

Formal Modeling and SMT-Based Parameterized Verification of Multi-Case Data-Aware BPMN

Diego Calvanese¹, Silvio Ghilardi², Alessandro Gianola¹,
Marco Montali¹, Andrey Rivkin¹

¹Faculty of Computer Science, Free University of Bozen-Bolzano (Italy)

²Dipartimento di Matematica, Università degli Studi di Milano (Italy)

Abstract. We propose DAB – a data-aware extension of the BPMN de-facto standard with the ability of operating over case and persistent data (partitioned into a read-only catalog and a read-write repository), and that balances between expressiveness and the possibility of supporting parameterized verification of safety properties on top of it. In particular, we take inspiration from the literature on verification of artifact systems, and consider verification problems where safety properties are checked irrespectively of the content of the read-only catalog, possibly considering an unbounded number of active cases and tuples in the catalog and repository. Such problems are tackled using fully implemented array-based backward reachability techniques belonging to the well-established tradition of SMT model checking. We also identify relevant classes of DABs for which the backward reachability procedure implemented in the MCMT array-based model checker is sound and complete, and then further strengthen such classes to ensure termination.

1 Introduction

In recent years, increasing attention has been given to multi-perspective models of business processes that strive to capture the interplay between the process and data dimensions [26,25]. The corresponding development of formal models and foundational results on their verification has flourished, but mainly focusing on data-centric approaches [28,3] that are quite different in nature from conventional notations. In parallel, conventional approaches such as the de-facto BPMN standard have been extended with various forms of data, though with a main focus on conceptual modeling and enactment [22,11,8], without considering formal verification. More on the formal side, many approaches within this line considering Petri nets as the main control-flow backbone for capturing the process, and enriching them with data locally carried by tokens [27,19,23], case data with different datatypes [12], and/or persistent relational data manipulated with the full power of FOL/SQL [13,24]. While this latter family of approaches qualify well to capture BPMN enriched with persistent data (such as [22,8]), they place two assumptions on verification: first, that the process is studied by considering a fully-specified, initial instance of the underlying database; second, that only boundedly many tuples can be stored in the database. Such approaches have not

yet been applied to formalize BPMN with data, and do not lend themselves to capture parameterized verification problems that can ascertain the correctness of the process without imposing the aforementioned limitations.

In this work, we attack this open challenge and propose a data-aware extension of BPMN, called *data-aware BPMN* (DAB), equipped with the ability of querying and operating over case and persistent data (partitioned into a read-only catalog and a read-write repository). The approach balances expressiveness with the possibility of supporting parameterized verification of safety properties. In particular, we take inspiration from the literature on the verification of artifact systems [28], and in particular the frameworks in [20,6], considering verification problems where safety properties are checked irrespectively of the content of the read-only catalog, and possibly considering an unbounded number of active cases and tuples in the catalog and repository.

We study this problem by establishing a bridge between our approach to business process modeling, and the line of research where techniques based on Satisfiability Modulo Theories (SMT) are employed to attack the verification of infinite-state *array-based systems* – originally introduced in [16,17] to handle the verification of distributed systems (parameterized on the number of interacting processes). In particular, we rely on novel foundational results on the verification of artifact-centric systems via array-based systems [4,6], which have been fully implemented in the state-of-the-art MCMT SMT symbolic model checker [18].

By exploiting this formal basis, we identify two relevant classes of DABs for which the backward reachability procedure implemented in the MCMT array-based model checker is sound and complete. This guarantees that if the procedure terminates, it produces a correct judgement. We then introduce further conditions that, by carefully controlling the interplay between the process and data components, guarantee the termination of the procedure, in turn witnessing decidability. Such conditions are expressed as syntactic restrictions over the DAB under study, thus providing a concrete, BPMN-grounded counterpart of the conditions imposed in [20,4,6] towards decidability.

These termination results obtained by translating DABs into the array-based artifact systems studied in [4,6]. The translation can then be straightforwardly implemented making it possible to effectively verify DABs using MCMT.

This article builds upon [5], extending it in two respects. On the one hand, while [5] focuses on the verification of DABs considering a single, running case, we consider here the possibility of (unboundedly many) cases running concurrently. On the other hand, we provide full proofs of the technical results, including those from [5] and those specifically introduced in this extended version.

2 Data-aware BPMN

We start by describing our formal model of data-aware BPMN processes (DABs). We focus here on private, single-pool processes. Incoming messages are therefore handled as pure nondeterministic events. The model combines a wide range of (block-structured) BPMN control-flow constructs with task, event-reaction,

and condition logic that inspect and modify persistent as well as case data. The combination achieves a balanced trade off between the expressive power of the resulting integrated model, and the possibility of carrying out sophisticated forms of parameterized verification, which will be tackled in Section 3. When going through the modeling features of DAB, it is then important to remember that if something is not supported, it is because it would hamper soundness and completeness of SMT-based (parameterized) verification.

First, some preliminary notation. We consider a set $\mathcal{S} = \mathcal{S}_v \uplus \mathcal{S}_{id}$ of (semantic) *types*, consisting of *primitive types* \mathcal{S}_v accounting for data objects, and *id types* \mathcal{S}_{id} accounting for identifiers. We assume that each type $S \in \mathcal{S}$ comes with a (possibly infinite) domain \mathbb{D}_S , a special constant $\text{undef}_S \in \mathbb{D}_S$ to denote an undefined value in that domain, and a type-wise equality operator $=_S$. We omit the type and simply write undef and $=$ when clear from the context. We do not consider here additional type-specific predicates (such as comparison and arithmetic operators for numerical primitive types); these will be added in future work (cf. Section 5 for a discussion on this). In the following, we simply use *typed* as a shortcut for *\mathcal{S} -typed*. We also denote by \mathbb{D} the overall domain of objects and identifiers (i.e., the union of all domains in \mathcal{S}). We consider a countably infinite set \mathcal{V} of typed variables. Given a variable or object x , we may explicitly indicate that x has type S by writing $x : S$. We omit types whenever clear from the context, or irrelevant. We compactly indicate a possibly empty tuple $\langle x_1, \dots, x_n \rangle$ of variables as \vec{x} , and with slight abuse of notation, we write $\vec{x} \subseteq \vec{y}$ if all variables in \vec{x} also appear in \vec{y} .

2.1 The Data Schema

Consistently with the BPMN standard, we consider two main forms of data: *case data*¹, instantiated and manipulated on a per-case basis; *persistent data* (cf. data store references in BPMN), accounting for global data that are accessed by all cases. For simplicity, case data are defined at the whole process level, and are directly visible by all tasks and subprocesses (without requiring the specification of input-output bindings and the like).

To account for persistent data, we consider relational databases. We describe relation schemas by using the *named perspective*, i.e., by assigning a dedicated typed attribute to each component (i.e., column) of a relation schema. Also for an attribute, we use the notation $a : S$ to explicitly indicate its type.

Definition 1. A relation schema is a pair $R = \langle N, A \rangle$, where: (i) $N = R.\text{name}$ is the relation name; (ii) $A = R.\text{attrs}$ is a nonempty tuple of attributes. We call $|A|$ the arity of R . ◁

We assume that distinct relation schemas use distinct names, blurring the distinction between the two notions (i.e., we set $R.\text{name} = R$). We also use the predicate notation $R(A)$ to represent a relation schema $\langle R, A \rangle$.

¹ These are called *data objects* in BPMN, but we prefer to use the term *case data* to avoid name clashes with the formal notions.

