

Singapore Management University

## Institutional Knowledge at Singapore Management University

---

Research Collection School Of Computing and  
Information Systems

School of Computing and Information Systems

---

3-2020

### Automated synthesis of local time requirement for service composition

Étienne ANDRÉ

Tian Huat TAN

Manman CHEN

Shuang LIU

Jun SUN

Singapore Management University, junsun@smu.edu.sg

*See next page for additional authors*

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)



Part of the [Software Engineering Commons](#), and the [Systems Architecture Commons](#)

---

#### Citation

ANDRÉ, Étienne; TAN, Tian Huat; CHEN, Manman; LIU, Shuang; SUN, Jun; LIU, Yang; and DONG, Jin Song. Automated synthesis of local time requirement for service composition. (2020). *Software and Systems Modeling*. 1-31.

Available at: [https://ink.library.smu.edu.sg/sis\\_research/5112](https://ink.library.smu.edu.sg/sis_research/5112)

This Journal Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [cherylds@smu.edu.sg](mailto:cherylds@smu.edu.sg).

---

**Author**

Étienne ANDRÉ, Tian Huat TAN, Manman CHEN, Shuang LIU, Jun SUN, Yang LIU, and Jin Song DONG

## Automated synthesis of local time requirement for service composition

Étienne André · Tian Huat Tan · Manman  
Chen · Shuang Liu · Jun Sun · Yang Liu ·  
Jin Song Dong

**Abstract** Service composition aims at achieving a business goal by composing existing service-based applications or components. The response time of a service is crucial especially in time critical business environments, which is often stated as a clause in service level agreements between service providers and service users. To meet the guaranteed response time requirement of a composite service, it is important to select a feasible set of component services such that their response time will collectively satisfy the response time requirement of the composite service. In this work, we use the BPEL modeling language, that aims at specifying Web services. We extend it with timing parameters, and equip it with a formal semantics. Then, we propose a fully automated approach to synthesize the response time requirement of component services modeled using BPEL, in the form of a constraint

---

This is a pre-print of an article published in the [International Journal on Software and Systems Modeling \(SoSyM\)](#). The final authenticated version is available online at: [10.1007/s10270-020-00787-5](https://doi.org/10.1007/s10270-020-00787-5). Étienne André, Jin Song Dong and Yang Liu are partially supported by CNRS STIC-Asie project CATS (“Compositional Analysis of Timed Systems”). Étienne André is partially supported by the ANR national research program ANR-14-CE28-0002 PACS (“Parametric Analyses of Concurrent Systems”). Étienne André and Jun Sun are partially supported by the ANR-NRF French-Singaporean research program ProMiS (ANR-19-CE25-0015).

---

É. André  
Université de Lorraine, CNRS, Inria, LORIA, Nancy, France

T. H. Tan  
IBM, Singapore

M. Chen  
Autodesk, Singapore

S. Liu  
School of Software, Tianjin University, China

J. Sun  
Singapore Management University, Singapore

Y. Liu  
Nanyang Technological University, Singapore

J. S. Dong  
National University of Singapore  
Griffith University, Australia

on the local response times. The synthesized requirement will guarantee the satisfaction of the global response time requirement, statically or dynamically. We implemented our work into a tool, SELAMAT, and performed several experiments to evaluate the validity of our approach.

**Keywords** Web service composition, Parameter synthesis, Modeling Web services, Formal semantics, BPEL, Parametric model checking

## 1 Introduction and motivation

Service-oriented architecture is a paradigm where building blocks are used as services for software applications. Services encapsulate their functionalities, information, and make them available through a set of operations accessible over a network infrastructure using standards like SOAP [Gud+07] and WSDL [Chi+07]. To make use of a set of services to achieve a business goal, service composition languages such as BPEL (Business Process Execution Language) [Alv+07] have been proposed. A service that is composed by other services is called a *composite* service, and services that the composite service makes use of are called *component* services.

The requirement on the service response time is often an important clause in service-level agreements (SLAs) especially in business where timing is critical. An SLA is a contract between service consumers and service providers specifying the expected quality of service (QoS) level. Henceforth, we refer to the response time requirement of composite services as *global time requirement*, and to the set of constraints on the response times of the component services as *local time requirement*. The response time of a composite service is highly dependent on that of each component service. It is therefore crucial to derive local time requirements (*i. e.*, requirements for the component services) from the global time requirement, so that it will help in the selection of component services when building a composite service while satisfying the response time requirement.

An additional motivation for our work is that of micro-services. As pointed out by [Tan+16], many big players in the market (*e. g.*, Netflix, Amazon, and Microsoft Azure) have adopted microservice architecture [Ric18] by decomposing their existing monolithic applications into smaller, and highly decoupled services (also known as microservices). These services are then composed for fulfilling their business requirements. For example, Netflix decomposed their monolithic DVD rental application into services that work together, and that stream digital entertainment to millions of Netflix customers every day. Services of Netflix are hosted in a cloud provided by Amazon EC2 [Ama18], which offers about 40 instance types. The problem of composition of Web services with a large set of microservices is more and more relevant now, as the micro-services are getting more popular than ever (see *e. g.*, [Mid17; ST19]). This justifies the use of techniques for which different services can be compared to and eventually selected. Service-oriented architecture and micro-service architectures are conceptually similar: service-oriented architecture is a term that is used earlier and also widely used in literature. Micro-service architecture is more of a newer term that is used and practise widely in current industry, for the purpose of agile development. (For detailed comparison, see *e. g.*, [CDP17].) The methods developed here are applicable to both service-oriented architecture and micro-service architecture.

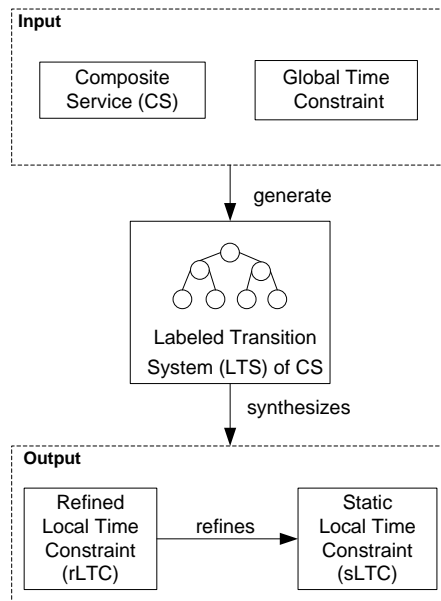


Fig. 1: General approach

Consider an example of a stock indices service, which has an SLA with the subscribed users requiring that the stock indices shall be returned within three seconds upon request. The stock indices service makes use of several component services, including a paid service, for requesting stock indices. The stock indices service provider would be interested in knowing the local time requirement of the component services, while satisfying the global response time requirement. To avoid discarding any service candidates that might be part of a feasible composition, the synthesized local time requirement needs to be as *weak* as possible, *i. e.*, to maintain as many combinations of local time requirements as possible. This is crucial as having a faster service might incur a higher cost.

### 1.1 Contribution

In this paper, we present a fully automated technique to perform a rigorous model-based analysis of Web services, in order to synthesize the local time requirement in composite services. Our approach performs an analysis of the composite service model behavior, using techniques inspired by parameter synthesis for timed systems. Our synthesis approach does not only avoid bad scenarios in the service composition, but also guarantees the fulfillment of the global time requirement.

We use as a formalism BPEL, which is a *de-facto* standard language for specifying service composition. BPEL supports control flow structures that involve complex timing constructs (*e. g.*, `<pick>` control structure) and concurrent execution of activities (*e. g.*, `<flow>` control structure). Due to the non-determinism in both time and control flow, it is unknown which execution path will be executed

at runtime. Such a combination of timing constructs, concurrent calls to external services, and complex control structures, makes it a challenge to synthesize the local time requirement. More precisely, response times of component services can be dependant; therefore, constraint between services may be of the form, *e.g.*,  $t_{fs} > t_{hs}$  (for two parametric component service times), rendering the problem quite delicate. In addition, there may be multiple possibilities of component services' response times that are satisfiable. This can be particularly delicate to find out with only manual human inspection.

Fig. 1 illustrates the main steps of our approach for synthesizing local time requirements. The required inputs are the specification of the composite service, and its global time requirement. The output will be local time requirements (at design time, and then at runtime) given in the form of a linear constraint.

We first propose a formal semantics for BPEL composite services augmented with timing parameters, *i.e.*, constants, the value of which is not known at design time; this symbolic semantics is given the form of a labeled transition system (LTS).

Based on the LTS resulting from the input composite service, we then propose an approach to synthesize local time requirements of component services, represented as a (linear) constraint, which we refer to as the *local time constraint*. During the design phase of a composite service, the local time constraint is synthesized based on *all* possible execution paths of the model, since it is unknown which execution path will be executed at runtime (this will depend on the dynamic behavior of the system). The local time constraint of a composite service that is synthesized during the design time is called the *static local time constraint* (hereafter sLTC).

The synthesized sLTC has several advantages. Firstly, when creating a new composite service, it allows the selection of feasible services from a large pool of services with similar functionalities but different local response times. Secondly, service designers can use the synthesized result to avoid over-approximations on the local response times, which may lead the service provider to purchase a service at a higher cost, while a service at a lower cost with a slower response time might have been sufficient to guarantee the global time requirement. Thirdly, the local time requirements serve as a safe guideline when component services need to be substituted or new services need to be introduced.

Due to the highly evolving and dynamic environment which the composite service is running in, the design time assumptions for Web service composition, even if they are initially accurate, may later change at runtime. For example, the execution time of a component service could violate the sLTC due to reasons such as network congestion. Nevertheless, this does not necessarily imply that the composite service will not satisfy the global time requirement. Indeed, the sLTC is synthesized based on all possible execution paths at design time, whereas only one path will be executed at runtime. At runtime, some of the execution paths can be eliminated. Therefore, we can use the runtime information to refine the sLTC to make it weaker—which results in a more relaxed constraint. We refer to the sLTC refined at runtime as the *refined local time constraint* (hereafter rLTC). The rLTC is used to decide whether the current composite service can still satisfy the global time requirement, despite some unplanned issues such as network congestion.

Our contributions are summarized as follows.

1. We augment the BPEL modeling language with timing parameters, and we equip it with a formal semantics in the form of a labeled transition system.
2. Given a composite service modeled using BPEL, we develop a sound method for synthesizing the local time requirement in the form of a set of constraints, which can be applied at the design stage of service composition.
3. We introduce a refinement procedure on the sLTC of a composite service based on the runtime information, which results in a more relaxed rLTC. The rLTC can be used to verify whether the composite service could still eventually satisfy the global time requirement at runtime.
4. We implement our algorithms into a tool SELAMAT. We then conduct experiments on several examples. The results show that the rLTC can indeed help to improve the accuracy of the sLTC. In addition, we show that the runtime adaptation does not incur much overhead in practice.

## 1.2 About this manuscript

This manuscript is an extended version of [Tan+13]. We in fact rewrote most of the manuscript for a better readability. The most notable differences between this manuscript and [Tan+13] are:

1. we replaced the formerly defined “AOLTS” with what we believe to be a simpler and more elegant presentation of labeled transition systems (LTS);
2. we added details on our implementation and used more service composition examples; and,
3. most importantly, we added a refinement procedure that attempts to meet the global time requirement at runtime even when the constraint computed statically is violated (Section 6).

## 1.3 Outline

The rest of this paper is structured as follows. Section 2 introduces a timed BPEL running example. Section 3 provides the necessary definitions and terminologies. Section 4 introduces our formal semantics for BPEL extended with timing parameters. Section 5 presents the synthesis algorithms for sLTC. Section 6 introduces rLTC, and its usage for runtime adaptation of a service composition. Section 7 evaluates our approach with four service composition examples. Section 8 reviews related works. Finally, Section 9 concludes the paper, and outlines future work.

## 2 A BPEL example with timed requirements

Let us introduce a *Stock Market Indices Service* (SMIS) that will be used as a running example. SMIS is a paid service and its goal is to provide updated stock indices to the subscribed users. It provides a service level agreement (SLA) to the subscribed users stating that it always responds within three seconds upon request.

SMIS has three component Web services, *i. e.*, a database service (*DS*), a free news feed service (*FS*) and a paid news feed service (*PS*). The strategy of the

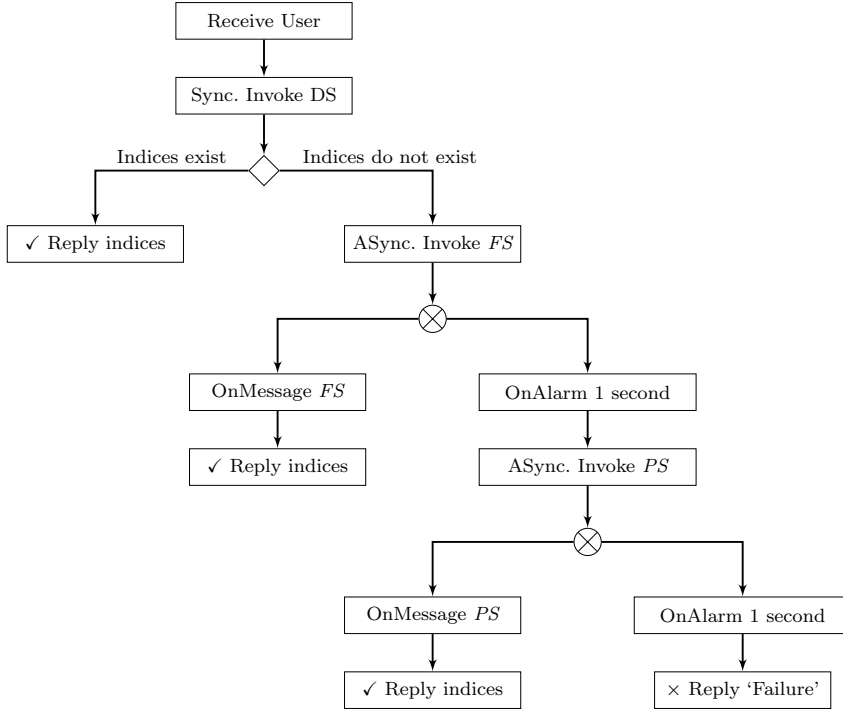


Fig. 2: Stock Market Indices Service

SMIS is calling the free service *FS* before calling the paid service *PS* in order to minimize the cost. Upon returning the result to the user, the SMIS also stores the latest results in an external database service provided by *DS* (storage of the results is omitted here). The workflow of the SMIS is sketched in Fig. 2 in the form of a tree. When a request is received from a subscribed customer (**Receive User**), it synchronously invokes (*i. e.*, invoke and wait for reply) the database service (**Sync. Invoke DS**) to request stock indices stored in the past minute. Upon receiving the response from *DS*, the process reaches an **<if>** branch (denoted by  $\diamond$ ). If the indices are available (**Indices exist**), then they are returned to the user (**Reply indices**). Otherwise, *FS* is invoked asynchronously (*i. e.*, the system moves on after the invocation without waiting for the reply). A **<pick>** construct (denoted by  $\otimes$ ) is used here to await an incoming response (**<onMessage>**) from previous asynchronous invocation or timeout (**<onAlarm>**), whichever occurs. If the response from *FS* (**OnMessage FS**) is received within one second, then the result is returned to the user (**Reply indices**). Otherwise, the timeout occurs (**OnAlarm 1 second**), then SMIS stops waiting for the result from *FS* and calls *PS* instead (**ASync. Invoke PS**). Similarly to *FS*, the result from *PS* is returned to the user, if the response from *PS* is received within one second. Otherwise, it notifies the user regarding the failure of getting stock indices (**Reply 'Failure'**). The states marked with a  $\checkmark$  (resp.  $\times$ ) represent desired (resp. undesired) end states.



The global time requirement for SMIS is that SMIS should respond within three seconds upon request. It is of particular interest to know the local time requirements for services  $PS$ ,  $FS$ , and  $DS$ , so as to fulfill the global time requirement. This information can also help to choose a paid service  $PS$  which is both cheap and responds quickly enough.

In this example, an activity to avoid (which will be referred to as a “bad activity” in the following) is the reply activity that is triggered after the component service  $PS$  fails to respond within one second, which is marked with  $\times$  in Fig. 2.

### 3 A formal model for parametric composite services

#### 3.1 Variables, clocks, parameters, and constraints

Given a finite set  $\mathcal{V}$  of finite-domain *variables*, a *variable valuation* for  $\mathcal{V}$  is a function assigning to each variable a value in its domain. We denote by  $\text{Valuations}(\mathcal{V})$  the set of all variable valuations of  $\mathcal{V}$ . Given a variable  $y \in \mathcal{V}$  and a variable valuation  $v \in \text{Valuations}(\mathcal{V})$ , we denote by  $v(y) = \perp$  the fact that variable  $y$  is uninitialized in valuation  $v$ .

The clocks, parameters and constraints that we use in this work are similar to the ones used in the formalisms of (parametric) timed automata [AD94; AHV93] and (parametric) stateful timed CSP [Sun+13; And+14]. Let  $X = \{x_1, \dots, x_h\}$  (for some integer  $h$ ) be a finite set of clocks, *i. e.*, real-valued variables evolving at the same rate. A *clock valuation* is a function  $w : X \rightarrow \mathbb{R}_{\geq 0}$ , that assigns a non-negative real value to each clock.

Let  $\Lambda = \{\lambda_1, \dots, \lambda_m\}$  (for some integer  $m$ ) be a finite set of *parameters*, *i. e.*, rational-valued constants that will be used here to represent the unknown response time of a component service. A *parameter valuation* is a function  $\pi : \Lambda \rightarrow \mathbb{Q}_{\geq 0}$  assigning a non-negative rational value to each parameter.

Henceforth, we use  $w$  (resp.  $\pi$ ) to denote a clock (resp. parameter) valuation.

A *linear term* over  $X \cup \Lambda$  is an expression of the form  $\sum_{1 \leq i \leq N} \alpha_i z_i + d$  for some  $N \in \mathbb{N}$ , with  $z_i \in X \cup \Lambda$ ,  $\alpha_i \in \mathbb{Q}_{\geq 0}$  for  $1 \leq i \leq N$ , and  $d \in \mathbb{Q}_{\geq 0}$ . We denote by  $\mathcal{L}_{X \cup \Lambda}$  the set of all linear terms over  $X$  and  $\Lambda$ . Similarly, we denote by  $\mathcal{L}_\Lambda$  the set of all linear terms over  $\Lambda$ . An *inequality* over  $X$  and  $\Lambda$  is of the form  $e \bowtie e'$  where  $\bowtie \in \{<, \leq\}$ , and  $e, e' \in \mathcal{L}_{X \cup \Lambda}$ .

A *convex constraint* (or *constraint*) is a conjunction of inequalities. We denote by  $\mathcal{C}_{X \cup \Lambda}$  the set of all convex constraints over  $X$  and  $\Lambda$ . Similarly, we denote by  $\mathcal{C}_\Lambda$  the set of all convex constraints over  $\Lambda$ .

Let  $C \in \mathcal{C}_{X \cup \Lambda}$ ,  $C[\pi]$  denotes the valuation of  $C$  with  $\pi$ , *i. e.*, the constraint over  $X$  obtained by replacing each  $\lambda \in \Lambda$  with  $\pi(\lambda)$  in  $C$ . Note that  $C[\pi]$  can be written as  $C \wedge \bigwedge_{\lambda_i \in \Lambda} \lambda_i = \pi(\lambda_i)$ . We say that  $w$  satisfies  $C[\pi]$  if the expression obtained by replacing each  $x \in X$  in  $C[\pi]$  with  $w(x)$  evaluates to true.

Given  $C \in \mathcal{C}_{X \cup \Lambda}$ , we define  $C^\uparrow$  as the *time elapsing* of  $C$ , *i. e.*, the constraint over  $X$  and  $\Lambda$  obtained from  $C$  by delaying all clocks by an arbitrary amount of time. That is:

$$C^\uparrow = \{(w', \pi) \mid w \text{ satisfies } C[\pi] \wedge \forall x \in X : w'(x) = w(x) + d, d \in \mathbb{R}_{\geq 0}\}.$$

Given  $C \in \mathcal{C}_{X \cup \Lambda}$  and  $X' \subseteq X$ , we denote by  $\text{prune}_{X'}(C)$  the constraint in  $\mathcal{C}_{X \cup \Lambda}$  that is obtained from  $C$  by pruning the clocks in  $X'$ ; this can be achieved using

variable elimination techniques such as Fourier-Motzkin (see, *e. g.*, [Sch86]). More generally, given  $C \in \mathcal{C}_{X \cup A}$ , we denote by  $C \downarrow_A$  the *projection* of constraint  $C$  onto  $A$ , *i. e.*, the constraint obtained from  $C$  by pruning all clock variables. Again, such a projection can be computed using Fourier-Motzkin elimination.

A *non-necessarily convex constraint* (or NNCC) is a conjunction of disjunction of inequalities<sup>1</sup>; NNCCs are used to represent the synthesized local time constraint obtained via the methods proposed in this paper. Note that the negation of an inequality remains an inequality; however, the negation of a convex constraint becomes (in the general case) an NNCC. We denote by  $\mathcal{NC}_A$  the set of all NNCCs over  $A$ .

Given  $C \in \mathcal{NC}_A$ , we say that  $\pi$  *satisfies*  $C$ , denoted by  $\pi \models C$ , if  $C[\pi]$  evaluates to true.  $C$  is *empty* if there does not exist a parameter valuation  $\pi$  such that  $\pi \models C$ ; otherwise  $C$  is *non-empty*. Given two constraints  $C_1, C_2 \in \mathcal{NC}_A$ , we say that  $C_2$  is *weaker* (or *more relaxed*) than  $C_1$ , denoted by  $C_1 \subseteq C_2$ , if  $\forall \pi : \pi \models C_1 \Rightarrow \pi \models C_2$ .

