

## Data-based scheduling framework

**Abstract.** Based on the analysis of the differences and relations between traditional and data-based scheduling methods for complex manufacturing systems, we propose a data-based scheduling framework. Then we discuss how to implement it for a semiconductor manufacturing system. Finally, we introduce the state-of-the-art research on the key technologies of data-based scheduling and point out their development trends.

**Streszczenie.** Zaproponowano szkielet opracowania terminarza działań w oparciu o bazę danych. Metodę sprawdzono na przykładzie przemysłu półprzewodnikowego. (Szkielet terminarza działań oparty o bazę danych)

**Keywords:** Data-based; scheduling; complex manufacturing system; data mining.

**Słowa kluczowe:** baza danych, terminarz działań

### Introduction

Production scheduling is an important way to increase a factory's productivity and enhance its competitiveness. Generally, it determines the machines to process jobs and their process orders, batch styles, and the assignments of other key resources to them to optimize the operational performance of the factory while meeting the process and resource constraints. Since 1950s, the research on scheduling has made much progress and some results have been applied to industries successfully.

However, the production styles of manufacturing factories have been dramatically changed along with the development of information technologies. The manufacturing processes become increasingly complex, such as large-scale tasks, complicated constraints, coupling of the operational performances and uncertain scheduling environments. Consequently, considerable existing scheduling modeling and optimization methods are no longer applicable.

The development of information technologies brings not only the challenges on the modeling and optimization of complex manufacturing system scheduling problems, but also many opportunities. There are rich data in information systems of a factory, such as Enterprise Resource Planning (ERP), Manufacturing Executive System (MES), Advanced Production Control (APC) and Supervisory Control and Data Acquisition (SCADA). These data contain a plentiful of scheduling relevant knowledge. Naturally, a new idea to solve complex scheduling problems is emerging, i.e., to extract useful knowledge from related on-line and off-line data to improve the operational performance of the factory.

The remainder of this paper is organized as follows. In Section 2, we design a data-based scheduling framework based on the discussions on the differences and relations between traditional and data-based scheduling. Then we discuss how to implement the data-based scheduling framework for a semiconductor manufacturing system in Section 3. In Section 4, we introduce the state-of-the-art research related to the key technologies of data-based scheduling and point out their future trends. Section 5 gives conclusions and future works.

### Data-based scheduling framework

#### (1) The challenges facing to traditional scheduling

The motivation of research on data-based scheduling is the limits of traditional scheduling facing to complex manufacturing systems.

In view of modeling, there are some difficulties to describe, e.g., the time constraints between some steps of some process flows and re-entrant process flows, which make a scheduling model inaccurate or difficult to be optimized and analyzed. Inaccurate models lead to the

solutions that cannot work well in the real scheduling environment. Especially, the uncertain events cannot be included in a model in a real-time way. The consequence is that the model lacks accurate and adaptive parameters and cannot respond dynamically to uncertain environments.

In view of optimization, most of the scheduling optimization problems are NP-hard. It is impossible to obtain an optimal solution in a polynomial time. The satisfactory solutions can be obtained in a reasonable computation time at the sacrifice of the performance. On the other hand, the knowledge-based scheduling methods (such as one-step heuristic rules) have the ability to obtain a solution in a short time. However, such knowledge is difficult to obtain. It needs a large amount of simulation experiments. The generalization ability of the knowledge base obtained is relatively low.

On one hand, the more complex a manufacturing process, the clearer the limit of traditional scheduling. On the other hand, there is a huge amount of data containing scheduling related knowledge in ERP, MES, APC, SCADA and other information systems. Meanwhile, the development on Wireless Sensor Network (WSN) and Radio Frequency Identification (RFID) enables real-time and accurate retrieval of on-line data. It has drawn great attentions from both academia and industry to apply data-based methods to the scheduling problems of complex manufacturing systems.

#### (2) The relations between data-based and traditional scheduling

Actually, data-based scheduling is not a complete departure from traditional scheduling. On the contrary, there are close relations between them. The basic tasks of data-based scheduling are the same as those of traditional scheduling, i.e., modeling and optimization. The model is still an important part of data-based scheduling. In addition, the solutions obtained with traditional methods can serve as study samples of data-based scheduling.

There are differences between data-based scheduling and traditional scheduling. The former pays more attention to the role of knowledge in scheduling, and emphasizes the adaptivity to real environments, the operability of the scheduling solutions and powerful real-time response to the uncertainties in complex manufacturing systems. In addition, its models are not used to directly guide the operations in a production line. Its solutions may be the ideal scheduling performance indicators characterizing its capacity or property to provide an optimal basis for optimization methods.

