

Partial-Order Reduction for Performance Analysis of Max-Plus Timed Systems

Bram van der Sanden*, Marc Geilen*, Michel Reniers*, Twan Basten*[§]

*Eindhoven University of Technology, Eindhoven, The Netherlands

[§]ESI, Eindhoven, The Netherlands

Email: b.v.d.sanden@tue.nl, m.c.w.geilen@tue.nl, m.a.reniers@tue.nl, a.a.basten@tue.nl

Abstract—This paper presents a partial-order reduction method for performance analysis of max-plus timed systems. A max-plus timed system is a network of automata, where the timing behavior of deterministic system tasks (events in an automaton) is captured in $(\max,+)$ matrices. These tasks can be characterized in various formalisms like synchronous data flow, Petri nets, or real-time calculus. The timing behavior of the system is captured in a $(\max,+)$ state space, calculated from the composition of the automata. This state space may exhibit redundant interleaving with respect to performance aspects like throughput or latency. The goal of this work is to obtain a smaller state space to speed up performance analysis. To achieve this, we first formalize state-space equivalence with respect to throughput and latency analysis. Then, we present a way to compute a reduced composition directly from the specification. This yields a smaller equivalent state space. We perform the reduction on-the-fly, without first computing the full composition. Experiments show the effectiveness of the method on a set of realistic manufacturing system models.

I. INTRODUCTION AND RELATED WORK

Performance is one of the key aspects in the design of complex systems. Besides having to meet functional requirements, systems need to adhere to timing constraints, and optimize productivity, typically expressed with throughput or latency metrics. Throughput describes the system performance in the long run, for instance the number of products produced by the system per hour. Latency describes the temporal distance between certain events, for instance the time between the start and end of processing a product. Usually, system performance can only be measured at a later stage in the development process, once the system is assembled. A model-based design approach to performance engineering [1] can be used to address this issue. In such an approach, formal models capture the system behavior under the various scenarios of execution. Moreover, by adding timing information, performance analysis techniques can be used to predict the system performance at an early stage in the design process. However, in many industrial applications, the underlying timed state space of these models becomes large, and performance analysis quickly becomes a bottleneck.

In this work, we introduce a new partial-order reduction technique to speed up performance analysis of timed systems. The reduction explores only a restricted number

of interleavings of concurrently enabled system operations that use different sets of system resources. The ample conditions on the reduction guarantee that the performance properties of the original model are preserved in the reduced model. Traditional partial-order techniques typically consider local properties such as deadlocks, and temporal properties formulated in logics like $LTL_{\setminus \circ}$ (next-time-free Linear-time Temporal Logic) [2] and $CTL_{\setminus \circ}^*$ (next-time-free Computation Tree Logic) [3].

There has been some initial work in applying partial-order reduction techniques to timed systems. Bengtsson et al. [4] apply standard partial-order reduction on timed automata for reachability analysis. These automata execute asynchronously, in their own local time scale, and synchronize their time scales on communication transitions. This work has been extended by Minea [5] to perform model checking for an extension of LTL, that can express timing relations between events. Yoneda et al. [6] investigated partial-order reduction for timed Petri nets, that allows the verification of similar timing relations. Theelen et al. [7] apply ideas from partial-order reduction on Scenario-Aware Data Flow models, where they use an independence relation among actions to resolve non-deterministic choices that have no impact on the performance metrics.

In this paper, we consider max-plus timed systems as a formal model. Such systems are described by a set of $(\max,+)$ automata [8] and a composition operator. A $(\max,+)$ automaton is a conventional automaton, where the timing semantics of each system task (event in an automaton) is described by a $(\max,+)$ matrix. Such a matrix captures the corresponding timing behavior, induced by the corresponding action execution times and action dependencies. Max-plus timed systems can express the timing semantics of a broad range of specification formalisms, such as Network Calculus [9], Real-Time Calculus [10], Synchronous Data Flow [11], Scenario-Aware Data Flow [12], and Timed-Event Graphs [13], [14], an important subclass of timed Petri nets. A broad range of industrial systems can be expressed using these formalisms.

To illustrate the concepts in this paper we use a specification framework [15] suitable for performance analysis of manufacturing systems. It allows analysis of both throughput and latency as performance metrics. System tasks are described by *activities*, which consist of a set of *actions* that execute on *resources* and dependencies

among those actions. The timing behavior of each activity is captured by a $(\max,+)$ matrix. A parallel composition of $(\max,+)$ automata describes the order in which activities can be executed. As composition operator, multi-party synchronization [16] is used. From the system specification, a timed $(\max,+)$ state space is derived that captures the system behavior and the necessary timing information to evaluate system throughput and latency. We use partial-order reduction to compute a reduced composition of $(\max,+)$ automata directly from the specification. From the reduced composition a reduced state space is computed that preserves performance properties.

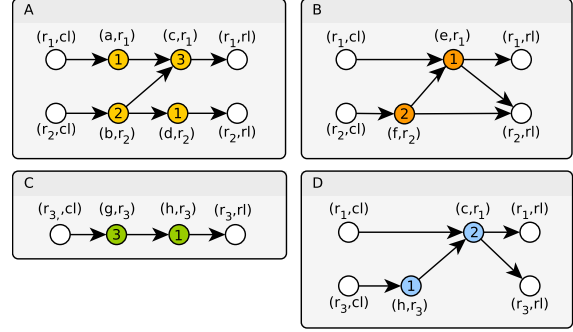
Well-known partial-order reduction techniques include the *stubborn sets* method of Valmari [17], the *persistent sets* method of Godefroid [18], and the *ample sets* method of Peled [19]. The idea in these methods is to exploit information about the independence of certain activities to reduce the size of the state space, while preserving the properties of interest. In each state, only a subset of all possible transitions is selected. In this work, we use ample sets and the extension of cluster-based ample sets [20].

In the remainder of this paper, we first formally introduce max-plus timed systems and define the properties of interest. Then, we introduce a reduction function that preserves the specified properties on the level of the state space. Next, we introduce local conditions to compute a reduced composition automaton directly from the network of $(\max,+)$ automata. The reduced state space can be computed from the reduced composition automaton. An experimental evaluation shows the effectiveness of the reduction technique. Due to space limits, the proofs have been omitted. Proofs of the theoretical results presented in this paper can be found in [21].

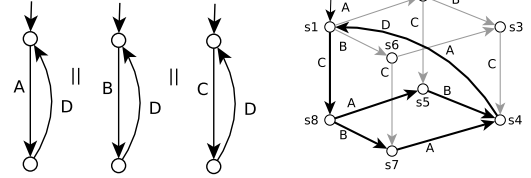
II. MAX-PLUS TIMED SYSTEMS

In the specification framework that we use [15], a system is modeled in terms of *resources* that provide the *actions* that the system can execute. Deterministic system tasks are described by *activities*. An activity consists of a fixed set of action instances and dependencies among those action instances. A resource must be claimed before its actions can be used. After execution of the actions, the resource must be released. The timing information of each activity is captured in a $(\max,+)$ matrix. This matrix describes the release time of each system resource in terms of when the resources are available at the start of executing the activity. The availability times of resources are captured in a $(\max,+)$ vector. Given such a *resource availability vector*, we obtain the new resource availability vector after execution of some activity by multiplying with the corresponding $(\max,+)$ matrix.

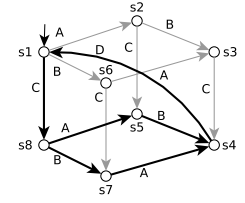
We use $(\max,+)$ algebra (see for instance [22]) to capture the timing semantics of activities in a concise way. Two essential characteristics of the execution of an activity are *synchronization*, when an action inside an activity waits until all preceding actions are finished, and *delay*, when an action execution takes an amount of time before it completes. These characteristics correspond well to



(a) Activities A, B, C and D. Actions are from a till h, and the action timings are given inside the nodes. Resources r_1, r_2 and r_3 are claimed (*cl*) and released (*rl*).



(b) Max-plus timed system. Activities A, B, D have reward 0; C has reward 1.



(c) Composition of the $(\max,+)$ automata. The reduction is shown with thick transitions.