### 3.2 Syntax of composite service processes

BPEL [Alv+07] is an industrial standard for implementing composition of existing Web services by specifying an executable workflow using predefined activities. In this work, we assume the composite service is specified using the BPEL language. Basic BPEL activities that communicate with component Web services are `<receive>`, `<invoke>`, and `<reply>`, which are used to receive messages, invoke an operation of component Web services and return values respectively. These activities are *communication activities*. The control flow of the service is defined using structural activities such as `<flow>`, `<sequence>`, `<pick>` and `<if>`.

A composite service CS makes use of a finite number of component services to accomplish a task. Let  $E = \{S_1, \dots, S_n\}$  be the set of all component services that are used by CS. In this work, we assume that the response time of a composite service is based on the time spent on individual communication activities, and the time incurred by internal operations of the composite service is negligible.<sup>2</sup>

Composite services are expressed using *processes*. We define a formal syntax definition in the following.

**Definition 1** *Processes* are defined using the following grammar:

$P \hat{=} \text{rec}(S)$	receive activity
$\text{reply}(S)$	reply activity
$s\text{Inv}(S)$	synchronous invocation
$a\text{Inv}(S)$	asynchronous invocation
$P \parallel Q$	concurrent activity
$P ; Q$	sequential activity
$P \triangleleft b \triangleright Q$	conditional activity
$\text{pick}(\biguplus_{i=1}^n S_i \Rightarrow P_i, \biguplus_{j=1}^k \text{alarm}(a_j) \Rightarrow Q_j)$	pick activity

where  $S$  is a component service,  $P$  and  $Q$  are composite service processes,  $b$  is a Boolean expression, and  $a_j \in \mathbb{Q}_{>0}$  are positive rational numbers, for  $1 \leq j \leq k$ .

<sup>1</sup> Without loss of generality, we assume here that all NNCCs are in conjunctive normal form (CNF).

<sup>2</sup> We discuss the time incurred for internal operations in Section 6.6.

Let us describe below the BPEL syntax notations introduced in [Definition 1](#):

- $rec(S)$  and  $reply(S)$  are used to denote “receive from” and “reply to” a service  $S$ , respectively;
- $sInv(S)$  (resp.  $aInv(S)$ ) denotes the synchronous (resp. asynchronous) invocation of a component service  $S$ ;
- $P \parallel Q$  denotes the concurrent composition of BPEL activities  $P$  and  $Q$ ;
- $P ; Q$  denotes the sequential composition of BPEL activities  $P$  and  $Q$ ;
- $P \triangleleft b \triangleright Q$  denotes the conditional composition, where  $b$  is a guard condition on the process variables. If  $b$  evaluates to true, BPEL activity  $P$  is executed, otherwise activity  $Q$  is executed;
- $pick(\biguplus_{i=1}^n S_i \Rightarrow P_i, \biguplus_{j=1}^k alm(a_j) \Rightarrow Q_j)$  denotes the BPEL *pick* composition, which contains two types of activities, *i. e.*, *onMessage* activity and *onAlarm* activity. An *onMessage* activity  $S_i \Rightarrow P_i$  is activated when the message from service  $S_i$  arrives and BPEL activity  $P_i$  is subsequently executed; an *onAlarm* activity  $alm(a_j) \Rightarrow Q_j$  is activated at  $a_j$  time units, and BPEL activity  $Q_j$  is subsequently executed. The *pick* activity contains  $n$  *onMessage* activities and  $k$  *onAlarm* activities. Exactly one activity from these  $n + k$  activities will be executed. If multiple activities are activated at the same time, one of the activities will be chosen non-deterministically for execution. Given a *pick* activity  $P$ , we use  $P.onMessage$  and  $P.onAlarm$  to denote the *onMessage* and *onAlarm* branches of  $P$  respectively.

A *structural activity* is an activity that contains other activities. Concurrent, sequential, conditional, and pick activities are examples of structural activities. An activity that does not contain other activities is called an *atomic activity*, which includes receive, reply, synchronous invocation and asynchronous invocation activities.

Note that the communication activities can implicitly make use of variables for passing information. For example, let  $S$  be a component service that calculates the stock indices for a particular date. For synchronous invocation  $sInv(S)$ , it requires an input variable  $v_i$  that specifies the date information, and an output variable  $v_o$  to hold the return value from  $sInv(S)$ . To keep the notations concise, we abstract the usage and assignment of variables for communication activities.

We make the following assumption throughout this manuscript:

**Assumption 1** *All loops have a bound on the number of iterations and on the execution time.*

This assumption is necessary to ensure termination of our approach. We believe it is reasonable in practice (see [Section 6.6](#) for a discussion).

### 3.3 Parametric composite service models

Let us now formally define composite service models and parametric composite service models. Let  $\mathcal{P}_{np}$  denote the set of all possible (non-parametric) composite service processes.

**Definition 2 (Composite service model)** *A composite service model CS is a tuple  $(\mathcal{V}, v_0, N_0)$ , where  $\mathcal{V}$  is a finite set of variables,  $v_0 \in \text{Valuations}(\mathcal{V})$  is an initial*

valuation that maps each variable to its initial value, and  $N_0 \in \mathcal{P}_{np}$  is a composite service process (defined according to the grammar of [Definition 1](#)) making use of the variables in  $\mathcal{V}$ .

Each service comes with a *response time*, which is a rational-valued constant, and can be seen as an upper bound on the time that a service needs to successfully return its answer.

Given a composite service CS, let  $t_i \in \mathbb{R}_{\geq 0}$  be the response time of component service  $S_i$  for  $i \in \{1, \dots, n\}$ , and let  $E_t = \{t_1, \dots, t_n\}$  be a set of component service response times that fulfill the global time requirement of service CS. Because  $t_i$ , for  $i \in \{1, \dots, n\}$ , is a rational number, there are infinitely many possible values, even in a bounded interval (and even if one restricts these values to rational numbers). A method to tackle this problem is to reason *parametrically*, by considering these response times as unknown constants, or *parameters*.

We now extend the definitions of services, composite service processes and composite service model to the parametric case. First, a parametric service is a service  $S_i$ , the response time of which is now a parameter  $\lambda_i \in A$ , instead of a rational-valued constant. Then, a parametric composite service process is a service process defined according to the grammar of [Definition 1](#), where services (“S” in [Definition 1](#)) are now parametric services. We denote by  $\mathcal{P}$  the set of all possible parametric composite service processes. Finally, parametric composite service models are defined similarly to composite service models, except that the composite service processes are now parametric composite service processes.

**Definition 3 (Parametric composite service model)** A *parametric composite service model* CS is a tuple  $(\mathcal{V}, v_0, A, P_0, C_0)$ , where  $\mathcal{V}$  is a finite set of variables;  $v_0 \in \text{Valuations}(\mathcal{V})$  is an initial valuation that maps each variable to its initial value;  $A$  is a finite set of parameters;  $P_0 \in \mathcal{P}$  is a parametric composite service process making use of the variables in  $\mathcal{V}$  and  $C_0 \in \mathcal{C}_A$  is the initial parametric constraint.

*Example 1* Let  $\mathcal{V} = \{y_1\}$ . Let  $v_0$  be such that  $v_0(y_1) = 0$ . Let  $A = \{\lambda_1, \lambda_2\}$ . Let  $P_0 = \text{pick}(S \Rightarrow \text{sInv}(S_1), \text{alm}(1) \Rightarrow \text{sInv}(S_2)) \triangleleft y_1 > 0 \triangleright \text{Stop}$ . Let  $C_0 = \lambda_1 < \lambda_2$ . Let  $\lambda_i$  denote the response time of  $\text{sInv}(S_i)$ ,  $i \in \{1, 2\}$ .

Then CS =  $(\mathcal{V}, v_0, A, P_0, C_0)$  is a parametric composite service model.

*Process and model valuation* Given a parametric composite service process  $P$  with a parameter set  $A = \{\lambda_1, \dots, \lambda_m\}$  and given a parameter valuation  $(\pi(\lambda_1), \dots, \pi(\lambda_m))$ ,  $P[\pi]$  denotes the *valuation* of  $P$  with  $\pi$ , *i. e.*, the process where each occurrence of a parameter  $\lambda_i$  is replaced with its valuation  $\pi(\lambda_i)$ .

Given a parametric composite service model CS with a parameter set  $A = \{\lambda_1, \dots, \lambda_m\}$ , and given a parameter valuation  $(\pi(\lambda_1), \dots, \pi(\lambda_m))$ ,  $\text{CS}[\pi]$  denotes the *valuation* of CS with  $\pi$ , *i. e.*, the model  $(\mathcal{V}, v_0, A, P_0, C)$ , where  $C$  is  $C_0 \wedge \bigwedge_{i=1}^m (\lambda_i = \pi(\lambda_i))$ . Note that  $\text{CS}[\pi]$  can be seen as a non-parametric service model  $(\mathcal{V}, v_0, P_0[\pi])$ .

*Example 2* Consider the parametric composite service model CS defined in [Example 1](#). Assume  $\pi$  such that  $\pi(\lambda_1) = 1$  and  $\pi(\lambda_2) = 2$ . Then  $P_0[\pi] = \text{pick}(S \Rightarrow \text{sInv}(S_1), \text{alm}(1) \Rightarrow \text{sInv}(S_2)) \triangleleft y_1 > 0 \triangleright \text{Stop}$ , where the response time of  $\text{sInv}(S_1)$  is 1, and the response time of  $\text{sInv}(S_2)$  is 2.

### 3.4 Bad activities

Given a BPEL service  $CS$ , we define a *bad activity* as an atomic activity such that its execution leads the composite service  $CS$  to violate the global time requirement. To distinguish bad activities, we allow the user to annotate a BPEL activity  $A$  as a bad activity. The annotation can be achieved, for example, by using extension attributes of BPEL activities. This work can be performed manually or using semi-automated procedures.

*Example 3* Consider again the example in Section 2. Then “Reply ‘Failure’” is a bad activity, denoted in Fig. 2 by  $\times$ .

## 4 A formal semantics for parametric composite services

In this section, we provide our parametric composite service model with a formal semantics, defined in the form of a labeled transition system (LTS). The semantics we use is inspired by the one proposed for (parametric) stateful timed Communicating Sequential Processes (CSP) [Sun+13; And+14], that makes use of implicit clocks.

We first recall LTSs (Section 4.1) and define symbolic states (Section 4.2). Following that, we define implicit clocks and the associated functions, *i. e.*, activation and idling (Section 4.3). We then introduce our formal semantics (Section 4.4), and apply it to an example (Section 4.5). We finally prove a technical result relating the non-parametric and the parametric service models (Section 4.6).

### 4.1 Labeled transition systems

**Definition 4 (Labeled transition system)** A *labeled transition system (LTS)* is a tuple  $LTS = (S, s_0, \Sigma, \delta)$ , where

- $S$  is a set of states;
- $s_0 \in S$  is the initial state;
- $\Sigma$  is a set of actions; and
- $\delta \subseteq S \times \Sigma \times S$  is a transition relation.

Given  $LTS = (S, s_0, \Sigma, \delta)$ , a state  $s \in S$  is a *terminal state* if there does not exist a state  $s' \in S$  and an action  $a \in \Sigma$  such that  $(s, a, s') \in \delta$ ; otherwise,  $s$  is said to be a *non-terminal state*. There is a *run* from a state  $s$  to state  $s'$ , where  $s, s' \in S$ , if there exists an alternating sequence of states and actions  $\langle s_1, a_1, s_2, \dots, a_{n-1}, s_n \rangle$ , where  $s_i \in S$  for  $1 \leq i \leq n$ ,  $a_i \in \Sigma$  for  $1 \leq i \leq n-1$ ,  $s_1 = s$ ,  $s_n = s'$ , and  $\forall i \in \{1, \dots, n-1\}, (s_i, a_i, s_{i+1}) \in \delta$ . A *complete run* is a run that starts in the initial state  $s_0$  and ends in a terminal state. Given a state  $s \in S$ , we use  $\text{succ}(s)$  to denote the set of states reachable in one step from  $s$ ; formally,  $\text{succ}(s) = \{s' \mid \exists a \in \Sigma, \exists s' \in S : (s, a, s') \in \delta\}$ .

In the following, we introduce the notion of LTS starting from a state  $s$  which is defined as the LTS containing  $s$  and all its successor states and transitions.

**Definition 5 (sub-LTS)** Let  $LTS = (S, s_0, \Sigma, \delta)$  be an LTS, and let  $s$  be a state of  $S$ . The *sub-LTS* of  $LTS$  starting from  $s$  is  $(S', s, \Sigma', \delta')$ , where

1.  $S' \subseteq S$  is the set of states reachable from  $s \in S$  in  $LTS$ ;
2.  $\delta' \subseteq \delta$  is the transition relation satisfying the following condition:  $(s_1, a, s_2) \in \delta'$  if  $s_1, s_2 \in S'$  and  $(s_1, a, s_2) \in \delta$ ; and
3.  $\Sigma' \subseteq \Sigma$  is the set of all actions used in  $\delta'$ , *i. e.*,  $\{a \mid \exists s_1, s_2 \in S' : (s_1, a, s_2) \in \delta'\}$ .

## 4.2 Symbolic states

In the following, we equip our parametric composite service models with a symbolic semantics, *i. e.*, a semantics, a run of which will capture a (possibly infinite) set of runs, for a (possibly infinite) set of parameter valuations.

Let us first define the notion of (symbolic) state of a parametric composite service model.

**Definition 6 (State)** Given a parametric composite service model  $CS = (\mathcal{V}, v_0, \Lambda, P_0, C_0)$ , a (symbolic) *state* of  $CS$  is a tuple  $s = (v, P, C, D)$ , where  $v \in \text{Valuations}(\mathcal{V})$  is a valuation of the variables,  $P$  is a composite service process,  $C$  is a constraint over  $\mathcal{C}_{X \cup \Lambda}$ , and  $D \in \mathcal{L}_A$  is the (parametric) elapsed time from the initial state  $s_0$  to state  $s$ , excluding the idling time in state  $s$ .

Given a state  $s = (v, P, C, D)$ , we use the notation  $s.v$  to denote the field  $v$  of  $s$ , and similarly for  $s.P$ ,  $s.C$  and  $s.D$ . When a parametric composite service model  $CS$  has no variable, we denote each state  $s \in S$  by  $(P, C, D)$  for the sake of brevity.

## 4.3 Implicit clocks

In order to provide parametric composite service models with a symbolic semantics, we use *clocks* to record the elapsing of time. Recall from [Section 3.1](#) that clocks are real-valued variables initially equal to 0, and evolving all at the same rate; some clocks may be reset to 0. Clocks are used to record the time elapsing in several formalisms, in particular in timed automata (TAs) [\[AD94\]](#). In TAs, the clocks are defined as part of the models and state space. It is known that the state space of the system may grow exponentially with the number of clocks and that the fewer clocks, the more efficient real-time model checking is [\[BY03\]](#). In (P)TAs, it is possible to dynamically reduce the number of clocks [\[DY96; And13\]](#). An alternative approach is to define a semantics that create clocks on the fly when necessary, and prune them when they are no longer needed. This approach was initially proposed for stateful timed CSP [\[Sun+13\]](#), and shares similarities with *firing times* in time Petri nets [\[Mer74\]](#). This allows a smaller state space compared to the explicit clock approach. We refer to this second approach [\[Sun+13\]](#) as the *implicit clock approach*, and adopt this implicit clock approach in our work.

### 4.3.1 Clock activation

Clocks are implicitly associated with processes. For instance, given a communication activity  $sInv(S)$ , a clock starts measuring time once the activity becomes activated. To introduce clocks on the fly, we define an activation function  $Act$  in the following definition, in the spirit of the one defined in [\[Sun+13; And+14\]](#).

In short, this definition explains how to associate a new clock to a process: this clock will only be associated to the new processes with timing constraints, while it will not be associated to untimed processes nor to processes to which another implicit clock is already associated.

**Definition 7** Given a process, we define the activation function  $Act$  using the following set of recursive rules:

$$\begin{array}{lll}
Act(A(S), x) & = & A(S)_x \quad \text{A1} \\
Act(mpick, x) & = & mpick_x \quad \text{A2} \\
Act(A(S)_{x'}, x) & = & A(S)_{x'} \quad \text{A3} \\
Act(mpick_{x'}, x) & = & mpick_{x'} \quad \text{A4} \\
Act(P \oplus Q, x) & = & Act(P, x) \oplus Act(Q, x) \quad \text{A5} \\
Act(P; Q, x) & = & Act(P, x); Q \quad \text{A6}
\end{array}$$

where  $A \in \{rec, sInv, aInv, reply\}$ ,  $\oplus \in \{|||, \langle b \rangle\}$ , and  $mpick = pick(\biguplus_{i=1}^n S_i \Rightarrow P_i, \biguplus_{j=1}^k alrm(a_j) \Rightarrow Q_j)$

Let us explain [Definition 7](#). Given a process  $P$ , we denote by  $P_x$  the corresponding process that has been associated with clock  $x$ . When a new state  $s$  is reached, the activation function is called to assign a new clock for each newly activated communication activity.

- Rules A1 and A2 state that a new clock is associated with a BPEL communication activity  $A$  if  $A$  is newly activated.
- Rules A3 and A4 state that if a BPEL communication activity has already been assigned a clock, it will not be reassigned one.
- Rules A5 and A6 state that function  $Act$  is applied recursively to activate the child activities for BPEL structural activities.
- For rule A6, function  $Act$  is applied only to activity  $P$ , but not to activity  $Q$ , since activity  $P$  is the immediate subsequent activity (activity  $Q$  will be executed only after the completion of activity  $P$ ).

*Example 4* Let  $P = sInv(S_1) ||| aInv(S_2)$ . Then, applying rules A5 and A1,  $Act(P, x) = sInv(S_1)_x ||| aInv(S_2)_x$ . Note that  $x$  is associated with both processes, as they are both simultaneously activated.

*Example 5* Let  $P = sInv(S_1)_{x'}; aInv(S_2)$ . Then, applying rules A6 and A3,  $Act(P, x) = sInv(S_1)_{x'}; aInv(S_2)$ . Indeed, the first invocation  $sInv(S_1)_{x'}$  is already associated to another clock  $x'$  (rule A3) while the right-hand part of the sequence is not yet activated (rule A6).

Given a process  $P$ , we denote by  $aclk(P)$  the set of *active clocks* associated with  $P$ .

*Example 6* Assume process  $P = sInv(S_1)_{x_0} ||| sInv(S_2)_{x_1}$ . The set of active clocks associated with  $P$  is  $aclk(P) = \{x_0, x_1\}$ .

#### 4.3.2 Idling function

We define in [Definition 8](#) below the function *idle* that, given a state  $s$ , returns a constraint that specifies how long an activity can idle at state  $s$ . The result is a

constraint over  $X \cup \Lambda$ . This idling function is similar in essence to the *time elapsing* on symbolic states (zones or parametric zones) defined for TAs or PTAs [BY03; Hun+02].

**Definition 8** Given a process, we define the idling function *idle* using the following set of recursive rules:

$$\begin{array}{lll}
idle(A(S)_x) & = & x \leq \lambda_S \quad \text{I1} \\
idle(B(S)_x) & = & x = 0 \quad \text{I2} \\
idle(P \oplus Q) & = & idle(P) \wedge idle(Q) \quad \text{I3} \\
idle(P ; Q) & = & idle(P) \quad \text{I4} \\
idle(mpick_x) & = & x \leq \lambda_S \wedge \bigwedge_{j=0}^k x \leq a_j \quad \text{I5}
\end{array}$$

where  $A \in \{rec, sInv\}$ ,  $B \in \{aInv, reply\}$ ,  $\oplus \in \{|||, \langle b \rangle\}$ ,  $mpick = pick(\bigoplus_{i=1}^n S_i \Rightarrow P_i, \bigoplus_{j=1}^k alrm(a_j) \Rightarrow Q_j)$ , and  $\lambda_S$  is the parametric response time of service  $mpick_x$ .

Let us explain [Definition 8](#).

- Rule I1 considers the situation when the communication requires waiting for the response of a component service  $S$ , and the value of clock  $x$  must not be larger than the response time parameter  $\lambda_S$  of the service: that is, one can only remain in this state while  $x \leq \lambda_S$  remains valid.
- Rule I2 considers the situation when no waiting is required: therefore, the clock constraint  $x = 0$  implies that this state should be left within 0-time, as these actions are instantaneous.
- Rules I3 and I4 state that the function *idle* is applied recursively to activate the child activities of a BPEL structural activity.
- Similar to rule A6, for rule I4, function *Act* is applied only to activity  $P$ , but not to activity  $Q$ , since only activity  $P$  is executed next. Therefore, given a state  $s$  and activity  $P ; Q$ , we only need to consider how long the activity  $P$  can idle at state  $s$ .
- Rule I5 states that the activity can idle only until  $\lambda_S$  or any of the alarms  $a_j$  is reached. The conjunction comes from the fact that, as soon as any alarm reaches its time-out, then it will be triggered, therefore leading the system to leave this symbolic state.

*Example 7* Let  $P = sInv(S_1) ||| aInv(S_2)$ . Assume the response time of  $S_i$  is  $\lambda_i$  for  $i \in \{1, 2\}$ . Recall from [Example 4](#) that  $Act(P, x) = sInv(S_1)_x ||| aInv(S_2)_x$ . Let us apply *idle* to  $Act(P, x)$ . Applying rules I3, I1 and I2, we get  $x \leq \lambda_1 \wedge x = 0$ .