Data schema. We start by defining the *catalog*, i.e., a read-only, persistent storage of data that is not modified during the execution of the process.

Definition 2. A catalog Cat is a set of relation schemas satisfying the following requirements:

- (**single-column primary key**) Every relation schema R is such that the first attribute in $R.attrs$ has type in \mathcal{S}_{id} , and denotes the primary key of the relation; we refer to such attribute using the dot notation $R.id$.
- (**non-ambiguity of primary keys**) for every pair R_1 and R_2 of distinct relation schemas in Cat , we have that the types of $R_1.id$ and $R_2.id$ are different.
- (**foreign keys**) for every relation schema $R \in Cat$ and non-*id* attribute $a \in R.attrs \setminus R.id$ with type $S \in \mathcal{S}_{id}$, there exists a relation schema $R_2 \in \mathcal{R}$ such that the type of $R_2.id$ is S ; a is hence a foreign key referring to R_2 . \triangleleft

We now define the data schema of a BPMN process, which combines a catalog with: (i) a persistent data *repository*, consisting of updatable relation schemas possibly referring to the catalog; (ii) a set of *case variables*, constituting local data carried by each process case.

Definition 3. A data schema \mathcal{D} is a tuple $\langle Cat, CType, Repo, X \rangle$, where (i) $Cat = \mathcal{D}.cat$ is a catalog, (ii) $CType = \mathcal{D}.ctype \in \mathcal{S}_{id}$ is a special case identifier type, (iii) $Repo = \mathcal{D}.repo$ is a set of relation schemas called repository, and (iv) $X = \mathcal{D}.cvars \subset \mathcal{V}$ is a finite set of typed variables called case variables, such that:

- $CType$ is disjoint from all identifier types used in Cat ;
- for every relation schema $R \in Repo$ and every attribute $a \in R.attrs$ whose type is $S \in \mathcal{S}_{id}$, there exists $R \in Cat \cup \{CType\}$ such that the type of $R.id$ is S ;
- for every case variable $\mathbf{x} \in X$ whose type is $S \in \mathcal{S}_{id}$, there exists $R \in Cat \cup \{CType\}$ such that the type of $R.id$ is S ;
- \mathcal{D} contains an special case variable **self**: $CType$ that is never modified, and that keeps track, for a case, of the corresponding case identifier. \triangleleft

We use bold-face to distinguish a case variable \mathbf{x} from a “normal” variable x . It is worth noting that relation schemas in the repository are not equipped with an explicit primary key, and thus they cannot reference each other, but may contain foreign keys pointing to the catalog or the case identifiers. *This is essential towards soundness and completeness of SMT-based verification of DABs.* It will be clear how tuples can be inserted and removed from the repository once we will introduce updates.

At runtime, a *data snapshot* of a data schema consists of three components:

- An immutable *catalog instance*, i.e., a fixed set of tuples for each relation schema contained therein, so that the primary and foreign keys are satisfied.
- A *case map* whose keys are the identifiers of active or completed cases (i.e., elements of the case identifier type), and whose values are assignments of the case variables to corresponding values (satisfying the foreign keys when pointing to identifiers in the catalog). Each entry then indicates, for a given case, which are the current values for the case variables of that case.
- A *repository instance*, i.e., a set of tuples for for each relation schema contained therein, so that the foreign key constraints pointing to the catalog or the case

map keys are satisfied. Each tuple is associated to a distinct primary key that is not explicitly accessible.

Example 1. We consider a simplified example of a job hiring process in a company. We describe here the data schema \mathcal{D}^h used to store data about job hirings and their corresponding applications. The catalog $\mathcal{D}^h.\text{cat}$ consists of the following relation schemas:

- *JobCategory*(*Jcid*:jobcatID) contains the different job categories available in the company (e.g., programmer, analyst, and the like) - we just store here the identifiers of such categories;
- *User*(*Uid*:userID, *Name*:StringName, *Age*:NumAge) stores data about users registered to the company website, and who are potentially interested in job positions offered by the company.

Each case of the process is about a job. Jobs are identified by the type `jobId`.

To manage key information about the applications submitted for the various job hirings, including data on users, the score they receive after having been interviewed and their eligibility, the company employs the repository $\mathcal{D}.\text{repo}$ that consists of one relation schema

Application(*Jid*:jobId, *Jcid*:jobcatID, *Uid*:userID, *Name*:StringName, *Age*:NumAge, *Score*:NumScore, *Eligible*:Bool)

Notice that `NumScore` is a finite-domain type containing 100 values in the range $[1, 100]$, and it is used to assign an overall score to each candidate application. For readability, we use the usual predicates $<$, $>$, and $=$ to compare variables of type `NumScore`: this is syntactic sugar and does not require to introduce rigid predicates in our framework.

Since each posted job is created using a dedicated portal, its corresponding data do not have to be stored persistently and thus can be maintained just for a given case. At the same time, some specific values have to be moved from a specific case to the repository and vice-versa. This is done by resorting to the following case variables $\mathcal{D}.\text{cvars}$: (i) **jcId** : `jobcatID` references a job type from the catalog, matching the type of job associated to the case; (ii) **uid** : `userID`, **name** : `StringName` and **age** : `NumAge` respectively reference the identifier, name, and age of a user who is applying for the job associated to the case; (iii) **result** : `Bool` indicates whether the user identified by **uid** is eligible for winning the position or not; (iv) **result** : `Bool` indicates whether the user identified by **uid** qualifies for directly getting the job (without the need of carrying out a comparative evaluation of all applicants); (v) **winner** : `userID` contains the identifier of the applicant winning the position; (vi) **tPassed** : `StringDate` contains special strings symbolically indicating the current temporal phase of the case in relation with its creation time. The last variable is not essential for the progression of job hirings through the process, but it is useful to formulate verification properties.<

Querying the data schema. To inspect the data contained in a snapshot, we need suitable query languages operating over the data schema of that snapshot. In the following, we assume that queries are well-typed, i.e., sorts of their elements are duly matched (this can be easily checked by scanning the query). We

start by considering boolean *conditions* over (case) variables. These conditions will be attached to choice points in the process.

Definition 4. A condition is a formula of the form $\varphi ::= (x = y) \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2$, where x and y are variables from \mathcal{V} or constant objects from \mathbb{D} . If in φ negation is restricted to be only in front of atoms, φ is called a cubical condition. \triangleleft

We make use of the standard abbreviation $\varphi_1 \vee \varphi_2 = \neg(\neg\varphi_1 \wedge \neg\varphi_2)$.

We now extend conditions to also access the data stored in the catalog and repository, and to ask for data objects subject to constraints. We consider the well-known language of unions of conjunctive queries with atomic negation, which correspond to unions of select-project-join SQL queries with table filters.

Definition 5. A conjunctive query with filters over a data component \mathcal{D} is a formula of the form $Q ::= \varphi \mid R(x_1, \dots, x_n) \mid \neg R(x_1, \dots, x_n) \mid Q_1 \wedge Q_2$, where φ is a cubical condition, $R \in \mathcal{D}.cat \cup \mathcal{D}.repo$ is a relation schema of arity n , and x_1, \dots, x_n are variables from \mathcal{V} (including $\mathcal{D}.cvars$) or constant objects from \mathbb{D} . We denote by $free(Q)$ the set of variables occurring in Q that are not case variables in $\mathcal{D}.cvars$. \triangleleft

Definition 6. A guard G over a data component \mathcal{D} is an expression of the form $q(\vec{x}) \leftarrow \bigvee_{i=1}^n Q_i$, where: (i) $q(\vec{x})$ is the head of the guard with answer variables \vec{x} ; (ii) each Q_i is a conjunctive query with filters over \mathcal{D} ; (iii) for each $i \in \{1, \dots, n\}$, $\vec{x} \subseteq free(Q_i)$. We denote by $casevars(G) \subseteq \mathcal{D}.cvars$ the set of case variables used in G , and by $normvars(G) = \bigcup_{i \in \{1, \dots, n\}} free(Q_i)$ the other variables used in G . \triangleleft

Definition 7. A guard G over a data component \mathcal{D} is repo-free if none of its atoms queries a relation schema from $\mathcal{D}.repo$. \triangleleft

Notice that *going beyond this guard query language* (e.g., by introducing universal quantification) *would hamper the soundness and completeness of SMT-based verification over the resulting DABs*. We will come back to this important aspect in the conclusion.