Fig. 1. Running Example.

$$\begin{aligned}
 M_A &= \begin{bmatrix} 4 & 5 & -\infty \\ -\infty & 3 & -\infty \\ -\infty & -\infty & 0 \end{bmatrix} & M_B &= \begin{bmatrix} 1 & 3 & -\infty \\ 1 & 3 & -\infty \\ -\infty & -\infty & 0 \end{bmatrix} \\
 M_C &= \begin{bmatrix} 0 & -\infty & -\infty \\ -\infty & 0 & -\infty \\ -\infty & -\infty & 4 \end{bmatrix} & M_D &= \begin{bmatrix} 2 & -\infty & 3 \\ -\infty & 0 & -\infty \\ 2 & -\infty & 3 \end{bmatrix}
 \end{aligned}$$

Fig. 2. $(\max,+)$ matrices of activities A, B, C and D.

the $(\max,+)$ operators *maximum* (\max) and *addition* ($+$), defined over the set $\mathbb{R}^{-\infty} = \mathbb{R} \cup \{-\infty\}$. Operators \max and $+$ are defined as usually in algebra, with the additional convention that $-\infty$ is the unit element of \max : $\max(-\infty, x) = \max(x, -\infty) = x$, and the zero-element of $+$: $-\infty + x = x + -\infty = -\infty$.

Since $(\max,+)$ algebra is a linear algebra, it can be extended to matrices and vectors in the usual way. Given matrix \mathbf{A} and matrix \mathbf{B} , we use $\mathbf{A} \otimes \mathbf{B}$ to denote the $(\max,+)$ matrix multiplication. Given $m \times p$ matrix \mathbf{A} and $p \times n$ matrix \mathbf{B} , the elements of the resulting matrix $\mathbf{A} \otimes \mathbf{B}$ are determined by: $[\mathbf{A} \otimes \mathbf{B}]_{ij} = \max_{k=1}^p ([\mathbf{A}]_{ik} + [\mathbf{B}]_{kj})$. For any vector \mathbf{x} , $\|\mathbf{x}\| = \max_i [\mathbf{x}]_i$ denotes the vector norm of \mathbf{x} . For vector \mathbf{x} , with $\|\mathbf{x}\| > -\infty$, we use $norm(\mathbf{x})$ to denote $\mathbf{x} - \|\mathbf{x}\|$, the normalized vector, such that $\|norm(\mathbf{x})\| = 0$. We use $\mathbf{0}$ to denote a vector with only zero-valued entries.

In this paper, we use the running example shown in Fig. 1. This max-plus timed system consists of activities A, B, C and D (see Fig. 1a), and three $(\max,+)$ automata (see Fig. 1b). Fig. 1c shows the composition of the $(\max,+)$ automata. Each activity has a corresponding $(\max,+)$

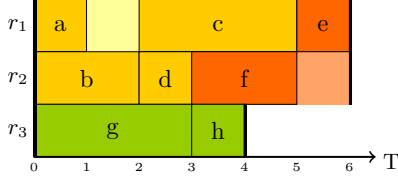


Fig. 3. Gantt chart of activity sequence ABC when all resources are initially available.

matrix that captures the timing behavior, shown in Fig. 2. Each matrix row represents the symbolic release time of a resource in terms of all system resources, which is denoted by set \mathcal{R} . Let R be a function that maps each activity to the set of resources the activity uses. As an example, consider the first row of matrix M_A . This row describes the release time of resource $r_1 \in \mathcal{R}$, expressed in terms of when resources r_1, r_2 and r_3 are available at the start of executing A . In the execution of activity A , there is a timing delay of 4 time units between the claiming of resource r_1 and the subsequent release of resource r_1 . Similarly, a delay of 5 is present between the claiming of resource r_2 and the release of r_1 . There is no dependency between the availability times of resource r_3 and the release of resource r_1 , indicated by $-\infty$. This can also be seen in the structure of activity A (Fig. 1a), since resource r_3 is not involved in the execution. For details on how to compute the $(\max,+)$ matrices of activities, see [15].

The timing evolution of the system is expressed using $(\max,+)$ matrix multiplication. Assume that all resources are initially available, captured in vector $\mathbf{0}$. The new availability times of the resources after executing activity A with corresponding matrix M_A are computed as follows:

$$M_A \otimes \mathbf{0} = \begin{bmatrix} \max(4 + 0, 5 + 0, -\infty + 0) \\ \max(-\infty + 0, 3 + 0, -\infty + 0) \\ \max(-\infty + 0, -\infty + 0, 0 + 0) \end{bmatrix} = \begin{bmatrix} 5 \\ 3 \\ 0 \end{bmatrix}.$$

Resources r_1 and r_2 are available again after 5 and 3 time units respectively. The availability time of resource r_3 stays 0, since it is not used by activity A .

The timing semantics of an activity sequence is defined in terms of repeated matrix multiplication. As an example, consider the execution of activity sequence ABC , shown in Fig. 3. The new resource availability vector after execution of this sequence is computed as follows:

$$M_C \otimes M_B \otimes M_A \otimes \mathbf{0} = [6, 6, 4]^T.$$

To capture all possible activity orderings in the system, we use $(\max,+)$ automata.

Definition 1 ($(\max,+)$ automaton (adapted from [8])). *A $(\max,+)$ automaton \mathcal{A} is a tuple $\langle S, \hat{s}, Act, reward, M, T \rangle$ where S is a finite set of states, $\hat{s} \in S$ is the initial state, Act is a set of activities, function $reward : Act \rightarrow \mathbb{R}^{\geq 0}$ quantifies the amount of progress per activity, function M maps each activity to its associated $(\max,+)$ matrix of size $|\mathcal{R}| \times |\mathcal{R}|$, and $T \subseteq S \times Act \times S$ is the transition relation. Let $s \xrightarrow{A} s'$ be a shorthand for $\langle s, A, s' \rangle \in T$. It is assumed that*

\mathcal{A} is deterministic, which means that for any $s, s', s'' \in S$ and $A \in Act$, $s \xrightarrow{A} s'$ and $s \xrightarrow{A} s''$ imply $s' = s''$.

Let $\mathcal{A} = \langle S, \hat{s}, Act, reward, M, T \rangle$ be a $(\max,+)$ automaton. An activity $A \in Act$ is said to be *enabled* in a state $s \in S$ if $s \xrightarrow{A} s'$ for some $s' \in S$. Set $enabled(s) = \{A \in Act \mid \exists s' : s \xrightarrow{A} s'\}$ contains all activities enabled in s . State s is a *deadlock* state if $enabled(s) = \emptyset$. Since \mathcal{A} is deterministic, for any activity $A \in enabled(s)$, there is a unique A -successor of s , denoted by $A(s)$. For an activity sequence $A_1 \dots A_n$, the resulting state is defined inductively as $(A_1)(s) = A_1(s)$ if $A_1 \in enabled(s)$, and $(A_1 \dots A_n A_{n+1})(s) = A_{n+1}((A_1 \dots A_n)(s))$ if $A_{n+1} \in enabled((A_1 \dots A_n)(s))$. Otherwise, $(A_1 \dots A_{n+1})(s)$ is undefined.

A $(\max,+)$ automaton is an ω -automaton [23] that accepts infinite ω -words over Act . There are no specific acceptance conditions on these words, so any infinite word that conforms to a sequence of transitions starting in the initial state is accepted. A possible behavior of the $(\max,+)$ automaton is described in a run. An *infinite* run ρ of \mathcal{A} is an infinite, alternating sequence of states and activities:

$$\rho = s_0 A_1 s_1 A_2 s_2 A_3 \dots \text{ such that } s_0 = \hat{s} \\ \text{and } s_{i+1} = A_{i+1}(s_i) \text{ for all } i \geq 0.$$

Given run ρ , let $\rho[..i]$ denote the prefix $s_0 A_1 \dots s_i$, and let $\rho[i, j] = s_i A_{i+1} \dots s_j$ denote a run fragment from state s_i until s_j . Furthermore, let $\rho[i]$ denote state s_i in ρ .

We use the term $(\max,+)$ automaton to emphasize that the timing semantics of the automaton is expressed in $(\max,+)$ algebra. Note that the original definition in [8] does not consider rewards, but this extension is considered for instance in Weakly-Consistent Scenario-Aware Data Flow [24]. It allows for a refined, explicit specification of progress. A \max -plus timed system is defined as a composition of $(\max,+)$ automata. In our framework we use multi-party synchronization as composition operator, which is defined in the following way.