#### 4.4 Operational semantics

The operational semantics will be defined in the form of an LTS. The actions labeling the LTS will be sequences of rules; these rules will be a set of rules (similar to those of parametric stateful timed CSP [[And+14](#)]) defining the transitions of the semantics, and will be explained below. Let

$$\text{Rules} = \{rSInv, rRec, rReply, rAInv, rCond1, rCond2, rCond3, rCond4, rSeq1, rSeq2, rFlow1, rFlow2, \} \cup (rPickM \times \mathbb{N}) \cup (rPickA \times \mathbb{N})$$



be the set of rules that will be used by the LTS. Two rules ( $rPickM$  and  $rPickA$ ) are associated with a positive integer, so as to remember which subprocess is derived (this will be explained later on). Let  $\mathbf{Sequences}(\mathbf{Rules})$  denote the set of *sequences* of rules, *i. e.*, non-empty ordered elements of  $\mathbf{Rules}$  (possibly used several times). An example of a sequence of rule is  $\langle rRec, rReply, rRec, (rPickM, 2) \rangle$ . Sequence concatenation is denoted by operator  $+$ .

We can now define the semantics of a parametric composite service model in the form of an LTS. Let  $ClkSeq = \langle x_0, x_1, \dots \rangle$  be a sequence of clocks. We will need  $ClkSeq$  to pick a fresh clock when applying the clock activation function  $Act$  defined previously.

**Definition 9 (semantics of composite services)** Let  $\mathbf{CS} = (\mathcal{V}, v_0, A, P_0, C_0)$  be a parametric composite service model. The *semantics* of  $\mathbf{CS}$  (hereafter denoted by  $LTS_{\mathbf{CS}}$ ) is the LTS  $(S, s_0, \mathbf{Sequences}(\mathbf{Rules}), \delta)$  where

$$\begin{aligned} S &= \{(v, P, C, D) \in \text{Valuations}(\mathcal{V}) \times \mathcal{P} \times \mathcal{C}_{X \cup A} \times \mathcal{L}_A\}, \\ s_0 &= (v_0, P_0, C_0, 0) \end{aligned}$$

and the transition relation  $\delta$  is the smallest transition relation satisfying the following. For all  $(v, P, C, D) \in S$ , if  $x$  is the first clock in the sequence  $ClkSeq$  which is not in  $acclk(P)$ , and  $(v, Act(P, x), C \wedge x = 0, D) \xrightarrow{seq} (v', P', C', D')$  where  $C'$  is satisfiable, then we have:  $((v, P, C, D), seq, (v', P', prune_{X \setminus acclk(P')}(C'), D')) \in \delta$ .

The transition relation  $\hookrightarrow$  is specified by a set of rules, given in [Appendix A](#). Let us first explain these rules, after which we will go back to the explanation of [Definition 9](#). The transition relation is labeled by a sequence of rules, that allows one to remember by using which sequence of rules a process evolves into another one.

**Synchronous invocation.** Rule  $rSInv$  states that a state  $s = (v, sInv(S)_x, C, D)$  may evolve into the state  $s' = (v', Stop, (x = \lambda_S) \wedge C^\uparrow, D + \lambda_S)$ , where  $Stop$  is the activity that does nothing, and  $\lambda_S$  is the parametric response time of component service  $S$ . Note that, from [Definition 9](#), the condition  $(x = \lambda_S) \wedge C^\uparrow$  is necessarily satisfied (otherwise this evolution is not possible). Furthermore, the parametric duration from the initial state ( $D$ ) is incremented by  $\lambda_S$ . Rules  $rRec$ ,  $rReply$  and  $rAInv$  are similar.

**Pick activity.** Rule  $rPickM$  encodes the transition that takes place due to an *onMessage* activity, where  $\lambda_i$  denotes the parametric response time of  $P_i$ . Let us explain the constraint  $(x = \lambda_i) \wedge idle(mpick_x) \wedge C^\uparrow$ . First, after the transition, the current clock  $x$  needs to be equal to the parametric response time of service  $S_i$ , *i. e.*,  $x = \lambda_i$ . Second, the constraint  $idle(mpick_x)$  is added to ensure that  $x$  remains smaller or equal to the maximum duration of the  $mpick_x$  activity. Third, the constraint  $C^\uparrow$  denotes the time elapsing of  $C$ . Observe that the transition in  $\hookrightarrow$  is labeled using the pair  $(rPickM, i)$  so as to remember that the  $i$ th process (*i. e.*,  $P_i$ ) has been selected.

Rule  $rPickA$  (for an *onAlarm* activity) is similar; observe that, instead of using the parametric response time, we use the time stipulated by the alarm (*i. e.*,  $a_j$ ) of process  $Q_j$ .

**Conditional activity.** Given a conditional composition  $A \triangleleft b \triangleright B$ , the guard condition  $b$  is a Boolean, hence its values are in  $\{true, false\}$ . As a consequence, given a valuation  $v$  of the variables, then  $v(b) \in \{true, false, \perp\}$ . We have that

$v(b) = \perp$  when the evaluation of  $b$  is unknown, due to the fact that there may be uninitialized variables in  $b$ . Since  $b$  might be evaluated to either true or false at certain stages at runtime, we explore both activities  $A$  and  $B$  when  $v(b) = \perp$  so as to reason about all possible scenarios. The case of  $v(b) = \perp$  is captured by rules  $rCond1$  and  $rCond2$ , and the cases where  $v(b) \in \{true, false\}$  are captured by rules  $rCond3$  and  $rCond4$ .

**Sequential activity.**  $rSeq1$  states that if activity  $A'$  is not a *Stop* activity (*i. e.*, activity  $A'$  has not finished its execution), then a state containing activity  $A;B$  may evolve into a state containing activity  $A';B$ . Otherwise, if  $A$  is a *Stop* activity (*i. e.*, activity  $A$  has finished its execution), then the state may evolve into  $B$ . This is captured by  $rSeq2$ .

**Concurrent activity.** For concurrent activity  $A \parallel B$ , both activities  $A$  and activity  $B$  are executed. This is captured by  $rFlow1$  and  $rFlow2$  respectively.  $rFlow1$  states that if state  $(v, A, C, D)$  can evolve into  $(v', A', C', D')$ , then a state containing  $A \parallel B$  can evolve into a state containing  $A' \parallel B$ , if  $C' \wedge idle(B)$  holds. That is, the clock constraints in  $C'$  cannot exceed the duration activity  $B$  can last for. Rule  $rFlow2$  is dual.

Let us now explain [Definition 9](#). Starting from the initial state  $s_0 = (v_0, P_0, C_0, 0)$ , we iteratively construct successor states as follows. Given a state  $(v, P, C, D)$ , a fresh clock  $x$  which is not currently associated with  $P$  is picked from  $ClkSeq$ . The state  $(v, P, C, D)$  is transformed into  $(v, Act(P, x), C \wedge x = 0, D)$ , *i. e.*, timed processes which just become activated are associated with  $x$  and  $C$  is conjuncted with  $x = 0$ . Then, a firing rule is applied to get a target state  $(v', P', C', D')$ . Lastly, clocks which do not appear within  $P'$  are pruned from  $C'$ . More in details, the expression  $prune_{X \setminus aclk(P')}(C')$  denotes that we remove all clocks from the obtained constraint  $C'$  by existential quantification, except those which are still active in the successor  $P'$  of  $P$  (recall that  $prune_X(C)$  was defined in [Section 3.1](#)).

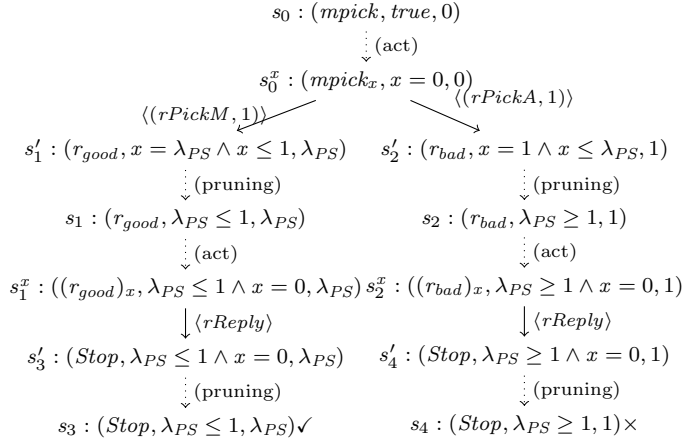
Observe that one clock is introduced and zero or more clocks may be pruned during a transition. In practice, a clock is introduced only when necessary; if the activation function does not activate any subprocess, no new clocks are created.

*Good and bad states* Let us define good and bad states in the LTS obtained from [Definition 9](#). The execution of a bad activity will make the execution of CS end in an undesired terminal state, which we refer to as a *bad state*. A terminal state which is not a bad state is called a *good state*.

#### 4.5 Application to an example

Consider a composite service CS starting from  $pick(PS \Rightarrow reply(User), alrm(1) \Rightarrow [reply(User)]_{bad})$ . Assume  $\lambda_{PS}$  is the parametric response time of service  $PS$ . (Note that CS is a part of the SMIS example from [Section 2](#).) The states of CS computed according to [Definition 9](#) are given in [Fig. 3](#), including intermediate states (detailed in the following). Since CS has no variable, then  $v = \emptyset$  in all states; therefore, we omit the component  $v$  from all states for sake of brevity.

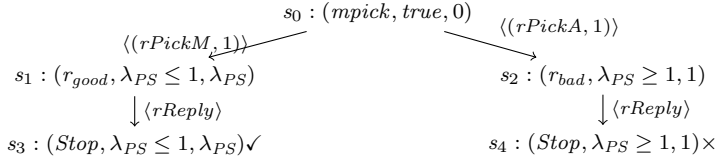
- At state  $s_0$ , the activation function assigns clock  $x$  to record time elapsing of pick activity  $mpick$ , with  $x$  initialized to zero. The tuple becomes the intermediate state  $s_0^x = (mpick_x, x = 0, 0)$ .



where  $mpick = pick(PS \Rightarrow r_{good}, alrm(1) \Rightarrow r_{bad})$ ,  $r_{good} = reply(User)$ ,  $r_{bad} = [reply(User)]_{bad}$ , and  $\lambda_{PS}$  is the parametric response time of service  $PS$ .

Fig. 3: Computing states of service CS (including intermediate states)

- From intermediate state  $s_0^x$ , the process may evolve into the intermediate state  $s_1'$  by applying the rule  $rPickM$ , if the constraint  $C_1 = ((x = \lambda_{PS}) \wedge idle(mpick_x) \wedge (x = 0)^\uparrow)$ , where  $idle(mpick_x) = (x \leq \lambda_{PS} \wedge x \leq 1)$  and  $(x = 0)^\uparrow$  (i. e.,  $x \geq 0$ ), is satisfiable. Intuitively,  $C_1$  denotes the constraint where  $\lambda_{PS}$  time units elapsed since clock  $x$  has started. In fact,  $C_1$  is satisfiable (for example with  $\lambda_{PS} = 0.5$  and  $x = 0.5$ ). Therefore, it may evolve into the intermediate state  $s_1' = (r_{good}, (x = \lambda_{PS}) \wedge idle(mpick_x) \wedge (x = 0)^\uparrow, \lambda_{PS}) = (r_{good}, (x = \lambda_{PS}) \wedge x \leq 1, \lambda_{PS})$ . Since clock  $x$  is not used anymore in  $s_1'.P$  which is  $r_{good}$ , it is pruned. After pruning of clock variable  $x$  and simplification of the expression, the intermediate state  $s_1'$  becomes the state  $s_1 = (r_{good}, \lambda_{PS} \leq 1, \lambda_{PS})$ .
- From intermediate state  $s_0^x$ , the process may also evolve into the intermediate state  $s_2'$ , by applying the rule  $rPickA$ , if the constraint  $C_2 = ((x = 1) \wedge idle(mpick_x) \wedge (x = 0)^\uparrow)$ , where  $idle(mpick_x) = (x \leq \lambda_{PS} \wedge x \leq 1)$  and  $(x = 0)^\uparrow$  (i. e.,  $x \geq 0$ ), is satisfiable. It is easy to see that  $C_2$  is satisfiable; therefore, the process may evolve into the intermediate state  $s_2' = (r_{bad}, (x = 1) \wedge x \leq \lambda_{PS}, 1)$ . After clock pruning from intermediate state  $s_2'$ , it becomes state  $s_2 = (r_{bad}, \lambda_{PS} \geq 1, 1)$ .
- From state  $s_1$ , activation function assigns clock  $x$  to the reply activity  $r_{good}$ , and the process evolves into intermediate state  $s_1^x$ . From  $s_1^x$ , the process may evolve into intermediate state  $s_3'$  by applying rule  $rReply$ , if the constraint  $C_3 = ((x = 0) \wedge (\lambda_{PS} \leq 1)^\uparrow)$  is satisfiable, where  $(\lambda_{PS} \leq 1)^\uparrow = \lambda_{PS} \leq 1$ . In fact it is, and therefore it evolves into state  $s_3' = (Stop, \lambda_{PS} \leq 1 \wedge (x = 0), \lambda_{PS})$ . After pruning of the non-active clock, it evolves into the terminal state  $s_3 = (Stop, \lambda_{PS} \leq 1, \lambda_{PS})$ . Since the terminal state is not caused by a bad activity,  $s_3$  is considered as a good state, denoted by  $\checkmark$  in Fig. 3.
- From state  $s_2$ , the process may also evolve into the terminal state  $s_4 = (Stop, \lambda_{PS} \geq 1, 1)$ . Since the terminal state is caused by a bad activity, it is considered as a bad state, denoted by  $\times$  in Fig. 3.



where  $mpick = pick(PS \Rightarrow r_{good}, alrm(1) \Rightarrow r_{bad})$ ,  $r_{good} = reply(User)$ ,  $r_{bad} = [reply(User)]_{bad}$ , and  $\lambda_{PS}$  is the parametric response time of service  $PS$ .

Fig. 4: LTS of service CS

Note that all states  $s_i^x$  and  $s_j^y$ , where  $i, j \in \mathbb{N}$  and  $0 \leq i \leq 4$ , are intermediate states. State  $s_i^x$  is the state  $s_i$  after clock assignment operations are applied. State  $s_j^y$  is the state  $s_j$  before clock pruning operations are applied. These intermediate states are given in Fig. 3 to illustrate in details the application of the semantics. The LTS of CS (without the intermediate states) is given in Fig. 4.

#### 4.6 A technical result: the reachability condition

We defined the operational semantics of parametric composite service models as an LTS, the states of which contain information on clocks and parameters in the form of a constraint  $C$ . We now show that, for any reachable state of this LTS along a run, a parameter valuation  $\pi$  satisfies  $C$  iff the model valuated with  $\pi$  has an equivalent run. This is called the *reachability condition*. Similar results have been proved for parametric timed automata [Hun+02], parametric time Petri nets [TLR09] or parametric stateful timed CSP [And+14].

We first need several definitions and intermediate results. Given a parametric service model CS and a parameter valuation  $\pi$ , let us relate runs of  $LTS_{CS}$  and  $LTS_{CS[\pi]}$ . We will say that two runs are equivalent if they share the same discrete support, *i. e.*, follow the same application of sequences of rules regardless of the actual timing values.

**Definition 10 (equivalent runs)** Let CS be a parametric service model, and let  $\pi$  be a parameter valuation.

Let  $\rho = \langle (v_0, P_0, C_0, D_0), seq_0, (v_1, P_1, C_1, D_1), \dots, seq_{n-1}, (v_n, P_n, C_n, D_n) \rangle$  be a run of  $LTS_{CS[\pi]}$ . Let  $\rho' = \langle (v'_0, P'_0, C'_0, D'_0), seq'_0, (v'_1, P'_1, C'_1, D'_1), \dots, seq'_{n-1}, (v'_n, P'_n, C'_n, D'_n) \rangle$  be a run of  $LTS_{CS}$ .

The two runs  $\rho$  and  $\rho'$  are *equivalent* if  $v_i = v'_i$  and  $P_i = P'_i[\pi]$  for  $0 \leq i \leq n$  and  $seq_i = seq'_i$  for  $0 \leq i \leq n-1$ .

The following lemma states that, given a run of  $LTS_{CS[\pi]}$ , there exists a unique equivalent run in  $LTS_{CS}$ .

**Proposition 1** Let CS be a parametric service model, and let  $\pi$  be a parameter valuation. Let  $\rho_\pi$  be a run of  $LTS_{CS[\pi]}$ .

Then there exists a unique run of  $LTS_{CS}$  equivalent to  $\rho_\pi$ .

*Proof* By induction on the length of the runs. We prove in fact a slightly stronger result: given a state  $(v, P, C, D)$  of a run  $\rho$  in  $LTS_{CS[\pi]}$ , and given a state  $(v', P', C', D')$  of the equivalent run  $\rho'$  in  $LTS_{CS}$ , we show that these two runs are not only equivalent, but also that  $C \subseteq C'$ .

**Base case.** From [Definition 9](#), the initial state of  $LTS_{CS}$  is  $(v_0, P_0, C_0, 0)$ . The initial state of  $LTS_{CS[\pi]}$  is  $(v_0, P_0[\pi], C_0[\pi], 0)$ . Since  $C_0[\pi] \subseteq C_0$ , then the result trivially holds.

**Induction step.** Assume  $\rho_\pi$  is a run of  $LTS_{CS[\pi]}$  of length  $m$  reaching state  $(v_1, P_1, C_1, D_1)$ ; assume there exists a unique run of  $LTS_{CS}$  equivalent to  $\rho_\pi$  and of length  $m$ , reaching state  $(v'_1, P'_1, C'_1, D'_1)$ . From [Definition 10](#), it holds that  $v_1 = v'_1$  and  $P_1 = P'_1[\pi]$ . From the induction hypothesis, it holds that  $C_1 \subseteq C'_1$ .

Let  $(v_2, P_2, C_2, D_2)$  be the successor state of  $(v_1, P_1, C_1, D_1)$  via a given sequence of rules  $seq$  in  $\rho_\pi$ .

Assume  $(v_2, P_2, C_2, D_2)$  is obtained from  $(v_1, P_1, C_1, D_1)$  by applying rule  $rSInv$  in [Appendix A](#). Since  $C_1 \subseteq C'_1$ , then rule  $rSInv$  can also be applied to  $(v'_1, P'_1, C'_1, D'_1)$ , yielding a state  $(v'_2, P'_2, C'_2, D'_2)$ . Now, we have:

$$\begin{aligned} C_1 \subseteq C'_1 &\implies C_1^\uparrow \subseteq C'_1{}^\uparrow \\ \implies (x = \pi(\lambda_S) \wedge C_1^\uparrow) &\subseteq (x = \lambda_S \wedge C'_1{}^\uparrow) \\ \implies C_2 \subseteq C'_2. \end{aligned}$$

In particular,  $C_2 \subseteq C'_2$  implies that  $C'_2$  is non-empty, hence the state  $(v'_2, P'_2, C'_2, D'_2)$  is a valid state. In addition, since  $P_1 = P'_1[\pi]$  and rule  $rSInv$  derives to  $Stop$ , then  $P_2 = P'_2[\pi]$ . Variables are updated in the same manner on both sides, hence  $v_2 = v'_2$ . The proof is similar for other rules in [Appendix A](#).

Finally, the successor state  $(v'_2, P'_2, C'_2, D'_2)$  is the unique successor state of  $(v'_1, P'_1, C'_1, D'_1)$  in  $LTS_{CS}$  via this sequence of rules. Hence there exists a unique run of  $LTS_{CS}$  equivalent to  $\rho_\pi$  and of length  $m + 1$ .

We now prove the dual result. [Proposition 2](#) states that, given a run  $\rho$  of  $LTS_{CS}$ , there exists a unique equivalent run in  $LTS_{CS[\pi]}$ , provided  $\pi$  satisfies the parametric constraint associated with the last state of  $\rho$ .

**Proposition 2** *Let CS be a parametric service model, and let  $\pi$  be a parameter valuation. Let  $\rho$  be a run of  $LTS_{CS}$  ending in a state  $(v_n, P_n, C_n, D_n)$ .*

*For any  $\pi \models C_n \downarrow_A$ , there exists a unique run of  $LTS_{CS[\pi]}$  equivalent to  $\rho$ .*

*Proof* By induction on the length of the runs. We prove in fact a slightly stronger result: given a state  $(v', P', C', D')$  of a run  $\rho'$  in  $LTS_{CS}$ , and given a state  $(v, P, C, D)$  of the equivalent run  $\rho$  in  $LTS_{CS[\pi]}$ , we show that these runs are not only equivalent, but also that  $C = C'[\pi]$ .

**Base Step:** From [Definition 9](#), the initial state of  $LTS_{CS}$  is  $(v_0, P_0, C_0, 0)$ . The initial state of  $LTS_{CS[\pi]}$  is  $(v_0, P_0[\pi], C_0[\pi], 0)$ . Since  $C_0 = true$  then  $C_0 = C_0[\pi]$ . Hence the result trivially holds in that case.

