

Validation of web service compositions

L. Baresi, D. Bianculli, C. Ghezzi, S. Guinea and P. Spoletini

Abstract: Web services support software architectures that can evolve dynamically. In particular, in this paper the focus is on architectures where services are composed (orchestrated) through a workflow described in the business process execution language (BPEL). It is assumed that the resulting composite service refers to external services through assertions that specify their expected functional and non-functional properties. On the basis of these assertions, the composite service may be verified at design time by checking that it ensures certain relevant properties. Because of the dynamic nature of web services and the multiple stakeholders involved in their provision, however, the external services may evolve dynamically, and even unexpectedly. They may become inconsistent with respect to the assertions against which the workflow was verified during development. As a consequence, validation of the composition must extend to run time. In this work, an assertion language, called assertion language for BPEL process interactions (ALBERT), is introduced; it can be used to specify both functional and non-functional properties. An environment which supports design-time verification of ALBERT assertions for BPEL workflows via model checking is also described. At run time, the assertions can be turned into checks that a software monitor performs on the composite system to verify that it continues to guarantee its required properties. A *TeleAssistance* application is provided as a running example to illustrate our validation framework.

1 Introduction

Service-oriented architectures (SoAs) recently emerged as a useful architectural paradigm in new and innovative computing domains, like ambient intelligence, context-aware applications and pervasive computing [1]. Many current technologies can be associated with SoAs, such as Jini [2], OSGi [3] and so on. However, because of substantial investments by important industrial players, such as BEA, IBM, Microsoft and Oracle, important open-source communities such as Apache and a very active research community, it is common to identify SoAs with the web-based implementations that go under the term ‘web services’.

The research we describe here focuses on web services. In particular, it deals with service compositions built using the business process execution language (BPEL) [4]. BPEL workflows (also called processes) may define new, composite services by coordinating (‘orchestrating’) external services, which are typically not under their jurisdiction. This leads to a distributed ownership of the composite system, which is ultimately responsible for its overall functionalities and quality of service (QoS). External partner services, which affect both functionalities and QoS of composite services, may in fact evolve independently, and even unexpectedly, even after the system is deployed.

When a composite service is designed, certain assumptions must be made on the external services which will be orchestrated. Not only must the syntax of their interface be considered, but also their semantics. In particular, the designer must decide which properties must be fulfilled by the external services and, in turn, based on these assumptions, which properties the composite service will guarantee to its own users.

In the dynamic world of SoAs, however, what is guaranteed at development time, unfortunately, may not be true at run time. The actual services to which the workflow is bound may change dynamically (In BPEL, only the implementation behind partner services can change, but there are many proposals [5] to complement BPEL with dynamic binding capabilities), possibly in an unexpected way that may cause the implemented composition to diverge from the assumptions made at design time. Traditional approaches, which restrict validation to being a design time activity, are no longer valid in this dynamic setting. Besides performing design-time validation, it is also necessary to perform continuous run-time validation to ensure that the required properties are maintained by the operational system. However, it is virtually impossible to predict all the evolutions and changes that might occur in the services we use, and the same is true for the environment. This leads us to consider monitoring as a defensive means. Since it is currently unrealistic to believe that external services will provide a formal and machine-readable specification of their functionality and QoS, all we can rely on are our process-side expectations of their functionality and QoS. Therefore monitoring allows us to take notice of infringements of such expectations.

In this article we propose a framework for validating BPEL processes that covers both design-time and run-time validations. Properties are expressed in assertion language

for BPEL process interactions (ALBERT). ALBERT assertions can be used for two purposes. First, they can formally specify the properties that partner services are required to fulfil. Such properties formalise the assumptions on the external services made at development time by the software developer while designing the workflow. These are called *assumed assertions* (AAs). Second, ALBERT assertions can be used to state properties that the workflow should satisfy, assuming that external services operate as specified. Such assertions, which characterise the behaviour the composite service should guarantee, are called *guaranteed assertions* (GAs). The use of assume-guarantee reasoning is a known technique in verification. It is used to support ‘divide and conquer’ and compositional reasoning [6]. AAs and GAs can state both functional and non-functional properties. Because of this use of assertions, ALBERT promotes ‘design by contract’, as advocated by Meyer [7].

Our validation environment supports the software designer at design time by verifying that GAs hold for a given workflow, assuming that AAs hold. This verification is achieved via model checking. We also provide a run-time monitoring facility that checks whether the external services satisfy the AAs and whether GAs also hold.

When designing our framework, we adhered to the following design principles:

- Use of standard technology. We decided to use standard technology (e.g. BPEL, XML, XPath and so on) to favour adoption of the proposed approach.
- Separation of concerns. The process designer should concentrate separately on the business logic implemented by the workflow and on the validation properties expressed in ALBERT. The two are kept in two separate documents. This is also true for the enactment of the workflow. The adoption of an aspect-oriented approach in the implementation allows the monitoring logic to be kept separate from the business logic.
- Defensive design. Because external services can evolve dynamically, the assumptions (AAs) made on the environment when the abstract workflow is statically validated may be violated at run time. ALBERT can be used to declaratively state the properties that must hold and then to support the monitoring of these properties at run time.

The main contribution of this article is the description of a complete and coherent framework for validating BPEL process compositions. It consolidates our previous work which investigated different aspects of SoAs, and in particular static and dynamic analysis of web services. The ALBERT language and the overall validation framework represent the novel contributions of the article. To the best of our knowledge, although many existing research efforts deal with different aspects of our approach, none provides a complete and coherent coverage of validation of web service compositions and the specification of both functional and non-functional properties.

The article is organised as follows. Section 2 gives an overview of BPEL to make the reader familiar with the language constructs which are then used throughout the paper. Section 3 introduces a running example, which will be used to illustrate the concepts and tools we provide. Section 4 introduces the specification language ALBERT. Section 5 describes the overall validation methodology and how it fits into a complete development process. Section 6 describes our model checking approach and Section 7 describes how we achieve run-time monitoring. Section 8 surveys the state of the art. Finally, Section 9 draws some conclusions and outlines future work directions.

Activity	Shape	Activity	Shape	Activity	Shape
<i>receive</i>		<i>wait</i>		<i>pick</i>	
<i>invoke</i>		<i>terminate</i>		<i>flow</i>	
<i>reply</i>		<i>sequence</i>		<i>fault handler</i>	
<i>assign</i>		<i>switch</i>		<i>event handler</i>	
<i>throw</i>		<i>while</i>		<i>compensation handler</i>	

Fig. 1 Graphical notation of BPEL

2 Overview of BPEL

BPEL [4] is a high-level XML-based language for the definition and execution of business processes by means of web service-based workflows. The definition of a process contains a set of global variables and the workflow logic expressed as a composition of *activities* (Fig. 1 shows the graphical notation we use in the rest of the paper); variables and activities can be defined at different visibility levels within the process using the *scope* construct.

Activities include primitives for communicating with other services (*receive*, *invoke*, *reply*), for executing assignments (*assign*), for signalling faults (*throw*), for pausing (*wait*) and for stopping the execution of the process (*terminate*). The *sequence*, *while* and *switch* constructs provide standard control structures to order activities, define loops and branches. The *pick* construct is peculiar to the domain of concurrent and distributed systems and waits either for the first out of several incoming messages to occur or for a time-out alarm to go off, to execute the activities associated with such event.

The *flow* construct supports the concurrent execution of activities. Synchronisation among the activities of a flow may be expressed using the *link* construct; a *link* can have a guard, which is called *transitionCondition*. Since an activity can be the target of more than one *link*, it may define a *joinCondition* for evaluating the *transitionCondition* of each incoming *link*. By default, if the *joinCondition* of an activity evaluates to false, a fault is generated. Alternatively, BPEL supports *dead path elimination* (DPE), to propagate a false condition rather than a fault over a path, thus disabling the activities along that path.

Each *scope* (including the top-level one) may contain the definition of the following handlers:

- An *event handler* reacts to an event by executing – concurrently with the main activity of the *scope* – the activity specified in its body. In BPEL there are two types of events: message events, associated with incoming messages, and alarms based on a timer.
- A *fault handler* catches faults in the local *scope*. If a suitable *fault handler* is not defined, the fault is propagated to the enclosing *scope*.
- A *compensation handler* restores the effects of a previously completed transaction. The *compensation handler* for a *scope* is invoked by using the *compensate* activity, from a *fault handler* or compensation handler associated with the parent *scope*.

3 Running example

TeleAssistance (TA) is a small company in the business of remote assistance to patients. Its server runs the BPEL

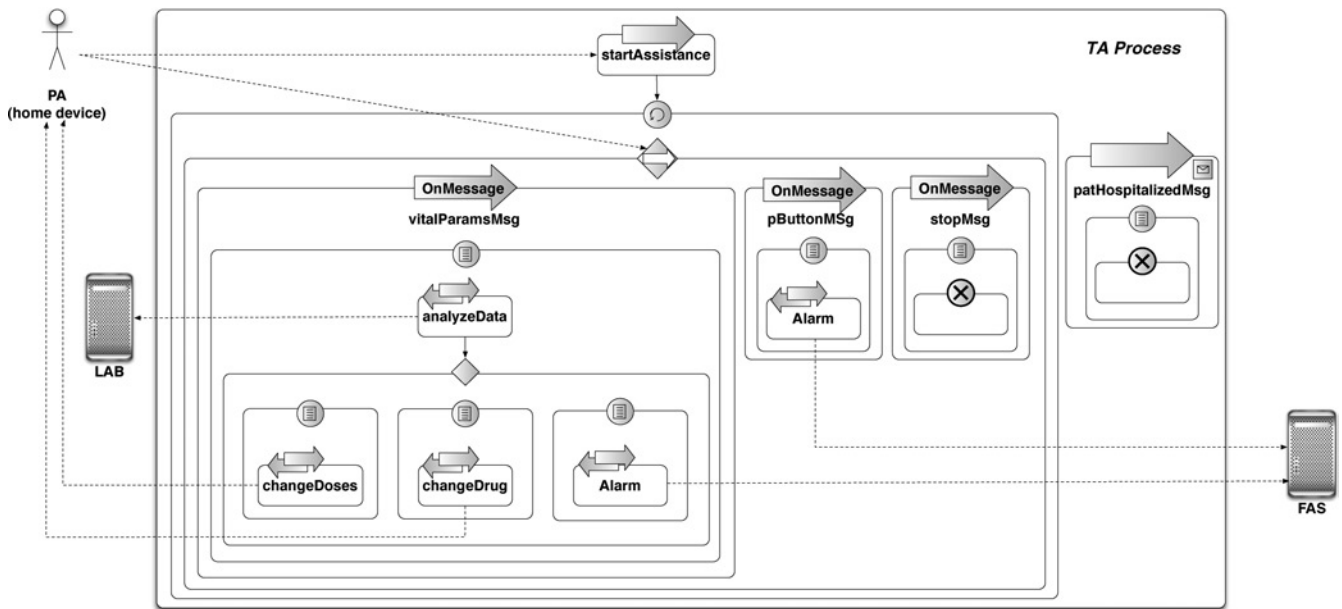


Fig. 2 The TeleAssistance service

process shown in Fig. 2 to assist its clients. The process starts as soon as a *Patient* (PA) enables the home device supplied by TA, which sends a message to the process' *receive* activity *startAssistance*. Then, it enters an infinite loop: every iteration is a *pick* activity that suspends the execution and waits for one of the following three messages:

- *vitalParamsMsg*. The home device (e.g. a glucometer) sends the patient's vital parameters (which are saved in the variable *vitalParams*, which has a field named *glucose*, containing the measured glucose level). This message enables the corresponding execution path, in which the vital parameters are sent to the service *Medical Laboratory* (LAB), by invoking operation *analyzeData*. The LAB is in charge of analysing the data and replies by sending a result value stored in a variable *analysisResult*. This variable contains a field *suggestion* whose value can be 'changeDrug', 'changeDoses' or 'sendAlarm'. In the last case, the TA process invokes service *First-aid Squad* (FAS). This service coordinates a group of doctors, nurses and paramedics who assist patients day-by-day and visit them at home in case of emergency. To alert the squad, the TA process invokes the operation *alarm* on the FAS, by passing the id of the patient and the severity level of the alarm ('mild' in this case).
- *pButtonMsg*. The patient can press a panic button, causing an alarm message to be sent to the TA process. If the patient feels sick, he or she can issue a request for immediate assistance. The TA process alerts the FAS with an alarm, whose severity level is 'high'.
- *stopMsg*. The patient may decide to cancel the TA service.

We assume that the TA process uses a variable *alarmNotif* to send an alarm notification to the FAS. This variable contains a field *level*, which can be set to 'mild' or 'high', and a field *pID* which represents the patient identification code.

If the patient needs special-purpose assistance, the FAS transfers the patient to the closest hospital. Upon arrival, the FAS notifies TA that the patient has been hospitalized by sending a message (*patHospitalizedMsg*), which

is received through an appropriate event handler, saving the parameter in variable *patHospitalized*. The current instance of the TA process then terminates.

4 ALBERT

ALBERT is an assertion language for BPEL processes. It is a reminiscent of assertion languages that were designed for specific programming languages such as ANNA [8], an annotation language for Ada, and the Java modelling language (JML) [9]. ALBERT supports the definition of AA and GA that state both functional and non-functional properties of BPEL processes.

For example, in our TA process we may define a number of assertions (AAs) that state the assumptions made on the partner services, upon which we base the design of the process. One of such properties (referred to as *VitalParams*) specifies that 'the glucose value sent by the patient's remote device to the process is between 40 and 300 mg/dL'. Another AA (referred to as *FASConfirmHospitalization*) specifies that 'if the FAS is invoked three times over a week, with a 'high' severity level alarm for a certain patient, it must notify the TA, within one day that the patient has been hospitalised.'

The language can also specify assertions (GAs) that must be satisfied by the workflow, assuming that external services behave as specified by AAs. One such property (referred to as *FASInvokeMildAlarm*) specifies that 'after receiving a message from the LAB, indicating that an alarm must be issued to the FAS, the TA process must invoke the FAS service within 4 h, passing a 'mild' alarm notification'.

Another example of a GA (referred to as *MDCheckUp*) specifies that 'if a certain patient sends the *pButtonMsg* three times during a span of a week, the FAS must hospitalise the patient within one day'. Intuitively, the truth of this assertion is assured by the AA *FASConfirmHospitalization*, combined with the structure of the workflow, shown in Fig. 2.

4.1 Variables

Variables used in ALBERT assertions can be of two kinds: internal and external. Internal variables consist of

elementary data (i.e. a number, a string or a boolean) that are extracted from BPEL variables. For example, we can refer to a patient's glucose level as (`$vitalParams/glucose`). In this expression, `$vitalParams` indicates the BPEL variable from which we are extracting an elementary value, and `/glucose` is the XPath expression used to extract the desired value.

ALBERT also provides means to predicate on variables whose values originate outside the process. This is useful when the correctness of a property can only be established by referring to contextual data provided by external data sources, such as the time and/or place of execution. These are called external variables. For example, an external variable can be used to define the following AA property: 'The new drug sent by the LAB service through variable `analysisResult` must be amongst those in the list of drugs approved by the Food and Drug Administration (FDA)'. To check if the drug is valid, one needs to refer to an external variable that is provided by querying the *FDA online registry*, an external service which is not part of the workflow. In our example the external variable could be defined as `FDA::inList(ins)/result`, where `FDA::inList` is the invoked remote method, `ins` the input message for the method (in our example it contains the name of the drug) and `/result` the XPath expression used to extract the desired datum.

4.2 Constructs

In this section we provide an introduction to the main constructs of the ALBERT language. The language defines formulae which specify invariant assertions for the workflow. Formulae are defined by the grammar shown in Fig. 3.

In the grammar `id` is an identifier, `var` an internal or external variable, `onEvent` an event predicate, `Becomes`, `Until`, `Between` and `Within` the temporal predicates, `count`, `elapsed`, `past` and all the functions derivable from the non-terminal `fun` the temporal functions of the language. Parameter μ identifies the start or the end of an *invoke* or *receive* activity, the reception of a message by a *pick* or an *event handler*, or the execution of any other BPEL activity. K is a positive real number, n a natural number and `const` a constant.

The grammar in Fig. 3 defines the core of the language. To improve its expressiveness, other constructs are also provided, including the obvious logical connectives, which can be trivially derived from \neg and \wedge and the temporal operators *Always* and *Eventually*, which can be derived from *Until*. Although, in principle, universal and existential quantifiers should not be part of the core of the language, because they predicate over finite sets of data values, we introduce for notational convenience and use them extensively in the rest of the paper.

Here we present the semantics of the language informally, assuming that the workflow process does not contain a *flow* activity. This is not a limitation of the language; it allows us to simplify our presentation. The formal semantics for the complete core language is reported in Appendix A.

The informal meaning of ALBERT formulae can be explained by referring to the sequences of (time-stamped) states of the BPEL process. A state is a triple (V, i, t) , where V is a set of \langle variable, value \rangle pairs, i a label of a BPEL instruction and t a time instant in the domain of positive real numbers. V is the set of \langle variable, value \rangle pairs that hold after executing the BPEL activity i , and t is the instant at which the execution of the activity is completed. Two states $s_j = (V_j, i_j, t_j)$ and $s_{j+1} = (V_{j+1}, i_{j+1}, t_{j+1})$ are adjacent in the sequence if i_{j+1} is the activity that follows i_j in the control flow and the execution of i_{j+1} on the variables in V_j terminates at time t_{j+1} , yielding V_{j+1} .

Boolean, relational and arithmetic operators have the conventional meaning; the same is true for quantifiers. Because ALBERT assertions are implicitly assumed to be invariant for the BPEL process, they express properties that must hold in all states. To express the fact that they must hold when the execution reaches a given point of the workflow, we need to use the predicate *onEvent*, which is true when the corresponding event occurs. For example, property `VitalParams` can be expressed as

$$\begin{aligned} & \text{onEvent}(\text{vitalParamsMsg}) \rightarrow \\ & (\text{\$vitalParams/glucose} \geq 40 \wedge \\ & \text{\$vitalParams/glucose} \leq 300) \end{aligned}$$

More precisely, in the case of *assign*, *pick*, *event handler* and the end of *invoke* or *receive* activities, it is true in a state whose label identifies the corresponding activity. In the case of the start of an *invoke* or *receive* activity, it is true in a state if the label of the next state in the sequence identifies the corresponding activity. In the case of a *while* or a *switch* activity, it is true in the state where the condition is evaluated.

Temporal predicate *Becomes* is evaluated on two adjacent elements of the sequence of states. The formula is true when its argument is true in the current state and false in the previous. The temporal predicate *Until*(ϕ, ξ) is true in a given state if ξ is true in the current state, or eventually in a future state, and ϕ holds in all the states from the current (included) until that state (excluded). The temporal predicate *Between*(ϕ, ξ, K) is true in a given state $s_j = (V_j, i_j, t_j)$ if ϕ is true, for the first time, in a state $s_k = (V_k, i_k, t_k)$ such that $t_k \geq t_j$, and ξ is true in a further subsequent state $s_w = (V_w, i_w, t_w)$ such that $t_w - t_k \leq K$, and for the successor state $s_{w+1} = (V_{w+1}, i_{w+1}, t_{w+1})$, $t_{w+1} - t_j > K$. The temporal predicate *Within*(ϕ, K) is evaluated on a finite sequence of states, built from (and including) the current state $s_j = (V_j, i_j, t_j)$ until we reach a subsequent state $s_k = (V_k, i_k, t_k)$ for which $t_k - t_j \leq K$, and for the successor state $s_{k+1} = (V_{k+1}, i_{k+1}, t_{k+1})$, $t_{k+1} - t_j > K$. The predicate is true if ϕ is true at least in one of these states.