#### (3) The application scope of data-based scheduling

Obviously, the scheduling problems have no requirements on data-based scheduling if they can be optimally solved with traditional scheduling methods. The

systems optimized with data-based scheduling should at least meet one of following conditions. Firstly, the study samples (e.g., on-line, off-line or simulation data) required by data-based scheduling can be obtained from the systems. Secondly, the scheduling problems of a system cannot be accurately modeled, or can be modeled, without expressions with high precision on the parameters. Thirdly, the scheduling problems of a system can be modeled with

high precision, but its optimum or satisfaction solution cannot be obtained in a reasonable time.

#### (4) Data-based scheduling framework

According to the above analysis, this work proposes a data-based scheduling framework that is composed of data layer, model layer and scheduling layer as illustrated in Fig.1.

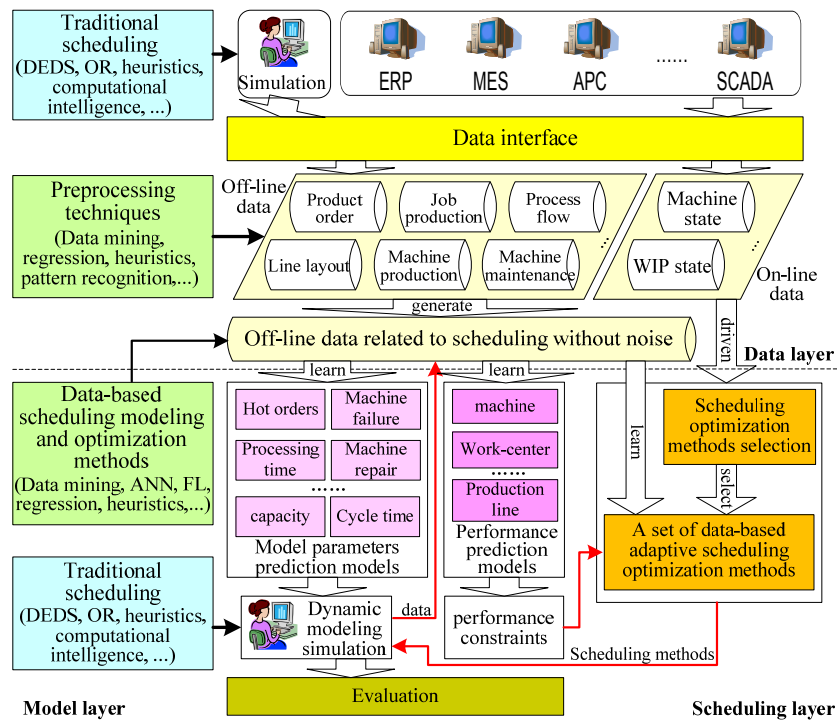


Fig.1. Data-based scheduling framework

- Data layer

The precondition to implement data-based scheduling is the rich data sources related to scheduling, which forms the data layer of the framework. There are lots of data related to scheduling in information systems of a company, such as ERP, MES, ACP and SCADA. Meanwhile, we can use traditional scheduling methods (such as discrete event simulation system, heuristics, and mathematical programming) to generate a plentiful of data related to scheduling. These data can be classified into off-line and on-line data. The former include the information about production orders, job production, process flows, machine layout, machine production, machine maintenance, etc. The latter are the real-time information reflecting the working condition of production environments, such as information about machine state (e.g., idle, busy, and occupied) and WIP state (being processed, queuing and hold). Because these data are commonly incomplete and contain noise, the preprocessing techniques are often required to execute filtration, purification, denoising, and optimization. Then these preprocessed data can be taken as study samples to extract useful scheduling-relevant knowledge.

- Model layer

There are two kinds of prediction models in the model layer, i.e., parameter and performance prediction models, which are built by learning from the preprocessed data from the data layer. The former is to predict the parameters of scheduling models, such as the possibility of hot orders' arrival, machine failure and machine maintenance, processing time of one step, capacity of a machine and cycle time of a product. These parameters represent the

occurrence probability of uncertain events (such as the first four), or the production performance (such as the last two) of a manufacturing system. If they are integrated within a dynamic modeling simulation system, the accuracy and precision of the performance prediction can be improved. Then the dispatching rules mined from these simulation data can be expected with better performance. The latter is to predict the operational performance and constraints of a machine, a work-center or production line in a period, which can provide useful guidance for the scheduling layer to select a suitable dispatching rule and find a near-optimal or satisfactory solution quickly.