Definition 2 (Multi-party synchronization). *Given $(\max,+)$ automata $\mathcal{A}_1 = \langle S_1, \hat{s}_1, Act_1, reward_1, M_1, T_1 \rangle$ and $\mathcal{A}_2 = \langle S_2, \hat{s}_2, Act_2, reward_2, M_2, T_2 \rangle$, we define the multi-party synchronization $\mathcal{A}_1 \parallel \mathcal{A}_2 = \langle S_1 \times S_2, \langle \hat{s}_1, \hat{s}_2 \rangle, Act_1 \cup Act_2, reward_1 \cup reward_2, M_1 \cup M_2, T_{12} \rangle$, where*

$$T_{12} = \begin{cases} \langle s_1, s_2 \rangle \xrightarrow{A}_{12} \langle s'_1, s'_2 \rangle & \text{if } A \in Act_1 \cap Act_2, \\ & s_1 \xrightarrow{A}_1 s'_1, s_2 \xrightarrow{A}_2 s'_2 \\ \langle s_1, s_2 \rangle \xrightarrow{A}_{12} \langle s'_1, s_2 \rangle & \text{if } A \in Act_1 \setminus Act_2, s_1 \xrightarrow{A}_1 s'_1 \\ \langle s_1, s_2 \rangle \xrightarrow{A}_{12} \langle s_1, s'_2 \rangle & \text{if } A \in Act_2 \setminus Act_1, s_2 \xrightarrow{A}_2 s'_2. \end{cases}$$

It is assumed that functions $reward_i$ and M_i with $i \in \{1, 2\}$ are equal on $Act_1 \cap Act_2$.

Definition 3 (Max-plus timed system). *A max-plus timed system \mathcal{M} is described by $M = \mathcal{A}_1 \parallel \dots \parallel \mathcal{A}_n$ with $(\max,+)$ automata \mathcal{A}_i and multi-party synchronization operator \parallel . It is assumed that all matrices have the same dimensions*

of $|\mathcal{R}| \times |\mathcal{R}|$, and that functions reward_i and M_i with $i \in \{1, \dots, n\}$ are equal on the same activities.

The composition of all the individual automata is again an automaton. Fig. 1c shows the composition of the automata shown in Fig. 1b. Each $(\max, +)$ automaton can be interpreted as a normalized $(\max, +)$ state space that captures all the accepted runs, and contains all the necessary information to evaluate performance properties.

Definition 4 (Normalized $(\max, +)$ state space (adapted from [12])). *Given $(\max, +)$ automaton $\mathcal{A} = \langle S, \hat{s}, \text{Act}, \text{reward}, M, T \rangle$ with matrices of size $|\mathcal{R}| \times |\mathcal{R}|$, we define the normalized $(\max, +)$ state space $\mathcal{S} = \langle C, \hat{c}, \text{Act}, \Delta, M, w_1, w_2 \rangle$ as follows:*

- set $C = S \times \mathbb{R}^{-\infty^{|\mathcal{R}|}}$ of configurations that consists of a state and a normalized (resource availability) vector;
- initial configuration $\hat{c} = \langle \hat{s}, \mathbf{0} \rangle$;
- a labeled transition relation $\Delta \subseteq C \times \text{Act} \times C$ that consists of the transitions in the set $\{ \langle \langle s, \gamma \rangle, A, \langle s', \text{norm}(\gamma') \rangle \rangle \mid s \xrightarrow{A} s' \wedge \gamma' = M(A) \otimes \gamma \}$;
- function w_1 that assigns a weight $w_1(c, A, c') = \text{reward}(A)$ to each transition $\langle c, A, c' \rangle \in \Delta$;
- function w_2 that assigns a weight $w_2(c, A, c') = \|M(A) \otimes \gamma\|$ to each transition $\langle c, A, c' \rangle \in \Delta$. This weight indicates the total added execution time to the complete schedule.

We define the set of enabled activities and runs in a $(\max, +)$ state space in a similar way as in a $(\max, +)$ automaton. The state space of $\mathcal{A}_1 \parallel \dots \parallel \mathcal{A}_n$ is computed in two steps. First, we compute the composition, and subsequently we compute the corresponding state space. Fig. 4 shows the state space of the max-plus timed system of Fig. 1. In some cases, this state space might be infinite [12]. The state space is guaranteed to be finite, if for every activity sequence u allowed by the $(\max, +)$ automaton and any $k \geq 0$, there is some $m > k$ such that the matrix $M_{u(k)} \otimes \dots \otimes M_{u(m-1)}$ contains no entries $-\infty$ [12]. This means that a complete resource-to-resource synchronization has been achieved. A sufficient condition for this is that the claim of each resource is linked to the claim of each other resource via the resource dependencies over each cycle in the automaton. In the $(\max, +)$ automaton shown in Fig. 1c, each cycle involves activities A, B, C , and D , and the resource claims of r_1, r_2 and r_3 are mutually dependent; r_1 to r_2 by activities A and B , and r_1 (and through it also r_2) to r_3 by activity D . Therefore, the corresponding state space shown in Fig. 4 is finite.

The behavior of a $(\max, +)$ automaton \mathcal{S} is captured by set $\mathcal{R}(\mathcal{S})$ of all allowed runs. A run $\rho \in \mathcal{R}(\mathcal{S})$ is an infinite, alternating sequence of configurations and activities:

$$\rho = c_0 A_1 c_1 A_2 c_2 A_3 \dots \text{ such that } c_0 = \hat{c} \\ \text{ and } c_{i+1} = A_{i+1}(c_i) \text{ for all } i \geq 0.$$

Given run ρ , we define run prefix $\rho[..i] = c_0 A_1 \dots c_i$, run fragment $\rho[i, j] = c_i A_{i+1} \dots c_j$ from configuration c_i until c_j , and $\rho[i] = c_i$. We also define vector $\bar{\gamma}_n =$

$(\otimes_{k=1}^n M(A_k)) \otimes \mathbf{0}$, which is the resulting resource availability vector after executing activities $A_1 \dots A_n$ without normalization. These vectors can be derived from the normalized $(\max, +)$ state space.

Theorem 5. *Let \mathcal{S} be a $(\max, +)$ state space, and $\rho = c_0 A_1 c_1 A_2 c_2 A_3 \dots$ be a run in \mathcal{S} . Then, for each $n \geq 0$ it holds that $\bar{\gamma}_n = \sum_{k=0}^{n-1} w_2(c_k, A_{k+1}, c_{k+1}) + \gamma_n$.*

Example 6. *Consider the normalized $(\max, +)$ state space shown in Fig. 4, and the execution of activity sequence ABC starting from the initial state. This corresponds to some run ρ that starts with run fragment $\rho[0, 3] =$*

$$\langle s_1, \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rangle \xrightarrow{A, 5} \langle s_2, \begin{bmatrix} 0 \\ -2 \\ -5 \end{bmatrix} \rangle \xrightarrow{B, 1} \langle s_3, \begin{bmatrix} 0 \\ 0 \\ -6 \end{bmatrix} \rangle \xrightarrow{C, 0} \langle s_4, \begin{bmatrix} 0 \\ 0 \\ -2 \end{bmatrix} \rangle.$$

The vector in each configuration without normalization can now be computed using Theorem 5; $\bar{\gamma}_0 = \gamma_0 = \mathbf{0}$ and

$$\bar{\gamma}_1 = \mathbf{0} + 5 + [0, -2, 5]^\top = [5, 3, 0]^\top \\ \bar{\gamma}_2 = \mathbf{0} + 5 + 1 + [0, 0, -6]^\top = [6, 6, 0]^\top \\ \bar{\gamma}_3 = \mathbf{0} + 5 + 1 + 0 + [0, 0, -2]^\top = [6, 6, 4]^\top.$$

Fig. 3 shows the same availability times of 6, 6, and 4 for resources r_1, r_2 , and r_3 after executing ABC .

III. PARTIAL-ORDER REDUCTION FOR PERFORMANCE ANALYSIS

In this section we present a partial-order reduction technique that preserves throughput and latency properties. We first define throughput and latency.

Throughput: We quantify the throughput of a run as the ratio between the total reward (sum of w_1 weights) and the total execution time (sum of w_2 weights).