Function *past*($\psi, \text{onEvent}(\mu), n$) is computed on a historical sequence of states, built backwards from (and excluding) the current state. The sequence must contain n states in which *onEvent*(μ) is true. The function returns the value of ψ in the state of the sequence with the smallest t . The function is undefined if such a sequence cannot be

```

 $\phi ::= \psi \text{ relop } \psi \mid \neg\phi \mid \phi \wedge \phi \mid (\text{op id in var} ; \phi) \mid \text{onEvent}(\mu) \mid \text{Becomes}(\phi) \mid \text{Until}(\phi, \phi) \mid \text{Between}(\phi, \phi, K) \mid \text{Within}(\phi, K)$ 
 $\psi ::= \text{var} \mid \psi \text{ arop } \psi \mid \text{const} \mid \text{past}(\psi, \text{onEvent}(\mu), n) \mid \text{count}(\phi, K) \mid \text{count}(\phi, \text{onEvent}(\mu), K) \mid \text{fun}(\psi, K) \mid \text{fun}(\psi, \text{onEvent}(\mu), K) \mid \text{elapsed}(\text{onEvent}(\mu))$ 
relop ::= <|≤|=|≥|>
op ::= forall|exists
arop ::= +|-|×|÷
fun ::= sum|avg|min|max|...

```

Fig. 3 Grammar of ALBERT formulae

found. Function $count(\phi, K)$ is also computed on a finite historical sequence of states, built backwards from (and including) the current state $s_j = (V_j, i_j, t_j)$ until we reach a state $s_k = (V_k, i_k, t_k)$ for which $t_j - t_k \leq K$, and for the predecessor state $s_{k-1} = (V_{k-1}, i_{k-1}, t_{k-1})$, $t_j - t_{k-1} > K$. The function returns the number of elements in this sequence in which ϕ holds. The overloaded version of the function ($count(\phi, onEvent(\mu), K)$) only considers states in which $onEvent(\mu)$ is true (The overloaded version does not add expressive power. Indeed, $count(\phi, onEvent(\mu), K)$ is equivalent to $count(\phi \wedge onEvent(\mu), K)$. The overloaded version is kept to simplify the specification of formulae).

The placeholder fun stands for any function (e.g. average, sum, minimum, maximum ...) that can be applied to sets of numerical values. As for $count$, there are two overloaded cases. $fun(\psi, K)$ is computed on a finite historical sequence of states, built backwards from (and including) the current state $s_j = (V_j, i_j, t_j)$ until we reach a state $s_k = (V_k, i_k, t_k)$ for which $t_j - t_k \leq K$, and for the predecessor state $s_{k-1} = (V_{k-1}, i_{k-1}, t_{k-1})$, $t_j - t_{k-1} > K$. The function returns the value resulting from the application of function fun to all values of the expression ψ in all states of the sequence. The overloaded version, which only considers states in which $onEvent(\mu)$ is true, as before, does not add expressive power to the language. Function $elapsed(onEvent(\mu))$ is computed on a finite historical sequence of states, built backwards from (and including) the current state $s_j = (V_j, i_j, t_j)$ until we reach the first state $s_k = (V_k, i_k, t_k)$ in which $onEvent(\mu)$ is true. The function returns $t_j - t_k$. The function is undefined if such a sequence cannot be found.

ALBERT can be used to specify both AAs and GAs for BPEL processes. However, when defining AAs: formulae should only refer to the BPEL activities that are responsible for interacting with external services. Typically, AAs express properties that must hold after the workflow has completed an interaction with an external service. The following formulae are common templates for AAs.

- $onEvent(\mu) \rightarrow \phi$
- $past(\psi', onEvent(\mu), n) = \psi \rightarrow \phi$
- $Becomes(count(\phi', onEvent(\mu), K) = \psi) \rightarrow \phi$
- $Becomes(fun(\psi', onEvent(\mu), K) = \psi) \rightarrow \phi$

where ϕ and ϕ' are ALBERT formulae, μ identifies the start or the end of an *invoke* or *receive* activity or the reception of a message by a *pick* or an *event handler*. ψ and ψ' are ALBERT expressions, n a natural number and K a positive real number.

These templates limit the states on which ϕ needs to be true to satisfy the formula. More precisely, the first template checks the truth value of ϕ only in the states in which $onEvent(\mu)$ is true; that is in the states preceding or following an interaction with an external service. In the second template, the fact that a property is checked depends on past interactions with the outside world. More precisely, we check ϕ only if, in association with a past interaction with an external service ($onEvent(\mu)$), ψ' was equal to ψ . In the third template, the property is checked if past interactions with the outside world have led to a certain number of specific events. More precisely, we are interested in checking ϕ in a state s , if, in s , it becomes true that, in the last K time instants, ϕ' is true ψ times. Notice that, when counting, we consider only states that are related to interactions with external services ($onEvent(\mu)$), meaning that we count the number of times in which, in relation to these interactions, ϕ' is true. In the last template, the property is checked if past interactions with the outside world have led to a certain value of an aggregate

function. In more detail, we are interested in checking ϕ in a state s , if, in s , it becomes true that fun , calculated over the values of ψ' , obtained from states associated with interactions with external services ($onEvent(\mu)$) over the last K time instants, is ψ . In all four cases, the decision to check ϕ depends on interactions with external services.

4.3 Examples

Referring to the example of Fig. 2, a non-functional assertion (hereafter referred to as **LabServiceTime**) could be 'After sending the patient's data to the LAB service, it should reply within 1 h'. This is an AA on the response time of the external LAB service that can be expressed as

```
onEvent(start_analyzeData) →
  Within(onEvent(end_analyzeData, 60)
```

Another non-functional assertion (hereafter referred to as **AverageLabServiceTime**) could be 'The average response time of all the invocations of operation `analyzeData` on service LAB completed in the past 10 h should be less than 45 min'. This AA can be expressed as

```
avg(elapsed(onEvent(start_analyzeData)),
  onEvent(end_analyzeData, 600) ≤ 45
```

Property **FASInvokeMildAlarm** can be expressed as

```
onEvent(end_analyzeData) ∧
  $analysisResult/suggestion
  = 'sendAlarm' →
  Within($alarmNotif/level = 'mild' ∧
  onEvent(start_alarm), 240)
```

In this case we are defining a GA which states that upon ending the execution of the activity `analyzeData`, if field `suggestion` of the output variable is `sendAlarm`, then the process must guarantee that an `alarmNotif` is sent within 4 h, with the severity level equal to 'mild'.

Property **MDCheckUp** can be expressed as

```
∀x(Becomes(count($alarmNotif/level = 'high' ∧
  x = $alarmNotif/pID,
  onEvent(pButtonMsg), 10080)
  = 3) →
  Within(onEvent(patHospitalizedMsg) ∧
  $patHospitalized/pId = x, 1440))
```

In the example, we count the number of times the `pButtonMsg` is received within a week; if it is received three times, the TA should receive the confirmation of hospitalisation within 24 h, by processing a `patHospitalizedMsg` event.

5 ALBERT-aware development process

Fig. 4 describes how our validation framework fits into a complete development process. The first step consists in designing the service composition, represented as a BPEL process; this task is usually accomplished by a BPEL designer, that is, an expert in business processes modelling.

The BPEL document produced by the design phase is passed to the verification and validation (V&V) engineer, who annotates it with ALBERT assertions. These assertions may be:

- AAs, which represent the assumptions made on the external services, that is, process-side expectations of the QoS and functionality the external services will provide;
- GAs, which state the properties that the workflow should satisfy, if the AAs hold.

The BPEL process, annotated with ALBERT assertions, is then provided to BPEL2BIR, which outputs a model suitable for verification through a model checker (Bogor, in this case). The design and design-time validation phases are continuously repeated until the GAs are satisfied.

Then, the BPEL process is deployed to an execution engine running Dynamo, our monitoring framework. Dynamo checks during run time if the assertions (both the GAs and the AAs) are satisfied.

If the assertions are violated, the process can be verified again, by relaxing the assumptions, or be redesigned and re-deployed.

The next two sections present the technical details of the two validation techniques we propose, including a qualitative evaluation of performance issues.

Both of the analysis techniques involve relevant complexity issues that could affect the overall performance of the proposed methodology. However, the technical details we present in the following sections allow us to deal with them efficiently, as we shall see from a qualitative stand point.

6 Design-time validation

In this section we describe our approach to design-time validation based on model checking. Model checking [6] is a completely automatic technique in which the state space

of the model representing the system under verification is exhaustively analysed. As a consequence, the model has to be finite in the number of states and can only contain variables that have finite sets of values. Since the state space must be completely explored, model checking suffers the well-known state explosion problem, which can be mitigated by the introduction of carefully crafted abstractions.

Several abstraction mechanisms have been proposed for programs written in common programming languages such as C (see, for example [10, 11]). Hereafter, we discuss how ALBERT and the assume-guarantee design methodology it enforces can support the definition of abstractions that can help in the verification of web service compositions.

According to our approach, we analyse workflow-based service compositions by abstracting the external environment (i.e. external partner services) through interfaces to external services, viewed as black boxes that fulfil certain desired properties, formalised through AAs. Often, these abstract interfaces are all, which one can know about their environment. Sometimes, however, some of the external partner services are shared within a restricted cooperating community. In such a case, it may be possible to inspect the black box. The design-and-conquer methodology supported by ALBERT allows the implementation of these partner services also to be analysed, to check what their users require as AAs are guaranteed by the implementations as GAs. Indeed, this is the essence of an assume-guarantee approach to design-time verification. ALBERT allows the approach to deal not only with functional requirements, but also with QoS agreements that bind service requesters and service providers, such as response time constraints.

