



**HAL**  
open science

## Concurrent WIP and an application to clearing functions for complex heterogenous systems

Pierre Lemaire, Kean Dequeant, Marie-Laure Espinouse, Philippe Vialletelle

► **To cite this version:**

Pierre Lemaire, Kean Dequeant, Marie-Laure Espinouse, Philippe Vialletelle. Concurrent WIP and an application to clearing functions for complex heterogenous systems. IFAC-PapersOnLine, 2022, 55 (10), pp.696-701. 10.1016/j.ifacol.2022.09.487 . hal-03863851

**HAL Id: hal-03863851**

**<https://hal.science/hal-03863851>**

Submitted on 21 Nov 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Concurrent WIP and an application to clearing functions for complex heterogeneous systems

Pierre Lemaire\* Kean Dequeant\* Marie-Laure Espinouse\*  
Philippe Vialletelle\*\*

\* Univ. Grenoble Alpes, CNRS, Grenoble INP, G-SCOP, 38000  
Grenoble, France (e-mail: pierre.lemaire@grenoble-inp.fr,  
marie-laure.espinouse@grenoble-inp.fr)

\*\* STMicroelectronics, 38926, Crolles Cedex, FRANCE (e-mail:  
philippe.vialletelle@st.com)

**Abstract:** Complex manufacturing systems are challenging to study because of the high level of information required and the inaccessibility of most of it. Their tractability is however essential for the efficiency of state-of-the-art industries. This is particularly the case in the semiconductor industry that faces high mix and low volume conditions, and for which traditional methods fail to capture the high complexity and require continuous actions and corrections to adjust to heterogeneous toolsets and product-mix.

We present the Concurrent WIP (CWIP), a new way of studying such systems at the level of a process-cluster by identifying each job's queue from its own perspective. CWIP is designed to be practical, with a low level of resource investments, yet informative. We explain how CWIP can be computed based on historical data and then used to derive capacity estimates and clearing functions without any assumptions on the system or on the form of the functions. In the process, we derive not only an average workload-dependent capacity, but also a confidence interval on this capacity. The relevance and efficiency of the proposed estimates are experimentally tested on a simulated system mimicking a small but complex process-cluster of the semiconductor industry. The estimates are used to predict WIP absorption times and we show how they characterize well not only the average behavior but also the full range of possible behaviors of the system. Finally, we discuss further applications of CWIP, that could be used to compute refined clearing functions or to monitor complex systems.

Copyright © 2022 The Authors. This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

*Keywords:* complex logistic systems, production planning and control, modelling and decision making in complex systems, methodologies and tools for analysis of complexity, data-driven decision making

## 1. INTRODUCTION

The shift of many industries from single-line single-product to complex reconfigurable and flexible multi-product manufacturing systems (ElMaraghy, 2008) is seen in semiconductors (Dequeant et al., 2016b), photonics (Frazee and Standridge, 2016), and many more “traditional” industries (Brettel et al., 2014) and is the natural result of increasing demands for specialized products and the necessity for industries to keep economies of scale. In these industries, a process-cluster (a set of process-units that share a flow of products to perform a particular manufacturing step) shifts from doing one specific task with identical process-units to doing multiple heterogeneous tasks with multiple heterogeneous process-units. The consequence of this shift to High Mix Low Volume production, or more generally to the context of Industry 4.0, is that an accurate understanding and a precise, continuously-

adapted control of the manufacturing systems become central yet hard to achieve.

Traditional approaches (e.g., queuing theory, simulation...), which attempt to model systems, become awfully difficult to carry out: the amount of information required to just describe accurately complex systems becomes huge, let alone putting it all together. Besides, in many cases, one does not need a full representation of the system; then the cost to set up and to maintain such a well-tailored model is prohibitive.

Indeed, in the scope of this paper, we are not so much interested by a manufacturing system itself, but rather by its behavior; our focus is on better understanding how an existing manufacturing system actually responds to usual conditions of use so as to be able to control it effectively. A typical need is to challenge production plans to anticipate bottlenecks, delays, etc., so that operational corrections can be implemented in time and in place. To this extend, we propose a new way of studying manufacturing systems

\* This work benefited from a funding of the French National Agency of Technical Research (ANRT).

that is generic, easy to implement, and nevertheless allows to derive useful characteristics such as clearing functions of a set of process-units.