Definition 7 (Ratio value of a run). *The ratio of a run $\rho = c_0 A_1 c_1 A_2 c_2 A_3 \dots$ is the ratio of the sums of weights w_1 and w_2 , defined as follows*

$$\text{Ratio}(\rho) = \limsup_{l \rightarrow \infty} \frac{\sum_{i=0}^l w_1(c_i, A_{i+1}, c_{i+1})}{\sum_{i=0}^l w_2(c_i, A_{i+1}, c_{i+1})}.$$

We define the ratio value of a run fragment $\rho[i, j]$ as

$$\text{Ratio}(\rho[i, j]) = \frac{\sum_{k=i}^j w_1(c_k, A_{k+1}, c_{k+1})}{\sum_{k=i}^j w_2(c_k, A_{k+1}, c_{k+1})}.$$

The system throughput is determined by the possible ratio values over all infinite runs on some state space \mathcal{S} . We can quantify a guarantee on the throughput of the system by the *minimum ratio value* achieved by any of those runs:

$$\tau_{\min}(\mathcal{S}) = \min_{\rho \in \mathcal{R}(\mathcal{S})} \text{Ratio}(\rho).$$

If \mathcal{S} is finite, each infinite run eventually reaches a recurrent configuration. Each reachable simple cycle in this state space allows for a periodic execution of the system. Since \mathcal{S} has a finite number of simple cycles (where no repetition of transitions is allowed), we can determine the minimum ratio value of the graph from a *minimum cycle*

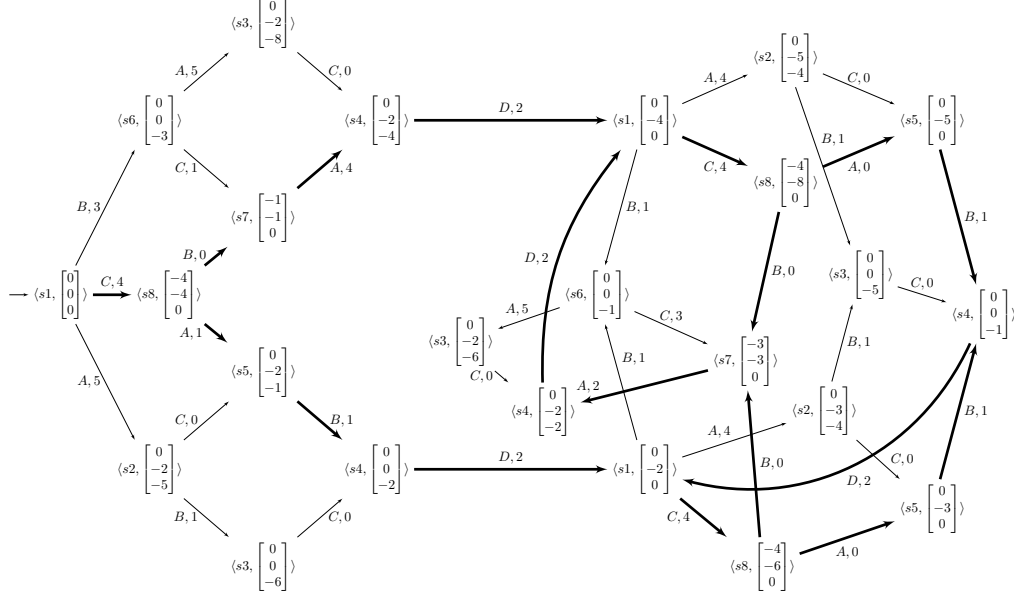


Fig. 4. Normalized (max,+) state space of the (max,+) automaton shown in Fig. 1c. The reduced state space is shown with thick transitions. Transitions are annotated with the corresponding activity and w_2 value. Activity C has reward 1, and activities A , B , and D have reward 0.

ratio (MCR) analysis [25]. The MCR over all cycles in \mathcal{S} , say $\text{cycles}(\mathcal{S})$, is defined in the following way:

$$\text{MCR}(\mathcal{S}) = \min_{c \in \text{cycles}(\mathcal{S})} \text{Ratio}(c) = \tau_{\min}(\mathcal{S}).$$

Example 8 (Cycle ratio). Consider the normalized (max,+) state space \mathcal{S} shown in Fig. 4. Recall that activity C has a reward of 1, and activities A , B , and D have a reward of 0. In this way, the ratio relates to the number of C occurrences per time unit. The minimum cycle ratio $\text{MCR}(\mathcal{S}) = 3/8$, which can for instance be found in the following cycle corresponding to the execution of $(CBAD)^\omega$:

$$\langle s_1, \begin{bmatrix} 0 \\ -4 \\ 0 \end{bmatrix} \rangle \xrightarrow{C,4} \langle s_8, \begin{bmatrix} -4 \\ -8 \\ 0 \end{bmatrix} \rangle \xrightarrow{B,0} \langle s_7, \begin{bmatrix} -3 \\ -3 \\ 0 \end{bmatrix} \rangle \xrightarrow{A,2} \langle s_4, \begin{bmatrix} 0 \\ -2 \\ -2 \end{bmatrix} \rangle \xrightarrow{D,2} \langle s_1, \begin{bmatrix} 0 \\ -4 \\ 0 \end{bmatrix} \rangle.$$

The other periodic executions where B precedes A , i.e. $(BACD)^\omega$ and $(BCAD)^\omega$, have the same minimum cycle ratio value.

Latency: In general, latency is the time delay between a stimulus and its effect. In the context of max-plus timed systems, we define the *latency* in terms of the temporal distance that separates the resource availability times of a resource at the start of two activities A_{src} and A_{snk} . In the state space, consider some run $\rho = c_0 A_1 c_1 A_2 \dots$ with $c_i = \langle s_i, \gamma_i \rangle$ containing run fragment $\rho[i, j+1] = c_i A_{i+1} \dots c_j A_{j+1} c_{j+1}$, with $A_{i+1} = A_{src}$ and $A_{j+1} = A_{snk}$. Then we define the start-to-start latency λ between the resource availability times of resource r in γ_i and γ_j as

$$\lambda(\rho, i, j, r) = [\bar{\gamma}_j]_r - [\bar{\gamma}_i]_r.$$

Example 9 (Latency). Consider again the execution of activity sequence ABC starting from configuration c_0 in

the (max,+) state space shown in Fig. 4. Suppose we want to compute the start-to-start latency between the resource availability times of r_1 in $\bar{\gamma}_0$ (start of activity A) and $\bar{\gamma}_2$ (start of activity C). Recall from Example 6 that $\bar{\gamma}_0 = \mathbf{0}$ and $\bar{\gamma}_2 = [6, 6, 0]^\top$. The latency is now computed as

$$\lambda(\rho, 0, 2, r_1) = [\bar{\gamma}_2]_{r_1} - [\bar{\gamma}_0]_{r_1} = \begin{bmatrix} 6 \\ 6 \\ 0 \end{bmatrix}_{r_1} - \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{r_1} = 6.$$

The above focuses on resources. The temporal distance between the start of two actions can also be determined from the resource availability times in the state space, by slightly adapting the corresponding activities that contain the actions. Assume action instances a and f in activities A_{src} and A_{snk} . To determine the start-to-start latency between a and f , we slightly transform activities A_{src} and A_{snk} . We illustrate the approach for A_{src} , the transformation for A_{snk} is analogous. First, we add a new action s_a to A_{src} , using the same resource as a , that does not take time. We remove the incoming dependencies from a and add them to s_a . Then we add a dependency from s_a to a . Next, we add a new resource r_l and a dependency from the claim node of resource r_l to s_a , and from s_a to the release node of r_l . In this way, we encode the start time of action a in terms of the resource availability of resource r_l . The transformation is illustrated graphically in Fig. 5.

We assume that the occurrences of A_{src} and A_{snk} activities are related. In any run, for any $k > 0$, the k -th occurrence of A_{src} is paired with the k -th occurrence of A_{snk} . We refer to such a pair of related activities as a source-sink pair. Let $\text{getOccurrence}(\rho, A, k)$ be a function that returns the index of the k -th occurrence of activity A in run ρ . The start-to-start latency for resource r in ρ with source-sink pair $A_{i+1} = A_{src}$ and $A_j = A_{snk}$ in

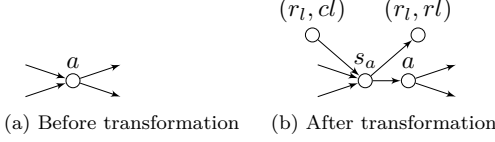


Fig. 5. Adaptation of an activity to measure the start time of action a in terms of the availability time of resource r_l .

run fragment $\rho[i, j + 1]$, is now equal to $\lambda(\rho, i, j, r)$. The maximum start-to-start latency in a run is now obtained by looking at all source-sink pairs:

$$\lambda_{max}(\rho, A_{src}, A_{snk}, r) = \sup_{k > 0} \lambda_k(\rho) \text{ where}$$

$$\lambda_k(\rho) = \lambda(\rho, i, j, r),$$

$$i = \text{getOccurrence}(\rho, A_{src}, k), \text{ and}$$

$$j = \text{getOccurrence}(\rho, A_{snk}, k).$$

Definition 10 (Latency). *Given normalized (max, +) state space \mathcal{S} , the maximum start-to-start latency of resource r with source-sink pair A_{src}, A_{snk} in \mathcal{S} is found by taking the maximum latency over all runs in the state space:*

$$\lambda_{max}(\mathcal{S}) = \sup_{\rho \in \mathcal{R}(\mathcal{S})} \lambda_{max}(\rho, A_{src}, A_{snk}, r).$$

Ratio Independence: In the state space, there can be redundancy with respect to multiple runs that have the same ratio value. Part of this redundancy is caused by the interleaving of activities that have no mutual influence. We reduce the size of the state space by removing redundant interleaving of ratio-independent activities.