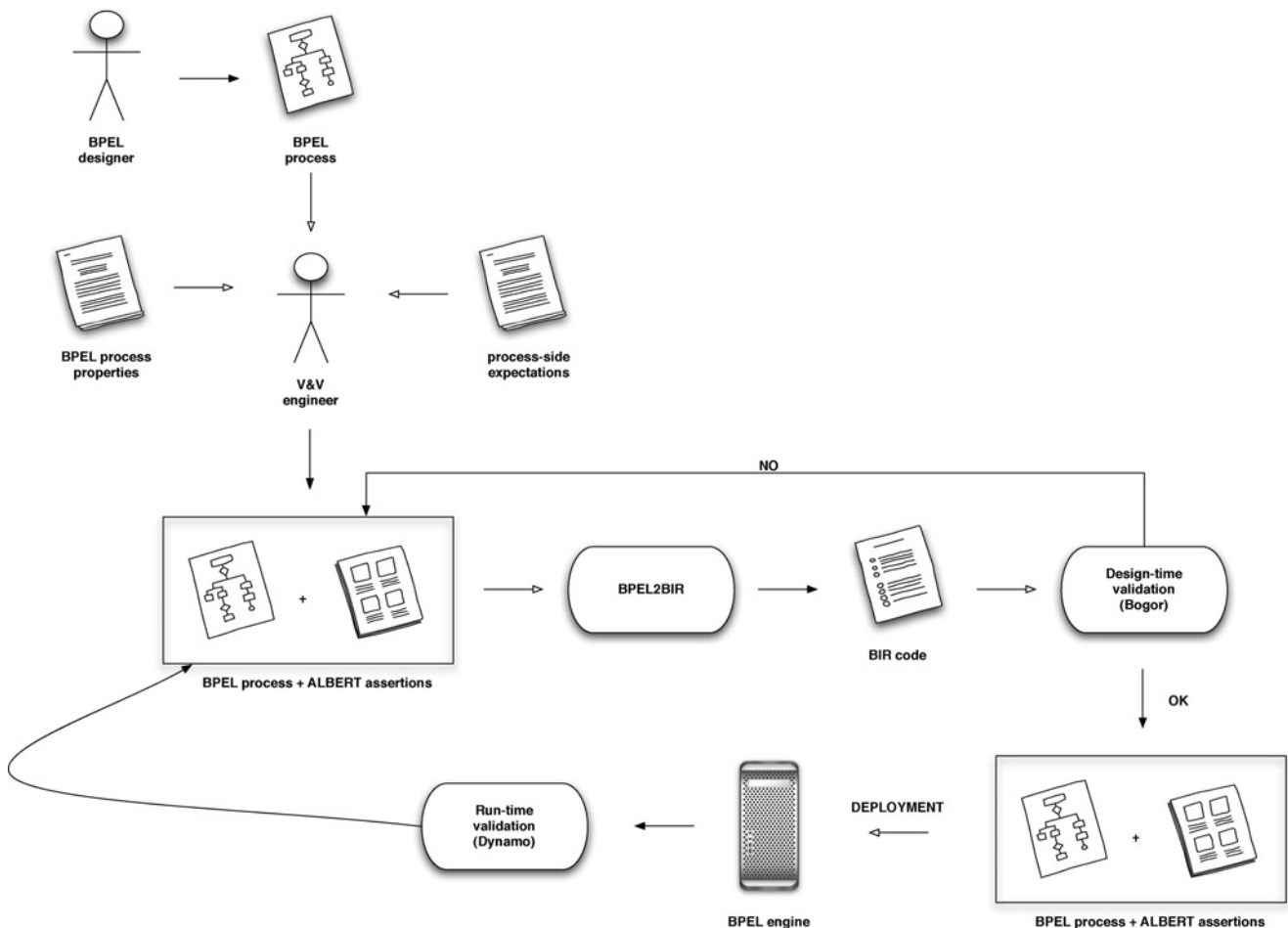


Fig. 4 ALBERT-aware development process

To support analysis, we developed BPEL2BIR [BPEL and ALBERT to Bogor's intermediate representation (BIR)], a model checking framework based on Bogor [12]. Our validation framework and the ALBERT language are theoretically independent of the selected model checking and monitoring technologies. In particular, Bogor was chosen because most BPEL constructs can be easily mapped onto BIR (Bogor's input language). Furthermore, Bogor's modular architecture supports the introduction of different model checking algorithms and customisations for particular domains. This will be exploited in the future evolution of this work, as outlined in Section 9.

To develop BPEL2BIR, we had to solve three main problems. First, we build a model from the BPEL workflow by modelling the interactions with the external world through the generation of random values for the relevant variables. Second, we exploit the abstractions that can be derived from AAs, which allow for a better representation of the interactions with the external world. Then, we translate GAs into BIR properties that must be verified. The model-checking problem is described in the following by presenting how we encode BPEL activities, AAs and GAs into BIR constructs.

The approach to model check web service compositions presented in this article differs from other proposals that appeared in the literature because:

- It is supported by an assume-guarantee verification methodology that fosters design-and-conquer.
- It covers the whole set of BPEL constructs, including those dealing with time.
- ALBERT allows one to describe a rich set of functional and non-functional properties.

Preliminary experimental results on the use of BPEL2BIR have shown that Bogor can provide a better support for model checking than other similar tools [13].

6.1 Bogor

Bogor [12] is a model-checking framework developed at Kansas State University. The input language for Bogor, called BIR, provides constructs found in modern programming languages, such as dynamic threads and object creation, exception handling, virtual functions, recursive functions and garbage collection. A low-level version of the intermediate representation, named low-level BIR, is a language for the description of transition systems based on guarded commands with explicit locations, explicit guards, sequences of statements comprising transition actions and explicit transitions.

The BIR data model contains primitive types (`boolean`, `int`, `enum`) and non-primitive types (`null`, `record`, `array`, `lock`). For record types, sub-typing declarations and virtual methods are also supported. The language is statically strongly typed. The memory model forbids pointer arithmetic and achieves object reclamation through garbage collection. To overcome the state explosion problem, Bogor implements some well-known optimisation and reduction strategies, such as data and thread symmetry [14], collapse compression [15] and partial order reduction [16].

Bogor is extensible, both in terms of input language and model-checking algorithms. It has an open, modular architecture that allows the development of extensions via new algorithms and optimisations, to improve core tasks, such as state encoding and state exploration. Several extensions have been implemented; among these, the ones related to partial-order reduction, state-encoding, search strategies, support for different property languages (regular

expressions, LTL and CTL, JML) and support for several domains (multithreading and Swing, event-based Java programs, CORBA-based avionics systems). For example, in our group we designed Bogor extensions to validate event-based service architectures based on the publish/subscribe paradigm [17].

6.2 Modelling BPEL in Bogor

A BPEL process is mapped onto a BIR system composed of threads that model the main control flow of the process and its flow activities. Data types are defined using an intuitive mapping between WSDL message/XML Schema types and BIR primitive/record types. Basic activities are trivially mapped onto their equivalent in BIR. Instead, a *receive* (resp., *invoke*) activity is translated as an assignment to its input (resp., output) variable, since the behaviour of external services is not modelled. If no assumptions are made on the external partner services, the assignment is performed with a non-deterministically generated value, ranging all over its domain. Otherwise, if AAs are provided to constraint the expected behaviour of external services, they generate ad-hoc abstractions as discussed next. Activities contained within a *flow* are translated into threads, preserving both transition and join conditions.

For each scope, *fault handlers* are translated as `try/catch` statements, with `catch(var)` clauses matching exception variables corresponding to faults. *Event handlers* are modelled using a dedicated thread, which non-deterministically consumes the events produced by a helper function.

A *pick* activity waits for the occurrence of one out of several messages delivered by external services. In this case, it is translated by invoking a function that models the occurrence of one of the messages being awaited; the message is then treated as if it had been received through a *receive*. Optionally, a *pick* activity can specify a timeout and a *wait* activity can specify a suspension. They both indicate how long (*for*) or *until* the process needs to wait.

To deal with time constraints in BIR, we need to include in the model an execution time for all BPEL activities represented in BIR and introduce a time counter for each sequential path of activities generated during the execution of the BPEL workflow.

For each activity, we insert a code block in BIR that randomly generates the duration of the activity within a certain interval. If there is an AA that constrains the duration of the activity, we use it to generate the BIR code as discussed in the next section. For a *flow* activity, the time consumed by the flow is the maximum time spent along all paths. Notice that if two activities in different paths are linked, we synchronise the time on this link by assigning to its time counter the maximum value of the counters on the incoming links.

6.3 Assumed assertions

In the previous section, we generated random values to model both the values generated through interactions with an external service and their duration. Now we show how AAs can provide a better abstraction of external environment by reducing the range of the generated values.

As an example, in the TA service, the AA associated with a receipt of message `VitalParamsMsg` (property `VitalParams`) states that the glucose value sent by the remote device must be between 40 and 300 mg/dL. This constraint is used to reduce the size of the domain from which the values are generated by the model checker for variable `$vitalParams/glucose`.

In some cases the formula representing an AA may refer to the value of an expression in the past, such as $past(\psi, onEvent(\mu), n)$. In such a case, it is necessary to log the historical values of the variables appearing in ψ when a location, in which $onEvent(\mu)$ is true, is reached during the execution. Logged data must be retrieved to evaluate $past$. Quantified ALBERT formulae are also used in the optimised generation of input values. The universal quantification is translated into a generation of values of all the quantified variables performed according to the quantified formula, and the existential quantifier is translated into a non-deterministic choice of the variable, whose value is then generated according to the quantified formula.

AAs can express constraints on the duration of activities executed by the invocation of external services. As an example, consider assertion **LabServiceTime**. In such a case, the time bound expressed by the *Within* operator can be assumed as an upper-bound for the duration of the activity whose *start* and *end* events are referred in the formula.

6.4 Guaranteed assertions

GAs specify invariant properties of the workflow. Thus each assertion must be evaluated after the execution of each activity of the BPEL process. The model checker executes the check by instantiating an evaluator for each assertion after executing the BIR code block corresponding to each activity in the workflow.

The evaluation of each assertion can be completed in a single step if it only refers to the present or past states of the computation. If, instead, it contains future temporal operators (such as *Within*), the instance of the evaluator carries on in all subsequent states until it terminates by producing the truth value of the assertion.

ALBERT assertions can refer to present and past values of variables and past sequences of events. To support their evaluation, the BIR system that models the BPEL process should also include book-keeping actions that collect the historical values and other auxiliary variables needed for verification.

Let us describe in more details how an evaluator for an ALBERT assertion works. The evaluation of arithmetic and relational operators is straightforward. Existential (universal) quantifiers over formulae are interpreted as disjunctions (respectively, conjunctions) of formulae. The evaluation of function $past(\psi, onEvent(\mu), n)$ requires accessing the historical values of expression ψ . The values of the variables referred by ψ are stored in an array of n elements, which is kept updated by the book-keeping actions we mentioned above.