As anticipated before, this language can be seen as a standard query language to retrieve data from a snapshot, but also as a mechanism to identify the allowed combinations of data objects that can be injected into the process from the external environment. For example, considering a case variable \mathbf{x} of type `string`, a simple guard $Input(y:string, z:string) \rightarrow y \neq \mathbf{x} \wedge y \neq z$ returns all pairs of strings that are different from each other, and so that the second string is different from that stored in the case variable \mathbf{x} . Picking an answer in this (infinite) set of pairs can be consequently seen as a (constrained) user decision on which values to input for y and z . We elaborate more on this in Section 2.2.

2.2 Tasks, Events, and Impact on Data

We now formalize how the process can access and update the data component when executing a task or reacting to the trigger of an external event.

The update logic. We start by discussing how data maintained in a snapshot can be subject to change while executing the process.

Definition 8. Given a data schema \mathcal{D} , an update specification α is a pair $\langle G, E \rangle$, where: (i) $G = \alpha.\text{pre}$ is a guard over \mathcal{D} of the form $q(\vec{x}) \leftarrow Q$, called precondition; (ii) $E = \alpha.\text{eff}$ is an effect rule that changes the content of case variables or that of the repository, as described next. Each effect rule has one of the following forms:

(Insert&Set) INSERT \vec{u} INTO R AND SET $\mathbf{x}_1 = v_1, \dots, \mathbf{x}_n = v_n$, where: (i) \vec{u}, \vec{v} are variables in \vec{x} or constant objects from \mathbb{D} ; (ii) $\vec{x} \in \mathcal{D}.\text{cvars} \setminus \{\text{self}\}$ are distinct case variables different from **self**; (iii) R is a relation schema from $\mathcal{D}.\text{repo}$ whose arity (and types) match \vec{u} . Either the INSERT or AND SET parts may be omitted, obtaining a pure (repository) **Insert rule** or (case variable) **Set rule**.

(Delete&Set) DEL \vec{u} FROM R AND SET $\mathbf{x}_1 = v_1, \dots, \mathbf{x}_n = v_n$, where: (i) \vec{u}, \vec{v} are variables in \vec{x} or constant objects from \mathbb{D} ; (ii) $\vec{x} \in \mathcal{D}.\text{cvars} \setminus \{\text{self}\}$; (iii) R is a relation schema from $\mathcal{D}.\text{repo}$ whose arity (and types) match \vec{u} . As in the previous rule type, the AND SET part may be omitted, obtaining a pure (repository) **Delete rule**.

(Conditional update) UPDATE $R(\vec{v})$ IF $\psi(\vec{u}, \vec{v})$ THEN η_1 ELSE η_2 , where: (i) \vec{u} is a tuple containing variables in \vec{x} or constant objects from \mathbb{D} ; (ii) ψ is a repo-free guard (called filter) (iii) \vec{v} is a tuple of new variables, i.e., such that $\vec{v} \cap (\vec{u} \cup \mathcal{D}.\text{cvars}) = \emptyset$; (iv) η_i is either an atomic formula of the form $R(\vec{u}')$ with \vec{u}' a tuple of elements from $\vec{x} \cup \mathbb{D} \cup \vec{v}$, or a nested IF... THEN... ELSE. \triangleleft

As stated in the definition, **self** is never explicitly set by any of the effect rules.

We now comment on the semantics of update specifications. An update specification α is executable in a given data snapshot if there is at least one answer to the precondition $\alpha.\text{pre}$ in that snapshot. If this is the case, then the process executor(s) can nondeterministically decide which answer to pick so as to *bind* the answer variables of $\alpha.\text{pre}$ to corresponding data objects in \mathbb{D} . This confirms the interpretation discussed in Section 2.1 for which the answer variables of $\alpha.\text{pre}$ can be seen as *constrained user inputs* in case multiple bindings are available.

Once a specific binding for the answer variables is selected, the corresponding effect rule $\alpha.\text{eff}$, instantiated using that binding, is issued. How this affects the current data snapshot depends on which effect rule is adopted.

If $\alpha.\text{eff}$ is an insert&set rule, the binding is used to *simultaneously* insert a tuple in one of the repository relations, and update some of the case variables – with the implicit assumption that those not explicitly mentioned in the SET part maintain their current values. Since repository relations do not have an explicit primary key, two possible semantics can be attached to the insertion of a tuple \vec{u} in the instance of a repository relation R :

(multiset insertion) Upon insertion, \vec{u} gets implicitly assigned to a fresh primary key. The insertion then always results in the genuine addition of the tuple to the current instance of R , even in the case where the tuple already exists there.

(set insertion) In this case, R comes not only with its implicit primary key, but also with an additional, genuine key constraint defined over a subset $K \subseteq R.\text{attrs}$ of its attributes. Upon insertion, if there already exists a tuple

in the current instance of R that agrees with \vec{u} on K , then that tuple is *updated* according to \vec{u} . If no such tuple exists, then as in the previous case \vec{u} gets implicitly assigned to a fresh primary key, and inserted into the current instance of R . By default, if no explicit key is defined over R , then the entire set of attributes $R.attrs$ is considered as a key, consequently enforcing a *set semantics* for insertion.

Example 2. We continue the job hiring example, by considering two update specifications of type insert&set. When a new case is created, the first update is about indicating what is the category of job associated to the case. This is done through the update specification **InsJobCat**, where:

- **InsJobCat.pre** $\triangleq GetJobType(jt) \leftarrow JobCategory(jt)$ selects a job category from the corresponding catalog relation;
- **InsJobCat.eff** $\triangleq SET \mathbf{jcid} = jt$ assigns the selected job category jt to the case variable **jcid**.