Definition 11 (Ratio independent). *Let $\mathcal{S} = \langle C, \hat{c}, Act, \Delta, M, w_1, w_2 \rangle$ be a (max, +) state space, $c \in C$ be a configuration, and $A, B \in \text{enabled}(c)$ be activities enabled in c . Activities A and B are ratio independent in c iff they satisfy the following conditions:*

- 1) if $A, B \in \text{enabled}(c)$, then $B \in \text{enabled}(A(c))$, $A \in \text{enabled}(B(c))$, and $AB(c) = BA(c)$;
- 2) $w_i(c, A, A(c)) + w_i(A(c), B, AB(c)) = w_i(c, B, B(c)) + w_i(B(c), A, BA(c))$ for $i \in \{1, 2\}$;
- 3) $R(A) \cap R(B) = \emptyset$.

Two activities are ratio dependent if they are not ratio independent.

The first property is the classical notion of *independence*: in every configuration where A and B are both enabled, the execution of one activity cannot disable the other activity, and the resulting configuration after executing both activities in any order is the same. The second property requires that the sum of the weights w_1 and w_2 of the corresponding transitions of A and B is the same. The third property requires that activities A and B do not share resources.

Reduced state space: In the remainder of this section, we formalize an ample reduction on a (max, +) state space that preserves throughput and latency.

Definition 12 (State space reduction function). *A reduction function reduce for a (max, +) state space $\mathcal{S} = \langle C, \hat{c}, Act, \Delta, M, w_1, w_2 \rangle$ is a mapping from C to 2^{Act} such that $\text{reduce}(c) \subseteq \text{enabled}(c)$ for each configuration $c \in C$. We define the reduction of \mathcal{S} induced by reduce as the smallest (max, +) state space $\mathcal{S}' = \langle C', \hat{c}', Act', \Delta', M', w'_1, w'_2 \rangle$ that satisfies the following conditions:*

- $C' \subseteq C$, $\hat{c}' = \hat{c}$, $Act' = Act$, $\Delta' \subseteq \Delta$, $M' = M$;
- for every $c \in C'$ and $A \in \text{reduce}(c)$, $(c, A, A(c)) \in \Delta'$, $w'_1(c, A, A(c)) = w_1(c, A, A(c))$, and $w'_2(c, A, A(c)) = w_2(c, A, A(c))$.

Definition 13 (Ample conditions state space). *Let ample be a reduction function on a (max, +) state space that satisfies the following conditions:*

- (R1) **Non-emptiness condition:** if $\text{enabled}(c) \neq \emptyset$, then $\text{ample}(c) \neq \emptyset$.
- (R2) **Ratio-dependency condition:** For any configuration $c_0 \in C'$ and run $c_0 A_1 c_1 A_2 \dots A_m c_m$ with $m \geq 1$ in \mathcal{S} , if activity A_m and some activity in $\text{ample}(c_0)$ are ratio dependent in c_0 , then there is an index i with $1 \leq i \leq m$ with $A_i \in \text{ample}(c_0)$.

Example 14. *Consider initial configuration $c_0 = \langle s1, \emptyset \rangle$ in the state space shown in Fig. 4. Activities A and B are ratio dependent in c_0 (they do not satisfy condition 1 in Def. 11) and ratio independent with activity C . Ample sets $\{A, B\}$ and $\{C\}$ both satisfy conditions (R1) and (R2).*

Conditions (R1) and (R2) ensure that for each run in the (max, +) state space we can find an equivalent run in the reduced (max, +) state space. First we define equivalence of run prefixes. Two run prefixes are equivalent iff their corresponding activity sequences can be obtained from each other by repeatedly commuting adjacent ratio-independent activities.

Definition 15. *Activity sequences $v, w \in Act^*$ are considered equivalent [26], denoted $v \equiv w$, iff there exists a list of activity sequences u_0, u_1, \dots, u_n , where $u_0 = v$, $u_n = w$, and for each $0 \leq i < n$, $u_i = \bar{u} A \hat{u}$ and $u_{i+1} = \bar{u} B A \hat{u}$ for some $\bar{u}, \hat{u} \in Act^*$ and ratio-independent activities $A, B \in Act$.*

Given prefix $\rho[.m] = c_0 A_1 \dots A_m c_m$ of some run ρ , let ρ^m denote the activity sequence $A_1 \dots A_m$.

Definition 16. *Prefixes $\rho[.m]$ and $\sigma[.m]$ of runs ρ and σ are equivalent, denoted $\rho[.m] \equiv \sigma[.m]$, iff $\rho^m \equiv \sigma^m$.*

Throughput is defined as a limit on prefix ratios of infinite runs. To define equivalence of runs in terms of throughput, we need to consider run prefixes with a bounded difference in activities following those prefixes.

Definition 17 (Equivalence of runs). *Let ρ and σ be two runs. We define $\rho \succeq \sigma$ iff there exists a $c \in \mathbb{N}$ such that for all $n \geq 0$ it holds that $\rho \succeq_n \sigma$, where $\rho \succeq_n \sigma$ is defined iff there exists some $k \geq n$, run prefixes $\rho[.k]$ and $\hat{\rho}[.k]$ with $\hat{\rho}^k \equiv \rho^k$ such that $\hat{\rho}^k = \sigma^n \cdot \tau$ for some τ , and $k - n \leq c$. Runs ρ and σ are equivalent, denoted $\rho \equiv \sigma$, iff $\rho \succeq \sigma$ and $\sigma \succeq \rho$.*

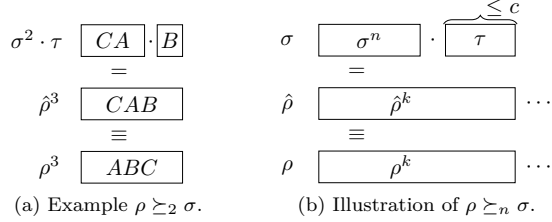


Fig. 6. Illustration of Def. 17.

Example 18. Consider run ρ with activity sequence $(ABCD)^\omega$ in the full state space. We want to construct an equivalent run σ in the reduced state space, i.e., a run that satisfies $\rho \succeq_n \sigma$ for all $n \geq 0$. Consider the case for $n = 2$, shown in Fig. 6a. Then, we need to match two (n) activities of ρ to the first two activities of σ in the reduced space. The prefix consisting of the first two activities of ρ , AB , is not a valid trace after reduction. In the reduced space the independent activity C must be performed first. Run $\hat{\rho}$, equivalent to ρ and identical to σ for $n = 2$ activities can be constructed by moving C to the front. ρ and σ must be such that this can be done for any $n \geq 0$. The general case is shown in Fig. 6b. Moreover, it is crucial for the preservation of throughput that the length k that one needs to consider in ρ to find the first n activities of σ exceeds n by maximally a finite amount c , independent of n .

If a reduction satisfies conditions (R1) and (R2), then for each run in the full ($\max, +$) state space we can find an equivalent run in the reduced ($\max, +$) state space.

Theorem 19 (Equivalent runs). Let $\mathcal{S} = \langle C, \hat{c}, Act, \Delta, M, w_1, w_2 \rangle$ be a finite normalized ($\max, +$) state space, and \mathcal{S}' be the reduced ($\max, +$) state space induced by reduction function *ample*. If *ample* satisfies conditions (R1) and (R2), then for each run $\rho \in \mathcal{R}(\mathcal{S})$, there exists a run $\sigma \in \mathcal{R}(\mathcal{S}')$ with $\rho \equiv \sigma$.

Theorem 20 (Equivalent runs have the same throughput). Let ρ and σ be runs. If $\rho \equiv \sigma$, then $\text{Ratio}(\rho) = \text{Ratio}(\sigma)$.

Theorem 21 (Equivalent runs have the same latency). Let $\rho, \sigma \in \mathcal{R}(\mathcal{S})$. Let A_{src} and A_{snk} be any source-sink pair, and let r be the resource for which we want to calculate the start-to-start latency. If $\rho \equiv \sigma$, then $\lambda_{max}(\rho, A_{src}, A_{snk}, r) = \lambda_{max}(\sigma, A_{src}, A_{snk}, r)$.