- Scheduling layer

The scheduling layer is responsible for the job dispatching of a manufacturing system. The adaptive scheduling optimization methods in this layer are learned from the preprocessed data in the data layer. The characteristics of the scheduling environments adapting to and concerned scheduling performance issues of each method may be different. In a real production environment, a choice of a method is jointly determined by the on-line data (such as machine state and WIP state) and the ideal performance objectives and constraints obtained by the data-based scheduling models.

#### The realization of the data-based scheduling framework

Taking a semiconductor manufacturing system as an example, the realization scheme of the proposed data-based scheduling framework has the following steps as shown in Fig.2.

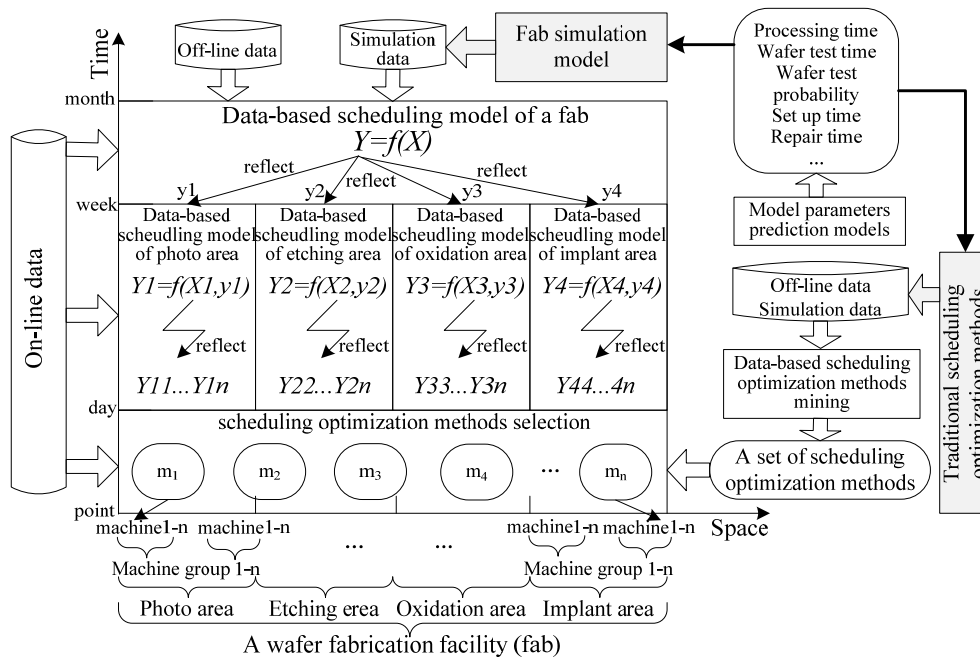


Fig.2. A realization scheme of the proposed data-based scheduling framework for a fab

Step 1: The off-line data in a semiconductor wafer fabrication facility (fab) and the simulation data generated by its dynamic modeling simulation model are taken as study samples. Then we use various methods (such as data mining, artificial neural-network, support vector machine and evolutionary algorithm) to learn from these data to build the scheduling models of the fab, its work areas (such as photo area, etching area, oxidation area and implant area), machine groups and machines. Next, the scheduling model of the fab is used to predict the expected performance of the fab in a certain period (e.g., a month, week or day), which are further decomposed to work areas, machine groups and machines. Their performances provide useful guidance to real-time dispatchers.

Step 2: The study samples for mining adaptive scheduling optimized methods include the off-line data of the fab and the scheduling plans generated by traditional scheduling methods (such as computational intelligence, simulation, and heuristic rules). It is notable that the concerned performance issues of each adaptive scheduling optimized method may be different due to its study samples with different characteristics.

Step 3: At the dispatch point (i.e., a machine becomes idle and there are queuing jobs before it), the machine selects the most suitable one from the set of data-based scheduling optimized methods according to its expected performance and on-line data relevant to scheduling.

### Key technologies of data-based scheduling

#### (1) Data pre-processing techniques

The existing research results on data preprocessing can be classified into attribute selection and data clustering.

Attribute selection is to select important ones from the condition attributes. The common methods for it include fuzzy set, feature selection, and computational intelligence. For example, Kusiak[1] proposed a fuzzy set based method to obtain an attribute selection rule by learning from samples for semiconductor manufacturing quality problems. Chen *et al.*[2] shrunk the search space with the concept of feature nuclear, and then used an ant colony algorithm to obtain the reduction of a set of attributes to improve the efficiency of knowledge reduction. Shiue *et al.* [3] developed a two stage decision tree based adaptive scheduling

system, and proposed an artificial neural network (ANN) based feature selection algorithm and a genetic algorithm (GA) for attribute selection, respectively.