When the case receives an application for its associated job, the user-related case variables are filled with the data of the user submitting the application - picked from the corresponding *User* catalog relation. This is done via the update specification **InsUser**, where:

$$\begin{aligned} \mathbf{InsUser.pre} &\triangleq GetUser(u, n, a) \leftarrow User(u, n, a) \\ \mathbf{InsUser.eff} &\triangleq SET \mathbf{uid} = u, \mathbf{name} = n, \mathbf{age} = a \end{aligned}$$

A different usage of precondition, resembling a pure external choice, is the update specification **CheckQual** to handle a quick evaluation of the candidate and check whether she has such a high profile qualifying her to directly get an offer:

$$\begin{aligned} \mathbf{CheckQual.pre} &\triangleq IsQualified(q : Bool) \leftarrow \mathbf{true} \\ \mathbf{CheckQual.eff} &\triangleq SET \mathbf{qualif} = q \end{aligned}$$

As an example of insertion rule, we consider the situation where the candidate whose data are currently stored in the case variables has not been directly judged as qualified. She is consequently subject to a more fine-grained evaluation of her application, resulting in a score that is then registered in the repository (together with the applicant data). This is done via the **Reg<eval** specification:

$$\begin{aligned} \mathbf{EvalApp.pre} &\triangleq GetScore(s : NumScore) \leftarrow 1 \leq s \wedge s \leq 100 \\ \mathbf{EvalApp.eff} &\triangleq INSERT \langle \mathbf{self}, \mathbf{jcid}, \mathbf{uid}, \mathbf{name}, \mathbf{age}, s, \mathbf{undef} \rangle \text{ INTO } Application \end{aligned}$$

Here, the insertion uses the applicant data currently stored in the corresponding case variables, the selected score, and **undef** eligibility (which is then assessed in a consequent step of the process). These objects are correlated to the case identifier, so as to keep track of the relationship between the application and the job to which the application has been submitted. This is essential, as the same user may apply for different jobs. Notice that, by adopting the *multiset insertion semantics*, the same user may even apply multiple times for the same job. With a *set insertion semantics*, one could instead ensure that each user can apply at most once to the same job, by indicating that the first two components of *Application* form a key. \triangleleft

If $\alpha.eff$ is a delete&set rule, then the executability of the update is subject to the fact that the tuple \vec{u} selected by the binding and to be removed from

R , is actually present in the current instance of R . If so, the binding is used to *simultaneously* delete \vec{u} from R and update some of the case variables – with the implicit assumption that those not explicitly mentioned in the **SET** part maintain their current values.

Finally, a conditional update rule applies, tuple by tuple, a bulk operation over the content of R . For each tuple in R , if it passes the filter associated to the rule, then the tuple is updated according to the **THEN** part, whereas if the filter evaluates to false, the tuple is updated according to the **ELSE** part.

Example 3. Continuing with our running example, we now consider the update specification **MarkE** handling the situation where no candidate has been directly considered as qualified, and so the eligibility of all received (and evaluated) applications has to be assessed. Here we consider that each application is eligible if and only if its evaluation resulted in a score greater than 80. Technically, **MarkE.pre** is a true precondition, and:

```

MarkE.eff  $\triangleq$  UPDATE Application( $j, jc, u, n, a, s, e$ )
    IF  $j = \mathbf{self} \wedge s > 80$  THEN Application( $j, jc, u, n, a, s, \mathbf{true}$ )
    ELSE IF  $j = \mathbf{self} \wedge s \leq 80$  THEN Application( $j, jc, u, n, a, s, \mathbf{false}$ )
    ELSE Application( $j, jc, u, n, a, s, e$ )

```

The update logic realized by **MarkE.eff** is the following: (i) applications sent for the considered job and with a score > 80 are marked as eligible; (ii) other applications sent for the considered job are marked as not eligible; (iii) applications sent for other jobs are left unaltered.

If there is at least one eligible candidate, she can be selected as a winner using the **SelWinner** update specification, which deletes the selected winner tuple from *Application*, and transfers its content to the corresponding case variables (also ensuring that the **winner** case variable is set to the applicant id). Technically:

```

SelWinner.pre  $\triangleq$  GetWinner( $j, jc, u, n, a, s, e$ )  $\leftarrow$  Application( $j, jc, u, n, a, s, e$ )
     $\wedge j = \mathbf{self} \wedge e = \mathbf{true}$ 

SelWinner.eff  $\triangleq$  DEL  $\langle j, jc, u, n, a, s, e \rangle$  FROM Application
    AND SET jcid =  $jc$ , uid =  $u$ , name =  $n$ ,
    age =  $a$ , winner =  $jc$ , result =  $e$ 

```

Deleting the tuple is useful in the situation where the selected winner may refuse the job, and consequently should not be considered again if a new winner selection is carried out. To keep such tuple in the repository, one would just need to remove the **DEL** part from **EvalApp.eff**. \triangleleft

The task/event logic. We now substantiate how the update logic is used to specify the task/event logic within a DAB process. The first important observation, which does not relate to our specific design choice for the update logic, but is inherently present whenever the process control flow is enriched with relational data, is that update effects manipulating the repository must be executed in an atomic, non-interruptible way. This is essential to ensure that insertions/deletions into/from the repository are applied on the same data snapshot where the precondition is checked. This cannot be guaranteed if the precondition and effect occur in different moments, as they may nondeterministically interleave

with other update specifications potentially operating over the same portion of the repository. This is why in our approach we consider two types of task: *atomic* and *nonatomic*. This goes beyond the BPMN standard, where generic tasks are implicitly assumed to be nonatomic.

Each atomic task/catching event is associated to a corresponding update specification. In the case of tasks, the specification precondition indicates under which circumstances the task can be enacted, and the specification effect how enacting the task impacts on the underlying data snapshot. In the case of events, the specification precondition constrains the data payload that comes with the event (possibly depending on the data snapshot, which is global and therefore accessible also from the perspective of an external event trigger), and the specification effect how reacting to a triggered event impacts on the underlying data snapshot. More concretely, this is realized according to the following lifecycle.

The task/event is initially **idle**, i.e., quiescent. When the progression of a case reaches an **idle** task/event, such a task/event becomes **enabled**. An **enabled** task/event may nondeterministically fire depending on the choice of the process executor(s). Upon firing, a binding satisfying the precondition of the update specification associated to the task/event is selected, consequently grounding and applying the corresponding effect. At the same time, the lifecycle moves from **enabled** to **compl**. Finally, a **compl** task/event triggers the progression of its case depending on the process-control flow, simultaneously bringing the task/event back to the **idle** state (which would then make it possible for the task to be executed again later on within the same case, if the process control-flow dictates so).

The lifecycle of a nonatomic task diverges in two crucial respects. First of all, upon firing it moves from **enabled** to **active**, and later on nondeterministically from **active** to **compl** (thus having a duration). The precondition of its update specification is checked and bound to one of the available answers when the task becomes **active**, while the corresponding effect is applied when the task becomes **compl**. Since these two transitions occur asynchronously, to avoid the aforementioned transactional issues we assume that the effect operates, in this context, only on case variables (and not on the repository).

2.3 Process Schema

A process schema consists of a block-structured BPMN diagram, enriched with conditions and update effects expressed over a given data schema, according to what described in the previous sections. As for the control flow, we consider a wide range of block-structured patterns compliant with the standard, taking inspiration and expanding those in [21]. We focus on private BPMN processes, thereby handling incoming messages in a pure nondeterministic way. So we do for timer events, nondeterministically accounting for their expiration without entering into their metric temporal semantics. Focusing on block-structured components helps us in obtaining a direct, execution semantics, and a consequent modular and clean translation of various BPMN constructs (including boundary

events and exception handling). However, it is important to stress that our approach would seamlessly work also for non-structured processes where each case introduces boundedly many tokens.

As usual, blocks are recursively decomposed into sub-blocks, the leaves being task or empty blocks. Depending on its type, a block may come with one or more nested blocks, and be associated with other elements, such as conditions, types of the involved events, and the like. We consider a wide range of blocks, covering basic, flow, and exception handling patterns. They are reported in Appendix A. Figure 1 gives an idea about what is covered by our approach. With these blocks at hand, we finally obtain the full definition of a DAB.