**Induction step:** Assume  $\rho$  is a run of  $LTS_{CS}$  of length  $m$  reaching state  $(v'_1, P'_1, C'_1, D'_1)$ . Let  $(v'_2, P'_2, C'_2, D'_2)$  be the successor state of  $(v'_1, P'_1, C'_1, D'_1)$  via a sequence of rules  $seq$  in  $\rho$ . Let  $\pi \models C'_2 \downarrow_A$ . Assume there exists a unique run of  $LTS_{CS[\pi]}$  equivalent to  $\rho$  and of length  $m$ , reaching state  $(v_1, P_1, C_1, D_1)$ . From [Definition 10](#), it holds that  $v_1 = v'_1$  and  $P_1 = P'_1[\pi]$ . From the induction hypothesis, it holds that  $C_1 = C'_1[\pi]$ .

Assume  $(v'_2, P'_2, C'_2, D'_2)$  is obtained from  $(v'_1, P'_1, C'_1, D'_1)$  by applying rule  $rSInv$  in [Appendix A](#). Recall that  $C_1 = C'_1[\pi]$ ; since  $P_1 = P'_1[\pi]$  (from

(Definition 10), we can apply rule  $rSInv$  to  $(v_1, P_1, C_1, D_1)$ , yielding a state  $(v_2, P_2, C_2, D_2)$ . From Appendix A, we know that  $C_2 = (x = \lambda_S \wedge (C_1^\dagger))$  and  $C_2' = (x = \lambda_S \wedge (C_1'^\dagger))$ . Now, we have:

$$\begin{aligned}
C_2 &= (x = \pi(\lambda_S) \wedge (C_1^\dagger)) \\
&= (x = \pi(\lambda_S) \wedge (C_1'[\pi])^\dagger) && \text{(induction hypothesis)} \\
&= (x = \pi(\lambda_S) \wedge (C_1' \wedge \bigwedge_{\lambda_i \in \Lambda} \lambda_i = \pi_i)^\dagger) && \text{(definition of valuation)} \\
&= (x = \pi(\lambda_S) \wedge (C_1'^\dagger) \wedge \bigwedge_{\lambda_i \in \Lambda} \lambda_i = \pi_i) && \text{(property of time elapsing)} \\
&= (x = \lambda_S \wedge (C_1'^\dagger) \wedge \bigwedge_{\lambda_i \in \Lambda} \lambda_i = \pi_i) && \text{(definition of valuation)} \\
&= C_2' \wedge \bigwedge_{\lambda_i \in \Lambda} \lambda_i = \pi_i && \text{(definition of } C_2') \\
&= C_2'[\pi] && \text{(definition of valuation)}
\end{aligned}$$

Note that adding  $x = \lambda_S$  while keeping satisfiability of the expression is only true because  $\pi \models C_2' \downarrow \Lambda$ . This implies that  $C_2'[\pi]$  is non-empty, hence the state  $(v_2, P_2, C_2, D_2)$  is a valid state. In addition, since  $P_1 = P_1'[\pi]$  and rule  $rSInv$  derives to  $Stop$ , then  $P_2 = P_2'[\pi]$ . Similarly, variables are updated in the same manner on both sides, hence  $v_2 = v_2'$ . The proof is similar for other rules in Appendix A.

The proof of uniqueness is identical to that of Proposition 1.

Propositions 1 and 2 give the following theorem.

**Theorem 1 (reachability condition)** *Let CS be a parametric service model, and let  $\pi$  be a parameter valuation. Let  $\rho$  be a run of  $LTS_{CS}$  ending in a state  $(v_n, P_n, C_n, D_n)$ .*

*There exists a run of  $LTS_{CS[\pi]}$  equivalent to  $\rho$  iff  $\pi \models C_n \downarrow \Lambda$ .*

## 5 Synthesizing the static LTC

Given  $CS = (\mathcal{V}, \Lambda, P_0, C_0)$ , the *global time requirement* for CS requires that, for every state  $(v, P, C, D)$  reachable from the initial state  $(v_0, P_0, C_0, 0)$  in its LTS, the constraint  $D \leq T_G$  is satisfied, where  $T_G \in \mathbb{R}_{\geq 0}$  is the *global time constraint*. The *local time requirement* requires that if the response times of all component services of CS satisfy the *local time constraint* (LTC)  $C_L \in \mathcal{C}_\Lambda$ , then the service CS satisfies the global time requirement.

In this section, given a global time constraint  $T_G$  for a service CS, we present an approach to synthesize the static LTC (sLTC)  $C_L$ . The sLTC will be given in the form of an NNCC over  $\Lambda$ . We show that if the response times of all component services of CS satisfy the local time requirement, then the service CS will end in a good state within  $T_G$  time units.

### 5.1 Motivation

Let  $\lambda_i \in \mathbb{Q}_{\geq 0}$  be the parametric response time of component service  $S_i$  for  $i \in \{1, \dots, n\}$ , and let  $\Lambda = \{\lambda_1, \dots, \lambda_n\}$  be the set of component service parametric response times. Using constraints over  $\Lambda$ , we can represent an infinite number of possible response times symbolically. The local time requirement of component services of CS is specified as a constraint over  $\Lambda$ . An example of a local time requirement is  $(\lambda_1 \leq 6) \wedge (\lambda_2 \leq 5)$ . This local time requirement specifies that, in

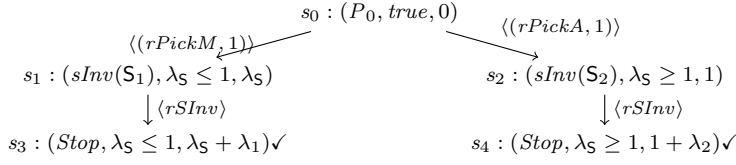


Fig. 5: LTS of composite service CS

order for CS to satisfy the global time requirement, service  $S_1$  needs to respond within 6 time units, and service  $S_2$  needs to respond within 5 time units. A local time requirement can also be in the form of a dependency between parametric response times, *e. g.*,  $(\lambda_2 \leq \lambda_1 \Rightarrow \lambda_1 + \lambda_2 \leq 6) \wedge (\lambda_1 \leq \lambda_2 \Rightarrow \lambda_1 \leq 6)$ .

In the following, we will propose a technique to synthesize the static LTC in the form of a convex over  $\Lambda$ . We first give an intuition concerning how to handle the good states (Section 5.2) and the bad states (Section 5.3); then, we give the full synthesis algorithm (Section 5.4), apply it to an example (Section 5.5) and prove its soundness (Section 5.6).

## 5.2 Addressing the good states

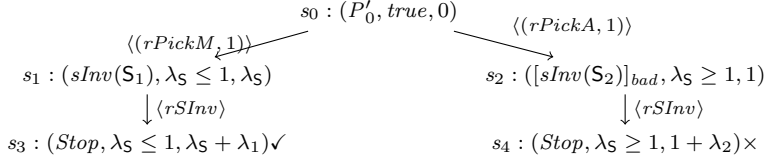
We assume a composite service CS and its LTS  $LTS_{CS} = (S, s_0, \text{Sequences}(\text{Rules}), \delta)$ ; let  $S_{good}$  be the set of all good states of  $LTS_{CS}$ . We make two observations here. First, from Theorem 1, a good state  $s_g = (v_g, P_g, C_g, D_g) \in S_{good}$  is reachable from the initial state  $s_0$  iff  $C_g$  is satisfiable. Second, whenever the good state  $s_g$  is reached, we require that the total delay from initial state  $s_0$  to state  $s_g$  must be no larger than the global time constraint  $T_G$ , *i. e.*,  $D_g \leq T_G$ . To sum up, given a good state  $s_g = (v_g, P_g, C_g, D_g)$  where  $s_g \in S_{good}$ , we require the constraint  $(C_g \downarrow \Lambda \Rightarrow (D_g \leq T_G))$  to hold. The constraint means that whenever  $s_g$  is reachable from  $s_0$ , the total (parametric) delay from  $s_0$  to  $s_g$  must be less than the global time constraint  $T_G$ . The synthesized sLTC for CS must include the conjunction of such constraints for each good state  $s_g \in S_{good}$ , that is:

$$\bigwedge_{(v_g, P_g, C_g, D_g) \in S_{good}} (C_g \downarrow \Lambda \Rightarrow (D_g \leq T_G)).$$

*Example 8* Let us consider a composite service CS whose process component is  $P_0 = pick(S \Rightarrow sInv(S_1), alrm(1) \Rightarrow sInv(S_2))$ , where  $S$  is a component service. Assume that  $sInv(S_j)$  is a component service with parametric response time  $\lambda_j$ , for  $j \in \{1, 2\}$ , and  $S$  has a response time  $\lambda_S$ . Suppose the global time requirement of the composite service CS is to respond within five seconds. Fig. 5 shows the LTS of CS.

For composite service CS in Fig. 5, we have two good states (states  $s_3$  and  $s_4$ ), and the synthesized local time requirement for composite service CS is:

$$(\lambda_S \leq 1) \Rightarrow (\lambda_S + \lambda_1 \leq 5) \wedge (\lambda_S \geq 1) \Rightarrow (1 + \lambda_2 \leq 5)$$

Fig. 6: LTS of composite service  $CS'$ **Algorithm 1:** synthSLTC(CS)

**input** : Composite service model CS with LTS  $LTS_{CS}$  of initial state  $s_0$   
**output**: The sLTC  $C_L \in \mathcal{NC}_A$

1 **return** synthRec( $s_0$ );

## 5.3 Addressing the bad states

Another goal we want to achieve is to avoid all bad states in  $LTS_{CS}$ . Let  $S_{bad}$  be the set of all bad states of service  $LTS_{CS}$ . Given a bad state  $s_b = (v_b, P_b, C_b, D_b) \in S_{bad}$ , this bad state must not be reachable from the initial state  $s_0$ . Hence, in order to prevent  $C_b$  to be satisfiable, we require that the parameters be taken in the negation of the projection of  $C_b$  onto  $A$ , *i. e.*, we require that  $\neg(C_b) \downarrow_A$  be satisfiable because of the reachability condition ([Theorem 1](#)). In addition to the good state constraint given in [Section 5.2](#), the synthesized sLTC for CS must also include the conjunction of such constraints for each bad state  $s_b \in S_{bad}$ , that is:

$$\bigwedge_{(v_b, P_b, C_b, D_b) \in S_{bad}} (\neg(C_b) \downarrow_A).$$

*Example 9* Consider a variant  $CS'$  of [Example 8](#), where  $sInv(S_2)$  is now treated as a bad activity, denoted by  $[sInv(S_2)]_{bad}$ . This service results in the LTS shown in [Fig. 6](#), where state  $s_4$  is a bad state. From [Theorem 1](#), a way to avoid the reachability of  $s_4$  is to negate its associated constraint  $C$ . Therefore, the local time requirement for composite service  $CS'$  is  $(s_3.C \downarrow_A \Rightarrow (s_3.D \leq T_G)) \wedge \neg(s_4.C \downarrow_A)$ : the first term guarantees the reachability of  $s_3$  while the second term guarantees the non-reachability of  $s_4$ . Therefore, this NNCC ensures that any complete run of the service ends in a good state. (This will be proved in [Section 5.6](#).)

## 5.4 Synthesis algorithms

[Algorithm 1](#) presents the entry algorithm for synthesizing the sLTC for a given service CS, by traversing the LTS of CS. [Algorithm 1](#) simply calls synthRec( $s$ ) applied to the initial state  $s_0$ ; this latter algorithm synthRec is given in [Algorithm 2](#).

Given a state  $s = (v, P, C, D)$  in the LTS of service CS, synthRec( $s$ ) returns a parameter constraint as follows. If state  $s$  is a good state ([line 1](#)), then it returns the constraint  $s.C \downarrow_A \Rightarrow (s.D \leq T_G)$  ([line 2](#)), where  $T_G$  is the given global time constraint of the service CS. If state  $s$  is a bad state ([line 3](#)), then the negation of the current constraint  $s.C \downarrow_A$  is returned ([line 4](#)). Finally, if  $s$  is a non-terminal



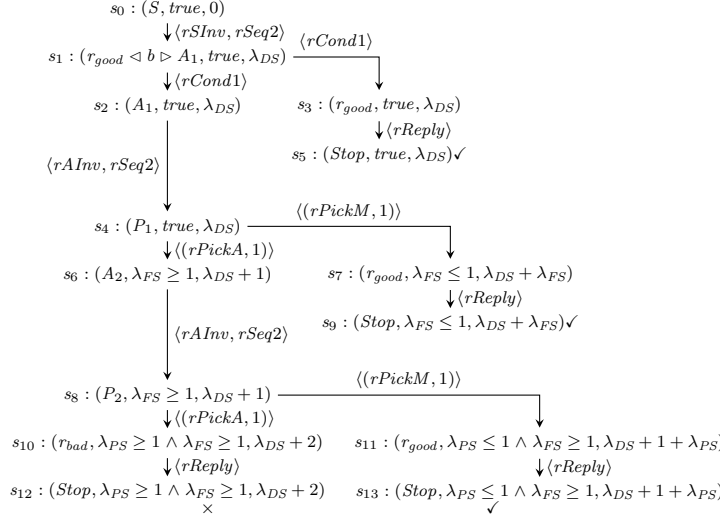
**Algorithm 2:** synthRec( $s$ )

---

**input** : State  $s$  of LTS  
**output**: The constraint for LTS that starts at  $s$

- 1 **if**  $s$  is a good state **then**
- 2 |   **return**  $(s.C \downarrow_A \Rightarrow (s.D \leq T_G))$ ;
- 3 **else if**  $s$  is a bad state **then**
- 4 |   **return**  $\neg (s.C \downarrow_A)$ ;
- 5 **else**
- 6 |   //  $s$  is a non-terminal state  
    **return**  $\bigwedge_{s' \in succ(s)} \text{synthRec}(s')$  ;

---



$S = (sInv(DS); r_{good} \triangleleft b \triangleright A_1)$   
 $A_1 = (aInv(FS); P_1)$   
 $P_1 = (pick(FS \Rightarrow r_{good}, alrm(1) \Rightarrow A_2))$   
 $A_2 = (aInv(PS); P_2)$   
 $P_2 = (pick(PS \Rightarrow r_{good}, alrm(1) \Rightarrow r_{bad}))$   
 $r_{good} = (reply(User))$   
 $r_{bad} = ([reply(User)]_{bad})$

Fig. 7: LTS of the SMIS

state (line 5), the algorithm returns the conjunction of the result of the algorithm recursively applied on the successors of  $s$  (line 6).

## 5.5 Application to the running example

Consider again the running example SMIS introduced in Section 2. Assume the parametric response times of  $FS$ ,  $PS$  and  $DS$  are  $\lambda_{FS}$ ,  $\lambda_{PS}$  and  $\lambda_{DS}$ , respectively. Recall that  $T_G = 3$ .

Fig. 7 shows the LTS of SMIS. The sLTC resulting from the application of synthSLTC is:

$$\begin{aligned} & ((\lambda_{DS} \leq 3) \wedge (\lambda_{FS} \leq 1) \Rightarrow (\lambda_{DS} + \lambda_{FS} \leq 3)) \wedge \\ & ((\lambda_{FS} \geq 1 \wedge \lambda_{PS} \leq 1) \Rightarrow (\lambda_{DS} + \lambda_{PS} \leq 2)) \wedge \neg (\lambda_{FS} \geq 1 \wedge \lambda_{PS} \geq 1) \end{aligned}$$

After simplification<sup>3</sup> using Z3 [DB08], a state-of-the-art Satisfiability Modulo Theories (SMT) solver developed by Microsoft Research, we get the following sLTC:

$$\begin{aligned} & (\lambda_{FS} < 1 \wedge \lambda_{DS} + \lambda_{FS} \leq 3) \vee \\ & (\lambda_{PS} < 1 \wedge \lambda_{FS} > 1 \wedge \lambda_{DS} + \lambda_{PS} \leq 2) \vee \\ & (\lambda_{PS} < 1 \wedge \lambda_{DS} + \lambda_{FS} \leq 3 \wedge \lambda_{DS} + \lambda_{PS} \leq 2) \end{aligned}$$

This result provides us useful information on how the component services collectively satisfy the global time constraint. That is useful when selecting component services. For the case of SMIS, one way to fulfill the global time requirement of SMIS is to select component service  $FS$  with response time that is less than 1 second, and component services  $DS$  and  $FS$  where the summation of their response times should be less than or equal to 3 seconds. For example, a suitable valuation is  $\pi$  such that  $\pi(\lambda_{FS}) = 0.5$ ,  $\pi(\lambda_{DS}) = 1.5$  and  $\pi(\lambda_{PS}) = 0.8$ .

## 5.6 Termination and soundness of synthSLTC

### 5.6.1 Termination

**Lemma 1** *Let CS be a service model. Then  $LTS_{CS}$  is acyclic and finite.*

*Proof* From [Assumption 1](#) and from the fact that there are no recursive activities in BPEL.

**Proposition 3** *Let CS be a service model. Then synthSLTC(CS) terminates.*

*Proof* From [Lemma 1](#),  $LTS_{CS}$  is acyclic. [Algorithm 1](#) is obviously non-recursive. Now, [Algorithm 2](#) is recursive ([line 6](#)). However, due to the acyclic nature of  $LTS_{CS}$  and the fact that [Algorithm 2](#) is called recursively on the *successors* of the current state, then no state is explored more than once. This ensures termination.

*Remark 1 (Complexity of Algorithm 2)* First, note that all states of  $LTS_{CS}$  are explored by [Algorithm 2](#): indeed, the algorithm is recursively called on non-terminal states, and stops only on terminal states—that have no successors anyway. So, the algorithm time is constant in the number of states of  $LTS_{CS}$ . In addition, the number of conjuncts in the result of [Algorithm 2](#) is at most the number of states of  $LTS_{CS}$ , and less if not all states are terminal states.

<sup>3</sup> For readability, we give the constraint as output in disjunctive normal form (DNF), instead of the usual conjunctive normal form (CNF).

### 5.6.2 Soundness

Let us prove that for any parameter valuation satisfying the output of `synthSLTC`, any complete run ends in a good state, and all reachable good states are reachable within the global delay  $T_G$ .

In the following, given a run  $\rho_\pi$  of  $LTS_{CS[\pi]}$ , from [Proposition 1](#) we can safely refer to the run of  $LTS_{CS}$  equivalent to  $\rho_\pi$ .

The following lemmas will be used to prove the subsequent [Theorem 2](#).

**Lemma 2** *Let  $CS$  be a service model. Let  $\pi \models \text{synthSLTC}(CS)$ . Then no bad state is reachable in  $LTS_{CS[\pi]}$ .*

*Proof* Let  $K = \text{synthSLTC}(CS)$ .  $K$  is a conjunction of “good” parameter constraints (accumulated from [line 2](#) in [Algorithm 2](#)) and “bad” parameter constraints (accumulated from [line 4](#) in [Algorithm 2](#)). Hence,  $K$  contains at least the negated constraints of all bad states. Hence, from [Theorem 1](#), the bad states are unreachable for any  $\pi \models K$ .

**Lemma 3** *Let  $CS$  be a service model. Let  $\pi \models \text{synthSLTC}(CS)$ . Then any complete run of  $LTS_{CS[\pi]}$  ends in a good state.*

*Proof* First, note that the initial state  $s_0$  is reachable in  $LTS_{CS[\pi]}$  (since  $s_0.C = \text{true}$ ). If the initial state is the only state, then from [Lemma 2](#), it is also not a bad state; hence it is a good state. Now, if it is not the only state, from the fact that all runs of  $LTS_{CS}$  end either in a good state or in a bad state, from the absence of bad states ([Lemma 2](#)), and from [Theorem 1](#), then any run of  $LTS_{CS[\pi]}$  ends in a good state.

**Lemma 4** *Let  $CS$  be a service model. Let  $\pi \models \text{synthSLTC}(CS)$ . Then for all good state  $(v, P_g, C, d)$  of  $LTS_{CS[\pi]}$ ,  $d \leq T_G$ .*

*Proof* Let  $s_g = (v, P_g, C, D)$  be a reachable state in  $LTS_{CS}$  such that  $s_g$  is a good state. From [Definition 9](#),  $C$  is satisfiable (and hence  $C \downarrow_A$  too). Since  $s_g$  is a good state, [Algorithm synthRec](#) added a constraint  $C \downarrow_A \Rightarrow D \leq T_G$  to the result. Hence,  $\text{synthSLTC}(CS) \subseteq (C \downarrow_A \Rightarrow D \leq T_G)$ . Now, for any  $\pi \models \text{synthSLTC}(CS)$ , we have that  $\pi \models (C \downarrow_A \Rightarrow D \leq T_G)$ , and hence all reachable states in  $LTS_{CS[\pi]}$  are such that  $d \leq T_G$ .

We can now formally state the soundness of `synthSLTC`.

**Theorem 2** *Let  $CS$  be a service model. Let  $\pi \models \text{synthSLTC}(CS)$ . Then:*

1. *Any complete run of  $LTS_{CS[\pi]}$  ends in a good state.*
2. *For all good state  $(v, P_g, C, d)$  of  $LTS_{CS[\pi]}$ ,  $d \leq T_G$ .*

*Proof* From [Lemmas 3](#) and [4](#).

Given a composite service  $CS$ , and assume  $S_g = \{s_1, \dots, s_n\}$  be the set of all good states in  $LTS_{CS}$ . In the following proposition, we show that any  $\pi \models \text{synthSLTC}(CS)$  necessarily satisfies (at least) one of the good states’ constraints, *i. e.*,  $\pi \models s_i.C \downarrow_A$  for some  $s_i \in S_g$ .

Indeed, recall  $\text{synthSLTC}(CS)$  is a conjunction of good and bad constraints. In the following proposition, we show that the good constraints of the form  $(C_1 \Rightarrow r_1 \wedge \dots \wedge C_n \Rightarrow r_n)$  will not hold trivially by just having  $C_i = \text{false}$ , for all  $i \in \{1, \dots, n\}$ .