From Theorems 19-21, it immediately follows that a safe reduction preserves throughput and latency aspects.

Corollary 22. Let \mathcal{S} be a finite normalized ($\max, +$) state space, and \mathcal{S}' the reduced ($\max, +$) state space induced by reduction function *ample*. If *ample* satisfies conditions (R1) and (R2), then $\tau_{min}(\mathcal{S}) = \tau_{min}(\mathcal{S}')$ and $\lambda_{max}(\mathcal{S}) = \lambda_{max}(\mathcal{S}')$.

Reduced ($\max, +$) automata composition: We want to perform the partial-order reduction at the level of ($\max, +$) automata, rather than at the level of the ($\max, +$) state space. To achieve this goal, we introduce the notion of resource independence at the level of a ($\max, +$) automaton.

This notion is used in the ample conditions on reductions of a ($\max, +$) automaton.

Definition 23 (Resource-independent activities). Given ($\max, +$) automaton $\mathcal{A} = \langle S, \hat{s}, Act, reward, M, T \rangle$ and state $s \in S$, activities $A, B \in \text{enabled}(s)$ are resource independent in s if they satisfy the following conditions:

- $B \in \text{enabled}(A(s))$ and $A \in \text{enabled}(B(s))$;
- $AB(s) = BA(s)$;
- $R(A) \cap R(B) = \emptyset$.

Two activities are resource dependent, if they are not resource independent.

If two activities are resource independent in some state in the ($\max, +$) automaton, then they are also *ratio* independent in the corresponding configurations in the underlying state space. When two activities are resource independent, it holds that $R(A) \cap R(B) = \emptyset$. In this case their corresponding ($\max, +$) matrices commute, which means that $M_A \otimes M_B = M_B \otimes M_A$. As a result, the resulting normalized vector after multiplication is the same, and the sum of the weights w_1 and w_2 is the same, independent of the execution order.

Theorem 24. Given are a ($\max, +$) automaton $\mathcal{A} = \langle S, \hat{s}, Act, reward, M, T \rangle$ and state $s \in S$ with activities $A, B \in \text{enabled}(s)$. Consider any configuration $c = \langle s, \gamma \rangle$ in the underlying normalized ($\max, +$) state space. If A and B are resource-(in)dependent in s , then they are ratio-(in)dependent in c .

A reduction function on a ($\max, +$) automaton is defined in the following way.

Definition 25 (($\max, +$) automaton reduction function). A reduction function *reduce* for a ($\max, +$) automaton $\mathcal{A} = \langle S, \hat{s}, Act, reward, M, T \rangle$ is a mapping from S to 2^{Act} such that $\text{reduce}(s) \subseteq \text{enabled}(s)$ for each state $s \in S$. We define the reduction of \mathcal{A} induced by *reduce* as the smallest ($\max, +$) automaton $\mathcal{A}' = \langle S', \hat{s}', Act', reward', M', T' \rangle$ that satisfies the following conditions:

- $S' \subseteq S$, $\hat{s}' = \hat{s}$, $Act' = Act$, $T' \subseteq T$;
- for every $s \in S'$ and $A \in \text{reduce}(s)$, $\langle s, A, A(s) \rangle \in T'$, $reward'(A) = reward(A)$, and $M'(A) = M(A)$.

To preserve latency and throughput properties, we impose the following ample conditions on a reduction function of a ($\max, +$) automaton.

Definition 26 (Ample conditions ($\max, +$) automaton). Let *ample* be a reduction function on a ($\max, +$) automaton that satisfies the following conditions:

- (A1) **Non-emptiness condition:** if $\text{enabled}(s) \neq \emptyset$, then $\text{ample}(s) \neq \emptyset$.
- (A2) **Dependency condition:** For any state $s_0 \in S'$ and run $s_0 A_1 s_1 A_2 \dots A_m s_m$ with $m \geq 1$ in \mathcal{A} , if activity A_m and some activity in $\text{ample}(s_0)$ are resource dependent in s_0 , then there is an index i with $1 \leq i \leq m$ with $A_i \in \text{ample}(s_0)$.

Example 27. Consider state s_1 in the ($\max, +$) automaton shown in Fig. 1c. Activities A and B are resource dependent

in s_1 and resource independent with activity C . Ample sets $\{A, B\}$ and $\{C\}$ both satisfy conditions (A1) and (A2).

An ample reduction on a $(\max, +)$ automaton preserves the minimum throughput and maximum latency values.

Theorem 28. *Let \mathcal{A} be any $(\max, +)$ automaton with corresponding state space \mathcal{S} . Let ample be an ample reduction on \mathcal{A} satisfying conditions (A1) and (A2) that yields \mathcal{A}' with corresponding state space \mathcal{S}' . Then $\tau_{\min}(\mathcal{S}) = \tau_{\min}(\mathcal{S}')$ and $\lambda_{\max}(\mathcal{S}) = \lambda_{\max}(\mathcal{S}')$.*

IV. ON-THE-FLY REDUCTION

The previous section gives sufficient conditions for a reduction function on a $(\max, +)$ automaton to preserve performance properties in the corresponding $(\max, +)$ state space. For an efficient reduction, we avoid first computing the full composition of the $(\max, +)$ automata. Rather, we use sufficient local conditions on the network of $(\max, +)$ automata to compute a reduced composition on-the-fly.

The traditional on-the-fly method of Peled [19] selects the enabled activities $enabled_i(s) = enabled(s) \cap Act_i$ of some automaton \mathcal{A}_i as ample set, while exploring a state $s = \langle s_1, s_2, \dots, s_n \rangle$. The ample condition requires that all locally enabled activities in $enabled(s_i)$ of this automaton \mathcal{A}_i are resource-independent with all activities in the other automata. Otherwise, all enabled activities in s are selected. In an experimental evaluation, we found that this approach did not yield any reduction on the models described in Section V. Therefore, we consider a generalization of the approach to clusters, as in [20]. The cluster-inspired ample reduction selects a safe cluster in each state in the composition. Given max-plus timed system $\mathcal{M} = \mathcal{A}_1 \parallel \dots \parallel \mathcal{A}_n$ and $\mathcal{A} = \{\mathcal{A}_i \mid 1 \leq i \leq n\}$, we define a *cluster* $\mathcal{C} \subseteq \mathcal{A}$. Let $Act(\mathcal{C}) = \bigcup_{\mathcal{A}_i \in \mathcal{C}} Act_i$ denote the set of activities that occur in \mathcal{C} . The set of enabled activities that is selected in state s induced by a cluster \mathcal{C} is $enabled_{\mathcal{C}}(s) = enabled(s) \cap Act(\mathcal{C})$. We define a projection to consider the local state $\pi_{\mathcal{C}}(s)$ in a cluster $\mathcal{C} \subseteq \mathcal{A}$:

$$\begin{aligned} \pi_{\mathcal{A}_i}(s) &= s_i \\ \pi_{\mathcal{C}}(s) &= \langle \pi_{\mathcal{A}_{c_1}}(s), \dots, \pi_{\mathcal{A}_{c_k}}(s) \rangle \text{ where } \mathcal{C} = \{\mathcal{A}_{c_1}, \dots, \mathcal{A}_{c_k}\}, \\ &\text{and for all } 1 \leq j < k, c_j < c_{j+1}. \end{aligned}$$

Given local state $\pi_{\mathcal{C}}(s)$, $enabled(\pi_{\mathcal{C}}(s))$ denotes the set of activities that are enabled in the composition of precisely the automata in \mathcal{C} . Note that $enabled_{\mathcal{C}}(s) \subseteq enabled(\pi_{\mathcal{C}}(s))$, since the latter might contain activities that are enabled in the local state of the cluster-composition, but disabled in the global composition due to an automaton outside the cluster that disables the activity. We only consider independence of activities among automata, and not within the same automaton. The former can be checked locally, whereas the latter requires an exploration on the internal transition structure. We treat activities inside the same automaton as resource dependent.

Definition 29 (Cluster safety). *Let $\mathcal{C} \subseteq \mathcal{A}$ be any cluster, and s be a state in the composition. Cluster \mathcal{C} is safe in s if the following conditions are satisfied.*

- (C1) if $enabled(s) \neq \emptyset$, then $enabled_{\mathcal{C}}(s) \neq \emptyset$;
- (C2.1) for any $A \in enabled_{\mathcal{C}}(s)$ and $B \in Act(\mathcal{A}) \setminus Act(\mathcal{C})$, $R(A) \cap R(B) = \emptyset$;
- (C2.2) for any $A \in enabled(\pi_{\mathcal{C}}(s))$, if $A \in Act_i$ then $\mathcal{A}_i \in \mathcal{C}$.