Definition 9. *A DAB \mathcal{M} is a pair $\langle \mathcal{D}, \mathcal{P} \rangle$ where \mathcal{D} is a data schema, and \mathcal{P} is a root process block such that all conditions and update effects attached to \mathcal{P} and its descendant blocks are expressed over \mathcal{D} .* \triangleleft

Example 4. The full hiring job process is shown in Figure 1, using the update effects described in Examples 2 and 2. Intuitively, the process works as follows. A new case for the process is created whenever a new job is posted. The case enters into a looping subprocess where it expects candidates to apply. Specifically, the case waits for an incoming application, or for an external message signalling that the hiring has to be stopped (e.g., because too much time has passed from the posting). Whenever an application is received, the CV of the candidate is evaluated, with two possible outcomes. The first outcome indicates that the candidate directly qualifies for the position, hence no further applications should be considered. In this case, the process continues by declaring the candidate as winner, and making an offer to her. The second outcome of the CV evaluation is instead that the candidate does not directly qualify. A more detailed evaluation is then carried out, assigning a score to the application and storing the outcome into the process repository, then waiting for additional applications to come. When the application management subprocess is stopped (which we model through an error so as to test various types of blocks in the experiments reported in Section 4), the applications present in the working memory are all processed in parallel, declaring which candidates are eligible and which not depending on their scores. Among the eligible ones, a winner is then selected, making an offer to her. We implicitly assume here that at least one applicant is eligible. We can remove this assumption and also handle the case where no eligible applicant exists for a job, by simply introducing (and consequently using) a boolean case variable that, upon the application evaluation, is set to true if the obtained score makes the application eligible. \triangleleft

As customary, each block has a lifecycle that, case by case, indicates the current state of the block, and how it can be evolved depending on the specific semantics of the block, and the evolution of its inner blocks. In Section 2.2 we have already characterized the lifecycle of tasks and catch events. For the other blocks, we continue to use the standard states `idle`, `enabled`, `active` and `compl`. We use the very same rules of execution described in the BPMN standard to regulate the progression of blocks through such states, taking advantage from the fact that, being the process block-structured, only one instance of a block

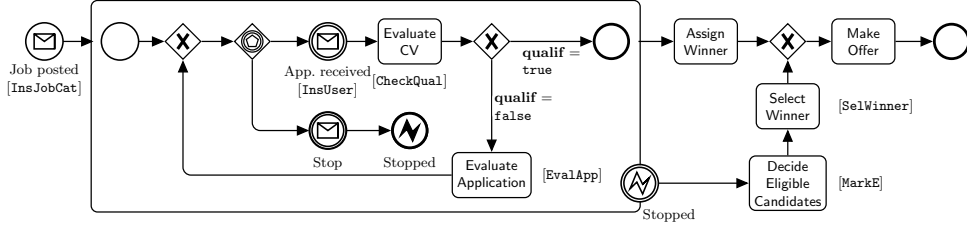


Fig. 1. The job hiring process. Elements in squared brackets indicate the update specifications attached to the corresponding tasks/events, and formalized as shown in Examples 2 and 3.

can be enabled/active at a given time for a given case. For example, the lifecycle of a sequence block S with nested blocks B_1 and B_2 can be described as follows (considering that the transitions of S from `idle` to `enabled` and from `compl` back to `idle` are inductively regulated by its parent block): (i) if S is `enabled`, then it becomes `active`, simultaneously inducing a transition B_1 from `idle` to `enabled`; (ii) if B_1 is `compl`, then it becomes `idle`, simultaneously inducing a transition of B_2 from `idle` to `enabled`; (iii) if B_2 is `compl`, then it becomes `idle`, simultaneously inducing S to move from `active` to `compl`. The lifecycle of other block types can be defined analogously.

2.4 Execution Semantics

We intuitively describe the execution semantics of a DAB $\mathcal{M} = \langle \mathcal{D}, \mathcal{P} \rangle$, using the update/task logic and progression rules of blocks as a basis. Upon execution, each state of \mathcal{M} is characterized by an \mathcal{M} -*snapshot*, in turn constituted by a data snapshot of \mathcal{D} (cf. Section 2.1), and a *control map* whose keys are the (active) case identifiers, and whose values are assignments from each block in \mathcal{P} to one of its lifecycle states.

Initially, the data snapshot fixes the immutable content of the catalog $\mathcal{D}.cat$, while the repository instance, the case map, and the control map, are all empty. At each moment in time, the \mathcal{M} -snapshot is then evolved by nondeterministically performing one of the two steps:

(new case creation) A new case is created, obtaining a corresponding fresh identifier id . The new \mathcal{M} -snapshot is then obtained by maintaining the catalog and repository unaltered, while updating: (i) the case map, creating a new entry with key id and setting all the case variables to `undef` in such a newly created entry, with the exception of `self`, which gets the fresh identifier id ; (ii) the control map, creating a new entry with key id and setting all the block lifecycles to `idle` in such a newly created entry.

(case progression) The identifier id of a case present in the control/case map is nondeterministically picked, nondeterministically evolving it of one step through the process, depending on the current \mathcal{M} -snapshot. The new \mathcal{M} -snapshot is then obtained by manipulating, accordingly, the id entries in the case and control maps, as well as possibly updating the repository if the selected step involves the application of an effect rule (cf. Section 2.2).

3 Parameterized Verification of Safety Properties

We now focus on parameterized verification of DABs using SMT-based techniques grounded in the theory of arrays.

3.1 Array-Based Artifact Systems and Safety Checking

We recall the key notions behind array-based systems, and the artifact variants recently studied in [4,6] to bridge the gap between SMT-based model checking of array-based systems [16,17] and verification of artifact-centric processes [28,15]. In general terms, an array-based system is described using a multi-sorted theory that contains two types of sorts, one accounting for the indexes of arrays, and the other for the elements stored therein. Since the content of an array changes over time, it is referred to using a *function* variable, whose interpretation in a state is that of a total function mapping indexes to elements (so that applying the function to an index denotes the classical *read* operation for arrays). The definition of an array-based system with array state variable a always requires (i) a formula $I(a)$ describing the *initial configuration* of the array a ; (ii) a formula $\tau(a, a')$ describing a *transition* that transforms the content of the array from a to a' . In such a setting, verifying whether the system can *reach* unsafe configurations described by a formula $K(a)$ amounts to check whether the formula $I(a_0) \wedge \tau(a_0, a_1) \wedge \dots \wedge \tau(a_{n-1}, a_n) \wedge K(a_n)$ is satisfiable for some n . Notably, several mature model checkers exist to ascertain safety of these type of systems, such as MCMT [18] and CUBICLE [9]. In [4,6], we have extended array-based systems towards an array-based version of the artifact-centric model, considering two main settings:

- (**simple artifact systems - SAS**) artifact systems operating over a read-only relational database that resembles our DAB catalog, and over a single tuple (or boundedly many tuples) of updatable elements;
- (**relational artifact systems - RAS**) systems that extend SAS with a relational storage for unboundedly many updatable elements.

Notably, safety properties are checked over such systems irrespectively of the content of the read-only database, following the tradition of [28,15]. Several soundness, completeness, and decidability results have been obtained by suitably controlling the expressiveness of these systems. In addition, from Version 2.8 MCMT has been extended to handle safety checking of SAS and RAS.

3.2 Verification Problems for DABs

First of all, we need a language to express undesired properties over a DAB $\mathcal{M} = \langle \mathcal{D}, \mathcal{P} \rangle$. To do so, we resort to the same *guard* language introduced in Definition 6, extended with two features. First, we support querying also the control map of each case, so as to express control-flow properties. This is done by simply extending \mathcal{D} with additional, special case control variables referring to the lifecycle state of the blocks in \mathcal{P} (where each block B gets variable **Blifecycle**). Each such variable is assigned, in a given snapshot, to the state of the corresponding

block lifecycle (i.e., `idle`, `enabled`, and the like). We call these additional case variables $F_{\mathcal{P}}$. Second, we allow the modeler to inspect the status of multiple cases at once, so as to check how different cases may implicitly influence each other via the global repository. This is done by “indexing” guards, with the assumption that different indexes denote different cases.