In the remainder of this paper, we review major existing methods for the study of manufacturing systems with a focus on cycle times and capacities (section 2). We introduce our new approach, Concurrent WIP (section 3) and explain how it can be used to derive capacity estimates and clearing functions (section 4). The approach is validated on a simulated complex heterogeneous system (section 5). A discussion and perspectives are proposed to conclude (section 6).

## 2. MODELLING COMPLEX MANUFACTURING SYSTEMS

Our work originates from the semiconductor industry, and more particularly from plants that display several challenging characteristics: flexibility, process-units' heterogeneity, batching, setups, re-entrancy, process-units' breakdowns, process-time variability, arrivals variability, job-dispatching, queuing rules, as well as other less conventional sources of variability, as reported by Dequeant et al. (2016b). However we intend our work to be widely applicable from simple cases to complex ones (by “complex”, we merely mean that the amount of information required to accurately describe the system is extremely high; usually part of this information is even unknown) and we do not assume the presence or absence of any characteristic. Besides, an important aspect of the tractability of a system is that not only the conjunction of many sources makes the system complex, but each individual source can be complex by itself (for instance real-life process-times are not independent and follow complex patterns based on the health of the process-unit). Another major component of the overall complexity is that the sum of all the small sources of complexity can be really impactful on the macroscopic scale: the conjunction of factors that are too complex to model individually cannot be neglected as a whole.

Many systems can be represented accurately with a variety of methodologies. Queueing theory, which relies on a classical server point-of-view, has been widely used. Whitt (1993) applied it to calling centers while Hopp and Spearman (2011) have set a milestone for all practitioners with their general queueing theory model; many other works have been done, partly summarized by Dequeant et al. (2016b), Shanthikumar et al. (2007), C and Appa Iyer (2013), and Wu (2014). Other approaches have also been used, for instances: markovian processes have been applied by Gurumurthi and Benjaafar (2004) to study flexible queueing systems; linear programming models have been proposed by Romauch and Hartl (2017) for capacity planning of cluster-tools restricted to multi-chamber tools in a static context; Vamsikrishna and Padmanabhan (2016) have reviewed the use of Petri nets in flexible manufacturing systems. To the best of our knowledge, the work closest to ours is that of Etman et al. (2011): based on a similar analysis (bottom-up modeling gets challenged due to the high level of information required), they have also proposed an aggregated, top-bottom, data-driven approach, but they have concentrated on the effective process

time, with an assumption of independency and with the purpose of estimating the parameters of an equivalent simple closed-form queueing model. On the practitioners' side and especially in the semiconductor manufacturing industry, a standard practice to study process-clusters is to build a capacity model from the theoretical capacity of all process-units, to individually measure all inefficiencies (down-times, idle-times, setup-times...) and to remove them to get an estimation of the overall capacity, as shown by Martin (1999).

All those approaches provide detailed and well-tailored representations of the system under study, but they require a precise and detailed knowledge of it. In the case of complex systems, in particular with heterogeneous jobs and heterogeneous process-units, that information may be not known, not available or, in the best case, amount to too high a level of information inducing too high costs when a detailed representation of the system is not needed. Our objective is therefore to allow local studies of the actual performances of a complex system at a low level of resource investments.

When it comes to planning purposes, cycle times and capacities are of particular importance, as stressed by Pahl et al. (2007).

Following Hopp and Spearman (2011), we define the cycle time of a job as the time between its arrival at a processing step and its departure from this step; it includes processing times but also waiting times and possibly transportation times (cycle times are often called lead times; however “cycle time” is the preferred term in microelectronics and some other industries, see Hopp and Spearman (2011)).

The capacity of a system is usually defined as the maximum throughput the system can achieve. This is well illustrated by Little's Law (Little, 1961, 2011) that states that the throughput ( $TH$ ) of a system is equal to the ratio of the average number of jobs in the system ( $WIP$ ) by the average cycle time ( $CT$ ):  $TH = WIP/CT$ . Then, the throughput corresponds to the capacity of the system if there is always at least one job waiting. Little's Law also underlines the relation between capacity and cycle time.

