Science*, 220 (4598): 671-680.

[33] Krucky, J. (1994), "Fuzzy family setup assignment and machine balancing," *Hewlett-Packard Journal*, June: 51-64.

[34] Lawler, E. and D. Wood (1966), "Branch and bound methods: a survey," *Operations Research*, 14, 699-719.

[35] Le Pape, C. (1995), "Scheduling as intelligent control of decision-making and constraint propagation,"

Intelligent Scheduling," edited by M. Zweben and M. Fox, San Francisco, California: Morgan Kaufman, 67-98.

[36] Lo, Z. and B. Bavarian (1991), "Scheduling with neural networks for flexible manufacturing systems,"

Proceedings of the IEEE International Conference on Robotics and Automation, Sacramento, California, 818-823.
[37] McKenzie, L. (1976), "Turnpike theory," *Econometrics*, 44: 841-864.

[38] Montazer, M. and L. Van Wassenhove (1990), "Analysis of scheduling rules for an FMS," *International Journal of Production Research*, 28: 785-802.

[39] Morton, E. and D. Pentico (1993), *Heuristic Scheduling Systems*, New York: John Wiley & Sons.

[40] Panwalker, S. and W. Iskander (1977), "A survey of scheduling rules," *Operations Research*, 25 (1): 45-61.

[41] Parunak, H., B. Irish, J. Kindrick and P. Lozo (1985), "Fractal actors for distributed manufacturing control,"

Proceedings of the Second IEEE Conference on Artificial Intelligence Applications, 653-660.

[42] Quinlan, J. (1986), "Induction of decision trees," *Machine Learning*, 1: 81-106.

[43] Rabelo, L. (1990), "Hybrid Artificial Neural Networks and Knowledge-Based Expert Systems Approach to Flexible Manufacturing System Scheduling," PhD. Dissertation, University of Missouri-Rolla.

[44] Rabelo, L., M. Sahinoglu and X. Avula (1994), "Flexible manufacturing systems scheduling using Q-Learning," *Proceedings of the World Congress on Neural Networks*, San Diego, California: I378-I385.

[45] Rumelhart, D., J. McClelland and the PDP Research Group (1986), *Parallel Distributed Processing:* 

Explorations in the Microstructure of Cognition, 1: Foundations, Cambridge, Massachusetts: MIT Press.

[46] Saleh, A. (1988), "Real-Time Control of a Flexible Manufacturing Cell," Ph.D. Dissertation, Lehigh University.

[47] Shapiro, J. (1979), "A survey of Lagrangian techniques for discrete optimization," *Annals of Discrete Mathematics*, 5: 113-138.

[48] Shaw, M., S. Park and N. Raman (1992), "Intelligent scheduling with machine learning capabilities: The induction of scheduling knowledge," *IEE Transactions on Design and Manufacturing*, 24: 156-168.

[49] Slany, W. (1994), "Scheduling as a fuzzy multiple criteria optimization problem." CD-Technical Report 94/62, Technical University of Vienna.

[50] Smith, S. (1995), "OPIS: A methodology and architecture for reactive scheduling," *Intelligent Scheduling*, edited by M. Zweben and M. Fox, San Francisco, California: Morgan Kaufman, 29-66.

[51] Srinivasan, V. (1971), "A hybrid algorithm for the one machine sequencing problem to minimize total tardiness," *Naval Research Logistics Quarterly*, 18: 317-327.

[52] Starkweather, T., D. Whitley, K. Mathias and S. McDaniel (1992), "Sequence scheduling with genetic algorithms," *Proceedings of the US/German Conference on New Directions for OR in Manufacturing*, 130-148.

[53] Starkweather, T., D. Whitley and B. Cookson (1993), "A Genetic Algorithm for scheduling with resource consumption.," *Proceedings of the Joint German/US Conference on Operations Research in Production Planning and Control*, 567-583.

[54] Subrahmanyam, S., M. Zentner and J. Pekny (1996), "Making the most out of corporate information assets: the next generation of process scheduling, planning, and design tool," *Proceedings of the Process Industry Technical Conference: Looking Toward the 21st Century*, June 26-27, Erie, Pennsylvania.

[55] Sutton, R. (1988), "Learning to predict by the methods of temporal differences," *Machine Learning*, 3: 9-44.

[56] Tesauro, G. (1992), "Practical issues in temporal difference learning," *Machine Learning*, 8: 257-277.

[57] Tsujimura, Y., S. Park, S. Chang and M. Gen (1993), "An effective method for solving flow shop scheduling problems with fuzzy processing times," *Computers and Industrial Engineering*, 25: 239-242.

[58] Vakharia A. and Y. Chang (1990), "A simulated annealing approach to scheduling a manufacturing cell," *Naval Research Logistics*, 37: 559-577.

[59] Watkins, C. (1989), "Learning from Delayed Rewards," Ph.D. Dissertation, King's College, Cambridge.
[60] Werbos, P. (1995), "Neurocontrol and supervised learning: An overview and evaluation," *Handbook of*

*Intelligent Control: Neural, Fuzzy, and Adaptive Approaches*, New York: Van Nostrand Reinhold Publication, 65-89. [61] Wilkerson, L. and J. Irwin (1971), "An improved algorithm for scheduling independent tasks," *AIIE* 

Transactions, 3: 239-245.

[62] Wu, D. (1987), "An Expert Systems Approach for the Control and Scheduling of Flexible Manufacturing Systems," Ph.D. Dissertation, Pennsylvania State University.

[63] Wysk, R., D. Wu and R. Yang (1986), "A multi-pass expert control system (MPECS) for flexible manufacturing systems," *NBS Special Publication*, 724: 251-278.

[64] Yih, Y. (1990), "Trace-driven knowledge acquisition (TDKA) for rule-based real-time scheduling systems," *Journal of Intelligent Manufacturing*, 1 (4): 217-230.

[65] Zentner, M., J. Pekny, G. Reklaitis and N. Gupta (1994), "Practical considerations in using model-based optimization for the scheduling and planning of batch/semicontinuous processes," *Journal of Process Cont*rol, 4 (4): 259-280.

[66] Zhang, M. and C. Zhang (1995), "The consensus of uncertainties in distributed expert systems," *Proceedings of the First International Conference on Multi-Agent Systems*, Cambridge, Massachusetts: MIT Press.

[67] Zhang, W. and T. Dietterich (1996), "High-performance job-shop scheduling with a time-delay TD(() network," *Advances in Neural Information Processing Systems*, Cambridge, Massachusetts: MIT Press: 1025-1030.

[68] Zhou, D., V. Cherkassky, T. Baldwin and D. Hong (1990), "Scaling neural networks for job shop scheduling," *Proceedings of the International Conference on Neural Networks*, 3: 889-894.

[69] Zweben, M., B. Daun, E. Davis and M. Deale (1995), "Scheduling and rescheduling with iterative repair," *Intelligent Scheduling*, edited by M. Zweben and M. Fox, San Francisco, California: Morgan Kaufman, 241-256.