To evaluate predicate $onEvent(\mu)$, we rely on a boolean auxiliary variable, which is set to true by the BIR system exactly when μ happens, and to false immediately afterwards. The evaluation of the predicate returns the value of this variable. The evaluation of predicate $Becomes(\phi)$ returns true if ϕ is true in the current state and false in the previous. This is implemented by using an auxiliary variable that contains the previous value of ϕ . Function $elapsed(onEvent(\mu))$ is computed by using a counter variable managed as auxiliary variable by the BIR system. This variable is set to 0 whenever $onEvent(\mu)$ is true in the current state and is incremented by the duration of the last-executed instruction, in each state in which $onEvent(\mu)$ does not hold. When the function is evaluated the value of this variable is returned. The evaluation of function $count$ and of the other functions derivable from **fun** requires an auxiliary array variable that keeps track

of the process state in the last K time units. All such function can be evaluated by using the values stored in the array; the size of this array is finite and limited by the number of activities performed in the last K units.

Let us now describe how future temporal predicates can be evaluated. The evaluation of $Until(\phi, \xi)$ returns false if ϕ is false in the current state and true if ξ is true in the current state; otherwise an evaluator for the formula is started and remains active in all future states until either ϕ is false (in which case the evaluator returns false) or ξ is true (in which case it terminates returning true). The evaluation of $Within(\phi, K)$ requires an evaluator to be started if ϕ is not true in the current state; otherwise it returns true. The evaluator remains active until either ϕ becomes true (in which case it returns true) or K time units elapsed (in which case it returns false). To evaluate predicate $Between(\phi, \xi, K)$, an evaluator is started to check whether ϕ occurs. If this is the case it remains active for K time units and when this time interval elapses, it checks if ξ is true. If it is the case it returns true, otherwise false.

6.5 Performance

The performance of a model checker is influenced by the dimension of its two inputs: the model, which in the context of this work is a BPEL process extended with AAs, and a formula, which in our case is a GA. Although the formula is generally small, the model can be huge and sometimes potentially infinite.

Model checking requires that the systems under analysis be finite, hence the main issue to consider when modelling is the involved data's size. In a BPEL process, variables can vary on a potentially infinite domain that needs to be made finite using abstraction techniques. In our approach we represent the model by randomly generating the data over their domain: hence there is a trade-off between abstraction and precision of data generation. Indeed, the data representation is dually crucial; if we select coarse-grained domains for the variables, the model itself loses generality and significance; on the other hand, a fine-grained domain leads to an exponential blow-up during verification.

AAs may help reduce the range of data exchanged with the environment, by considerably affecting the space required for verification. For instance, by annotating our TA process with the **VitalParams** assertion, we can reduce the size of the `$vitalParams/glucose` variable's domain, and as a consequence, the number of states visited during verification: in our experimentation, it decreased from 422 to 282.

The metric ALBERT induces over time is another main issue that affects performance: we support it by enriching the model with time annotations. The temporal domain we consider is discrete, but not finite. Hence, even though we explicitly consider time, we do not use a global temporal variable to represent it, since this would make verification infeasible, but local timers.

The timers' granularity is set using the state of the process with respect to the GA that must be verified, and to the time guards in the model. More precisely, whenever the passing of time does not affect the model's behaviour, it is abstracted and considered as a single time slot.

7 Run-time validation

When performing model checking, we distinguish between AAs and GAs. The former defines assumptions on the outside world, whereas the latter defines properties that, when statically verified, allow us to state the 'internal'

correctness of the workflow process. When moving from development time to run time, the workflow interacts with real external services. Such services, as we observed, are owned, run and evolved by independent authorities. Compliance of their behaviour with the properties assumed by the workflow at development time is not automatically ensured, and must be checked by monitoring. Moreover, to ensure system robustness, we may also be interested in monitoring the GAs.

Our validation framework proposes a rule-based monitoring approach where AAs and GAs are checked at run time by Dynamo, which is our monitoring infrastructure for dynamic monitoring.

7.1 Monitoring rules

Monitoring rules specify the directives for the monitoring framework, and are made up of two main parts: (1) a set of optional meta-level information called **Monitoring Parameters**, and (2) a **Monitoring Property** expressed in the ALBERT language.

Monitoring Parameters allow our approach to be flexible and adjustable with respect to the context of execution. Each rule is associated with a set of optional monitoring parameters. These are meta-level information used at run time to decide whether the rule should be taken into account or not. By default, if no parameters are given, monitoring is performed. Supported parameters are `priority`, `validity` and `trusted providers`. A `priority` associated with a monitoring rule defines a simple notion of ‘importance’. In our model we support five levels of priority: very low, low, medium, high and very high. When the rule is about to be evaluated, its priority is compared with a threshold value (The threshold value is set by the owner of the process); the rule is taken into account if its priority is less than or equal to the threshold value. By dynamically changing the threshold value we can dynamically set the intensity of probing. A `validity` parameter defines time constraints on when a supervision rule should be considered. The supervision designer can define two different kinds of constraints: time windows and periodicity. The former defines time-frames within which monitoring is performed. When outside of this frame, any new monitoring activities are ignored. Rule checking, however, when started within a valid time-frame, is always completed. Monitoring periodicity, on the other hand, can be specified by using the `every` keyword. Accepted values are durations and dates, for example, ‘every 3D’, meaning every 3 days, or every ‘01/01’, meaning every January 1st. Finally, `trusted providers` is a list of service providers for which supervision is not necessary. This is useful because, in abstract process definitions, the actual service to which the process binds could be chosen at development time or at run time. If a rule refers to at least one non-trusted provider, it is checked.

7.2 Dynamo

Fig. 5 presents the Dynamo execution and monitoring framework, by illustrating the dependencies existing between the various components and the technologies used in the implementation.

The **Configuration Manager** is a storage component for all the ALBERT properties that have been defined. The **ActiveBPEL engine** is a modified version of ActiveBPEL [18] in which we embed monitoring. This is achieved by following an aspect-oriented programming approach [19]. The engine is a Java program in which we

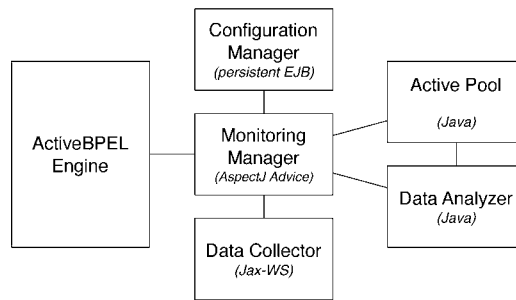


Fig. 5 A run time view of the monitoring framework

weave the cross-cutting monitoring features via AspectJ [20]. ActiveBPEL works by creating an internal tree representation of the process being executed. In this tree, each node represents a single BPEL activity in the process definition, and it is an appropriate extension of the `AEActivityDefinition` class. Each node contains the information necessary to perform the particular activity it is associated with. At run time, the tree is visited and the definition classes are used by the engine to instantiate appropriate `AEActivityImpl` classes, all of which implement a common interface. Amongst other things, this interface provides an `execute` method where the activity’s primary action is performed. For example, a `scope` activity will set up its internal variables, whereas an `invoke` activity will perform the appropriate external invocation. To perform monitoring, we intercept the process after the `execute` method is called for the various BPEL activities. These are the points where the **Monitoring Manager** (implemented as an AspectJ advice) is activated. Its main responsibility is data collection, both from within the process and from the outside world (in the case of external variables, through the **Data Collector**). The **Monitoring Manager** collects all the values of the variables defining the state of the process, and time-stamps and stores them in the **Active Pool**, together with the label of the BPEL activity after which they were collected. In practice, a complete process state is built and saved in the **Active Pool**.

The **Active Pool** is responsible for keeping track of historical sequences of process states. These states are needed by the **Data Analyzer** to perform analysis, and can be obtained in two ways. One way is to use an API, provided by the **Active Pool**, which can return the last state or a historical sequence of states. Another way is to subscribe to new states being generated (via the `publish/subscribe` paradigm), to receive them as they become available.

ALBERT formulae are evaluated by the **Data Analyzer**, which retrieves the data from the **Active Pool**, as described in the next section.

7.3 Evaluating ALBERT properties

This section describes how Dynamo evaluates ALBERT formulae in an intuitive manner. The **Active Pool** and the **Data Analyzer** perform optimisations that are ignored here for simplicity.

The evaluation of ALBERT formulae that do not contain references to the present state and/or to the past history (i.e. formulae that do not contain `Until`, `Between`, or `Within` operators) can be performed by the **Data Analyzer** by retrieving the relevant values from the **Active Pool**. We show that given a formula ϕ that describes an assertion to be monitored at run time, the size of the state history to be kept

by the **Active Pool** is limited. For the sake of simplicity, we first assume that the formulae do not contain nested past operators.

When considering $past(\psi, onEvent(\mu), n)$ functions, the **Active Pool** needs to keep n history states. In the case of $Becomes(\phi)$, the **Active Pool** needs to keep one history state. Assertion ϕ evaluated in such state must be false, and it must be true when evaluated in the current state. When considering $count(\phi, K)$ and $fun(\psi, K)$ functions, K represents a historical time window, going back from the current time. Thus, the **Active Pool** needs to keep all the history states that were collected within this time frame, whose number is of course bound. Therefore the number of states in the **Active Pool** will be the maximum among the maximum of the n_i states needed for the $past_i$ functions, 1 (if there is a *Becomes* predicate), and the maximum number of states needed for the various *count* and *fun* time windows.

To illustrate a case where the formula contains nested operators that refer to the past, we analyse the following expression: $past(past(\psi, onEvent(\mu), a), onEvent(\mu'), b)$ where a and b are two natural numbers. The size of the history sequence to be kept to evaluate the expression is $a + b$. In general, nested operators that refer to the past lead to a bounded growth of the history sequence.

Function *elapsed* needs special attention. For each event μ appearing in the argument of $onEvent(\mu)$, the **Active Pool** keeps one temporal distance variable. The variable is initially undefined, and is set to zero the first time $onEvent(\mu)$ becomes true. The time-stamp is updated upon receiving new states to store the amount of time passed from the last time $onEvent(\mu)$ was true.

The evaluation of formulae that contain *Until*, *Between* or *Within* predicates is more complex. From a theoretical viewpoint, it could be explained by referring to the well-known correspondence between linear temporal logic and alternating automata [21].

Conceptually, the evaluation of these formulae cannot be completed in the current state. Their value, in fact, depends on the values the variables will assume in future states. For this reason, as soon as the **Data Analyzer** needs to evaluate one such subformula, it spawns a new evaluation thread for that subformula. The thread terminates in some future time instant. If the formula comprises subformulae, its evaluation can only be completed when all threads spawned by the evaluation terminate.