Even though accurate cycle times or capacities are essential for accurate planning, they are often considered constant and independent of resource utilization (Kacar and Uzsoy, 2010); the MRP methodology is a typical example. Pahl et al. (2007) states that load dependent cycles times are still rare in the literature and that the work-around of using maximal cycle times leads to earlier job releases, higher WIP levels and thus even longer cycle times. It is thus essential to take into account that cycle times and capacities are influenced by workload, batching and sequencing decisions, WIP levels, etc.

The most significant effort to establish variable cycle times is probably provided by clearing functions. The notion originates in Graves (1986) (constant proportions, linear cycle times) and has been refined ever since. As illustrated by Figure 1, typical clearing functions relate throughput (or capacity) as a non-decreasing concave function of the WIP level. Several close-forms have been proposed on several queueing systems, e.g., by Karmarkar (1989) and Asmundsson et al. (2009) (see Pahl et al. (2007)

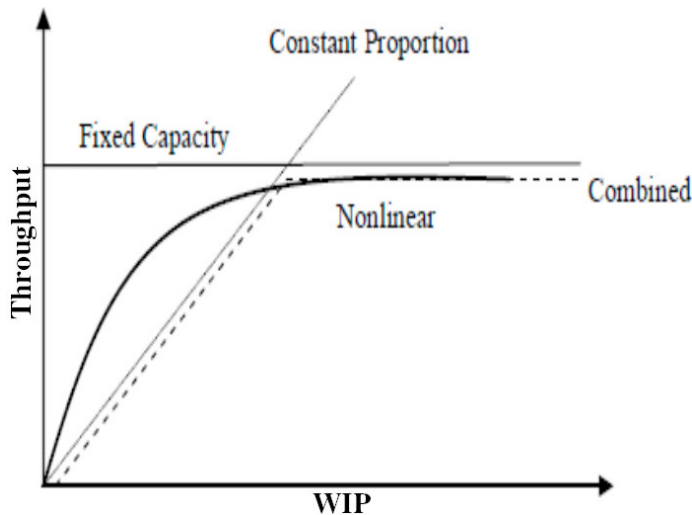


Fig. 1. Examples of clearing functions (adapted from Karmarkar (1989); Kacar and Uzsoy (2010)).

for an overview). A major advantage of those clearing functions is that their close-form allows for an efficient integration within mathematical programming model for optimal planning. However, even more recent versions as those of Kacar and Uzsoy (2010) remain problem-specific.

In what follows, we aim at deriving clearing functions that are representative of the actual behavior of complex systems without making assumptions regarding distributions, independences, or the existence of an equivalent closed-form model (we want to cope with real complex and fuzzy systems).

### 3. CONCURRENT WIP (CWIP)

Concurrent WIP (CWIP) has been first introduced in Dequeant et al. (2016a); Dequeant (2017) but is worth being described here. It all starts by changing the reference frame: instead of classically studying the system from the process-clusters point-of-view, we propose to study it from the jobs point-of-view.

The CWIP is a notion defined for every job that waited before being processed. Informally, it corresponds to the total amount of processing that has been achieved by the process-cluster while the job was waiting for being processed. Discarding jobs that did not wait is not an issue in our case as such jobs would not provide insightful information about the process-cluster; somehow, such jobs have witness nothing, so they have nothing to say. They may fall into two categories: either “emergency” jobs that are processed right-away at the cost of stopping a currently-processed job, but such jobs should be exceptional in an industrial context and talk about themselves, not about the normal behavior of the process-cluster; or jobs that arrived when some process-unit was idle and ready, and that only informs that the process-cluster was currently under-utilized.