Definition 10. *Let I be a set of n distinct indexes. A property over n cases of $\mathcal{M} = \langle \mathcal{D}, \mathcal{P} \rangle$ is a formula of the form $\bigwedge_{i \in I} G_i[i]$, where each G_i is a guard over \mathcal{D} and the case control variables $F_{\mathcal{P}}$ satisfying the following condition: every non-case variable $x \in \text{normvars}(G_i)$ appearing in G_i must also appear in an atom whose relation schema belongs to the repository $\mathcal{D}.\text{repo}$. \triangleleft*

Example 5. By calling HP the root, process block of Figure 1, we can test whether some case of the process can terminate through the property $(\text{HPlifecycle} = \text{completed})[i]$. If the property is unreachable, then no case can be progressed from the start to the end of the process. Since DAB processes are block structured, this is enough to ascertain whether the process is *unsound*. \triangleleft

We study (un)safety of these properties by considering general, unrestricted DABd, and also two interesting fragments:

(case-bounded) DABs that only introduce a bounded number of cases during their execution, where the bound is known a-priori - a typical setting being the one where a single case is studied, in the style of soundness verification for workflow nets;

(repo-bounded) DABs that introduce only boundedly many tuples in the repository during their execution, where the bound is known a-priori.

In this light, verification of safety properties over DABs present multiple dimensions of parameterization, depending on the type of DAB under analysis. The two extremes are: (i) case- and repo-bounded DABs, for which the only parameter is the instantiation of the catalog; (ii) unrestricted DABs, where verification is carried out parametrically w.r.t. the instantiation of the catalog, the number of tuples in the repository, and the overall number of cases.

3.3 Translating DABs into Array-Based Artifact Systems

To attack parameterized verification problems over DABs, we translate them into corresponding verification problems over SAS and RAS. We only provide here the main intuitions behind the translation, which is fully addressed in the Appendix. Let $\mathcal{M} = \langle \mathcal{D}, \mathcal{P} \rangle$ be a DAB. $\mathcal{D}.\text{cat}$ is maintained unaltered, as it is addressed in SAS and RAS in its full generality. The translation of $\mathcal{D}.\text{cvars}$ and $\mathcal{D}.\text{repo}$ depends instead on whether \mathcal{M} is studied unrestrictedly, or under case- and/or repo-boundedness. Consider $\mathcal{D}.\text{repo}$. If \mathcal{M} is repo-bounded with a bound of 1, then every relation schema in $\mathcal{D}.\text{repo}$ has just one tuple, and consequently can be represented using a set of (global) variables, one per relation attribute, in the style of a SAS. A k -bounded setting is handled similarly, just replicating the variables for k times. If instead no boundedness assumption on the report is placed, for each relation schema $R \in \mathcal{D}.\text{repo}$, and each attribute $a \in R.\text{attrs}$, a dedicated array is introduced. The index of the array represents the (implicit)

identifier of R , in line with our repository model. To reconstruct a specific tuple from R , one just needs to retrieve the objects present in the arrays corresponding to the different attributes of R , always using the same index i . The resulting model corresponds to that of a RAS and its notion of *artifact relation* [4,6].

A similar strategy is adopted for the case variables: if no bound on the number of cases is given, then each case variable is translated into a corresponding array, whose elements maintain the value that one case is assigning to that variable. Accessing all such arrays with the same index produces back the entire case variable assignments for the corresponding case. Finally, **self** is handled by introducing an array with the property that its elements are in bijection with the indexes (i.e., no element repeats twice in the array). Exactly the same approach is replicated to store the control information about blocks on a per-case basis.

All in all, depending on the boundedness assumptions on cases and/or repository, the translation produces a SAS or a RAS with different artifact relations. Each transition formula realizes one of the progression rules that collectively realize the execution semantics of the input DAB (cf. Section 2.4).

In [4,6], we focus on parameterized (un)safety of RAS, verifying whether there exists an instance of the read-only database such that the artifact system can reach an unsafe configuration. Since the cells of the arrays may point to identifiers in the catalog, in turn related to other catalog relations via foreign keys, the standard backward reachability procedure needs to be suitably revised [6]. In fact, when computing preimage formulae over RAS, existentially quantified “data” variables may be introduced, breaking the format of state formulae. To restore the key property that the preimage of a state is again represented symbolically as a state formula, such additional quantified variables must be eliminated. Suitable quantifier elimination techniques have been studied in [6,7] and implemented in the latest version 2.8 of MCMT, which can now natively handle the verification of RAS. In addition, while the unsafety verification is in general undecidable for RAS, several subclasses with decidable unsafety have been singled out. One of such classes corresponds to RAS operating over arrays whose maximum size is bounded a-priori, i.e. SAS. All in all, the RAS framework provides a natural foundational and practical basis to formally analyze DABs, which we tackle next.

From now on, we use **BackReach** to refer to the backward reachability procedure that:

- takes as input a DAB, a property to be verified, and a series of additional information related to the boundedness assumptions and the adopted semantics for insertion (set vs multiset);
- translates the input DAB into a corresponding SAS/RAS (according to the provided additional information), and the input property into a corresponding property over the target SAS/RAS;
- invokes the SAS/RAS backward reachability procedure described in [4,6] and implemented in MCMT;
- returns *yes* if and only if the property is reachable.

All the proofs of the following theorems are obtained by exploiting the translation, and by showing that the SAS/RAS produced from the translation enjoys the property stated in the theorem.

3.4 Soundness and Completeness Results

We start by considering case-bounded systems.

Theorem 1. *BackReach is sound and complete for case-bounded DABs that use the multiset or set insertion semantics.* \triangleleft

While for case-bounded DABs soundness and completeness are guaranteed without additional restrictions, this is not the case in the unrestricted setting. The problem is, in fact, the usage of **self**, which implicitly allows to create references across read-write relations (something that is not allowed in a DAB repository, nor in the corresponding model of RAS). We recover soundness and completeness by disallowing the explicit usage of **self**.

Definition 11. *A DAB \mathcal{M} is case-identifier-agnostic if none of the update specifications in \mathcal{M} mentions **self**.* \triangleleft

Theorem 2. *BackReach is sound and complete for case-identifier-agnostic DABs that use the multiset or set insertion semantics.* \triangleleft

We stress here that soundness and completeness indicate that whenever BackReach terminates, it produces a correct answer. Termination is not guaranteed in the general case (but may very well be obtained on the analyzed DAB), and consequently BackReach is a semi-decision procedure.

3.5 Termination Results

We now discuss how the previous results can be strengthened to ensure termination (thus witnessing decidability of parameterized verification). The first, unavoidable limitation that we have to impose is on the constraints used in the catalog: its foreign keys cannot form cycles. This is in line with [20,6]. To define acyclicity, we associate a catalog Cat to a characteristic graph $G(Cat)$ that captures the dependencies between relation schema components induced by primary and foreign keys. Specifically, $G(Cat)$ is a directed graph such that:

- for every $R \in Cat$ and every attribute $a \in R.attrs$, the pair $\langle R, a \rangle$ is a node of $G(Cat)$ (and nothing else is a node);
- $\langle R_1, a_1 \rangle \rightarrow \langle R_2, a_2 \rangle$ is and only if one of the two cases apply: (i) $R_1 = R_2$, $a_2 \neq a_1$, and $a_1 = R.id$; (ii) $a_2 = R_2.id$ and a_1 is a foreign key referring R_2 .