Whenever a new state is stored in the **Active Pool**, the **Data Analyzer** is notified. Consequently, each of the threads that was spawned for the evaluation of temporal subformulae is evaluated, according to the following rules:

- If the subformula is of the type $Until(\phi, \xi)$, ξ is evaluated by the thread in the state notified by the **Active Pool** and, if it is false, ϕ is evaluated. If ϕ is also false, the evaluation thread terminates by returning false. Otherwise, the thread continues to evaluate the subformula in future states.
- If the subformula is of the type $Between(\phi, \xi, K)$ a timer is associated with its evaluation thread. The timer is initialised to 0 the first time ϕ evaluates to true in a state notified by the **Active Pool**. As the timer reaches K , the thread terminates by returning the value of ξ in the current state.
- If the subformula is of the type $Within(\phi, K)$, the thread checks the truth of ϕ . If it is true, the thread terminates by returning true. If it is false, a timer (initialised to 0) is associated with the thread. If ϕ becomes true before the timer reaches K , the thread terminates by returning true; otherwise, it terminates by returning false.

Consider for example property **MDCheckUp**, which is recalled here for convenience.

$$\begin{aligned} & \forall x (Becomes(count(\$alarmNotif/level = 'high' \wedge \\ & \quad x = \$alarmNotif/pID, \\ & \quad \quad onEvent(pButtonMsg), 10080) \\ & = 3) \rightarrow \\ & \quad Within(onEvent(patHospitalizedMsg) \wedge \\ & \quad \quad \$patHospitalized/pID = x, 1440)) \end{aligned}$$

The formula includes elements of different nature: two state variables $\$alarmNotif/level$ and $\$alarmNotif/pID$, one *pick* message $pButtonMsg$, and one event $patHospitalizedMsg$ and its associated variable $\$patHospitalized/pID$.

The formula is universally quantified over the patients, hence it is rewritten as a conjunction of formulae of the same structure, where variable x is substituted each time with a different patient. More precisely, if the set of patients is $P = \{p_1, \dots, p_N\}$, the formula is internally considered as

$$\begin{aligned} & \bigwedge_{i=1}^N (Becomes(count(\$alarmNotif/level = 'high' \wedge \\ & \quad p_i = \$alarmNotif/pID, \\ & \quad \quad onEvent(pButtonMsg), 10080) \\ & = 3) \rightarrow \\ & \quad Within(onEvent(patHospitalizedMsg) \wedge \\ & \quad \quad \$patHospitalized/pID = p_i, 1440)) \end{aligned}$$

Each component of the conjunction is monitored by the **Data Analyzer** as follows, taking into account that since there is an external conjunction all the subformulae have to be true.

Operator *Becomes* is evaluated for a single patient by checking the value of its argument in the previous and in the current state. To evaluate such an argument in either state, function *count* must be evaluated, and this requires examining the historical sequence of states that occurred in the previous time-span, starting from the time-stamp of the state we are considering (either the previous or the current). The evaluation of the consequent of the formula is evaluated only when the *Becomes* is true and spawns a thread, which terminates, at the latest, after 1440 time units from the current time value. This is the latest time at which the value of the overall formula becomes known.

7.4 Performance

Many different aspects must be considered when studying the time **Dynamo** takes to verify an **ALBERT** expression. However, before starting our analysis, we must recall that monitoring is achieved mainly asynchronously. As we shall see, in a realm in which long ongoing processes are the norm, this will turn out not to be a real issue.

Monitoring can be broken down into various steps. In the first step the process execution is intercepted to build a new process state and to save it to the **Active Pool**. This is the only step taken synchronously, as the **Monitoring Manager** performs data collection from the process and from external services. Tests performed on an AMD Athlon(tm) XP 2600+(1.93 GHz) with 512 MB of RAM, running Windows XP, show that our system takes less than 2 ms to intercept the execution and to commence data collection. Since the **Monitoring Manager** lives in

the same application space of the executing process, internal variables are collected extremely rapidly leveraging ActiveBPEL's own APIs (i.e. in a time quantifiable in milliseconds). External variables, on the other hand, are a completely different matter. The time needed to obtain data from an external service mostly depends on issues we are not responsible for, such as network-related issues or middleware serialisation and de-serialisation. The more the designer decides to include external variables in his properties, the more this becomes obviously an issue. In general, however, the literature has taught us to expect the use of external variables to be limited with respect to internal variables. Nevertheless, many interesting properties can only be expressed with the help of external data, leaving the decision of how much external variables to use up to the designer.

Once data collection has been completed, the data are time-stamped and sent to the **Active Pool**, and the process is free to proceed. From this point all monitoring activities are performed asynchronously, meaning they have no further impact on performance.

On the other hand, the time needed to actually verify a property depends solely on the property itself. If a property does not contain any *Until*, *Between* or *Within* predicates, it can be evaluated immediately, and in the long running business processes this may very well occur before the next significant business step is taken. Moreover, if any one of these predicates appears, the analysis time will depend on the evolution of the process execution (i.e. the appearance of new process states in the **Active Pool**), and on the timers mentioned in Section 7.3. Regardless of the property, the result will be known in an amount of time that depends linearly on and is bounded by the time constants used in these predicates. When the process terminates, all properties still under analysis are immediately evaluated considering that internal variables can no longer evolve. Properties involving external variables remain pending, for a bounded amount of time, waiting for the external data to become available.

8 Related work

In this section we discuss some related approaches. We first review the existing work on the application of model checking to web service compositions. Then, we move to run-time monitoring. Besides introducing the different proposals, we also describe how our approach is different from the others.

8.1 Model checking

Research on model checking web service compositions is quite recent, but it has attracted considerable attention.

WSAT [22] is a framework for analysing the interactions among composite web services modelled as conversations. BPEL specifications of web services are translated into an intermediate representation, an XPath-guarded automaton augmented with unbounded queues for incoming messages. This model is then translated into Promela and LTL properties, which can also be derived from XPath expressions, are then checked with the SPIN model checker [23].

The Verbus verification framework [24] is based on an intermediate formalism, which decouples the approach from any particular process definition language or verification tool. The support for BPEL is incomplete: *compensation* and *event handlers* are not considered. The current version of the prototype performs reachability analysis and supports the verification of properties like invariants, goals, activity pre- and postconditions, as well as generic properties defined in temporal logic.

Table 1: Comparison of BPEL constructs support among model checking approaches

BPEL constructs	BPEL2BIR	WSAT	Verbus	Nakajima
basic + structured activities ^a	yes	yes	yes	yes
fault handler	yes	yes	yes	no
event handler	yes	no	no	no
compensation handler	yes	no	no	no

^aActivities described in Sections 11 and 12 of [4]

Nakajima [25] proposes a method to extract the behavioural specification from a BPEL process and to analyse it by using the SPIN model checker. A finite-state automaton extended with variable annotations (definitions and updates) is used as an intermediate representation. This approach provides only partial BPEL support, which does not deal with *fault/event/compensation handlers*. The tool checks for deadlocks and verifies user-defined LTL properties.

Table 1 summarises the results of comparing our approach with the three previous SPIN-based verification frameworks in terms of the support they provide to the BPEL language. BPEL2BIR is the only approach that supports all the constructs of the language; all the others have some limitations in dealing with handlers.

Other authors use different computational models for verifying BPEL processes. Schlingloff *et al.* [26] use Petri Nets to define the semantics of BPEL. Validation is performed by using the LoLA [27] model checking tool. Process algebras are used in [28] and [29]. Foster *et al.* [28] verify web service compositions against properties created from design specifications and implementation models. Specifications, in the form of message sequence charts, and implementations, in the form of BPEL processes, are translated into the Finite State Process notation, which is the input language for the Labelled Transition System Analyser (LTSA) model checker. Koshkina and van Breugel use a process algebra, the BPE-calculus, to abstract the BPEL control flow. This calculus is used as input for a process algebra compiler to produce a front-end for the concurrency workbench (CWB) [30], which performs equivalence checking, preorder checking and model checking.

8.2 Monitoring

Several works define specification languages for functional and non-functional properties [usually expressed in the form of a service level agreement (SLA)] and propose an associated monitoring architecture. Sahai *et al.* [31] describe an automated and distributed SLA-monitoring engine. The monitor acquires data from instrumented processes and – by analysing the execution of activities and message passing – then verifies the SLAs. Keller and Ludwig [32] propose a framework to define and monitor SLAs, focusing on QoS properties such as performance and costs. The language defines a type system for various SLA artefacts such as parties, obligations, parameters, metrics and functions. The monitoring component is composed of two services. The first (the measurement service) measures parameters defined in the SLA, by probing client invocations or by retrieving metrics from internal resources. The second (the condition evaluation service) tests measured values against the thresholds defined in the SLA and, in case of a violation, triggers corrective management actions. Skene *et al.* [33] propose the SLAng language for SLAs; in [34] they extend their work by providing a

Table 2: Comparison of monitoring approaches

Approach	Language		Abstraction		Properties	
	Logic	HL/VHL	Domain	Implementation	Safety	Temporal
Sahai <i>et al.</i> [31]	—	x	x	—	—	x
Keller and Ludwig [32]	—	x	x	—	—	x
Skene <i>et al.</i> [33, 34]	—	x	x	—	—	x
Robinson [35]	x	—	x	—	x	x
Mahbub and Spanoudakis [36]	x	—	x	—	x	x
Barbon <i>et al.</i> [37]	x	x	—	x	—	x
ALBERT	x	x	—	x	x	x
	Directives			Timeliness		
	Process	Activity	Event	Post-mortem	Synchronous	Asynchronous
Sahai <i>et al.</i> [31]	x	—	—	—	—	x
Keller and Ludwig [32]	x	—	—	—	—	x
Skene <i>et al.</i> [33, 34]	x	—	—	—	—	x
Robinson [35]	x	—	—	—	x	x
Mahbub and Spanoudakis [36]	x	—	—	x	—	—
Barbon <i>et al.</i> [37]	x	x	x	—	—	x
ALBERT	x	x	x	—	—	x

model and an analysis technique for reasoning about the monitorability of SLAs.

All these approaches focus on formally defining high-level contracts among parties (typically, between a service consumer and a service provider), hence they do not allow to specify properties that should hold on specific events occurring during the execution of the service, such as the completion of a certain activity.