**Proposition 4** *Let CS be a service model, and  $S_{good}$  be the set of all good states in  $LTS_{CS}$ . Let  $\pi \models \text{synthSLTC}(\text{CS})$ .*

*Then  $\exists s \in S_{good} : \pi \models s.C \downarrow_A$ .*

*Proof* From Algorithm 2,  $\text{synthSLTC}(\text{CS})$  is a conjunction of “good” constraints (accumulated at line 2) and “bad” constraints (accumulated at line 4). That is, assume  $\text{synthSLTC}(\text{CS}) = (C_g \wedge C_b)$ , where  $C_g = \bigwedge_{s_i \in S_{good}} (s_i.C \downarrow_A \Rightarrow (s_i.D \leq T_G))$ , and  $T_G$  be the global time constraint, and  $C_b = \bigwedge_{s_j \in S_{bad}} \neg(s_j.C \downarrow_A)$ . Hence, since  $\pi \models \text{synthSLTC}(\text{CS})$  then  $\pi \models C_g$ , hence  $\exists s \in S_{good} : \pi \models s.C \downarrow_A$ .

## 5.7 Incompleteness of $\text{synthSLTC}$

A limitation of  $\text{synthSLTC}$  is that it is incomplete, *i. e.*, it does not include all parameter valuations that could give a solution to the problem of the local time requirement. Given an expression  $A \triangleleft a = 1 \triangleright B$ , since  $a$  may be unknown at design time, we explore both branches (activities  $A$  and  $B$ ) for synthesizing the sLTC. Nevertheless, only exactly one of these activities will be executed at runtime. Including constraints from activities  $A$  and  $B$  will make the constraints stricter than necessary; therefore some of the feasible parameter valuations are excluded—this makes the synthesis by  $\text{synthSLTC}$  incomplete. This can be seen as a trade-off to make the synthesized local time requirement more general, *i. e.*, to hold in any composite service instance. In Section 6, we will introduce a method that leverages on runtime information to mitigate this problem.

## 6 Runtime refinement of local time requirement

In order to improve the local time requirement computed statically using the algorithms presented in Section 5, we introduce in this section a *refined* local time requirement, together with its usage for runtime adaptation of a service composition.

### 6.1 Motivation

Let us consider a composite service CS. Assume that we have selected a set of component services such that their stipulated response times fulfill the sLTC of CS. Since the composite service is executed under a highly evolving dynamic environment, the design time assumptions may evolve at runtime. For instance, the response times of component services could be affected by network congestion. This might result in the non-conformance of stipulated response times for some component services. However, the non-conformance of stipulated response times of component services does not necessary imply that the composite service will not satisfy its global time requirement. This is because the sLTC is synthesized at the design time to hold in *any* execution trace of CS; whereas at runtime, the runtime information can be used to synthesize a more relaxed constraint for CS.

More specifically, given a composite service CS, we have two pieces of runtime information that may help to synthesize a more relaxed constraint: the execution

path that has been taken by CS, and the elapsed time of CS. First, the execution path taken by CS can be used for LTS simplification. This is because in the midst of execution, some of the execution traces can be disregarded and therefore a weaker LTC, that includes more parameter valuations, may be synthesized. Second, the time elapsed of CS can be used to instantiate some of the response time parameters with real-valued constants; this makes the synthesized LTC contain less uncertainty and be more precise.

*Example 10* For example, consider the SMIS composite service, the LTS of which is depicted in Fig. 7. Assume a valuation  $\pi$  satisfying the *sLTC*. At runtime, after invocation of the component service *DS*, SMIS will be at state  $s_2$ . Assume that *DS* does not conform to its stipulated response time. Therefore, it is desirable to check whether invoking *FS* can still satisfy the global time requirement of CS. One can make use of *sLTC* for this purpose. Nevertheless, a more precise LTC may be synthesized at state  $s_2$ .

The first observation is that, from state  $s_2$ , we can safely ignore the constraints from the good state  $s_5$ , since it is not reachable from  $s_2$ . The second observation is that the delay from state  $s_0$  to state  $s_2$  (say  $r$  time units, with  $r \in \mathbb{R}_{\geq 0}$ ) is known. For this reason, we can substitute the delay component of state  $s_2$ , which is the response time  $\pi(\lambda_{DS})$ , with the actual time delay  $r$ . This motivates the use of runtime information of the composite service to refine the LTC. We refer to the runtime refined LTC as the runtime LTC (denoted by *rLTC*). In addition to this refinement, we can also simplify the LTS by pruning the states corresponding to past states (e.g.,  $s_0, s_1$  in Fig. 7), as well as the successors of these past states that were not met in practice (e.g.,  $s_3$  and  $s_5$  in Fig. 7), because another branch was taken at runtime. We show the LTS of SMIS before and after simplification in Figs. 8a and 8b respectively.

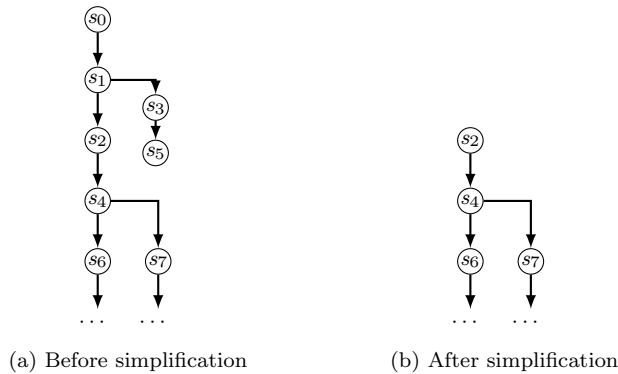


Fig. 8: LTS Simplification of SMIS

By incorporating the runtime information, the resulting *rLTC* at state  $s_2$  is:

$$\boxed{\begin{aligned} & ((\lambda_{FS} \leq 1) \Rightarrow (r + \lambda_{FS} \leq 3)) \wedge \\ & ((\lambda_{FS} \geq 1 \wedge \lambda_{PS} \leq 1) \Rightarrow (r + \lambda_{PS} \leq 2)) \wedge \\ & \neg (\lambda_{FS} \geq 1 \wedge \lambda_{PS} \geq 1) \end{aligned}}$$

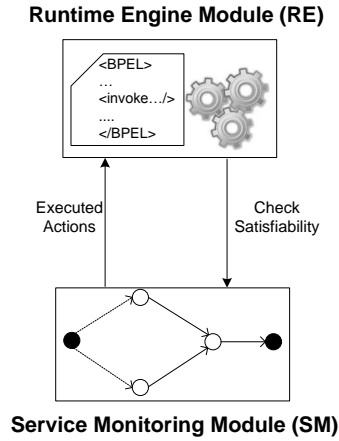


Fig. 9: Service adaptation framework

## 6.2 Runtime adaptation of a BPEL process

We now introduce a service adaptation framework to improve the conformance of global time requirement for a composite service. The architecture of the framework is shown in Fig. 9. There are two modules in the framework—the Runtime Engine Module (*RE*) and the Service Monitoring Module (*SM*). *RE* provides an environment for the execution of a BPEL service; here, we use ApacheODE [Fou07], an open source BPEL engine. We instrument the runtime component of Apache ODE to communicate with the service monitoring module.

*SM* is used to monitor the execution of a BPEL service. During the deployment of a service *CS*, *SM* generates the LTS of *CS* and stores it in the cache of *SM* so that it is available when *CS* is executing.

During the execution of the composite service *CS*, the sequences of rules from *RE* are used to update the active state  $s_a \in S$  of LTS stored in *SM*. The sequence of rules is also stored as part of the current execution run. *SM* also keeps track of the total execution time for this execution run, as well as the response time for each component service invocation.

Prior to the invocation of a component service *S*, *RE* will consult *SM* to check the satisfiability of rLTC. If the rLTC of  $s_a$  is satisfiable, then *SM* will instruct *RE* to continue invoking *S* as usual. Otherwise, some kind of mitigation procedure may be triggered. One of the possible mitigation procedures is to invoke a backup service of *S*,  $S_{bak}$ , which has a faster stipulated response time than *S* (that may come with a cost).

*Example 11* Consider again the running example SMIS in Section 2. An example of *S* and  $S_{bak}$ , are services *FS* and *PS* respectively.

In the following, we introduce the details on the synthesis of rLTC (Section 6.3) and satisfiability checking (Section 6.4).

### 6.3 Algorithm for runtime refinement

A way to calculate the rLTC could be to run `synthSLTC` (Algorithm 2) from a state  $s$  in the LTS. However, this requires traversing the state-space repeatedly for every calculation of the rLTC. To make it more efficient, we extend `synthSLTC` by calculating the rLTC for each state  $s$  during the synthesis of the LTC at the design time. Therefore, at runtime, we only need to retrieve the synthesized rLTC of the corresponding state for direct usage.

`synthRLTC` (given in Algorithm 3) synthesizes the rLTC for each state in the LTS. Before explaining the algorithm, let us introduce a few notations used in Algorithm 3. First, we assume that states in the LTS of CS are augmented with an additional “field” to store the computed rLTC. We use  $s.rLTC$  to denote the rLTC associated with state  $s$ . Additionally, we use the following shorthand to perform a conjunction of pairs of parametric constraints  $(cons_i.g, cons_i.b)$  such that the resulting pair is such that its left-hand (resp. right-hand) side is the conjunction of all left-hand (resp. right-hand) sides:  $\prod((cons_1.g, cons_1.b), \dots, (cons_n.g, cons_n.b))$  denotes  $((cons_n.g \wedge \dots \wedge cons_1.g), (cons_n.b \wedge \dots \wedge cons_1.b))$ .

Given a composite service CS together with its associated LTS, and a state in  $LTS_{CS}$ , `synthRLTC` returns a constraint pair  $C_s = (g, b)$ , where  $g, b \in \mathcal{C}_A$ . In this pair,  $g$  (resp.  $b$ ) denotes the constraint associated to a good (resp. bad) state. Given a constraint pair  $C_s$ , we use  $C_s.g$  (resp.  $C_s.b$ ) to refer to the first (resp. second) component of  $C_s$ . Variables  $d_f$  and  $r_f$  are *free variables*, which are variables to be substituted at runtime. In particular, given a state  $s$ , free variables  $d_f$  and  $r_f$  in  $s.rLTC$  are to be substituted by the delay component  $s.D \in \mathcal{L}_A$  and the actual delay  $r \in \mathbb{R}_{\geq 0}$  from the initial state to the state  $s$  respectively.

---

#### Algorithm 3: `synthRLTC(CS, LTSCS, s)`

---

```

input : Composite service CS
input : LTS  $LTS_{CS}$  of CS
input : State  $s$  in LTS of CS
output: Constraint pair for sub-LTS of CS starting with  $s$ 
1  $cons \leftarrow \emptyset$ ;
2 if  $s$  is a good state then
3    $cons \leftarrow (s.C \downarrow_A \Rightarrow (s.D - d_f + r_f \leq T_G), true)$ ;
4    $s.rLTC \leftarrow cons.g \wedge (d_f = s.D)$ ;
5 else if  $s$  is a bad state then
6    $cons \leftarrow (true, \neg (s.C \downarrow_A))$ ;
7    $s.rLTC \leftarrow cons.b$ ;
8 else
9   //  $s$  is a non-terminal state
9    $cons \leftarrow \prod_{s' \in succ(s)} synthRLTC(s')$ ;
10   $s.rLTC \leftarrow cons.g \wedge cons.b \wedge (d_f = s.D)$ ;
11 return  $cons$ ;

```

---

Let us now explain `synthRLTC` in details. Given a good state  $s$  (line 2),  $s.rLTC$  is assigned with value  $cons.g$ , with free variable  $d_f$  substituted with  $s.D$  (line 4); note that substitution is here achieved using conjunction of the constraint with the equality  $d_f = s.D$ . As an illustration, consider the good state  $s_{13}$  in the SMIS

example (the LTS of which is given in Fig. 7). At runtime, assume the active state is at state  $s_{13}$ , and assume that it takes  $r \in \mathbb{R}_{\geq 0}$  time units to execute from the initial state  $s_0$  to state  $s_{13}$ . Therefore, the previously unknown parametric response time in the delay component of state  $s_{13}$ , *i. e.*,  $t_{DS} + 1 + t_{PS}$ , can be substituted with the real value  $r$ . To achieve this, at line 3, we subtract away the free variable  $d_f$ , which is to be substituted with the response time parameter of state  $s_{13}$ , and add back the free variable  $r_f$ , which is to be substituted with the real value  $r$ . We substitute the free variable  $d_f$  at line 4. For free variable  $r_f$ , it is only substituted in Algorithm 4 at runtime when the delay is known. In the case of the SMIS example, the  $rLTC$  of state  $s_{13}$  after substituting free variable  $r_f$  with value  $r$  (*i. e.*,  $s_{13}.rLTC \wedge (r_f = r)$ ) is  $((t_{PS} \leq 1 \wedge t_{FS} \geq 2) \Rightarrow (r \leq 3))$ .

When  $s$  is a bad state (lines 5 to 7), we simply compute the negation of the associated constraint so as to keep the system reaching this bad state (just as in Algorithm 2).

When  $s$  is a non-terminal state (line 8),  $s.rLTC$  is assigned with the conjunction of all good and bad constraints computed by recursively calling `synthRLTC` on the successor states of  $s$ , where free variable  $d_f$  is substituted with  $s.D$  (line 10).

## 6.4 Satisfiability checking

We now introduce a satisfiability checking algorithm. This satisfiability checking is done before the invocation of a component service. Suppose that, before the invocation of a component service  $S_i$ , CS is at the active state  $s_a$ . The satisfiability of the  $rLTC$  at  $s_a$  will be checked before  $S_i$  is invoked. If it is satisfiable, then it will invoke  $S_i$  as usual. Otherwise, some mitigation procedures will be triggered. A mitigation procedure could consist of invoking a faster backup service  $S'_i$  instead of  $S_i$ .

---

### Algorithm 4: `checkSat(LTSCS, $s_a$ , $r$ , $\Lambda$ , $\pi$ )`

---

**input** : LTS of the parametric composite service CS, Active state  $s_a \in S$ , Elapsed time  $r \in \mathbb{R}_{\geq 0}$ , Set of parametric response times  $\Lambda$ , Parameter valuation  $\pi$   
**output**: True if the local time constraint at  $s_a$  is satisfiable, false otherwise  
**1 return**  $Is\_Sat((\bigwedge_{1 \leq i \leq n} \lambda_i \leq \pi(\lambda_i)) \Rightarrow (s_a.rLTC \wedge (r_f = r)))$ ;

---

We give in Algorithm 4 the algorithm checking the satisfiability of  $rLTC$  at state  $s_a \in Q$ . With the assumption that all component services will reply within their stipulated response times ( $\bigwedge_{1 \leq i \leq n} \lambda_i \leq \pi(\lambda_i)$ ), `checkSat` checks whether the  $rLTC$  at state  $s_a$  can be satisfied with free variables  $r_f$  substituted with the actual elapsed time  $r \in \mathbb{R}_{\geq 0}$ . The function `Is_Sat` returns true iff the input constraint is satisfiable.

## 6.5 Termination and soundness of `synthRLTC`

### 6.5.1 Termination

**Proposition 5** *Let CS be a service model,  $s$  be a state in  $LTS_{CS}$ .*



Then  $\text{synthRLTC}(\text{CS}, \text{LTS}_{\text{CS}}, s)$  terminates.

*Proof* Observe that [Algorithm 3](#) is recursive (on [line 9](#)). However, due to the acyclic nature of  $\text{LTS}_{\text{CS}}$  (from [Lemma 1](#)) and the fact that [Algorithm 3](#) is called recursively on the successors of the current state, then no state is explored more than once. This ensures termination.

### 6.5.2 Soundness

[Theorem 3](#) formally states the correctness of our runtime refinement algorithm.

**Theorem 3** *Let CS be a service model. Let  $\text{LTS}_{\text{CS}}$  be the LTS of CS. Let  $s$  be the current state in  $\text{LTS}_{\text{CS}}$  and  $r$  be the current elapsed time.*

*Fix  $\pi \models \text{synthRLTC}(\text{CS}, \text{LTS}_{\text{CS}}, s)$ . Then:*

1. *there exists a run in  $\text{LTS}_{\text{CS}[\pi]}$  ending in some state  $s_\pi$  such that this run is equivalent to a run of  $\text{LTS}_{\text{CS}}$  ending in  $s$ ;*
2. *any complete run of the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$  ends in a good state;*
3. *for all good states  $(v, P_g, C, d)$  in the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$ , then  $d \leq T_G$ .*

*Proof* 1. From [Proposition 2](#).

2. From [Definition 5](#), the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$  contains the successors of  $s_\pi$  in  $\text{LTS}_{\text{CS}[\pi]}$ , and hence any complete run of the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$  corresponds to the end of some complete run of  $\text{LTS}_{\text{CS}[\pi]}$ . From [Lemma 3](#), any complete run of  $\text{LTS}_{\text{CS}[\pi]}$  ends in a good state, which gives the result.
3. Any good state of the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$  is also a good state of  $\text{LTS}_{\text{CS}[\pi]}$ . From [Lemma 4](#), for all good state of  $\text{LTS}_{\text{CS}[\pi]}$ ,  $d \leq T_G$ , which gives the result.

*Remark 2 (Complexity of Algorithm 3)* First, note that all states of  $\text{LTS}_{\text{CS}}$  are explored by [Algorithm 3](#): indeed, the algorithm is recursively called on non-terminal states, and stops only on terminal states—that have no successors anyway. So, the algorithm time is constant in the number of states of the sub-LTS of  $\text{LTS}_{\text{CS}[\pi]}$  starting from  $s_\pi$ .

Let us now investigate the worst-case number of conjuncts in the result of [Algorithm 3](#). The algorithm returns the good conjuncts ( $\text{cons.g}$ ), the bad conjuncts ( $\text{cons.b}$ ) and a last term (“ $d_f = s.D$ ”) ([line 10](#) in [Algorithm 3](#)). Any good terminal state or bad terminal state adds exactly one conjunct to either  $\text{cons.g}$  or  $\text{cons.b}$ . Therefore, the number of conjuncts is exactly the number of terminal states, plus one due to the last term.

## 6.6 Discussion

*Termination* From [Proposition 5](#), our method terminates due to the fact that BPEL composite services do not support recursion, and [Assumption 1](#) on the loop activities ensuring that the upper bound on the number of iterations and the time of execution are known. We briefly discuss how to enforce this assumption

in the presence of loops in the composite service. The upper bound on the number of iterations could be either inferred by using loop bound analysis tool (*e. g.*, [Erm+07]), or could be provided by the user otherwise. In the worst case, an alternative option is to set up a bound arbitrary but “large enough”. Concerning the maximum time of loop executions, it could be enforced by using proper timeout mechanism in BPEL.

*Time for internal operations* For simplicity, we do not account for the time taken for the internal operations of the system. In reality, the time taken by the internal operations may become significant, especially when the process is large. We can provide a more accurate synthesis of the constraints by including an additional constraint  $t_{overhead} \leq b$ , where  $t_{overhead} \in \mathbb{R}_{\geq 0}$  is a time overhead for an internal operation, and  $b \in \mathbb{R}_{\geq 0}$  is a machine dependent upper bound for  $t_{overhead}$ . The method to obtain an estimation of  $b$  is beyond the scope of this work; interested readers may refer to, *e. g.*, [MRD08].

## 7 Evaluation

As a proof of concept, we applied our method to several examples. After briefly presenting our implementation, we describe the examples we use (Section 7.1). We then evaluate our methods for the synthesis of local time requirement at the design time (Section 7.2) and for the runtime refinement (Section 7.3).

*Implementation* We implemented our algorithms for synthesizing the sLTC and rLTC in SELAMAT, a tool developed in C#. We use an *ad-hoc* input syntax very close to that of Definition 1. Our prototype implementation uses basic state space reduction techniques, notably zone inclusion (see, *e. g.*, [HSW16; NPV18] for recent such techniques in the (parametric) timed setting), to prune whole branches of the state space. The front-end GUI relies on the PAT model checker [Sun+09].

The simplification of the final results of sLTC and rLTC is achieved using Microsoft Z3 [DB08]. For the runtime adaptation, we use Apache ODE 1.3.6 as runtime engine module (*RE*). The service monitoring module (*SM*) is developed in C#, which uses Microsoft Z3 for the satisfiability checking. The tool and examples can be downloaded at [Tan+19].

### 7.1 Examples

#### 7.1.1 Stock market indices service (*SMIS*)

This is the running example introduced in Section 2.