Data clustering is a technique to classify the samples according to their similarity. The samples with high similarity are included in the same class. This technique can be used to delete the noisy data. The common methods for data clustering include SOM, Fuzzy-C means, K-means and ANN. For example, Hu and Su[4] presented a hierarchical clustering method to find the machines related to the decreasing yield. Chen[5] attempted to integrate fuzzy-C means and K-means with backward propagation neuro-network (BPNN) to cluster the training samples and train an ANN for each clustering data. SOM was used to smooth the noisy data in the samples and improve the learning effects.

Powerful data pre-processing technologies are still required to find useful knowledge from complicated mass data with noise and deficiency.

#### (2) Data-based scheduling modeling methods

Data-based scheduling modeling of a manufacturing system is to reflect its off-line or on-line data in its information systems in the models describing its production scheduling process or predicting its operational performance. The existing work can be classified into data-based description and prediction models.

Data-based description models focus on mapping rules between their data sets and scheduling models of a manufacturing system. As a result, they can be changed conveniently and flexibly by modifying the related data. For example, Mueller [6] proposed a modeling method to reflect the production data of a fab into its object oriented Petri net simulation model. The factors, such as batch processing, failure time of machines and tools and rework jobs, were considered. The main shortage of the work was the great simplification on the real fab and lacked the consideration on its non-zero initial state as well. Ye *et al.*[7] developed a dynamic modeling technique to dynamically build a discrete event simulation model of a fab with off-line and on-line data in its information systems. The on-line data made the simulation model able to reflect the non-zero initial state of the fab. Since the mapping between the data and the simulation model is closely related to the simulation

software eM-Plant, the generality of the conversion method requires further improvement.

Data-based prediction models are used to determine the parameters of scheduling models and predict the occurrence probability of uncertain events (such as hot order arrival and machine failure) by learning from related historical data. For example, Bagchi *et al.*[8] pointed out that the processing time of a job in a fab was related to the recipes of machines, the number of the wafers in a lot (i.e., job), the set-up requirements of machines, and the number of parallel loading ports and slots (chambers). Then they provided a multiple regression analysis method to build the relations between these parameters and the processing time. Chen and Wang[9] applied a fuzzy ANN to building the prediction model of cycle time with high precision. Arrendo and Martinez[10] proposed a reinforcement learning strategy with local weighted regression to determine whether to accept the hot orders. Shukla *et al.*[11] developed a bidding based multi-agent system for flexible job shops, and used a fuzzy decision tree to mine the occurrence probability of tools failure.

Unfortunately, there still lacks the research on the data-based operational performance prediction models. Further study is needed to offer necessary guidance for data-based scheduling optimization methods.

### (3) Data-based scheduling optimization methods

Data-based scheduling optimization methods are used to learn useful dispatching rules from off-line data in information systems or the scheduling plans obtained by simulations or intelligent optimization methods.

For example, Lee[12] extracted fuzzy dispatching rules from training samples by using an ANN-based learning algorithm. Yang *et al.* [13] proposed the concept of MAS-based intelligent manufacturing systems, and developed a B-Q learning algorithm with reference to clustering technology to learn how to select dispatching rules. Chaudhuriz and De[14] obtained a black block based scheduler to dispatch jobs by making use of rough fuzzy multilayer perception neural networks to approximate optimal plans obtained by GA. Kumar and Rao[15] optimized the flow shop with an ant colony algorithm and used C4.5 to learn dispatching rules from those optimal plans. Olafsson and Li[16] proposed a GA-based dispatching rule selection method. It had two stages: the first stage was to select better plans by using a GA algorithm; and the second one was to learn the dispatching rules from these plans by using decision-tree based methods. Choi *et al.* [17] used a decision-tree based method to find the knowledge for selecting dispatching rules from off-line data in a re-entrant manufacturing system, with consideration of its real-time state.

Existing research results on data-based scheduling optimization methods have their common deficiency, i.e., lack of flexibility. They select a dispatching rule from a specified rules set at the dispatching point according to real-time running state or use a learned dispatching rule off-line during the whole decision process without adaptations. The performance of the former is dependent on the specified dispatching rules set, while that of the latter is dependent on the diversity of study samples. In addition, the manufacturing systems considered are small-scale job shops or flow shops. The research should move forward.