Definition 12. *A DAB is acyclic if the characteristic graph of its catalog is so.* \triangleleft

Theorem 3. *BackReach terminates when verifying properties over case- and repo-bounded, acyclic DABs using the multiset or set insertion semantics.* \triangleleft

This strong result is obtained due to the fact that case- and repo-bounded DABs get translated into SAS, where the read-write storage is constituted by a fixed set of variables. If instead we consider more sophisticated DABs that get translated into RAS with their sophisticated read-write relational storage, then termination requires to carefully control the interplay between the different components of

the DAB. While the required conditions are quite difficult to grasp at the syntactic level, they can be intuitively understood as follows: to ensure termination, whenever the progression of the DAB depends on the repository, it does so only via a single entry in one of its relations. This indicates that direct or indirect comparisons and joins of distinct tuples within the same or different repository relations cannot be used in an update. Towards avoiding indirect joins, queries cannot mix case variables and repository relations, nor update case variables with the content of other case variables. The following definition is instrumental to enforce this principle.

Definition 13. A guard $G \triangleq q(\vec{x}) \leftarrow \bigvee_{i=1}^n Q_i$ over a data component \mathcal{D} is separated if, for every i, j we have that $\text{normvars}(Q_i) \cap \text{normvars}(Q_j) = \emptyset$ and each Q_i is of the form $\chi \wedge R(\vec{y}) \wedge \xi$ (here, χ , $R(\vec{y})$ or ξ are optional), where: (i) χ is a conjunctive query with filters over $\mathcal{D}.\text{cat}$ only (that can employ case variables); (ii) $R \in \mathcal{D}.\text{repo}$ is a repository relation schema; (iii) \vec{y} is a tuple of variables and/or constant objects in \mathbb{D} , such that $\vec{y} \cap \mathcal{D}.\text{cvars} = \emptyset$, and $\text{normvars}(\chi) \cap \vec{y} = \emptyset$; (iv) ξ is a conjunctive query with filters over $\mathcal{D}.\text{cat}$ only, that possibly mentions variables in \vec{y} but does not include any case variable (i.e., $\text{casevars}(\xi) = \emptyset$), and such that $\text{normvars}(\chi) \cap \text{normvars}(\xi) = \emptyset$. A property is separated if all its inner guards are separated. \triangleleft

Intuitively, a separated guard consists of two isolated parts: one part χ inspecting the content of case variables and their relationship with the catalog, and another part $R(\vec{y}) \wedge \xi$ retrieving a single tuple \vec{y} in some repository relation R , possibly filtering it through inspection of the catalog via ξ .

Example 6. Consider the refinement $\text{EvalApp.pre} \triangleq \text{GetScore}(s : \text{NumScore}) \leftarrow \xi \wedge \chi$ of the guard EvalApp.pre from Example 2, where $\chi := \text{User}(\mathbf{uid}, \text{name}, \text{age})$ checks if the variables $\langle \mathbf{uid}, \text{name}, \text{age} \rangle$ form a tuple in User , and $\xi := 1 \leq s \wedge s \leq 100$. This guard is separated since χ and ξ match the requirements of the previous definition. \triangleleft

Theorem 4. Let \mathcal{M} be a case-bounded, acyclic DAB that uses the multiset insertion semantics, and is so that for each update specification \mathbf{u} of \mathcal{M} , the following holds:

- If $\mathbf{u}.\text{eff}$ is an insert&set rule (with an explicit INSERT part), then $\mathbf{u}.\text{pre}$ is repo-free;
- If $\mathbf{u}.\text{eff}$ is a set rule (not containing an INSERT part), then either (i) $\mathbf{u}.\text{pre}$ is repo-free, or (ii) $\mathbf{u}.\text{pre}$ is separated and all case variables $\mathcal{D}.\text{cvars} \setminus \{\mathbf{self}\}$ appear in the SET part of $\mathbf{u}.\text{eff}$;
- If $\mathbf{u}.\text{eff}$ is a delete&set rule, then $\mathbf{u}.\text{pre}$ is separated and all case variables $\mathcal{D}.\text{cvars} \setminus \{\mathbf{self}\}$ appear in the SET part of $\mathbf{u}.\text{eff}$;
- If $\mathbf{u}.\text{eff}$ is a conditional update rule, then $\mathbf{u}.\text{pre}$ is repo-free and boolean, so that $\mathbf{u}.\text{eff}$ only makes use of the new variables introduced in its UPDATE part (as well as constant objects in \mathbb{D}).

Then, BackReach terminates when verifying separated properties over \mathcal{M} . \triangleleft

Theorem 5. Let \mathcal{M} be a case-identifier-agnostic acyclic DAB that uses the multiset insertion semantics, and is so that for each update specification \mathbf{u} in \mathcal{M} ,

u satisfies the same conditions as of Theorem 4. Then, **BackReach** terminates when verifying separated properties over \mathcal{M} . \triangleleft

It is important to notice that the conditions in Theorems 4 and 5 represent a concrete, BPMN-like counterpart of the decidability results in [20] and [6].

Example 7. Our hiring job DAB makes use of **self** towards relating applications to the identifier of the case to which they were submitted. Hence, if we want to retain soundness and completeness of **BackReach**, we have to restrict the analysis to the case-bounded setting. By considering the data and process schema of the DAB, we can directly show that it obeys to all conditions in Theorem 1, in turn guaranteeing termination of **BackReach**. \triangleleft

4 First Experiments with MCMT

We have manually encoded the job hiring DAB described in the paper as an MCMT specification, using the same translation rules recalled in Section 3.3 and fully spelled out in the Appendix towards proving the main theorems in Sections 3.4 and 3.5. MCMT checks unsafety of a property through a symbolic, backward reachability procedure. The algorithm computes iterated preimages of the given property and applies to them quantifier elimination, until a fixpoint is reached or until a set intersecting the initial state (i.e., also characterized using a formula) is found. To do this efficiently, MCMT is equipped with dedicated quantifier elimination techniques, and discharge safety and fixpoint tests encountered during the backward search to state-of-the-art SMT solvers.

We have checked the encoding of the process against five safe and five unsafe properties. The first property ascertains whether a job hiring case can actually reach the end point of the process, in turn witnessing soundness. The **BPM_SAFE1** property checks whether it is possible to have a situation where the **Select Winner** task is **enabled** and the case variable **result** indicates that the winner is not eligible.

Table 1 shows the so-obtained, very encouraging results. Such initial results nicely complement those carried out on RAS in [4], indicating that the approach is promising not just foundationally, but also in terms of tool support. These experiments, together with the ones reported in this paper, are available as part of the last distribution 2.8 of MCMT.² Experiments were performed on a machine with Ubuntu 16.04, 2.6 GHz Intel Core i7 and 16 GB RAM.

5 Conclusion

In this paper, we have introduced a data-aware version of BPMN, called DAB, that achieves an interesting trade-off between expressiveness, and the possibility of applying sophisticated parameterized verification techniques to ascertain

² <http://users.mat.unimi.it/users/ghilardi/mcmt/>, subdirectory [/examples/dbdriven](#) of the distribution. The user manual contains a new section (pages 36–39) on how to encode RASs in MCMT specifications.