#### 7.1.2 Computer purchasing services (*CPS*)

The goal of a CPS is to allow a user to purchase a computer system online using credit cards. Our CPS makes use of five component services, namely Shipping Service (*SS*), Logistic Service (*LS*), Inventory Service (*IS*), Manufacture Service

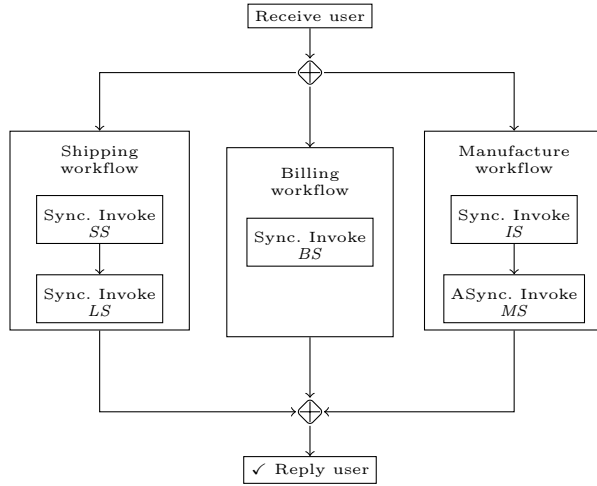


Fig. 10: Computer Purchasing Service (CPS)

(*MS*), and Billing Service (*BS*). The global time requirement of the CPS is to respond within three seconds. The CPS workflow is shown in Fig. 10. The CPS starts upon receiving the purchase request from the client with credit card information, and the CPS spawns three workflows (*viz.*, shipping workflow, billing workflow, and manufacture workflow) concurrently. In the shipping workflow, the shipping service provider is invoked synchronously for the shipping service on computer systems. Upon receiving the reply, *LS* (which is a service provided by the internal logistic department) is invoked synchronously to record the shipping schedule. In the billing workflow, the billing service (which is offered by a third party merchant) is invoked synchronously for billing the customer with credit card information. In the manufacture workflow, *IS* is invoked synchronously to check for the availability of the goods. Subsequently, *MS* is invoked asynchronously to update the manufacture department regarding the current inventory stock. Upon receiving the reply message from *LS* and *BS*, the result of the computer purchasing will be returned to the user.

### 7.1.3 Travel booking service (TBS)

The goal of a travel booking service (TBS) is to provide a combined flight and hotel booking service by integrating two independent existing services. TBS provides an SLA for its subscribed users, saying that it must respond within five seconds upon request. The travel booking system has four component services, namely Flight Service (*FS*), Backup Flight Service (*FS<sub>bak</sub>*), Hotel Service (*HS*) and Backup Hotel Service (*HS<sub>bak</sub>*). The TBS workflow is given in Fig. 11. Upon receiving the request from users, the variable *res* is assigned to true. After that, TBS spawns two workflows (*viz.*, a flight request workflow, and a hotel request workflow) concurrently. In the flight request workflow, it starts by invoking *FS*, which is a service provided by a flight service booking agent. If service *FS* does not respond within two seconds, then *FS* is abandoned, and another backup flight

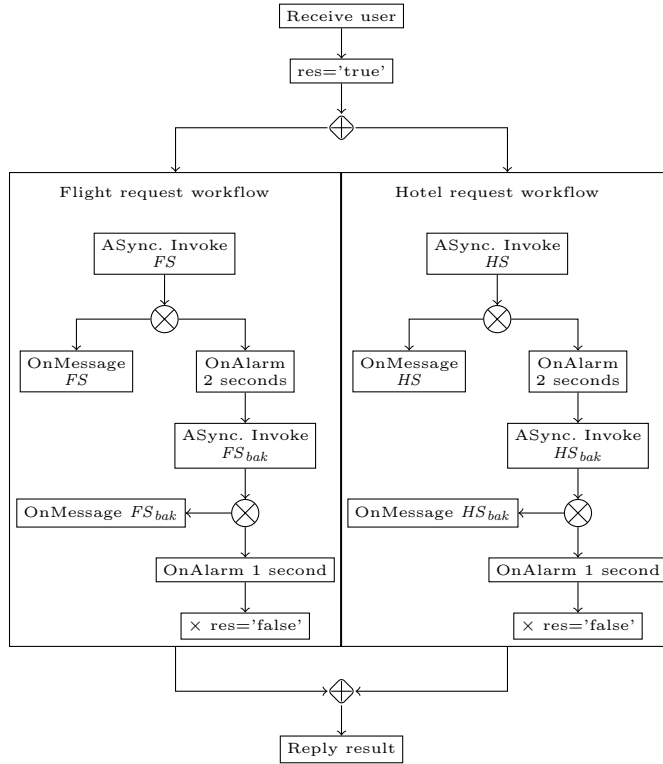


Fig. 11: Travel Booking Service (TBS)

service  $FS_{bak}$  is invoked. If  $FS_{bak}$  returns within one second, then the workflow is completed; otherwise the variable  $res$  is assigned to false. The hotel request workflow shares the same process as the flight request workflow, by replacing  $FS$  with  $HS$  and  $FS_{bak}$  with  $HS_{bak}$ . The booking result will be replied to the user if  $res$  is true; otherwise, the user will be informed of the booking failure.

#### 7.1.4 Rescue team service (RS)

The goal of a Rescue Team service (RS) is to identify the place, weather, and nearest rescue team, by the longitude and latitude on Earth. RS makes use of three component services, namely Terra Service ( $TS$ ), Weather Service ( $WS$ ) and Distance Service ( $DS$ ). The global requirement of the RS is to respond within five seconds. The RS workflow is given in Fig. 12. RS starts upon receiving longitude and latitude coordinates from the user. After that, it invokes Terra Service ( $TS$ ), Weather Service ( $WS$ ), and Distance Service ( $DS$ ) concurrently. Service  $TS$  (resp.  $WS$ ) will return the name of the place (resp. the weather information) that corresponds to the longitude and latitude.  $DS$  is used to calculate the distance between each rescue team and the event location. In particular,  $DS_{com}$  and  $DS_{sea}$  are used to calculate the distance between commander team and sergeant team to the event location. If the distance to the event of the commander team ( $d_{com}$ )

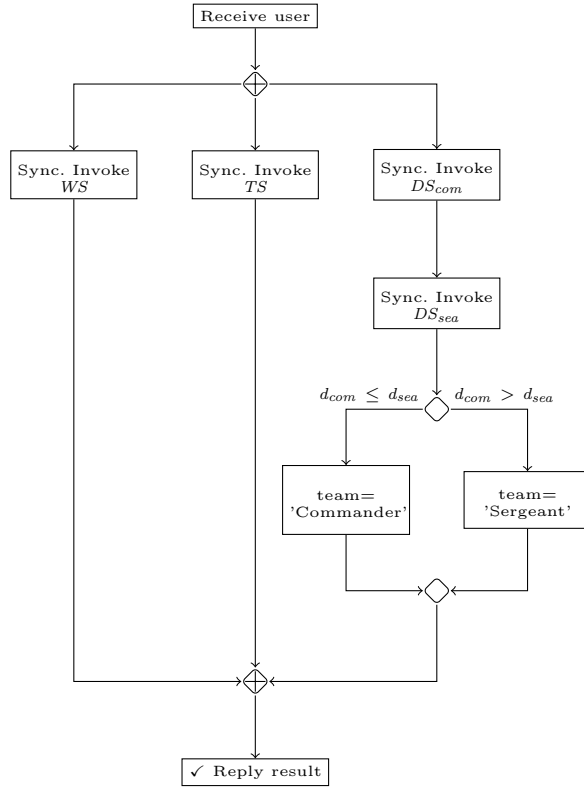


Fig. 12: Rescue Team Service (RS)

is not larger than the distance to the event of the sergeant team ( $d_{sea}$ ), then the commander team will be chosen. Otherwise, the sergeant team will be chosen. Subsequently, the place, weather and rescue team information is returned to the user.

## 7.2 Synthesis of local time requirement

### 7.2.1 Environment of the experiments

We run our algorithms to synthesize the sLTC and rLTC for the four examples on a computer with Intel Core I5 2410M CPU with 4 GiB RAM.

### 7.2.2 Evaluation results

The details of the synthesis are shown in Table 1. The **#states** and **#transitions** columns provide the information of number of states and transitions of the LTS, respectively. We repeated all experiments 30 times; we report here the average time for each experiment. The **sLTC** and **rLTC** columns provide the average time (in seconds) spent for synthesizing sLTC (for the entire LTS), and rLTC (for each

Example	#states	#transitions	sLTC (s)	rLTC (s)
SMIS	14	13	0.0076	0.0078
TBS	683	3677	1.8501	1.9000
CPS	120	119	0.0529	0.0559
RS	85	134	0.0701	0.0733

Table 1: Synthesis of sLTC and rLTC

$$(\lambda_{SS} + \lambda_{LS} + \lambda_{IS} + \lambda_{BS}) \leq 3$$

(a) sLTC for CPS

$$(\lambda_{TS} + \lambda_{WS} + 2 \cdot \lambda_{DS}) \leq 5$$

(b) sLTC for RS

$$\begin{aligned} & ((2 \cdot \lambda_{HSbak} < \lambda_{FSbak}) \wedge (2 \cdot \lambda_{FSbak} < \lambda_{HSbak}) \wedge (\lambda_{HSbak} < 1) \wedge (\lambda_{FSbak} < 1)) \\ & \vee ((\lambda_{HSbak} < 1) \wedge (\lambda_{FSbak} < 1) \wedge (\lambda_{FSbak} + \lambda_{HSbak} \leq 1)) \\ & \vee ((\lambda_{HSbak} < 1) \wedge (\lambda_{FS} < 2)) \vee ((\lambda_{HS} < 2) \wedge (\lambda_{FSbak} < 1)) \vee ((\lambda_{HS} < 2) \wedge (\lambda_{FS} < 2)) \end{aligned}$$

(c) sLTC for TBS

Fig. 13: Synthesized sLTC

state in the LTS), respectively. TBS takes a longer time than SMIS, CPS, and RS for synthesizing sLTC and rLTC, as it contains a larger number of states and transitions compared to SMIS, CPS, and RS. Nevertheless, since both sLTC and rLTC are synthesized offline, the time for synthesizing the constraints (less than two seconds) for TBS is considered to be reasonable.

The synthesized sLTC for SMIS has been given in Section 5.5, while the synthesized sLTC for CPS, TBS, and RS are shown in Fig. 13. Note that  $\lambda_{MS}$  does not appear in the sLTC of CPS. The reason is that  $MS$  is invoked asynchronously without expecting a response; therefore its response time is irrelevant to the global time requirement of CPS.

The synthesized rLTC are used for runtime adaptation at runtime. We will evaluate the runtime adaptation of a composite service with rLTC in the following section.

### 7.3 Runtime adaptation

We now conduct experiments to answer the following two questions:

- Q1.** What is the *overhead* of the runtime adaptation?
- Q2.** What is the *improvement* provided by the runtime adaptation?

#### 7.3.1 Environment of the experiments

The evaluation was conducted using two different physical machines, connected by a 100 Mbit LAN. One machine is running ApacheODE [Fou07] to host the *RE* module to execute the BPEL program, configured with Intel Core I5 2410M CPU with 4 GiB RAM. The other machine hosts the *SM* module, configured with Intel I7 3520M CPU with 8 GiB RAM.

To test the composite service under controlled situation, we introduce the notion of *execution configuration*. An execution configuration defines a particular

Example	Avg. #SAT	Avg. SAT runtime (s)
SMIS	1.74	13
TBS	2.25	17
CPS	4.00	27
RS	4.00	19

Table 2: Satisfiability checking

	$p_c$	$N_{se}$	$N_e$	Improvement (%)	Avg. backup service
SMIS	0.9	9441	8976	5.18	0.127
	0.8	9211	8374	10.00	0.352
	0.7	8109	6965	16.42	0.577
	0.6	7593	6348	19.61	0.702
TBS	0.9	10000	9743	2.64	0.384
	0.8	10000	9364	6.79	0.779
	0.7	10000	8460	18.20	0.948
	0.6	10000	7700	29.87	1.05
CPS	0.9	9523	8809	8.11	1.259
	0.8	9241	7156	29.14	1.509
	0.7	8504	6108	39.23	2.014
	0.6	8430	5650	49.20	2.578
RS	0.9	8181	7271	12.52	1.787
	0.8	7201	7011	2.71	1.589
	0.7	6590	5227	26.08	1.659
	0.6	5609	4146	35.29	1.54

Table 3: Improvement of runtime conformance

execution scenario for the composite service. Formally, an execution configuration  $E$  is a tuple  $(M, R)$ , where  $M$  decides which path to choose for an  $\langle \text{if} \rangle$  activity and  $R$  is a function that maps a component service  $S_i$  to a real value  $r \in \mathbb{R}_{\geq 0}$ , which represents the response time of  $S_i$ . We discuss how an execution configuration  $E = (M, R)$  is generated.  $M$  is generated by choosing one of the branches of an  $\langle \text{if} \rangle$  activity uniformly among all possible branches.

Let  $CS$  be a composite service model, where a component service  $S_i$  of  $CS$  has a stipulated response time  $\pi(\lambda_i) \in \mathbb{Q}_{\geq 0}$ . Then  $R(S_i)$  will be assigned with a response time within the stipulated response time  $\pi(\lambda_i)$  with a probability of  $p_c \in \mathbb{Q}_{\geq 0} \cap [0, 1]$ .  $p_c$  is the *response time conformance threshold*. More specifically,  $R(S_i)$  will be assigned with a value in  $[0, \pi(\lambda_i)]$  uniformly with a probability of  $p_c$ , and assigned to a value in  $(\pi(\lambda_i), \pi(\lambda_i) + t_e]$  uniformly with a probability of  $1 - p_c$ .  $t_e \in \mathbb{R}_{\geq 0}$  is the *exceeding threshold*; and assume after  $\pi(\lambda_i) + t_e$  seconds, the component service  $S_i$  will be automatically timeout by  $RE$  to prevent an infinite delay.

Given a composite service  $CS$ , and an execution configuration  $E$ , a *run* is denoted by  $\rho(CS, AM, E)$ , where the first argument is the composite service  $CS$  that is running, the second argument  $AM \in \{rr, \emptyset\}$  is the adaptive mechanism where  $rr$  denotes the runtime adaptation, and  $\emptyset$  denotes no runtime adaptation.

### 7.3.2 Evaluation results

We conducted two experiments Exp1 and Exp2, to answer the questions Q1 and Q2, respectively. Each experiment goes through 10,000 rounds of simulations,

and an execution configuration  $E$  is generated for each round of simulation. Given a composite service  $CS$ , we assume that for each component service  $S_i$  with a stipulated response time  $\pi(\lambda_i)$ , there exists a backup service  $S'_i$ , with a stipulated response time  $\pi(\lambda_i)/2$  and a conformance threshold of 1. Suppose that before the invocation of a component service  $S_i$ ,  $CS$  is at active state  $s_a$ . The satisfiability of the rLTC at  $s_a$  will be checked (using [Algorithm 4](#)) before  $S_i$  is invoked. If it is satisfiable, then it will invoke  $S_i$  as usual. Otherwise, the faster backup service  $S'_i$  will be invoked instead, as a mitigation procedure.

*Experiment Exp1* Given a composite service  $CS$ , in order to measure the overhead, we use an execution configuration  $E = (M, Q)$  for an adaptive run  $\rho(CS, rr, E)$ , and non-adaptive run  $\rho(CS, \emptyset, E)$ . We have modified the runtime adaptation mechanism for  $rr$  so that, if the rLTC of the active state is checked to be unsatisfiable, component service  $S_i$  will still be used (instead of  $S'_i$ ). The purpose for this modification is to make  $\rho(CS, rr, E)$  and  $\rho(CS, \emptyset, E)$  invoke the same set of component services, so that we can effectively compare the overhead of  $\rho(CS, rr, E)$ .

**Results.** Suppose at round  $k$ , the times spent for  $\rho(CS, rr, E)$  and  $\rho(CS, \emptyset, E)$  are  $r_{rr}^k \in \mathbb{R}_{\geq 0}$  time units and  $r_{\emptyset}^k \in \mathbb{R}_{\geq 0}$  time units respectively. The overhead  $O_k$  at round  $k$  is the time difference between  $r_{rr}^k$  and  $r_{\emptyset}^k$ , *i. e.*,  $O_k = r_{rr}^k - r_{\emptyset}^k$ . The average overhead at round  $k$  is calculated using [Eq. \(1\)](#).

$$Avg. \text{ overhead} = \left( \sum_{i=1}^k O_i \right) / k \quad (1)$$

The main source of overhead for runtime adaptation comes from the satisfiability checking with [Algorithm 4](#). We make use of Z3 [[DB08](#)] for this purpose. Other sources of overhead include update of active state in  $SM$ , and communications between  $SM$  and  $RE$ .

The average overheads of SMIS, CPS, TBS, and RS after 10,000 rounds are 15 ms, 21 ms, 30 ms, and 23 ms respectively. The results convey to us that the additional operations involved in the runtime adaptation, including the satisfiability checking, can be done efficiently.

We further evaluate the overhead on satisfiability checking. [Table 2](#) shows the results of satisfiability checking. The average number of satisfiability checking for each round (Avg. #SAT) is calculated using [Eq. \(2\)](#) where  $N_i$  is the total number of satisfiability checking for  $i$ -th round and  $r$  is the total number of running rounds. The average time (given in milliseconds) spent on satisfiability checking for each round (Avg. SAT runtime) is calculated using [Eq. \(3\)](#), where  $T_i$  is the time spent on satisfiability checking for  $i$ -th round. [Table 2](#) shows that the satisfiability checking has contributed most of the overhead of runtime adaptation.

$$Avg. \#SAT = \left( \sum_{i=1}^r N_i \right) / r \quad (2)$$

$$Avg. \text{ SAT runtime} = \left( \sum_{i=1}^r T_i \right) / r \quad (3)$$



*Experiment Exp2* In this second experiment, we measure the improvement for the conformance of global constraints due to  $rr$ . Given a composite service  $CS$ , an execution configuration  $E$ , two runs  $\rho(CS, rr, E)$  and  $\rho(CS, \emptyset, E)$  are conducted for each round of simulation.  $N_{se}$  is the number of executions that satisfy global constraints for composite service with  $rr$ , and  $N_e$  is the number of executions that satisfy global constraints for composite service without  $rr$ , the improvement is calculated by Eq. (4).

$$Improvement = \frac{(N_{se} - N_e) * 100}{N_e} \quad (4)$$

**Results.** The experiment results can be found in Table 3. The *Improvement (%)* column provides the information of improvement (in percentage) that is calculated using Eq. (4). The *Avg. Backup Service* column provides the average number of backup service used (calculated by summing up the number of backup services used for 10,000 rounds, and divided by 10,000).

The decrement of  $p_c$  represents the undesired situation where component services have a higher chance for not conforming to their stipulated response time. This may be due to situations such as poor network conditions. For each example, the improvement provided by the runtime adaptation increases when  $p_c$  decreases. This shows that runtime adaptation improves the conformance of global time requirement. In addition, the average number of backup service used increases when  $p_c$  decreases. This shows the adaptive nature of runtime adaptation with respect to different  $p_c$ —more corrective actions are likely to perform when the chances that component services do not satisfy their stipulated response time increase.

The results in Exp1 and Exp2 have shown that the runtime adaptation has a low overhead, and improves the runtime conformance, especially when the response time conformance threshold of the component services is low.

#### 7.4 Threats to validity

Our experiments show a good efficiency of our technique for the examples we considered; these are arguably on the smaller side, but we claim that they are non-trivial enough to not be analyzable by hand, and therefore our technique proposes what we believe to be a valuable contribution.

## 8 Related work

*Model-based analysis of Web services using LTSs* Our method is related to using LTSs for model-based analysis of Web services. In [BGP11], the authors propose an approach to obtain behavioral interfaces in the form of LTSs of external services by decomposing the global interface specification. It also has been used in model checking the safety and liveness properties of BPEL services. For example, Foster *et al.* [Fos06; Fos+06] transform BPEL process into FSP [MK06], subsequently using a tool named “WS-Engineer” for checking safety and liveness properties. Simmonds *et al.* [SBC10] propose a user-guided recovery framework for Web services based on LTSs. Our work uses LTSs in synthesizing local time requirement.

*Constraint synthesis for scheduling problems* Our work shares common techniques with work for constraint synthesis for scheduling problems. The use of models such as parametric timed automata (PTAs) [AHV93] and parametric time Petri nets (PTPNs) [TLR09] for solving such problems has received recent attention. In particular, in [CPR08; Le+10; Fri+12], parametric constraints are inferred, guaranteeing the feasibility of a schedule using PTAs extended with stopwatches (see, *e.g.*, [AM02]). In [And+14], we proposed a parametric, timed extension of CSP, to which we extended the “inverse method”, a parameter synthesis algorithms preserving the discrete behavior of the system (see, *e.g.*, [AS13]). Although PTAs or PTPNs might have been used to encode (part of) the BPEL language, our work is specifically adapted and optimized for synthesizing local timing constraint in the area of service composition.

*Finding suitable quality of service* Our method is related to the finding of a suitable quality of service (QoS) for the system [YZL07]. The authors of [YZL07] propose two models for the QoS-based service composition problem: a combinatorial model and a graph model. The combinatorial model defines the problem as a multidimension multichoice 0-1 knapsack problem. The graph model defines the problem as a multiconstraint optimal path problem. A heuristic algorithm is proposed for each model: the WS-HEU algorithm for the combinatorial model and the MCSP-K algorithm for the graph model. The authors of [AP05] model the service composition problem as a mixed integer linear problem where constraints of global and local component service can be specified. The difference with our work is that, in their work, the local constraint is specified, whereas in ours, the local constraint is synthesized. An approach of decomposing the global QoS to local QoS has been proposed in [AR09]. It uses the mixed integer programming (MIP) to find optimal decomposition of QoS constraint. However, the approach only concerns simple sequential composition of Web services method calls, without considering complex control flows and timing requirements.

*Response time estimation* Our approach is also related to response time estimation. In [Kra+09], the authors propose to use linear regression method and a maximum likelihood technique for estimating the service demands of requests based on their response times. [Men04] has also discussed the impact of slow services on the overall response time on a transaction that use several services concurrently. Our work is focused on decomposing the global requirement into local requirement, which is orthogonal to these works. Our work [LTC14] complements with this work by proposing a method on building LTCs that under-approximate the sLTC of a composite service. The under-approximated LTCs consisting of independent constraints over components, which can be used to improve the design, monitoring and repair of component-based systems under time requirements.