Hence, we place ourselves at an individual process-cluster and consider a given job  $p$  that waited for a non-null time before being processed. The reference frame of job  $p$  is the time interval between its arrival at the process-cluster and its actual process start: we call it its Effective

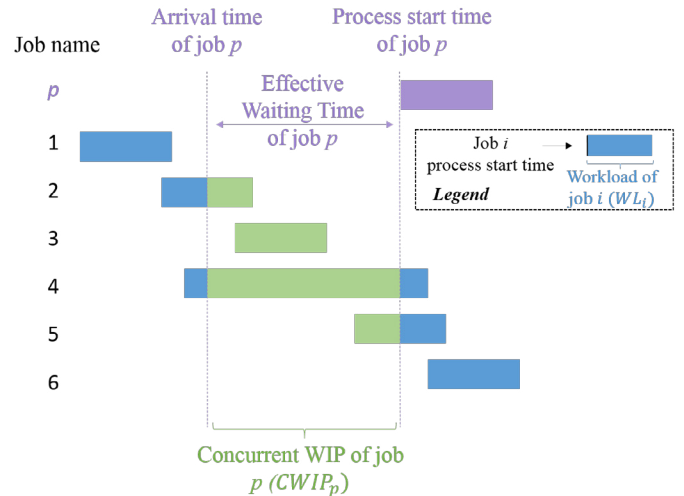


Fig. 2. The Concurrent WIP of a job  $p$  as the sum of the workloads of other jobs theoretically processed during job  $p$  Effective Wait Time window on the same process-cluster.

Waiting Time ( $EWT_p$ ). For every other job  $i$ , we call  $WL_i$  its Workload (several definitions are possible; see details below) and we then define the Concurrent WIP of a job  $p$  ( $CWIP_p$ ) as:

$$CWIP_p = \sum_i WL_i \times x_{i,p} \quad (1)$$

where  $x_{i,p}$  is the theoretical fraction of the workload of job  $i$  processed in the time-window defined by  $EWT_p$ .

The theoretical workload process of each job  $i$  that we consider here starts when the process of job  $i$  effectively starts, and lasts for a duration  $WL_i$ : if the process starts at time  $s_i$ , it is assumed to end at time  $c_i = s_i + WL_i$ . The theoretical fraction  $x_{i,p}$  is then the proportion of theoretical workload process that happens during  $EWT_p$ , so that  $WL_i \times x_{i,p}$  corresponds to the workload amount of job  $i$  that is expected to have been processed while  $p$  was waiting. Fig. 2 illustrates the 6 different cases. If  $a_p$  is the arrival date of job  $p$  and  $s_p$  the date its process starts, then the general formula is:

$$x_{i,p} = \max \left\{ 0, \frac{\min(c_i, s_p) - \max(s_i, a_p)}{c_i - s_i} \right\}. \quad (2)$$

Effectively,  $CWIP_p$  is the total amount of workload job  $p$  had to wait to be processed before starting its own process. The theoretical process used for this definition ensures the consistency between what is waiting to be processed and what has been processed. Note that in a simple FIFO single-family-job case, all jobs have the same workload and  $CWIP_p$  simply corresponds to the number of unprocessed jobs at the time of arrival (which corresponds to the traditional view of WIP).

We purposely speak of Effective Waiting Time as it is not necessarily composed of only waiting time. In practice, it can be composed of transportation time, loading time, and other non-waiting but non-processing times. However, from a logistical point-of-view, it can be seen as effective waiting time on the same reasoning as in (Etman et al., 2011).

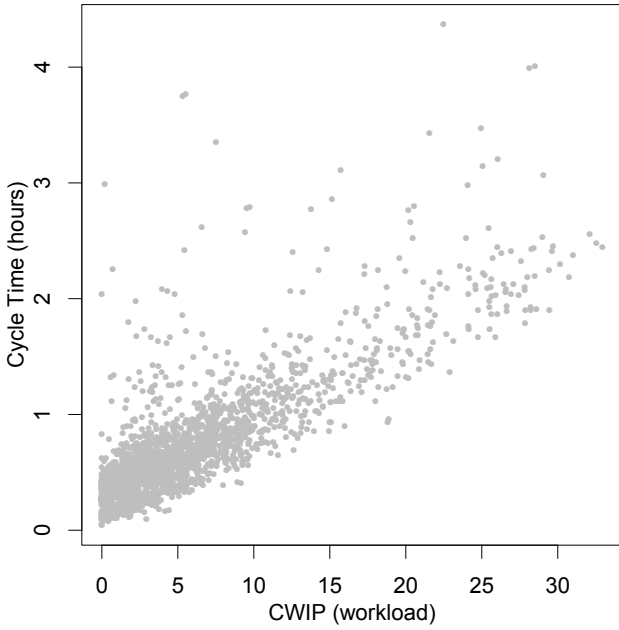


Fig. 3. Cycle Time vs Concurrent WIP for a simulated complex system.