Robinson [35] uses temporal logic and KAOS to express requirements, such as timeliness constraints. These requirements are then analysed to identify conditions under which they can be violated. If such conditions correspond to a pattern of events observable at run time, each of them is assigned to an agent for monitoring. At run time, an event adaptor translates SOAP messages into events and forwards them to the corresponding monitoring agent.

Mahbub and Spanoudakis [36] propose a framework for the run-time verification of requirements of service-based software systems. Requirements can be behavioural properties of a service composition, or assumptions on the behaviour of the different services composing the system. The first can be automatically extracted from the composite process, expressed in BPEL; the latter are specified by system providers using the event calculus. System events are collected at run time and stored in an event database; defined properties are checked by means of an algorithm based on integrity constraint checking in temporal deductive databases.

Barbon *et al.* [37] describe an approach to monitor BPEL compositions. Monitors can be attached to a single instance or to the whole class of process instances; they can check temporal, boolean, time-related and statistics properties, expressed in a run-time monitoring specification language. Business and monitoring logics are kept separated by executing in parallel the monitor engine and the BPEL execution engine; code-implementing monitors are automatically generated from high-level specifications.

A comparison of our approach against the others mentioned above is presented in Table 2. The classification of

the approaches follows the taxonomy presented by Delgado *et al.* [38], with some modifications/extensions of the metrics to adapt them to the service-oriented context. ‘Language’ indicates the type of specification used by the approach (logic or HL/VHL), ‘abstraction’ indicates the abstraction level at which properties are defined (domain or implementation), ‘properties’ is used to indicate the kind of properties definable by the language (safety or temporal), ‘directives’ indicates the level at which a property can be evaluated (process, activity, event), ‘timeliness’ indicates when the monitoring activity is performed (post-mortem, synchronous or asynchronous).

The comparison shows that ALBERT is one of the few logic-based specification languages to fully support BPEL. This means that we can define assertions that predicate both on the whole process and on the single activity or event.

9 Conclusions

SoAs provide unprecedented degrees of dynamism and flexibility to software systems. Independently developed and deployed services are made available dynamically and then composed by third parties to provide new useful features. Web services achieve these goals at the Internet level, supporting dynamic federations of business services. The intrinsically dynamic nature of these systems and the multiple stakeholders involved in their construction and composition, however, challenge our ability to provide dependable solutions. In particular, the traditional boundary between development time, during which applications are carefully validated, and run time, during which systems are operated in the real world, disappears.

Given such a premise, our goals are to provide designers with a coherent validation framework for composite services described in BPEL. We believe that designers can benefit from a language – ALBERT – built from the ground-up for specifying functional and non-functional properties, both for design-time and run-time validations. This is why we have also set out to provide appropriate

model-checking tools and a monitoring-aware execution environment. In conclusion, such a framework allows designers to produce more dependable solutions and to promptly discover whether their systems deviate from an expected and desirable QoS.

Our future work will focus on exploiting the results of run-time monitoring, by providing mechanisms and strategies to react to the detection of undesirable behaviours. Our goal will be to incorporate self-managing features in SoAs to allow service compositions to reorganise themselves to dynamically optimise the overall QoS. Moreover, regarding the design-time validation, we plan to exploit Bogor's extensibility for further development. The CEGAR (Counterexample Guided Abstraction Refinement) [11] loop and predicate abstraction [39] state space reduction techniques – which proved to be highly beneficial when applied to software model checking – may be implemented as Bogor plugins to improve verification efficiency.

10 Acknowledgments

Part of this work has been supported by the IST EU project SeCSE – contract number 511680, the IST EU project PLASTIC – contract number 026955, the Italian FIRB project ART DECO and the FAR project DISCORSO.

11 References

- Baresi, L., Di Nitto, E., and Ghezzi, C.: 'Towards open-world software: issues and challenges', *IEEE Comp.*, 2006, **39**, pp. 36–43
- Jini Network Technology [homepage on the Internet]. Sun Corporation; 2007. Available from: <http://sun.com/jini>
- OSGi [homepage on the Internet]. OSGi Alliance; 2007. Available from: <http://www.osgi.org>
- Andrews, T., Curbera, F., Dholakia, H., Golland, Y., Klein, J., Leymann, F., Liu, K., Roller, D., Smith, D., Thatte, S., Trickoric, I., and Weerawarana, S.: 'Business Process Execution Language for Web Services, Version 1.1, 2003. BPEL4WS specification
- Colombo, M., Di Nitto, E., and Mauri, M.: 'SCENE: a service composition execution environment supporting dynamic changes disciplined through rule'. *Service-Oriented Computing - ICSOC 2006, 4th Int. Conf., Proc. Lecture Notes in Computer Science* (Springer, 2006, vol. 4294), pp. 191–202
- Clarke, E.M., Grumberg, O., and Peled, D.A.: 'Model checking' (MIT Press, Cambridge, MA, 1999)
- Meyer, B.: 'Applying "design by contract"', *Computer*, 1992, **25**, (10), pp. 40–51
- Luckham, D.C., von Henke, F.W., Krieg-Brueckner, B., and Owe, O.: 'ANNA: a language for annotating Ada programs' ((Springer-Verlag, New York, NY, 1987)
- Leavens, G.T., Baker, A.L., and Ruby, C.: 'JML: a notation for detailed design' in Kilov, H., Rumpe, B., and Simmonds, I. (Eds.): 'Behavioral specifications of businesses and systems' (Kluwer Academic Publishers, Boston, MA, 1999), pp. 175–188
- Henzinger, T.A., Jhala, R., Majumdar, R., and Sutre, G.: 'Lazy abstraction'. *POPL 2002: Proc. 29th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 2002, pp. 58–70
- Clarke, E.M., Grumberg, O., Jha, S., Lu, Y., and Veith, H.: 'Counterexample-guided abstraction refinement'. *CAV 2000 Proc. 12th Int. Conf. on Computer Aided Verification*, Vol. 1855 of *Lecture Notes in Computer Science*, Springer, 2000, pp. 154–169
- Robby, Dwyer, M.B., and Hatcliff, J.: 'Bogor: an extensible and highly-modular software model checking framework'. *ESEC/FSE-11: Proc. 9th European Software Engineering Conference held jointly with 11th ACM SIGSOFT Int. Symp. Foundations of Software Engineering*. (ACM Press, 2003), New York, pp. 267–276
- Bianculli, D., Ghezzi, C., and Spoletini, P.: 'A model checking approach to verify BPEL4WS workflows'. *Proc. 2007 IEEE Int. Conf. Service-oriented Computing and Applications* (IEEE Computer Society Press, 2007), pp. 13–20
- Emerson, E.A., and Sistla, A.P.: 'Symmetry and model checking', *Form. Methods Syst. Des.*, 1996, **9**, (1–2), pp. 105–131
- Robby, Dwyer, M.B., Hatcliff, J., and Iosif, R.: 'Space-reduction strategies for model checking dynamic systems'. *Proc. 2003 Workshop on Software Model Checking*, Vol. 89 of *Electronic Notes in Theoretical Computer Science*, Elsevier, pp. 499–517
- Godefroid, P.: 'Using partial orders to improve automatic verification methods'. *CAV 1990 Proc. 2nd Int. Workshop on Computer Aided Verification*, Vol. 531 of *Lecture Notes in Computer Science*, Springer, 1991, pp. 176–185
- Baresi, L., Ghezzi, C., and Mottola, L.: 'Towards fine-grained automated verification of publish-subscribe architectures'. (FORTE06) *Proc. 26th Int. Conf. on Formal Methods for Networked and Distributed Systems*, Vol. 4229 of *Lecture Notes in Computer Science*, Springer, 2006, pp. 131–135
- ActiveBPEL Engine Architecture [homepage on the Internet]. Active Endpoints; 2006. Available from: <http://www.activebpel.org/docs/architecture.html>
- Kiczales, G., Lamping, J., Mendhekar, A., Maeda, C., Lopes, C.V., Loingtier, J.M., and Irwin, J.: 'Aspect-oriented Programming. ECOOP'97 – Object-Oriented Programming, 11th European Conf., Proc. Vol. 1241, *Lecture Notes in Computer Science*, Springer, 1997, pp. 220–242
- Kiczales, G., Hilsdale, E., Hugunin, J., Kersten, M., Palm, J., and Griswold, W.G.: 'An overview of AspectJ'. *ECOOP 2001 - Object-Oriented Programming, 15th European Conf., Proc. Vol. 2072 of Lecture Notes in Computer Science*, Springer, 2001, pp. 327–353
- Kupferman, O., and Vardi, M.: 'Weak alternating automata are not that weak'. *Fifth Israile Symp. Theory of Computing and Systems, ISTCS'97, Proc. (IEEE Computer Society Press, 1997)*, pp. 147–158
- Fu, X., Bultan, T., and Su, J.: 'Analysis of interacting BPEL web services'. *WWW 2004: Proc. 13th Int. Conf. World Wide Web*, New York, NY, (ACM Press, 2004), pp. 621–630
- Holzmann, G.J.: 'The model checker SPIN', *IEEE Trans. Softw. Eng.*, 1997, **23**, (5), pp. 279–295
- Arias-Fisteus, J., Fernández, L.S., and Kloos, C.D.: 'Formal verification of BPEL4WS business collaborations'. *E-Commerce and Web Technologies, 5th Int. Conf., EC-Web 2004, Proc. Vol. 3182 of Lecture Notes in Computer Science* Springer, 2004, pp. 76–85
- Nakajima, S.: 'Model-checking behavioral specification of BPEL applications', *Electron. Notes Theor. Comput. Sci.*, 2006, **151**, (2), pp. 89–105
- Schlingloff, B.H., Martens, A., and Schmidt, K.: 'Modeling and model checking web services', *Electron. Notes Theor. Comput. Sci.*, 2005 (Issue on Logic and Communication in Multi-Agent Systems), **126**, pp. 3–26
- Schmidt, K.: 'LoLA: a low level analyser', *Application and Theory of Petri Nets, 21st Int. Conf. (ICATPN 2000) Vol. 1825 of Lecture Notes in Computer Science*, Springer, 2000, pp. 465–474
- Foster, H., Uchitel, S., Magee, J., and Kramer, J.: 'Model-based verification of web service compositions'. *Proc. 18th IEEE Int. Conf. Automated Software Engineering (ASE 2003)*. IEEE Computer Society, 2003, pp. 152–163
- Koshkina, M., and van Breugel, F.: 'Verification of business processes for web services'. 4700 Keele Street, Toronto, M3J 1P3 Canada: (York, University - Department of Computer Science, 2003, CS-2003-11
- Cleaveland, R., Parrow, J., and Steffen, B.: 'The concurrency workbench: a semantics-based tool for the verification of concurrent systems', *ACM Trans. Program. Lang. Syst.*, 1993, **15**, (1), pp. 36–72
- Sahai, A., Machiraju, V., Sayal, M., Jin, L.J., and Casati, F.: 'Automated SLA Monitoring for Web Services'. *Proc. 13th IFIP/IEEE Int. Workshop on Distributed Systems: Operations and Management - Management Technologies for E-Commerce and E-Business Applications*, Vol. 2506 of *Lecture Notes in Computer Science*, Springer, 2002, pp. 28–41
- Keller, A., and Ludwig, H.: 'Defining and monitoring service-level agreements for dynamic e-business'. *Proc. 16th Conf. Systems Administration, (USENIX Association, Berkeley, CA, 2002)*, pp. 189–204
- Skene, J., Lamanna, D.D., and Emmerich, W.: 'Precise service level agreements'. *ICSE 2004: Proc. 26th Int. Conf. Software Engineering*, (IEEE Computer Society, Washington, DC, 2004), pp. 179–188
- Skene, J., Skene, A., Crampton, J., and Emmerich, W.: 'The monitorability of service-level agreements for application-service provision'. *WOSP 2007: Proc. 6th Int. Workshop on Software and performance*, (ACM Press, New York, NY, 2007), pp. 3–14
- Robinson, W.N.: 'Monitoring web service requirements'. *RE 2003: Proc. 11th IEEE Int. Conf. Requirements Engineering*, (IEEE Computer Society, Washington, DC, 2003), p. 65
- Mahbub, K., and Spanoudakis, G.: 'A framework for requirements monitoring of service based systems'. *ICSOC 2004: Proc. 2nd int. conf. Service Oriented Computing*, (ACM Press, New York, NY, 2004), pp. 84–93
- Barbon, F., Traverso, P., Pistore, M., and Trainotti, M.: 'Run-time monitoring of instances and classes of web service compositions'. *ICWS 2006: Proc. 2006 IEEE Int. Conf. Web Services*, (IEEE Computer Society, Washington, DC, 2006), pp. 63–71