*Service monitoring* Our method is related to service monitoring. Moser *et al.* [MRD08] present VieDAME, a non-intrusive approach to monitoring. VieDAME allows monitoring of BPEL composite service on quality of service attributes, and existing component services are replaced based on different replacement strategies. They make use of the aspect-oriented approach (AOP); therefore the VieDAME engine adapter could be interwoven into the BPEL runtime engine

at runtime. Baresi *et al.* [BG11] propose an idea of self-supervising BPEL processes by supporting both service monitoring and recovery for BPEL processes. They propose using Web Service Constraint Language (WScCoL) to specify the monitoring directives to indicate properties that need to hold during the runtime of composite service. They also make use of the AOP approach to integrate their monitoring adapters with the BPEL runtime engine. Our work is orthogonal to the aforementioned works, as we do not assume any particular service monitoring framework for monitoring the composite service, and those methods can be used to aid our monitoring approach, as discussed in Section 6.2. Our previous work [Tan+14] proposes an automated approach based on a genetic algorithm to calculate the recovery plan that can guarantee the satisfaction of functional properties of the composite service after recovery.

*Service selection* In [Zen+03; Zen+04], the authors present an approach that makes use of global planning to search dynamically for the best services component for service composition. Their approach involves the use of mixed integer programming (MIP) techniques to find the optimal selection of component services. Ardagna *et al.* [Ard+07] extend the MIP methods to include local constraints. Cardellini *et al.* [Car+09] propose a methodology to integrate different adaptation mechanisms for combining concrete services to an abstract service, in order to achieve a greater flexibility in facing different operating environments. Our work is orthogonal to aforementioned works, as it does not assume particular formulation of the MIP problems.

Although the method in aforementioned works efficiently for small case studies, it suffers from scalability problems when the size of the case studies becomes larger, since the time required grows exponentially with the size of problem. To address this problem, Yu *et al.* [YZL07] propose a heuristic algorithm that could be used to find a near-optimal solution. The authors proposed two QoS compositional models, a combinatorial model and a graph model. The time complexity for the combinatorial model is polynomial, while the time complexity for the graph model is exponential. However the algorithm does not scale with the increasing number of web services. To address this problem, Alrifai *et al.* present an approach that pruned the search space using skyline methods, and they make use of a hierarchical  $k$ -means clustering method [Llo82] for representative selection. The work of Alrifai *et al.* is the closest to ours. Despite its reasonable performance, a limitation for the MIP approach is that it cannot deal with non-linear objective functions or aggregated constraints. To address this problem, Canfora *et al.* [Can+05] have formulated the problem as a genetic algorithm problem. Genetic algorithms (GA) are algorithms based on stochastic search methods, that support non-linear objective functions. Two different GAs encodings are proposed in [Can+05; Zha+03]. In [Zha+03], the authors propose to encode the chromosome using binary strings, and every gene is a chromosome representing a service candidate with value 0 (respectively, 1) that represents the unselected (respectively, selected) service. Therefore the length of the genome can be very long, given a large number of service candidates. In [Can+05], the authors propose to encode the chromosome using an integer value which represents the index of the concrete services stored in an array. This coding scheme results in shorter chromosomes, and the length of a chromosome is independent of the number of service candidates. In [gao2007QoS], Gao *et al.* propose a tree coding scheme to represent the service

composition. They reported a 40% performance improvement with respect to the single dimension coding scheme used in [Can+05]. This is because the tree coding scheme does not need to recalculate the entire fitness value each time compare to the single-dimension encoding. Our work does not assume any particular encoding scheme and it can be used with any existing coding techniques. In [AT08], Ai *et al.* proposed an approach extending the GA methods for handling inter-service dependencies and conflicts using a penalty-based genetic algorithm. Our work does not assume a particular fitness function. In [MZ08], Ma *et al.* proposed an enhanced initial population policy and an evolution policy based on population diversity and a relation matrix coding scheme. They considered all concrete services for each service class starting from the initial population. Different from their approach, we only consider a subset of services with high local utility value from the start, and we progressively add more services later on. In [SN20] the problem of functionally equivalent service composition is considered.

*Verification of services* Concerning verification of services, Filieri *et al.* [FGT11] focus on checking the reliability of component (service)-based systems. They make use of Discrete Time Markov Chain (DTMC) to check the reliability of models at runtime. Our previous works [Che+13; Che+14] develop a tool to verify combined functional and non-functional requirements of Web service composition. In contrast, the current work focuses on response time: given the global response time of the composite service, we synthesize the response time requirement for component services at design time and refine it at runtime. Schmieders *et al.* [SM11] proposed the SPADE approach. SPADE invokes the BOGOR model checker to model check the SLAs at design time and at runtime. Our work is different from theirs in two aspects. First, we focus on the synthesis of the local time requirement, which is a formal requirement on the response time requirement of component services. Second, at runtime, [SM11] performs model checking on a given state to check whether an adaptation is needed. In contrast, we have precomputed the constraints for every state at design time. Therefore, we only require evaluation of constraints by substituting the free variables at runtime, and this allows a more efficient runtime-analysis.

*BPEL* In [GMJ08] a template-based system is used to reconfigure service composition, using BPEL. In [Pau09], service composition using RESTful (Representational State Transfer) is performed using the BPEL extension “BPEL for REST”. In [Eic+10], a tool based on Services Creation Environments and BPEL is proposed, that also allows translation to Java.

The work [TBM13] automates the formalization and verification process of BPEL. It extends the existing spring framework to represent BPEL activities with Java bean, which is subsequently transformed into XML-based formal model (like colored Petri nets) for verification. The work [Mi+16] also transforms BPEL services into Probabilistic Labeled Transition Systems, which are then used to conduct probabilistic model checking to verify reliability properties on the BPEL models. Their work did not consider timing requirements.

*Surveys* Finally, composition of Web services has been recently surveyed. In [Ord+15] the larger domain of composition of convergent services is surveyed;

however, BPEL is still surveyed in this work. In [LDB16], a taxonomy of Web service composition is provided with different directions surveyed such as language, knowledge reuse, automation, tool support, execution platform or target user.

## 9 Conclusion and future work

### 9.1 Conclusion

We have presented a model-based approach for synthesizing local time constraints for component services of a composite service CS, knowing its global time requirement. Our approach makes use of parameterized timed techniques.

We first proposed a design-time synthesis algorithm, that utilizes the parametric constraints from the LTS, to synthesize static local time constraint (sLTC) for component services. The sLTC is then used to select a set of component services that could collectively satisfy the global time requirement in design time.

Then, we use the runtime information to weaken the sLTC, which becomes the refined local time constraint (rLTC). In particular, two pieces of runtime information have been leveraged—the execution path that has been taken by the composite service, and the elapsed time of the composite service. The rLTC is then used to validate whether the composite service can still satisfy the global time requirement at runtime.

As a proof of concept, we have implemented our approach into a tool SELAMAT, and applied it to four service composition examples. Our experiments show that the runtime refinement leads to an improvement of the global time requirement, with limited overhead.

### 9.2 Future work

We plan to further improve and develop the technique presented in this paper.

*General and dedicated optimizations* First, the goal of our work is to propose a full framework for analyzing composition of Web services using parametric timings; therefore, integrating existing state space reduction techniques is perhaps a more practical work, orthogonal to the original goal of our approach. Nevertheless, in order to address huge sets of services, one could use efficient state-of-the-art techniques developed for timed systems or parametric timed systems such as (parametric) data difference bound matrices [Hum+02; QSW17], efficient L/U-zone abstractions [HSW16; NPV18], convex state merging [AFS13], integer-hull abstractions [JLR15; ALR15], or abstraction-refinement algorithms [RSM19].

*Soft deadlines* Second, we will investigate the usage of soft deadlines that allow to run a service with a delay, possibly with an acceptable penalty.

*Constraints satisfiability* Regarding our implementation, the bottleneck seems to be the satisfiability test using Z3; from our experience, switching to a polyhedra library (such as the Parma Polyhedra Library [BHZ08]) may give better results, and could help render our work scalable.

*Uncertain response times* Our work so far deals with exact response times. A different approach would be to consider that the response time should be fulfilled with some probability. In that setting, the goal would be to synthesize the values for the timing parameters such that the response time is indeed below the threshold with a given probability. To achieve this, we could reuse recent works involving probabilities and timing parameters (*e.g.*, [JK14; Ces+14]). An even more challenging problem would be to combine both kinds of parameters (timing parameters and probabilistic parameters), so as to infer the probability under which the response time can be fulfilled.

*Statistical model checking* Finally, when concurrent systems with or without timing constraints are too huge to be analyzed in an exact manner, a recent trend is to propose non-exact techniques, and notably statistical model-checking. This technique could be used for compositions of Web services arguably too large to be analyzed in an exact manner. Recent techniques developed in the timed setting (*e.g.*, [Dav+11; LL16; Med+18]) could be applied to our formalism.

## References

- [AD94] Rajeev Alur and David L. Dill. “A theory of timed automata”. In: *TC* 126.2 (Apr. 1994), pp. 183–235. issn: 0304-3975. DOI: [10.1016/0304-3975\(94\)90010-8](https://doi.org/10.1016/0304-3975(94)90010-8) (cit. on pp. 7, 12).
- [AFS13] Étienne André, Laurent Fribourg, and Romain Soulat. “Merge and Conquer: State Merging in Parametric Timed Automata”. In: *ATVA*. Ed. by Dang-Van Hung and Mizuhito Ogawa. Vol. 8172. Lecture Notes in Computer Science. Hanoi, Vietnam: Springer, Oct. 2013, pp. 381–396. DOI: [10.1007/978-3-319-02444-8\\_27](https://doi.org/10.1007/978-3-319-02444-8_27) (cit. on p. 43).
- [AHV93] Rajeev Alur, Thomas A. Henzinger, and Moshe Y. Vardi. “Parametric real-time reasoning”. In: *STOC*. Ed. by S. Rao Kosaraju, David S. Johnson, and Alok Aggarwal. San Diego, California, United States: ACM, 1993, pp. 592–601. ISBN: 0-89791-591-7. DOI: [10.1145/167088.167242](https://doi.org/10.1145/167088.167242) (cit. on pp. 7, 40).
- [ALR15] Étienne André, Didier Lime, and Olivier H. Roux. “Integer-Complete Synthesis for Bounded Parametric Timed Automata”. In: *RP*. Ed. by Mikołaj Bojańczyk, Sławomir Lasota, and Igor Potapov. Vol. 9328. Lecture Notes in Computer Science. Warsaw, Poland: Springer, Sept. 2015, pp. 7–19. DOI: [10.1007/978-3-319-24537-9\\_2](https://doi.org/10.1007/978-3-319-24537-9_2) (cit. on p. 43).
- [Alv+07] Alexandre Alves, Assaf Arkin, Sid Askary, Charlton Barreto, Ben Bloch, Francisco Curbera, Mark Ford, Yaron Goland, Alejandro Guízar, Neelakantan Kartha, Canyang Kevin Liu, Rania Khalaf, Dieter König, Mike Marin, IBM, Vinkesh Mehta, Satish Thatte, Danny van der Rijn, Prasad Yendluri, and Alex Yiu. *Web Services Business Process Execution Language Version, version 2.0*. Apr. 2007 (cit. on pp. 2, 8).
- [AM02] Yasmina Adbeddaïm and Oded Maler. “Preemptive Job-Shop Scheduling using Stopwatch Automata”. In: *TACAS*. Ed. by Joost-Pieter Katoen and Perdita Stevens. Vol. 2280. Lecture Notes in Computer Science. Grenoble, France: Springer-Verlag, Apr. 2002, pp. 113–126. DOI: [10.1007/3-540-46002-0\\_9](https://doi.org/10.1007/3-540-46002-0_9) (cit. on p. 40).
- [Ama18] Amazon. *Amazon Elastic Compute Cloud (Amazon EC2)*. 2018. URL: <https://aws.amazon.com/ec2/> (cit. on p. 2).
- [And+14] Étienne André, Yang Liu, Jun Sun, and Jin Song Dong. “Parameter Synthesis for Hierarchical Concurrent Real-Time Systems”. In: *Real-Time Systems* 50.5-6 (Sept. 2014), pp. 620–679. DOI: [10.1007/s11241-014-9208-6](https://doi.org/10.1007/s11241-014-9208-6) (cit. on pp. 7, 11, 12, 14, 18, 40).

- [And13] Étienne André. “Dynamic Clock Elimination in Parametric Timed Automata”. In: *FSFMA*. Ed. by Christine Choppy and Jun Sun. Vol. 31. OpenAccess Series in Informatics (OASICS). Singapore: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, July 2013, pp. 18–31. DOI: [10.4230/OASICS.FSFMA.2013.18](https://doi.org/10.4230/OASICS.FSFMA.2013.18) (cit. on p. 12).
- [AP05] Danilo Ardagna and Barbara Pernici. “Global and Local QoS Guarantee in Web Service Selection”. In: *Business Process Management Workshops*. Ed. by Christoph Bussler and Armin Haller. Vol. 3812. 2005, pp. 32–46. DOI: [10.1007/11678564\\_4](https://doi.org/10.1007/11678564_4) (cit. on p. 40).
- [AR09] Mohammad Alrifai and Thomas Risse. “Combining global optimization with local selection for efficient QoS-aware service composition”. In: *WWW*. Ed. by Juan Quemada, Gonzalo León, Yoëlle S. Maarek, and Wolfgang Nejdl. ACM, 2009, pp. 881–890. ISBN: 978-1-60558-487-4. DOI: [10.1145/1526709.1526828](https://doi.org/10.1145/1526709.1526828) (cit. on p. 40).
- [Ard+07] Danilo Ardagna, Marco Comuzzi, Enrico Mussi, Barbara Pernici, and Pierluigi Plebani. “PAWS: A Framework for Executing Adaptive Web-Service Processes”. In: *IEEE Software* 24.6 (2007), pp. 39–46. DOI: [10.1109/MS.2007.174](https://doi.org/10.1109/MS.2007.174) (cit. on p. 41).
- [AS13] Étienne André and Romain Soulat. *The Inverse Method*. FOCUS Series in Computer Engineering and Information Technology. 176 pages. ISTE Ltd and John Wiley & Sons Inc., 2013. ISBN: 9781848214477 (cit. on p. 40).
- [AT08] Lifeng Ai and Maolin Tang. “A penalty-based genetic algorithm for QoS-aware Web service composition with inter-service dependencies and conflicts”. In: ed. by Masoud Mohammadian. IEEE Computer Society, 2008, pp. 738–743. DOI: [10.1109/CIMCA.2008.104](https://doi.org/10.1109/CIMCA.2008.104) (cit. on p. 42).
- [BG11] Luciano Baresi and Sam Guinea. “Self-Supervising BPEL Processes”. In: *IEEE Transactions on Software Engineering* 37.2 (2011), pp. 247–263. DOI: [10.1109/TSE.2010.37](https://doi.org/10.1109/TSE.2010.37) (cit. on p. 41).
- [BGP11] Domenico Bianculli, Dimitra Giannakopoulou, and Corina S. Pasareanu. “Interface decomposition for service compositions”. In: *ICSE*. 2011, pp. 501–510. DOI: [10.1145/1985793.1985862](https://doi.org/10.1145/1985793.1985862) (cit. on p. 39).
- [BHZ08] Roberto Bagnara, Patricia M. Hill, and Enea Zaffanella. “The Parma Polyhedra Library: Toward a Complete Set of Numerical Abstractions for the Analysis and Verification of Hardware and Software Systems”. In: *Science of Computer Programming* 72.1–2 (2008), pp. 3–21. DOI: [10.1016/j.scico.2007.08.001](https://doi.org/10.1016/j.scico.2007.08.001) (cit. on p. 43).
- [BY03] Johan Bengtsson and Wang Yi. “Timed Automata: Semantics, Algorithms and Tools”. In: *Lectures on Concurrency and Petri Nets, Advances in Petri Nets*. Vol. 3098. Lecture Notes in Computer Science. Springer, 2003, pp. 87–124. ISBN: 3-540-22261-8. DOI: [10.1007/978-3-540-27755-2\\_3](https://doi.org/10.1007/978-3-540-27755-2_3) (cit. on pp. 12, 14).
- [Can+05] Gerardo Canfora, Massimiliano Di Penta, Raffaele Esposito, and Maria Luisa Villani. “An approach for QoS-aware service composition based on genetic algorithms”. In: *GECCO*. Ed. by Hans-Georg Beyer and Una-May O’Reilly. ACM, 2005, pp. 1069–1075. DOI: [10.1145/1068009.1068189](https://doi.org/10.1145/1068009.1068189) (cit. on pp. 41, 42).
- [Car+09] Valeria Cardellini, Emiliano Casalicchio, Vincenzo Grassi, Francesco Lo Presti, and Raffaella Mirandola. “QoS-driven runtime adaptation of service oriented architectures”. In: *FSE*. Ed. by Hans van Vliet and Valérie Issarny. ACM, 2009, pp. 131–140. DOI: [10.1145/1595696.1595718](https://doi.org/10.1145/1595696.1595718) (cit. on p. 41).
- [CDP17] Tomás Cerný, Michael J. Donahoo, and Jiri Pechanec. “Disambiguation and Comparison of SOA, Microservices and Self-Contained Systems”. In: *RACS*. (Sept. 20–23, 2017). Krakow, Poland: ACM, 2017, pp. 228–235. DOI: [10.1145/3129676.3129682](https://doi.org/10.1145/3129676.3129682) (cit. on p. 2).
- [Ces+14] Milan Ceska, Frits Dannenberg, Marta Z. Kwiatkowska, and Nicola Paoletti. “Precise Parameter Synthesis for Stochastic Biochemical Systems”. In: *CMSB*. Ed. by Pedro Mendes, Joseph O. Dada, and Kieran Smallbone. Vol. 8859. Lecture Notes in Computer Science. Springer, 2014, pp. 86–98. DOI: [10.1007/978-3-319-12982-2\\_7](https://doi.org/10.1007/978-3-319-12982-2_7) (cit. on p. 44).
- [Che+13] Manman Chen, Tian Huat Tan, Jun Sun, Yang Liu, Jun Pang, and Xiaohong Li. “Verification of Functional and Non-functional Requirements of Web Service



- Composition”. In: *ICFEM*. 2013, pp. 313–328. DOI: [10.1007/978-3-642-41202-8\\_21](https://doi.org/10.1007/978-3-642-41202-8_21) (cit. on p. 42).
- [Che+14] Manman Chen, Tian Huat Tan, Jun Sun, Yang Liu, and Jin Song Dong. “VeriWS: a tool for verification of combined functional and non-functional requirements of Web service composition”. In: *ICSE*. 2014, pp. 564–567. DOI: [10.1145/2591062.2591070](https://doi.org/10.1145/2591062.2591070) (cit. on p. 42).
- [Chi+07] Roberto Chinnici, Jean-Jacques Moreau, Arthur Ryman, and Sanjiva Weerawarana. *Web Services Description Language (WSDL) Version 2.0*. W3C Recommendation, available at <http://www.w3.org/TR/wsdl20/>. June 2007 (cit. on p. 2).
- [CPR08] Alessandro Cimatti, Luigi Palopoli, and Yusi Ramadian. “Symbolic Computation of Schedulability Regions Using Parametric Timed Automata”. In: *RTSS*. IEEE Computer Society, 2008, pp. 80–89. ISBN: 978-0-7695-3477-0. DOI: <http://dx.doi.org/10.1109/RTSS.2008.36> (cit. on p. 40).
- [Dav+11] Alexandre David, Kim Gulstrand Larsen, Axel Legay, Marius Mikucionis, and Zheng Wang. “Time for Statistical Model Checking of Real-Time Systems”. In: *CAV*. Vol. 6806. Lecture Notes in Computer Science. Snowbird, UT, USA: Springer, 2011, pp. 349–355. DOI: [10.1007/978-3-642-22110-1\\_27](https://doi.org/10.1007/978-3-642-22110-1_27) (cit. on p. 44).
- [DB08] Leonardo Mendonça De Moura and Nikolaj Bjørner. “Z3: An Efficient SMT Solver”. In: *TACAS*. Ed. by C. R. Ramakrishnan and Jakob Rehof. Vol. 4963. Lecture Notes in Computer Science. Springer, 2008, pp. 337–340. DOI: [10.1007/978-3-540-78800-3\\_24](https://doi.org/10.1007/978-3-540-78800-3_24) (cit. on pp. 24, 32, 38).
- [DY96] Conrado Daws and Sergio Yovine. “Reducing the number of clock variables of timed automata”. In: *RTSS*. IEEE Computer Society, 1996, pp. 73–81. ISBN: 0-8186-7689-2. DOI: [10.1109/REAL.1996.563702](https://doi.org/10.1109/REAL.1996.563702) (cit. on p. 12).
- [Eic+10] Thomas Eichelmann, Woldemar F. Fuhrmann, Ulrich Trick, and Bogdan V. Ghita. “Enhanced Concept of the TeamCom SCE for Automated Generated Services Based on JSLEE”. In: *INC*. Ed. by Udo Bleimann, Paul Dowland, Steven Furnell, and Oliver Schneider. University of Plymouth, 2010, pp. 75–84 (cit. on p. 42).
- [Erm+07] Andreas Ermedahl, Christer Sandberg, Jan Gustafsson, Stefan Bygde, and Björn Lisper. “Loop Bound Analysis based on a Combination of Program Slicing, Abstract Interpretation, and Invariant Analysis”. In: *WCET*. Ed. by Christine Rochange. Vol. 6. OASICS. Internationales Begegnungs- und Forschungszentrum fuer Informatik (IBFI), Schloss Dagstuhl, Germany, 2007 (cit. on p. 32).
- [FGT11] Antonio Filieri, Carlo Ghezzi, and Giordano Tamburrelli. “Run-time efficient probabilistic model checking”. In: *ICSE*. 2011, pp. 341–350. DOI: [10.1145/1985793.1985840](https://doi.org/10.1145/1985793.1985840) (cit. on p. 42).
- [Fos+06] Howard Foster, Sebastián Uchitel, Jeff Magee, and Jeff Kramer. “LTSA-WS: a tool for model-based verification of Web service compositions and choreography”. In: *ICSE*. 2006, pp. 771–774. DOI: [10.1145/1134408](https://doi.org/10.1145/1134408) (cit. on p. 39).
- [Fos06] Howard Foster. “A Rigorous Approach To Engineering Web Service Compositions”. PhD thesis. Imperial College of London, 2006 (cit. on p. 39).
- [Fou07] Apache Software Foundation. *Apache ODE*. 2007. URL: <http://ode.apache.org/> (cit. on pp. 28, 36).
- [Fri+12] Laurent Fribourg, David Lesens, Pierre Moro, and Romain Soulat. “Robustness Analysis for Scheduling Problems using the Inverse Method”. In: *TIME*. Ed. by Mark Reynolds, Paolo Terenziani, and Ben Moszkowski. Leicester, UK: IEEE Computer Society Press, Sept. 2012, pp. 73–80. DOI: [10.1109/TIME.2012.10](https://doi.org/10.1109/TIME.2012.10) (cit. on p. 40).
- [GMJ08] Kristof Geebelen, Sam Michiels, and Wouter Joosen. “Dynamic reconfiguration using template based Web service composition”. In: *MW4SOC*. Ed. by Karl M. Göschka, Schahram Dustdar, Frank Leymann, and Vladimir Tosic. ACM, 2008, pp. 49–54. DOI: [10.1145/1462802.1462811](https://doi.org/10.1145/1462802.1462811) (cit. on p. 42).
- [Gud+07] Martin Gudgin, Marc Hadley, Noah Mendelsohn, Jean-Jacques Moreau, Henrik Frystyk Nielsen, Anish Karmarkar, and Yves Lafon. *Simple Object Access Protocol (SOAP) Version 1.2*. W3C Recommendation, available at <http://www.w3.org/TR/soap12/>. Apr. 2007 (cit. on p. 2).