It is essential to define  $CWIP_p$  in terms of workload, and not in terms of number of jobs, because heterogeneous jobs are assumed, and different job families have different expected process times. The workload is therefore a measure of the expected occupancy weight that a job has on the process-cluster. In what follows, workloads can be assumed to be average per-family processing-times measured in (expected) hours-of-process, but other workload measures and units could be chosen, depending on the underlying need.

For an actual system, CWIP is easy to compute from historical data, provided one knows arrival times, start times and completion times of each job. As an example, Figure 3 pictures cycle times vs CWIP for a simulated complex system (this system is described in appendix A).

#### 4. CWIP-BASED CLEARING FUNCTIONS

We know, from section 2, that estimating capacity is useful but tricky for heterogeneous systems. In this section we show how Concurrent WIP enables to capture the behavior of the system and, in particular, to derive measures of capacity and clearing functions, without making any particular assumptions.

From Little’s Law, one can derive the capacity of a system as the ratio between the average amount of WIP and the average cycle time (provided there is always WIP present). However, this quantity is a long-term capacity of a stable system and it is thus improper for short-term planning taking system variability into account. To this extend, we adapt the ratio to define  $C_p$ , the capacity witnessed by job  $p$  while waiting:

$$C_p = CWIP_p / EWT_p \quad (3)$$

By definition of the CWIP, this quantity is well defined ( $EWT_p > 0$ ) and is indeed a capacity since there is always at least one job waiting (job  $p$ ). This capacity is

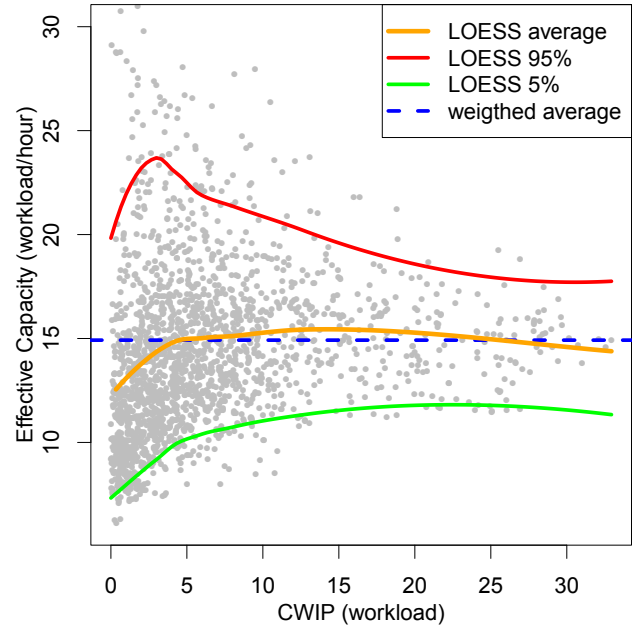


Fig. 4. Effective capacity vs Concurrent WIP for a simulated complex system.

expressed in terms of units-of-workload/hours (instead of a traditional number-of-jobs/hours), in order to enable to cope with heterogeneous jobs and different job families. Besides, note that this capacity is not *the* capacity of the process-cluster, in the sense of a maximum achievable performance; it is the *effective capacity* of the whole system defined by the process-cluster, dispatching rules, workload level, mix, *etc.*, all those particular conditions being caught by the CWIP. Doing so, we shift from a theoretical capacity barely seen in practice to a practical measure of the response of a process-cluster to particular conditions.

For each job we have a snapshot of the effective capacity, as pictured by Figure 4 (it is, of course, the same data as the data used for Figure 3). One can see that the capacity varies in mean and in variance with the CWIP. To better quantify this last assertion, several aggregations are possible, as we describe below.