Condition (C2.1) requires that each enabled activity in $enabled_{\mathcal{C}}(s)$ does not share resources with any activity outside $Act(\mathcal{C})$. Condition (C2.2) requires that each activity in $enabled(\pi_{\mathcal{C}}(s))$ does not occur outside of the cluster. Together, these two conditions ensure that no activity $A \in Act(\mathcal{A}) \setminus enabled_{\mathcal{C}}(s)$, dependent on some activity in $enabled_{\mathcal{C}}(s)$, becomes enabled by executing only activities outside the cluster. We define a cluster-inspired ample reduction through a safety condition M on a max-plus timed system:

Definition 30 (Cluster-inspired ample reduction). *A cluster-inspired ample reduction ample for a max-plus timed system $\mathcal{M} = \mathcal{A}_1 \parallel \dots \parallel \mathcal{A}_n$ is a mapping from $S = S_1 \times \dots \times S_n$ to 2^{Act} such that $ample(s) = enabled_{\mathcal{C}}(s)$ for some cluster $\mathcal{C} \subseteq \mathcal{A}$, and satisfies the following condition:*

- (M) **Cluster-safety condition:** for any state s , $ample(s) = enabled_{\mathcal{C}}(s)$ where \mathcal{C} is safe in s .

The reduction of \mathcal{M} is defined using Def. 25.

Theorem 31. *Let ample be a cluster-inspired ample reduction on $\mathcal{M} = \mathcal{A}_1 \parallel \dots \parallel \mathcal{A}_n$. Then ample satisfies conditions (A1) and (A2).*

In each state in the composition, there possibly exist multiple safe clusters that can be chosen. Heuristics can be used to select a cluster that likely leads to a large reduction. In our experiments, we use a heuristic that chooses a smallest safe cluster in each state. This heuristic often performs well [27], because it allows to prune most enabled transitions. Starting from each enabled activity in s as a candidate, we construct a safe cluster. After trying all candidates, we select a cluster \mathcal{C} with the smallest $enabled_{\mathcal{C}}(s)$ set. This approach is best suited when the enabled sets are small. If there are many states with large enabled sets, then a heuristic could be used to compute only one safe cluster starting from activities that occur in only few automata.

Algorithm 1 shows the algorithm to compute a safe cluster in a state. The algorithm checks for each activity enabled in the current cluster \mathcal{C} whether condition (C2.1) or (C2.2) is violated. The algorithm starts with initial *candidate* activity A . If A is enabled in the composition, we add all automata that contain A and add an automaton for each dependent activity outside the current cluster. This ensures that condition (C2.1) is satisfied for activity A and the cluster obtained after executing lines 5-9. If A is enabled in the composition of automata in the cluster, but not in the full composition, then we add an automaton that causes A to be disabled in the full composition. This ensures that condition (C2.2) is satisfied for A for the cluster obtained after executing lines 10-13. We use the notation $[\mathcal{A}_i \mid B \in Act_i]$ to denote a list comprehension,

Algorithm 1. Algorithm to compute a safe cluster.

```

1: proc COMPUTECLUSTER( $s, candidate$ )
2:    $A \leftarrow candidate$ ;  $C \leftarrow \emptyset$ ;  $processed \leftarrow \emptyset$ 
3:   while  $A \neq \perp$  do
4:      $processed \leftarrow processed \cup \{A\}$ 
5:     if  $A \in enabled(s)$  then
6:        $C \leftarrow C \cup \{A_i \mid A \in Act_i\}$ 
7:       for  $B \in \{D \in Act \mid R(D) \cap R(A) \neq \emptyset\}$  do
8:         if  $B \notin Act(C)$  then
9:            $C \leftarrow C \cup [A_i \mid B \in Act_i].first()$ 
10:      if  $A \notin enabled(s) \wedge A \in enabled(\pi_C(s))$  then
11:        for  $A_i \in A$  do
12:          if  $A \in Act_i \wedge A_i \notin C \wedge A \notin enabled(s_i)$  then
13:             $C \leftarrow C \cup \{A_i\}$ ; break
14:      if  $|enabled(\pi_C(s)) \setminus processed| > 0$  then
15:         $A \leftarrow [enabled(\pi_C(s)) \setminus processed].first()$ 
16:      else
17:         $A \leftarrow \perp$ 
18:      return  $C$ 

```

and function $first()$ picks the first element from a list. After handling the activity, we check whether there are other activities that are locally enabled in the new cluster and not yet processed (line 14-17). We continue until all locally enabled activities are processed. The algorithm is repeatedly called for each activity in $enabled(s)$ as a candidate, and afterwards the cluster with the smallest set of enabled activities is chosen.

Theorem 32. *Let s be a state in the composition, and A be the candidate activity. $COMPUTECLUSTER(s, A)$ returns a cluster that is safe in s .*

Example 33. *Consider the initial state s_1 in the composition $\mathcal{A}_1 \parallel \mathcal{A}_2 \parallel \mathcal{A}_3$ shown in Fig. 1b. We show how Algorithm 1 computes a valid ample set for s_1 . Assume that A is the candidate activity. Since A is enabled, we add all automata with A (line 6) to cluster C , i.e. \mathcal{A}_1 . Since activity B is resource-dependent with A in s_1 and B is not yet in the cluster alphabet, we add the first automaton with B , i.e. \mathcal{A}_2 , to C (line 9). There are no other activities enabled within the local state $\pi_C(s_1)$ of the cluster, so we skip lines 11-13. Since we processed all activities and $A = \perp$, cluster C is returned. This yields $ample(s_1) = enabled_C(s_1) = enabled(s_1) \cap Act(C) = \{A, B, C\} \cap \{A, B, D\} = \{A, B\}$.*

Using a similar reasoning, one can validate that $C_2 = \{\mathcal{A}_3\}$ (with $enabled_{C_2}(s_1) = \{C\}$) is also a safe cluster in s_1 . A reduced composition obtained with an ample reduction is shown in Fig. 1c. Since activity C is resource-independent with activities A and B , only one interleaving of C with A and B is explored. Both interleavings of A and B are still present, since activities A and B are resource-dependent. Fig. 4 shows the $(max, +)$ state space after reduction.

V. EXPERIMENTAL EVALUATION

In the previous section, we described an on-the-fly reduction on the level of $(max, +)$ automata to compute a reduced composition that preserves throughput and latency values. To test the effectiveness of the partial-order

| | automata composition | | | | normalized $(max, +)$ state space | | | | | |
|-----|----------------------|-------|---------|-------|-----------------------------------|------------|---------|------------|------|-------|
| | full | | reduced | | full | | reduced | | | |
| | $ S $ | $ T $ | $ S $ | $ T $ | $ C $ | $ \Delta $ | $ C $ | $ \Delta $ | | |
| TW1 | 1280 | 2170 | 1267 | 2144 | 3585 | 6120 | 3568 | 0.5% | 6086 | 0.6% |
| TW2 | 63 | 128 | 34 | 44 | 379 | 756 | 149 | 60.7% | 207 | 72.6% |
| TW3 | 343 | 759 | 151 | 209 | 4949 | 10655 | 1507 | 69.5% | 2261 | 78.8% |
| TW4 | 318 | 768 | 139 | 202 | 18792 | 40957 | 5280 | 71.9% | 7785 | 81.0% |

TABLE I
REDUCTIONS ACHIEVED ON THE TWILIGHT SYSTEM MODELS.

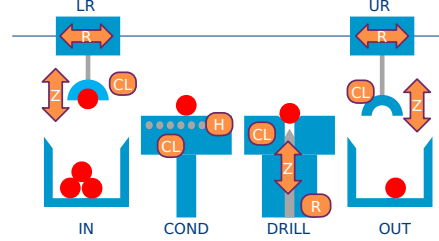


Fig. 7. Twilight manufacturing system [15].

reduction on models of manufacturing systems, we created a number of variants of the Twilight system [15] model.