- [HSW16] Frédéric Herbreteau, B. Srivathsan, and Igor Walukiewicz. “Better abstractions for timed automata”. In: *Information and Computation* 251 (2016), pp. 67–90. DOI: [10.1016/j.ic.2016.07.004](https://doi.org/10.1016/j.ic.2016.07.004) (cit. on pp. 32, 43).
- [Hun+02] Thomas Hune, Judi Romijn, Mariëlle Stoelinga, and Frits W. Vaandrager. “Linear parametric model checking of timed automata”. In: *Journal of Logic and Algebraic Programming* 52-53 (2002), pp. 183–220. DOI: [10.1016/S1567-8326\(02\)00037-1](https://doi.org/10.1016/S1567-8326(02)00037-1) (cit. on pp. 14, 18, 43).
- [JK14] Aleksandra Jovanović and Marta Z. Kwiatkowska. “Parameter Synthesis for Probabilistic Timed Automata Using Stochastic Game Abstractions”. In: *RP*. Ed. by Joël Ouaknine, Igor Potapov, and James Worrell. Vol. 8762. Lecture Notes in Computer Science. Springer, 2014, pp. 176–189. DOI: [10.1007/978-3-319-11439-2\\_14](https://doi.org/10.1007/978-3-319-11439-2_14) (cit. on p. 44).
- [JLR15] Aleksandra Jovanović, Didier Lime, and Olivier H. Roux. “Integer Parameter Synthesis for Real-Time Systems”. In: *IEEE Transactions on Software Engineering* 41.5 (2015), pp. 445–461. DOI: [10.1109/TSE.2014.2357445](https://doi.org/10.1109/TSE.2014.2357445) (cit. on p. 43).
- [Kra+09] Stephan Kraft, Sergio Pacheco-Sanchez, Giuliano Casale, and Stephen Dawson. “Estimating service resource consumption from response time measurements”. In: *VALUETOOLS*. 2009, p. 48. DOI: [10.4108/ICST.VALETOOLS2009.7526](https://doi.org/10.4108/ICST.VALETOOLS2009.7526) (cit. on p. 40).
- [LDB16] Angel Lagares Lemos, Florian Daniel, and Boualem Benatallah. “Web Service Composition: A Survey of Techniques and Tools”. In: *ACM Computing Surveys* 48.3 (2016), 33:1–33:41. DOI: [10.1145/2831270](https://doi.org/10.1145/2831270) (cit. on p. 43).
- [Le+10] Thi Thieu Hoa Le, Luigi Palopoli, Roberto Passerone, Yusi Ramadian, and Alessandro Cimatti. “Parametric analysis of distributed firm real-time systems: A case study”. In: *ETFA*. IEEE, 2010, pp. 1–8. ISBN: 978-1-4244-6848-5. DOI: [10.1109/ETFA.2010.5641315](https://doi.org/10.1109/ETFA.2010.5641315) (cit. on p. 40).
- [LL16] Kim Guldstrand Larsen and Axel Legay. “Statistical Model Checking: Past, Present, and Future”. In: *ISoLA, Part I*. Ed. by Tiziana Margaria and Bernhard Steffen. Vol. 9952. Lecture Notes in Computer Science. Imperial, Corfu, Greece, 2016, pp. 3–15. DOI: [10.1007/978-3-319-47166-2\\_1](https://doi.org/10.1007/978-3-319-47166-2_1) (cit. on p. 44).
- [Llo82] Stuart P. Lloyd. “Least squares quantization in PCM”. In: *IEEE Transactions on Information Theory* 28.2 (1982), pp. 129–136. DOI: [10.1109/TIT.1982.1056489](https://doi.org/10.1109/TIT.1982.1056489) (cit. on p. 41).
- [LTC14] Yi Li, Tian Huat Tan, and Marsha Chechik. “Management of Time Requirements in Component-Based Systems”. In: *FM*. 2014, pp. 399–415. DOI: [10.1007/978-3-319-06410-9\\_28](https://doi.org/10.1007/978-3-319-06410-9_28) (cit. on p. 40).
- [Med+18] Braham Lotfi Mediouni, Ayoub Nouri, Marius Bozga, Mahieddine Dellabani, Axel Legay, and Saddek Bensalem. “S BIP 2.0: Statistical Model Checking Stochastic Real-Time Systems”. In: *ATVA*. Ed. by Shuvendu K. Lahiri and Chao Wang. Vol. 11138. Lecture Notes in Computer Science. Los Angeles, CA, USA: Springer, 2018, pp. 536–542. DOI: [10.1007/978-3-030-01090-4\\_33](https://doi.org/10.1007/978-3-030-01090-4_33) (cit. on p. 44).
- [Men04] Daniel A. Menascé. “Response-Time Analysis of Composite Web Services”. In: *IEEE Internet Computing* 8.1 (2004), pp. 90–92. DOI: [10.1109/MIC.2004.1260710](https://doi.org/10.1109/MIC.2004.1260710) (cit. on p. 40).
- [Mer74] Philip Meir Merlin. “A study of the recoverability of computing systems.” PhD thesis. University of California, Irvine, CA, USA, 1974 (cit. on p. 12).
- [Mi+16] Chengyang Mi, Huaikou Miao, Jinyu Kai, and Honghao Gao. “Reliability modeling and verification of BPEL-based Web services composition by probabilistic model checking”. In: *SERA*. Ed. by Yeong-Tae Song. IEEE Computer Society, 2016, pp. 149–154. DOI: [10.1109/SERA.2016.7516140](https://doi.org/10.1109/SERA.2016.7516140) (cit. on p. 42).
- [Mid17] Middleware. *The State of Microservices Survey 2017 – Eight trends you need to know*. 2017. URL: <https://middlewareblog.redhat.com/2017/12/05/the-state-of-microservices-survey-2017-eight-trends-you-need-to-know/> (cit. on p. 2).
- [MK06] Jeff Magee and Jeff Kramer. *Concurrency – State models and Java programs (2. ed.)*. Wiley, 2006, pp. I–XVIII, 1–413. ISBN: 978-0-470-09355-9 (cit. on p. 39).

- [MRD08] Oliver Moser, Florian Rosenberg, and Schahram Dustdar. “Non-intrusive monitoring and service adaptation for WS-BPEL”. In: *WWW*. 2008, pp. 815–824. DOI: [10.1145/1367497.1367607](https://doi.org/10.1145/1367497.1367607) (cit. on pp. 32, 40).
- [MZ08] Yue Ma and Chengwen Zhang. “Quick convergence of genetic algorithm for QoS-driven Web service selection”. In: *Computer Networks* 52.5 (2008), pp. 1093–1104. DOI: [10.1016/j.comnet.2007.12.003](https://doi.org/10.1016/j.comnet.2007.12.003) (cit. on p. 42).
- [NPV18] Hoang Gia Nguyen, Laure Petrucci, and Jaco Van de Pol. “Layered and Collecting NDFS with Subsumption for Parametric Timed Automata”. In: *ICECCS*. Ed. by Anthony Widjaja Lin and Jun Sun. Melbourne, Australia: IEEE Computer Society, Dec. 2018, pp. 1–9. DOI: [10.1109/ICECCS2018.2018.00009](https://doi.org/10.1109/ICECCS2018.2018.00009) (cit. on pp. 32, 43).
- [Ord+15] Armando Ordóñez, Vidal Alcázar, Oscar Mauricio Caicedo Rendon, Paolo Falcarin, Juan Carlos Corrales, and Lisandro Zambenedetti Granville. “Towards automated composition of convergent services: A survey”. In: *Computer Communications* 69 (2015), pp. 1–21. DOI: [10.1016/j.comcom.2015.07.025](https://doi.org/10.1016/j.comcom.2015.07.025) (cit. on p. 42).
- [Pau09] Cesare Pautasso. “RESTful Web service composition with BPEL for REST”. In: *Data & Knowledge Engineering* 68.9 (2009), pp. 851–866. DOI: [10.1016/j.datak.2009.02.016](https://doi.org/10.1016/j.datak.2009.02.016) (cit. on p. 42).
- [QSW17] Karin Quaas, Mahsa Shirmohammadi, and James Worrell. “Revisiting reachability in timed automata”. In: *LICS*. Reykjavik, Iceland: IEEE Computer Society, 2017, pp. 1–12. DOI: [10.1109/LICS.2017.8005098](https://doi.org/10.1109/LICS.2017.8005098) (cit. on p. 43).
- [Ric18] Chris Richardson. *Microservices*. 2018. URL: <http://microservices.io/patterns/microservices.html> (cit. on p. 2).
- [RSM19] Victor Roussanaly, Ocan Sankur, and Nicolas Markey. “Abstraction Refinement Algorithms for Timed Automata”. In: *CAV, Part I*. Ed. by Işıl Dillig and Serdar Tasiran. Vol. 11561. Lecture Notes in Computer Science. New York City, USA: Springer, 2019, pp. 22–40. DOI: [10.1007/978-3-030-25540-4\\_2](https://doi.org/10.1007/978-3-030-25540-4_2) (cit. on p. 43).
- [SBC10] Jocelyn Simmonds, Shoham Ben-David, and Marsha Chechik. “Guided recovery for Web service applications”. In: *SIGSOFT FSE*. Ed. by André Roman Gruia-Catalinand van der Hoek. ACM, 2010, pp. 247–256. DOI: [10.1145/1882291.1882328](https://doi.org/10.1145/1882291.1882328) (cit. on p. 39).
- [Sch86] Alexander Schrijver. *Theory of linear and integer programming*. John Wiley and Sons, 1986 (cit. on p. 8).
- [SM11] Eric Schmieders and Andreas Metzger. “Preventing Performance Violations of Service Compositions Using Assumption-Based Run-Time Verification”. In: *ServiceWave*. Poznań, Poland, 2011, pp. 194–205. DOI: [10.1007/978-3-642-24755-2\\_19](https://doi.org/10.1007/978-3-642-24755-2_19) (cit. on p. 42).
- [SN20] Sujata Swain and Rajdeep Niyogi. “FESC: Functionally Equivalent Service Composition”. In: *Internet of Things* 9 (2020), p. 100151. ISSN: 2542-6605. DOI: <https://doi.org/10.1016/j.iot.2019.100151> (cit. on p. 42).
- [ST19] Zheng Song and Eli Tilevich. “Equivalence-Enhanced Microservice Workflow Orchestration to Efficiently Increase Reliability”. In: *ICWS*. (July 8–13, 2019). Ed. by Elisa Bertino, Carl K. Chang, Peter Chen, Ernesto Damiani, Michael Goul, and Katsunori Oyama. Milan, Italy: IEEE, 2019, pp. 426–433. DOI: [10.1109/ICWS.2019.00076](https://doi.org/10.1109/ICWS.2019.00076) (cit. on p. 2).
- [Sun+09] Jun Sun, Yang Liu, Jin Song Dong, and Jun Pang. “PAT: Towards Flexible Verification under Fairness”. In: *CAV*. Ed. by Ahmed Bouajjani and Oded Maler. Vol. 5643. Lecture Notes in Computer Science. Grenoble, France: Springer, 2009, pp. 709–714. ISBN: 978-3-642-02657-7. DOI: [10.1007/978-3-642-02658-4\\_59](https://doi.org/10.1007/978-3-642-02658-4_59) (cit. on p. 32).
- [Sun+13] Jun Sun, Yang Liu, Jin Song Dong, Yan Liu, Ling Shi, and Étienne André. “Modeling and Verifying Hierarchical Real-time Systems using Stateful Timed CSP”. In: *ACM Transactions on Software Engineering and Methodology* 22.1 (2013), pp. 3.1–3.29. DOI: [10.1145/2430536.2430537](https://doi.org/10.1145/2430536.2430537) (cit. on pp. 7, 11, 12).
- [Tan+13] Tian Huat Tan, Étienne André, Jun Sun, Yang Liu, Jin Song Dong, and Manman Chen. “Dynamic Synthesis of Local Time Requirement for Service Composition”. In: *ICSE*. Ed. by Betty H.C. Cheng and Klaus Pohl. San Francisco, USA: IEEE, May 2013, pp. 542–551. DOI: [10.1109/ICSE.2013.6606600](https://doi.org/10.1109/ICSE.2013.6606600) (cit. on p. 5).

- [Tan+14] Tian Huat Tan, Manman Chen, Étienne André, Jun Sun, Yang Liu, and Jin Song Dong. “Automated runtime recovery for QoS-based service composition”. In: *WWW*. 2014, pp. 563–574. DOI: [10.1145/2566486.2568048](https://doi.org/10.1145/2566486.2568048) (cit. on p. 41).
- [Tan+16] Tian Huat Tan, Manman Chen, Jun Sun, Yang Liu, Étienne André, Jin Song Dong, and Yinxing Xue. “Optimizing Selection of Competing Services with Probabilistic Hierarchical Refinement”. In: *ICSE*. Ed. by Willem Visser and Laurie Williams. Austin, Texas, USA: ACM, 2016, pp. 85–95. ISBN: 978-1-4503-3900-1. DOI: [10.1145/2884781.2884861](https://doi.org/10.1145/2884781.2884861) (cit. on p. 2).
- [Tan+19] Tian Huat Tan, Étienne André, Manman Chen, Jun Sun, Yang Liu, and Jin Song Dong. *Selamat: binary and experiment data*. 2019. URL: <https://sites.google.com/site/automatedsynthesis/home/> (cit. on p. 32).
- [TBM13] Zahir Tari, Peter Bertok, and Anshuman Mukherjee. “Framework for modeling, simulation and verification of a BPEL specification”. In: *Verification of Communication Protocols in Web Services*. John Wiley & Sons, Inc, 2013, pp. 205–244. ISBN: 9781118720103. DOI: [10.1002/9781118720103.ch8](https://doi.org/10.1002/9781118720103.ch8) (cit. on p. 42).
- [TLR09] Louis-Marie Traonouez, Didier Lime, and Olivier H. Roux. “Parametric Model-Checking of Stopwatch Petri Nets”. In: *Journal of Universal Computer Science* 15.17 (2009), pp. 3273–3304. DOI: [10.3217/jucs-015-17-3273](https://doi.org/10.3217/jucs-015-17-3273) (cit. on pp. 18, 40).
- [YZL07] Tao Yu, Yue Zhang, and Kwei-Jay Lin. “Efficient algorithms for Web services selection with end-to-end QoS constraints”. In: *ACM Transactions on the Web* 1.1 (2007), p. 6. DOI: [10.1145/1232722.1232728](https://doi.org/10.1145/1232722.1232728) (cit. on pp. 40, 41).
- [Zen+03] Liangzhao Zeng, Boualem Benatallah, Marlon Dumas, Jayant Kalagnanam, and Quan Z. Sheng. “Quality driven Web services composition”. In: *WWW*. Ed. by Gusztáv Hencsey, Bebo White, Yih-Farn Robin Chen, László Kovács, and Steve Lawrence. ACM, 2003, pp. 411–421. DOI: [10.1145/775152.775211](https://doi.org/10.1145/775152.775211) (cit. on p. 41).
- [Zen+04] Liangzhao Zeng, Boualem Benatallah, Anne H. H. Ngu, Marlon Dumas, Jayant Kalagnanam, and Henry Chang. “QoS-aware middleware for Web services composition”. In: *IEEE Transactions on Software Engineering* 30.5 (2004), pp. 311–327. DOI: [10.1109/TSE.2004.11](https://doi.org/10.1109/TSE.2004.11) (cit. on p. 41).
- [Zha+03] L.J. Zhang, B. Li, T. Chao, and H. Chang. “On demand Web services-based business process composition”. In: *ICSMC*. Vol. 4. IEEE. 2003, pp. 4057–4064 (cit. on p. 41).

## A Operational semantics

Set of rules for the transition relation  $\hookrightarrow$

Let  $mpick = pick(\biguplus_{i=1}^n S_i \Rightarrow P_i, \biguplus_{j=1}^k alrm(a_j) \Rightarrow Q_j)$

$$\frac{}{(v, sInv(S)_x, C, D) \xrightarrow{\langle rSInv \rangle} (v', Stop, (x = \lambda_S) \wedge C^\uparrow, D + \lambda_S)} [rSInv]$$

$$\frac{}{(v, rec(S)_x, C, D) \xrightarrow{\langle rRec \rangle} (v', Stop, (x = \lambda_S) \wedge C^\uparrow, D + \lambda_S)} [rRec]$$

$$\frac{}{(v, reply(S)_x, C, D) \xrightarrow{\langle rReply \rangle} (v', Stop, (x = 0) \wedge C^\uparrow, D)} [rReply]$$

$$\frac{}{(v, aInv(S)_x, C, D) \xrightarrow{\langle rAInv \rangle} (v', Stop, (x = 0) \wedge C^\uparrow, D)} [rAInv]$$

$$\frac{}{(v, mpick_x, C, D) \xrightarrow{\langle (rPickM, i) \rangle} (v', P_i, (x = \lambda_i) \wedge idle(mpick_x) \wedge C^\uparrow, D + \lambda_i)} [rPickM]$$

$$\frac{}{(v, mpick_x, C, D) \xrightarrow{\langle (rPickA, j) \rangle} (v', Q_j, (x = a_j) \wedge idle(mpick_x) \wedge C^\uparrow, D + a_j)} [rPickA]$$

$$\frac{v(b) = \perp}{(v, A \triangleleft b \triangleright B, C, D) \xrightarrow{\langle rCond1 \rangle} (v', A, C, D)} [rCond1]$$

$$\frac{v(b) = \perp}{(v, A \triangleleft b \triangleright B, C, D) \xrightarrow{\langle rCond2 \rangle} (v', B, C, D)} [rCond2]$$

$$\frac{v(b) = true}{(v, A \triangleleft b \triangleright B, C, D) \xrightarrow{\langle rCond3 \rangle} (v', A, C, D)} [rCond3]$$

$$\frac{v(b) = false}{(v, A \triangleleft b \triangleright B, C, D) \xrightarrow{\langle rCond4 \rangle} (v', B, C, D)} [rCond4]$$

$$\frac{(v, A, C, D) \xrightarrow{seq} (v', A', C', D'), A' \neq Stop}{(v, A ; B, C, D) \xrightarrow{seq+\langle rSeq1 \rangle} (v', A' ; B, C', D')} [rSeq1]$$

$$\frac{(v, A, C, D) \xrightarrow{seq} (v', Stop, C', D')}{(v, A ; B, C, D) \xrightarrow{seq+\langle rSeq2 \rangle} (v', B, C', D')} [rSeq2]$$

$$\frac{(v, A, C, D) \xrightarrow{seq} (v', A', C', D')}{(v, A \parallel B, C, D) \xrightarrow{seq+\langle rFlow1 \rangle} (v', A' \parallel B, C' \wedge idle(B), D')} \quad [ rFlow1 ]$$

$$\frac{(v, B, C, D) \xrightarrow{seq} (v', B', C', D')}{(v, A \parallel B, C, D) \xrightarrow{seq+\langle rFlow2 \rangle} (v', A \parallel B', C' \wedge idle(A), D')} \quad [ rFlow2 ]$$

Last modification to this document: March 19, 2020.