To capture capacity as a single measure, we propose the following formula:

$$\bar{C} = \frac{\sum_p CWIP_p \times C_p}{\sum_p CWIP_p} \quad (4)$$

This formula computes the average effective capacity seen by all jobs. Jobs are weighted by their Concurrent WIP so as to give a proportionally higher importance to jobs that witnessed the system’s capacity over a higher workload. This measure is shown as a blue dashed line on Figure 4.

To capture the variability of the capacity, we use a LOESS regression procedure (Cleveland et al., 1992), to get the average capacity as a function of the workload (orange curve on Figure 4). This curve corresponds to the clearing function of the system under study. Remark that, as expected for a clearing function, it is concave and non-decreasing (except at the end where data is sparse).

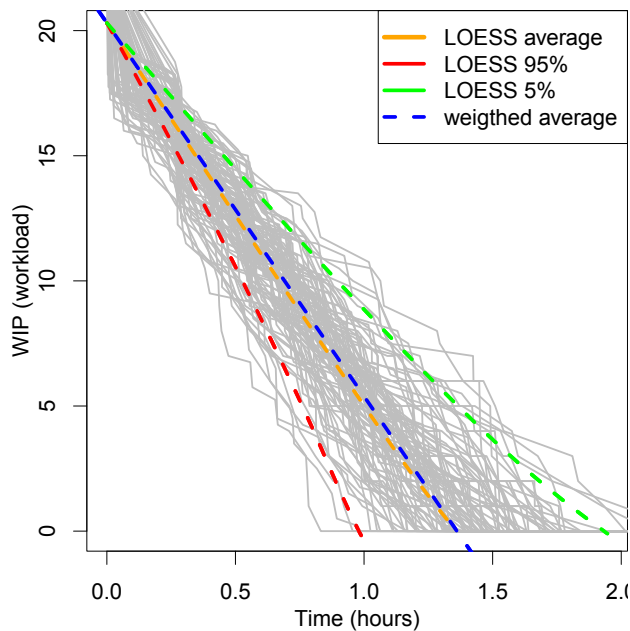


Fig. 5. WIP absorption for a simulated complex system.

In addition to the average capacity, we also propose to derive regressions for extreme quantiles (in our case, 5% and 95%, green and red curves on Figure 4) in order to get a more precise description of the capacity. The capacity can now be seen as a random variable for which we know the average behavior and a 90% confidence interval, both depending on the level of CWIP.

The relevance of the proposed measures is experimentally tested on the next section.

## 5. AN APPLICATION TO WIP ABSORPTION

To ensure the relevance of what we have just discussed, we challenge our capacity measures and clearing functions on our simulated system. More particularly, we test whether WIP absorption can be well-anticipated.

WIP absorption corresponds to the time needed to process all waiting jobs, providing there are no more arrivals. In the process, the amount of WIP decreases, which should impact the effectiveness of the system.

For our simulated system, historical data is known from a first simulation and have already been presented (Figure 4). From this data, we know that a WIP of 20 hours-of-workload is a large yet realistic level of WIP; it corresponds to the upper-bound level of WIP experienced by 95% of the jobs. Then, the weighted average capacity and the LOESS regressions allow to anticipate the time required to absorb such an amount of WIP.

Figure 5 depicts this “corridor” of WIP absorption. Absorbing a level of WIP of 20 hours-of-workload should take at least 0.99 hour and at most 1.95 hours (90% confidence level, red and green curves) and, on average 1.36 hours (given by both LOESS average and weighted average, orange and blue curves). Those values are derived a priori from the corresponding effective capacities of Fig 4.

Now that WIP absorption has been anticipated, let see how it actually goes. An experiment runs as follows:

we stop arrivals when 20 hours-of-workload of WIP are waiting in front of our simulated process-cluster, and we measure how long it takes for the system to process it all. Since the process is stochastic, we run the same experiment 100 times to get a full view of what could happen. On Figure 5, each grey line corresponds to the trajectory (remaining WIP level as a function of time) for one experiment.