The Twilight system, shown in Fig. 7, processes balls that need to be drilled. First, a ball is picked up by the load robot (LR) from the input buffer (IN). Then, the ball is put on the conditioner (COND), to heat the ball to a desired temperature. The conditioner has a heater action (H) to heat a ball. Once heated, the ball is transported to the (DRILL) by the load robot or unload robot (UR) to drill a hole in the ball. The drill has an R-motor (R) to rotate the drill bit, and a Z-motor (Z) to move the drill bit up and down. After drilling, the unload robot puts the finished ball in the output buffer (OUT). Both robots have a clamp (CL) to pick up and hold a ball, an R-motor (R) to move along the rail, and a Z-motor (Z) to move the clamp up and down. Both processing stations have a clamp (CL) to ensure that the ball stays positioned correctly. The system specification contains activities that describe the conditioning and drilling step, as well as activities that describe how a robot transports a product. To ensure safe movements, we model collision areas above the conditioner and the drill as virtual resources. The automata describe constraints on valid activity orderings. The throughput and latency values are analyzed on the $(max, +)$ state space of the composition of these automata.

We examine four variants of the system. The first variant (TW1) is the model described above, of which the full model is given in [15]. The life cycle and location of each product is explicitly modeled. In TW2, we remove these product-location and life-cycle automata, and instead use a set of automata that ensure that products are always moved forward in the production process. In TW3, we extend TW2 with a polish station, where each product undergoes a polish and drill step after the condition step but not in a fixed order. In TW4 we fix the order so that a product is always first conditioned, then drilled, and then polished.

Table I shows that a reduction is achieved in all models. The reduction for TW1 is very small (0.5%) compared to the reductions of the other models (60.7%, 69.5%, and 71.9% for TW2, TW3, and TW4). The reduction for model TW1 is very small, because there is a lot of event synchronization by the product-location and life-cycle automata. Recall condition (C2.2), that requires that each enabled activity in the local state of a safe cluster must be independent with activities outside the cluster. During state-space exploration of model TW1, the algorithm often needs to add product-location or life-cycle automata to the cluster to satisfy this condition (C2.2); this limits reduction possibilities. The reduction for TW2, TW3, and TW4 is much larger, since we do not explicitly model the product-location and life-cycle automata.

VI. CONCLUSION

In this paper, we presented a new partial-order reduction technique to speed up performance analysis of max-plus timed systems. The technique is inspired by an existing cluster-based ample set reduction for non-timed systems. It tries to compute a smaller state space by exploiting the structure of the concurrent (max,+) automata, and information about resource sharing. We derived conditions to compute a reduced composition, from which a reduced state space can be calculated. In this reduced state space, performance properties (throughput and latency) of the system model are preserved. The experimental evaluation shows that the partial-order reduction technique can be successfully used to reduce the size of the state space for a set of example models of manufacturing systems. The reductions that can be achieved are highly dependent on the structure of the input model, the amount of synchronization on activities among automata, and the extent to which activities use the same resources. In our models, the partial-order reduction technique successfully exploits redundant interleaving related to processing stations that can perform operations on products in parallel, and movements of the robots that can be executed simultaneously.

ACKNOWLEDGMENT

This research is supported by the Dutch NWO-TTW, carried out as part of the Robust Cyber-Physical Systems (RCPS) program, project number 12694.

REFERENCES

- [1] M. Hendriks, J. Verriet, T. Basten *et al.*, “Performance engineering for industrial embedded data-processing systems,” in *Proc. of Int. Conf. on Product-Focused Software Process Improvement (PROFES)*, 2015, pp. 399–414.
- [2] Z. Manna and A. Pnueli, *The Temporal Logic of Reactive and Concurrent Systems*. New York, NY, USA: Springer-Verlag New York, Inc., 1992.
- [3] E. M. Clarke, E. A. Emerson, and A. P. Sistla, “Automatic verification of finite-state concurrent systems using temporal logic specifications,” *ACM Trans. on Programming Languages and Systems (TOPLAS)*, vol. 8, no. 2, pp. 244–263, 1986.
- [4] J. Bengtsson, B. Jonsson, J. Lilius, and W. Yi, “Partial order reductions for timed systems,” in *Proc. of Int. Conf. on Concurrency Theory (CONCUR)*, D. Sangiorgi and R. de Simone, Eds. Berlin, Heidelberg: Springer, 1998, pp. 485–500.
- [5] M. Minea, “Partial order reduction for model checking of timed automata,” in *Proc. of Int. Conf. on Concurrency Theory (CONCUR)*. London, UK: Springer, 1999, pp. 431–446.
- [6] T. Yoneda and B.-H. Schlingloff, “Efficient verification of parallel real-time systems,” *Formal Methods in System Design*, vol. 11, no. 2, pp. 187–215, 1997.
- [7] B. Theelen, M. Geilen, and J. Voeten, “Performance model checking scenario-aware dataflow,” in *Int. Conf. on Formal Modeling and Analysis of Timed Systems (FORMATS)*, U. Fahrenberg and S. Tripakis, Eds. Berlin, Heidelberg: Springer, 2011, pp. 43–59.
- [8] S. Gaubert, “Performance evaluation of (max,+) automata,” *IEEE Trans. on Automatic Control*, vol. 40, no. 12, Dec 1995.
- [9] J.-Y. Le Boudec and P. Thiran, *Network Calculus: A Theory of Deterministic Queuing Systems for the Internet*. Berlin, Heidelberg: Springer-Verlag, 2001.
- [10] L. Thiele, S. Chakraborty, and M. Naedele, “Real-time calculus for scheduling hard real-time systems,” in *2000 IEEE International Symposium on Circuits and Systems. Emerging Technologies for the 21st Century. Proceedings (IEEE Cat No.00CH36353)*, vol. 4, 2000, pp. 101–104 vol.4.
- [11] M. Geilen, “Synchronous dataflow scenarios,” *ACM Trans. Embed. Comput. Syst.*, vol. 10, no. 2, pp. 16:1–16:31, Jan. 2011.
- [12] M. Geilen and S. Stuijk, “Worst-case performance analysis of Synchronous Dataflow scenarios,” *Int. Conf. on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, pp. 125–134, 2010.
- [13] J. Campos, G. Chiola, J. M. Colom, and M. Silva, “Properties and performance bounds for timed marked graphs,” *IEEE Trans. on Circuits and Systems I: Fundamental Theory and Applications*, vol. 39, no. 5, pp. 386–401, May 1992.
- [14] B. Heidergott, G. J. Olsder, and J. Van Der Woude, *Max Plus at Work*. Princeton University Press, 2014.
- [15] B. van der Sanden, J. Bastos, J. Voeten *et al.*, “Compositional specification of functionality and timing of manufacturing systems,” in *Forum on specification & Design Languages (FDL) 2016*, 2016, pp. 1–8.
- [16] B. van der Sanden, M. Reniers, M. Geilen *et al.*, “Modular model-based supervisory controller design for wafer logistics in lithography machines,” in *Int. Conf. on Model Driven Engineering Languages and Systems (MODELS)*, 2015, pp. 416–425.
- [17] A. Valmari, “A stubborn attack on state explosion,” *Formal Methods in System Design*, vol. 1, no. 4, pp. 297–322, 1992.
- [18] P. Godefroid, *Using partial orders to improve automatic verification methods*. Berlin, Heidelberg: Springer, 1991, pp. 176–185.
- [19] D. Peled, “Combining partial order reductions with on-the-fly model-checking,” *Formal Methods in System Design*, vol. 8, no. 1, pp. 39–64, 1996.
- [20] T. Basten, D. Bošnački, and M. Geilen, “Cluster-based partial-order reduction,” *Automated Software Engineering*, vol. 11, no. 4, pp. 365–402, 2004.
- [21] B. van der Sanden, M. Geilen, M. Reniers, and T. Basten, Eindhoven University of Technology, Department of Electrical Engineering, Electronic Systems, Tech. Rep. ESR-2018-01, 2018.
- [22] F. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat, *Synchronization and linearity: an algebra for discrete event systems*. John Wiley and Sons, 1992.
- [23] W. Thomas, “Automata on infinite objects,” *Handbook of theoretical computer science, Volume B*, pp. 133–191, 1990.
- [24] M. Geilen, J. Falk, C. Haubelt *et al.*, “Performance analysis of weakly-consistent scenario-aware dataflow graphs,” *Journal of Signal Processing Systems*, vol. 87, no. 1, pp. 157–175, Apr 2017.
- [25] A. Dasdan, “Experimental analysis of the fastest optimum cycle ratio and mean algorithms,” *ACM-TODAES*, vol. 9, no. 4, pp. 385–418, 2004.
- [26] A. W. Mazurkiewicz, “Trace theory,” in *Petri Nets: Central Models and Their Properties, Advances in Petri Nets 1986, Part II, Proceedings of an Advanced Course, Bad Honnef, 8.-19. September 1986*, 1986, pp. 279–324.
- [27] J. Geldenhuys, H. Hansen, and A. Valmari, *Exploring the Scope for Partial Order Reduction*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 39–53.