- 38 Delgado, N., Gates, A.Q., and Roach, S.: 'A taxonomy and catalog of runtime software-fault monitoring tools', *IEEE Trans. Softw. Eng.*, 2004, **30**, (12), pp. 859–872
- 39 Graf, S., and Saidi, H.: 'Construction of abstract state graphs with PVS'. CAV 1997: Proc. 9th Int. Conf. Computer Aided Verification Vol. 1254 of Lectures Notes in Computer Science, Springer, 1997, pp. 72–83

12 Appendix

12.1 Formal semantics of ALBERT

In this appendix we provide the formal definition of the semantics of the core of ALBERT. We start by formalising the notions of state and sequence of states.

A state is defined as a triple (V, I, t) , where V is a set of \langle variable, value \rangle pairs, I a location of the process and t a time-stamp. A state completely describes the system in the particular time instant indicated by the time-stamp. States can be considered snapshots of the process. A location is defined as a set of labels of BPEL activities; in the case of a *flow* activity, it contains, for each branch of the flow, the last instruction executed in that branch. The set is needed to deal with *flow* activities; it is a singleton if the workflow processes do not contain *flow* activities. Sequences of states are often called timed state words, defined as follows.

Definition: A timed state word is an infinite sequence $s = s_1, s_2, \dots$ such that $s_i = (V_i, I_i, t_i)$.

As a consequence of how states are defined, timed state words are strictly monotonic. In fact, between subsequent states there is always at least one time-consuming interaction with the outside world or one internal activity execution (e.g. an *assign* activity).

The formal semantics of ALBERT operators can be defined as follows. Given a timed state word $s = s_1, s_2, \dots, s_i, \dots$, we introduce a helper function $\text{numState}(i, K)$ where i is an index of a state in the word and $K > 0$ is a real value denoting a time interval. The function returns the number of states in the word, encountered in the time window of size K by moving backwards in past from the i -th state:

- $\text{numState}(i, K) = 1$ iff $t_i - t_{i-1} > K$
- $\text{numState}(i, K) = 1 + \text{numState}(i-1, K - (t_i - t_{i-1}))$ iff $t_i - t_{i-1} \leq K$

An overloaded version is $\text{numState}(K, \text{onEvent}(\mu), s_i)$ which operates exactly as before, with the only difference that it counts only the states in which $\text{onEvent}(\mu)$ is true.

We introduce now the function $\text{eval}(\psi, s_i)$, which takes as parameters an ALBERT expression ψ and the state s_i in a word s and returns the value of ψ in s_i .

- $\text{eval}(\text{const}, s_i) = \text{const}$;
- $\text{eval}(\text{var}, s_i) = \text{value}$ iff $(\text{var}, \text{value}) \in V_i$;
- $\text{eval}(\psi_1 \text{ arop } \psi_2, s_i) = \text{eval}(\psi_1, s_i) \text{ arop } \text{eval}(\psi_2, s_i)$;
- $\text{eval}(\text{past}(\psi, \text{onEvent}(\mu), n), s_i) = \text{value}$ iff $\exists j < i \mid \text{eval}(\psi, s_j) = \text{value}$ and $w, j \models \text{onEvent}(\mu)$ (the satisfiability relation \models is defined below) and \exists exactly $n-1$ disjoint values $h_1, \dots, h_m \mid \forall m \in \{1, \dots, n-1\}, j < h_m < i$ and $w, h_m \models \text{onEvent}(\mu)$. If such a value of j cannot be found, $\text{eval}(\text{past}(\psi, \text{onEvent}(\mu), n), s_i)$ is undefined;

- $\text{eval}(\text{elapsed}(\text{onEvent}(\mu)), s_i) = \text{value}$ iff $\exists j \leq i \mid w, j \models \text{onEvent}(\mu)$ and $\neg \exists h \mid j < h < i$ and $w, h \models \text{onEvent}(\mu)$ and $t_i - t_j = \text{value}$;
- Let j be such that $j \leq i, t_i - t_j \leq K$ and $t_i - t_{j-1} > K$. Then $\text{eval}(\text{sum}(\psi, K), s_i)$ is defined as follows:
 - if $i = j$ then $\text{eval}(\psi, s_i)$,
 - if $i \neq j$ then $\text{eval}(\psi, s_i) + \text{eval}(\text{sum}(\psi, K - (t_i - t_{i-1})), s_{i-1})$.
- $\text{eval}(\text{avg}(\psi, K), s_i) = \text{eval}(\text{sum}(\psi, K)) / \text{numState}(i, K)$
- Let j be such that $j \leq i, t_i - t_j \leq K$ and $t_i - t_{j-1} > K$. Then $\text{eval}(\text{max}(\psi, K), s_i)$ is defined as follows:
 - if $i = j$ then $\text{eval}(\psi, s_i)$,
 - if $i > j$ then
 - * if $\text{eval}(\psi, s_i) < \text{eval}(\text{max}(\psi, K - (t_i - t_{i-1})), s_{i-1})$ then $\text{eval}(\text{max}(\psi, K - (t_i - t_{i-1})), s_{i-1})$
 - * if $\text{eval}(\psi, s_i) \geq \text{eval}(\text{max}(\psi, K - (t_i - t_{i-1})))$ then $\text{eval}(\psi, s_i)$
- Let j be such that $j \leq i, t_i - t_j \leq K$ and $t_i - t_{j-1} > K$. Then $\text{eval}(\text{count}(\phi, K), s_i)$ is defined as follows:

- if $i = j$ then
 - * if $w, i \models \phi$ then 1,
 - * if $w, i \not\models \phi$ then 0;
- if $i > j$ then
 - * if $w, i \models \phi$ then $1 + \text{eval}(\text{count}(\phi, K - (t_i - t_{i-1})), s_{i-1})$,
 - * if $w, i \not\models \phi$ then $\text{eval}(\text{count}(\phi, K - (t_i - t_{i-1})), s_{i-1})$.

Function $\text{eval}(\text{min}(\psi, K), s_i)$ can be computed similarly to $\text{eval}(\text{max}(\psi, K), s_i)$. The overloaded version $\text{fun}(\psi, \text{onEvent}(\mu), K)$ of functions fun is performed similarly to the evaluation of the original versions, but only considering the states in which $\text{onEvent}(\mu)$ holds.

For all timed word w , for all $i \in \mathbb{N}$, the satisfaction relation \models is defined as

- $w, i \models \psi \text{ relop } \psi'$ iff $\text{eval}(\psi, s_i) \text{ relop } \text{eval}(\psi', s_i)$.
- $w, i \models \neg \phi$ iff $w, i \not\models \phi$.
- $w, i \models \phi \wedge \xi$ iff $w, i \models \phi$ and $w, i \models \xi$.
- $w, i \models \text{onEvent}(\mu)$ iff
 - if μ is a start event, $\mu \in I_{i+1}$,
 - otherwise, $\mu \in I_i$
- $w, i \models \text{Becomes}(\phi)$ iff $i > 0$ and $w, i \models \phi$ and $w, i-1 \not\models \phi$
- $w, i \models \text{Until}(\phi, \xi)$ iff $\exists j \geq i \mid w, j \models \xi$ and $\forall k$, if $i \leq k < j$ then $w, k \models \phi$;
- $w, i \models \text{Between}(\phi, \xi, K)$ iff $\exists j \geq i \mid w, j \models \phi$ and $\forall l$ if $i \leq l < j$ then $w, l \not\models \phi$ and $\exists h \mid t_h \leq t_j + K, t_{h+1} > t_j + K$ and $w, h \models \xi$
- $w, i \models \text{Within}(\phi, K)$ iff $\exists j \geq i \mid t_j - t_i \leq K$ and $w, j \models \phi$

Notice that even if the definition of the satisfaction relation recalls the function eval and viceversa, the two definitions are not cyclic.