One can observe that those trajectories indeed fall around the mean absorption curve, and within the whole range between the optimistic and the pessimistic guesses. On the 100 runs, the 5%, 50% and 95% quantiles for absorption times are 1.01, 1.36 and 1.72 hours respectively, very close to the anticipated values.

This experiment shows that CWIP-derived clearing functions allow to capture not only the average behavior of the system, but also its variability and the whole range of possible behaviors.

## 6. CONCLUSION

We have introduced the notion of Concurrent WIP (CWIP) as a novel way to describe the queue of each job from their unique perspective. The usefulness of CWIP lies in its ability to describe the short-term response of the system under variable conditions. This has been illustrated for WIP absorption and a variable amount of WIP. More generally, it allows to derive relevant data-driven clearing-functions, even in the context of complex manufacturing systems. Contrary to other approaches found in the literature, no assumptions have to be made on the system itself, on the influencing factors or on the form of the clearing function. Moreover, CWIP can be computed at a very low cost, since only standard historical data is required.

Practically, CWIP and the derived clearing functions allow to provide accurate short-term estimations of cycle times, with measures of possible deviations. It should integrate well into finite capacity planning heuristics such as the one proposed by Mhiri et al. (2015).

Among the perspectives, CWIP and derived estimates could be refined, similarly to what is discussed in Etman et al. (2011). Indeed, they can be adapted to integrate any aspect of choice: for instance, if one wants to know the effective capacity for a given job-family or as a function of priority, one just has to compute the aggregation (e.g., formula (4)) accordingly; the only requirement is that enough data is available to compute reliable values. For example, computing the effective capacity for different priorities would allow to describe the impact of priorities on cycle times.

CWIP could also be used to monitor a system. Computing capacities and clearing functions (similar to Figure 4) for the same system on different periods would allow to compare the effectiveness of the system on the different periods and reveal capacity losses.

## REFERENCES

- Asmundsson, J., Rardin, R.L., Turkseven, C.H., and Uzsoy, R. (2009). Production planning with resources subject to congestion. *Naval Research Logistics (NRL)*, 56(2), 142–157. doi:10.1002/nav.20335.