### Conclusion

The development of information technology helps advance the automation level of complex manufacturing systems. There are much more on-line and off-line data in their information systems, which provide useful knowledge, rules and optimal decisions for solving their scheduling

problems. We propose a data-based scheduling framework, introduce how to use it by taking a semiconductor manufacturing system as an example and summarize key technologies related to data-based scheduling. Our future work is to further improve them and apply the research results to the actual production environments.

### REFERENCES

- [1] Kusiak A. Feature Transformation Methods in Data Mining. *IEEE Transactions on Electronics Packing Manufacturing*, 24(2001), No.3, 214-221
- [2] Chen Y M, Miao D Q, Wang R Z. A Rough Set Approach to Feature Selection Based on Ant Colony Optimization, *Pattern Recognition Letters*, 31(2010), No.3, 226-233
- [3] Shiue Y R. Development of Two-Level Decision Tree-Based Real-Time Scheduling System under Product Mix Variety Environment, *Robotics and Computer-Integrated Manufacturing*, 25(2009), No.4-5,709-720
- [4] Hu C H, Su S F. Hierarchical Clustering Methods for Semiconductor Manufacturing Data, *Proceedings of the 2004 IEEE International Conference on Networking, Sensing Control*, 2004, Taiwan, 1063-1068
- [5] Chen T. Predicting Wafer-Lot Output Time with a Hybrid FCM-FBPN Approach, *IEEE Transactions on System, Man and cybernetics-Part B: Cybernetics*, 37(2007), No.4, 784-793
- [6] Mueller R, McGinnis L F. Automatic Generation of Simulation Models for Semiconductor Manufacturing, *Proceedings of the 2007 Winter Simulation Conference*, Washington, DC, United States, 2007, 648-657
- [7] Ye K, Qiao F, Ma Y M. General Structure of the Semiconductor Production Scheduling Model, *Applied Mechanics and Materials*, 20-23(2010), 465-469
- [8] Bagchi S, Baseman R J, Davenport A, Natarajan R, Slonim N, Weiss S. Data Analytics and Stochastic Modeling in a Semiconductor Fab, *Applied Stochastic Models in Business and Industry*, 26(2010), No.1, 1-27
- [9] Chen T, Wang Y C. A Nonlinear Scheduling Rule Incorporating Fuzzy-Neural Remaining Cycle Time Estimator for Scheduling a Semiconductor Manufacturing Factory—a Simulation Study, *International Journal of Advanced Manufacturing Technology*, 45(2009), No.1-2, 110-121
- [10] Arredondo F, Martinez E. Learning and Adaptation of a Policy for Dynamic Order Acceptance in Make-To-Order Manufacturing, *Computers and Industrial Engineering*, 58(2010), No.1, 70-83
- [11] Shukla K S, Tiwari M K, Son Y J. Bidding-Based Multi-Agent System for Integrated Process Planning and Scheduling a Data-Mining and Hybrid Tabu-SA Algorithm-Oriented Approach, *International Journal of Advanced Manufacturing Technology*, 38(2008), No.1-2, 163-175
- [12] Lee K K. Fuzzy Rule Generation for Adaptive Scheduling in a Dynamic Manufacturing, *Applied Soft Computing*, 8(2008), No.4, 1295-1304
- [13] Yang H B, Yan H S. An Adaptive Approach to Dynamic Scheduling in Knowledgeable Manufacturing Cell, *International Journal of Advanced Manufacturing Technology*, 42(2009), No.3-4, 312-320
- [14] Chaudhuri A, De K. Job Scheduling Problem Using Rough Fuzzy Multilayer Perception Neural Networks, *Journal of Artificial Intelligence: Theory and Application*, 1(2010), No.1, 4-19
- [15] Kumar S, Rao C S P. Application of Ant Colony, Genetic Algorithm and Data Mining-Based Techniques for Scheduling, *Robotics and Computer-Integrated Manufacturing*, 25(2009), No.6, 901-908
- [16] Olafsson S, Li X N. Learning Effective New Single Machine Dispatching Rules from Optimal Scheduling Data, *International Journal of Production Economics*, 128(2010), No.1, 118-126
- [17] Choi H S, Kim J S, Lee D H. Real-time scheduling for reentrant hybrid flow shops: A decision tree based mechanism and its application to a TFT-LCD line, *Expert Systems with Applications*, 38(2011), No.4, 3514-3521

**Authors:** Dr. Li Li is an associate professor in School of Electronics and Information Engineering, Tongji University, No.1239, Siping Road, Shanghai, China, E-mail: lili@tongji.edu.cn. The correspondence address is: lili@tongji.edu.cn