- Brettel, M., Friederichsen, N., Keller, M., and Rosenberg, M. (2014). How Virtualization, Decentralization and Network Building Change the Manufacturing Landscape: An Industry 4.0 Perspective. *International Journal of Mechanical, Industrial Science and Engineering*, 8(11), 37–44.
- C, L. and Appa Iyer, S. (2013). Application of queueing theory in health care: A literature review. *Operations Research for Health Care*, 2(1), 25–39. doi:10.1016/j.orhc.2013.03.002.
- Cleveland, W.S., Grosse, E., and Shyu, W.M. (1992). Local Regression Models. In *Statistical Models in S*. Routledge.
- Dequeant, K. (2017). *Workflow variability modeling in microelectronic manufacturing*. phdthesis, Université Grenoble Alpes.
- Dequeant, K., Lemaire, P., Espinouse, M.L., and Vialletelle, P. (2016a). Le wip concurrent : une proposition de file d'attente du point de vue du produit pour caractériser le temps de cycle. In *11th International Conference on Modeling, Optimization & SIMulation*, 9 pages. Montreal, Canada.
- Dequeant, K., Vialletelle, P., Lemaire, P., and Espinouse, M.L. (2016b). A literature review on variability in semiconductor manufacturing: The next forward leap to Industry 4.0. In *2016 Winter Simulation Conference (WSC)*, 2598–2609. doi:10.1109/WSC.2016.7822298.
- ElMaraghy, H.A. (2008). *Changeable and Reconfigurable Manufacturing Systems*. Springer Science & Business Media.
- Etman, L., Veeger, C., Lefeber, E., Adan, I., and Rooda, J. (2011). Aggregate modeling of semiconductor equipment using effective process times. In *Proceedings of the 2011 Winter Simulation Conference (WSC)*, 1790–1802. doi:10.1109/WSC.2011.6147894.
- Frazer, T. and Standridge, C.R. (2016). CONWIP versus POLCA: A comparative analysis in a high-mix, low-volume (HMLV) manufacturing environment with batch processing. *Journal of Industrial Engineering and Management*, 9(2), 432–449. doi:10.3926/jiem.1248.
- Graves, S.C. (1986). A Tactical Planning Model for a Job Shop. *Operations Research*, 34(4), 522–533. doi:10.1287/opre.34.4.522.
- Gurumurthi, S. and Benjaafar, S. (2004). Modeling and analysis of flexible queueing systems. *Naval Research Logistics*, 51(5), 755–782. doi:10.1002/nav.20020.
- Hopp, W.J. and Spearman, M.L. (2011). *Factory Physics: Third Edition*. Waveland Press.
- Kacar, N.B. and Uzsoy, R. (2010). Estimating clearing functions from simulation data. In *Proceedings of the 2010 Winter Simulation Conference*, 1699–1710. doi:10.1109/WSC.2010.5678899.
- Karmarkar, U.S. (1989). Capacity loading and release planning with work-in-progress (WIP) and leadtimes. *Journal of Manufacturing and Operations Management*, 2(105-123).
- Lemaire, P. (2019). Concurrent-WIP: an illustrative example with application to WIP absorption. Technical report, <https://hal.archives-ouvertes.fr/hal-02084050>.
- Little, J.D.C. (1961). A Proof for the Queueing Formula:  $L = \lambda W$ . *Operations Research*, 9(3), 383–387. doi:10.1287/opre.9.3.383.
- Little, J.D.C. (2011). OR FORUM—Little's Law as Viewed on Its 50th Anniversary. *Operations Research*, 59(3), 536–549. doi:10.1287/opre.1110.0940.
- Martin, D. (1999). Capacity and cycle time-throughput understanding system (CAC-TUS) an analysis tool to determine the components of capacity and cycle time in a semiconductor manufacturing line. In *10th Annual IEEE/SEMI. Advanced Semiconductor Manufacturing Conference and Workshop. ASMC 99 Proceedings (Cat. No.99CH36295)*, 127–131. doi:10.1109/ASMC.1999.798198.
- Mhiri, E., Jacomino, M., Mangione, F., Vialletelle, P., and Lepelletier, G. (2015). Finite capacity planning algorithm for semiconductor industry considering lots priority. *IFAC-PapersOnLine*, 48(3), 1598–1603. doi:10.1016/j.ifacol.2015.06.314.
- Pahl, J., Voß, S., and Woodruff, D.L. (2007). Production planning with load dependent lead times: an update of research. *Annals of Operations Research*, 153(1), 297–345. doi:10.1007/s10479-007-0173-5.
- Romauch, M. and Hartl, R.F. (2017). Capacity planning for cluster tools in the semiconductor industry. *International Journal of Production Economics*, 194, 167–180. doi:10.1016/j.ijpe.2017.01.005.
- Shanthikumar, J.G., Ding, S., and Zhang, M.T. (2007). Queueing Theory for Semiconductor Manufacturing Systems: A Survey and Open Problems. *IEEE Transactions on Automation Science and Engineering*, 4(4), 513–522. doi:10.1109/TASE.2007.906348.
- Vamsikrishna, C. and Padmanabhan, G. (2016). Role of petri nets in flexible manufacturing system: a review. *International Journal of Engineering Trends and Technology*, 41(2), 90–100.
- Whitt, W. (1993). Approximations for the GI/G/m queue. *Production and Operations Management*, 2(2), 114–161. doi:10.1111/j.1937-5956.1993.tb00094.x.
- Wu, K. (2014). Classification of queueing models for a workstation with interruptions: a review. *International Journal of Production Research*, 52(3), 902–917. doi:10.1080/00207543.2013.843799.

## Appendix A. A SIMULATED COMPLEX SYSTEM

Throughout this paper, we use a simulated system, rather small in size, but complex in the settings. It corresponds to the processing of jobs from three families on 5 process-units.

Each family represents, in average, one third of the jobs; priorities and process-times follow family-based non-standard distributions; arrivals follow a Poisson process. The jobs are scheduling according to a priority-based policy. The process-units have different speeds, different batching-capabilities and there are incompatibilities among job-families and process-units (see Case J4-M5 in Lemaire (2019) for details).

2000 jobs have been generated and scheduled to form our “historical data”.