Exp.	Res.	Time (s)	Exp.	Res.	Time (s)
BPM_end_process	UNSAFE	0.43			
BPM_SAFE1	SAFE	0.20	BPM_UNSAFE1	UNSAFE	0.18
BPM_SAFE2	SAFE	5.85	BPM_UNSAFE2	UNSAFE	1.17
BPM_SAFE3	SAFE	3.56	BPM_UNSAFE3	UNSAFE	4.45
BPM_SAFE4	SAFE	0.03	BPM_UNSAFE4	UNSAFE	1.43
BPM_SAFE5	SAFE	0.27	BPM_UNSAFE5	UNSAFE	1.14

Table 1. Time spent by MCMT to check different properties over the job hiring DAB. Names of experiments coincide with those of the MCMT files from the Ancillary files of arXiv:1905.12991.

safety of the produced models. In particular, we have identified classes of DABs for which backward reachability techniques coming from the SMT tradition are correct, further strengthening them to also guarantee termination of backward reachability. From the foundational point of view, we are interested in equipping DABs with datatypes and corresponding rigid predicates, including arithmetic operators, as done in [15] for artifact-centric systems. This is promising especially considering that there are plenty of state-of-the-art SMT techniques to handle arithmetics. At the same time, we want to attack the main limitation of our approach, namely that guards and conditions are actually existential formulae, and the only (restricted) form of universal quantification available in the update language is that of conditional updates. Universal guards in transition formulae could be very useful in specifications: for example, they would allow us to specify a branch in a job hiring process that is followed only if no applicant satisfies a certain condition. This is reminiscent to what happens in the verification of distributed systems, where universal guards frequently occur in specifications. The question has been debated since longtime in the literature and the most effective solution adopted to cope with this problem so far is the introduction of suitable “monotonic abstractions” (see [1] for a survey). Notably, this solution is currently implemented in MCMT. Monotonic abstractions could introduce spurious unsafe traces, and in fact MCMT warns the user about this (in practice, not so frequent) possibility.

An orthogonal, challenging question is how, and to what extent, some of the most recent techniques developed for temporal model checking of artifact-centric systems [15] can be incorporated in our approach, allowing us to prove more sophisticated properties beyond safety.

From the experimental point of view, while a systematic evaluation is out of scope of this paper, the initial experiments carried out in this paper and [4] indicate that the approach is promising. We intend to fully automate the translation from DABs to array-based systems, and to set up a benchmark to evaluate the performance of verifiers for data-aware processes, starting from the examples collected in [20] (which are inspired from BPMN processes, and consequently should be straightforwardly encoded as DABs).

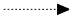




References

1. F. Alberti, S. Ghilardi, and N. Sharygina. Monotonic abstraction techniques: from parametric to software model checking. In *Proc. MOD**, EPTCS, pages 1–11, 2014.
2. M. Bojańczyk, L. Segoufin, and S. Toruńczyk. Verification of database-driven systems via amalgamation. In *Proc. PODS*, pages 63–74, 2013.
3. D. Calvanese, G. De Giacomo, and M. Montali. Foundations of data aware process analysis: A database theory perspective. In *Proc. PODS*, pages 1–12, 2013.
4. D. Calvanese, S. Ghilardi, A. Gianola, M. Montali, and A. Rivkin. Verification of data-aware processes via array-based systems (extended version). Technical Report arXiv:1806.11459, arXiv.org, 2018.
5. D. Calvanese, S. Ghilardi, A. Gianola, M. Montali, and A. Rivkin. Formal modeling and SMT-based parameterized verification of data-aware BPMN (extended version). Technical Report arXiv:1906.07811, arXiv.org, 2019.
6. D. Calvanese, S. Ghilardi, A. Gianola, M. Montali, and A. Rivkin. From model completeness to verification of data-aware processes. In *Description Logic, Theory Combination, and All That*, LNCS. Springer, 2019.
7. D. Calvanese, S. Ghilardi, A. Gianola, M. Montali, and A. Rivkin. Model completeness, covers and superposition. In *Proc. of CADE*, 2019.
8. C. Combi, B. Oliboni, M. Weske, and F. Zerbato. Conceptual modeling of processes and data. In *Proc. ER*, volume 11157 of LNCS, pages 236–250. Springer, 2018.
9. S. Conchon, A. Goel, S. Krstic, A. Mebsout, and F. Zaïdi. Cubicle: A parallel SMT-based model checker for parameterized systems - Tool paper. In *Proc. CAV*, pages 718–724, 2012.
10. E. Damaggio, A. Deutsch, and V. Vianu. Artifact systems with data dependencies and arithmetic. *ACM TODS*, 37(3):22, 2012.
11. G. De Giacomo, X. Oriol, M. Estañol, and E. Teniente. Linking data and BPMN processes to achieve executable models. In *Proc. CAISE*, 2017.
12. M. de Leoni, P. Felli, and M. Montali. A holistic approach for soundness verification of decision-aware process models. In *Proc. ER*, LNCS, pages 219–235. Springer, 2018.
13. R. De Masellis, C. Di Francescomarino, C. Ghidini, M. Montali, and S. Tessaris. Add data into business process verification: Bridging the gap between theory and practice. In *Proc. AAAI*, pages 1091–1099. AAAI Press, 2017.
14. A. Deutsch, R. Hull, F. Patrizi, and V. Vianu. Automatic verification of data-centric business processes. In *Proc. ICDT*, pages 252–267, 2009.
15. A. Deutsch, Y. Li, and V. Vianu. Verification of hierarchical artifact systems. In *Proc. PODS*, pages 179–194, 2016.
16. S. Ghilardi, E. Nicolini, S. Ranise, and D. Zucchelli. Towards SMT model checking of array-based systems. In *Proc. IJCAR*, pages 67–82, 2008.
17. S. Ghilardi and S. Ranise. Backward reachability of array-based systems by SMT solving: Termination and invariant synthesis. *Logical Methods in Computer Science*, 6(4), 2010.
18. S. Ghilardi and S. Ranise. MCMT: A model checker modulo theories. In *Proc. IJCAR*, 2010.
19. S. Lasota. Decidability border for petri nets with data: WQO dichotomy conjecture. In *Proc. PETRI NETS*, volume 9698 of LNCS, pages 20–36. Springer.
20. Y. Li, A. Deutsch, and V. Vianu. VERIFAS: A practical verifier for artifact systems. *PVLDB*, 11(3):283–296, 2017.

21. G. Meroni, L. Baresi, M. Montali, and P. Plebani. Multi-party business process compliance monitoring through iot-enabled artifacts. *Inf. Syst.*, 73:61–78, 2018.
22. A. Meyer, L. Pufahl, D. Fahland, and M. Weske. Modeling and enacting complex data dependencies in business processes. In *Proc. BPM*, 2013.
23. M. Montali and A. Rivkin. Model checking petri nets with names using data-centric dynamic systems. *Form. Asp. Comp.*, pages 1–27, 2016.
24. M. Montali and A. Rivkin. DB-Nets: on the marriage of colored Petri Nets and relational databases. *ToPNoC*, 28(4), 2017.
25. M. Reichert. Process and data: Two sides of the same coin? In *Proc. OTM*, 2012.
26. C. Richardson. Warning: Don’t assume your business processes use master data. In *Proc. BPM*, pages 11–12, 2010.
27. F. Rosa-Velardo and D. de Frutos-Escrig. Decidability and complexity of petri nets with unordered data. *Theor. Comput. Sci.*, 412(34):4439–4451, 2011.
28. V. Vianu. Automatic verification of database-driven systems: a new frontier. In *Proc. ICDT*, pages 1–13, 2009.

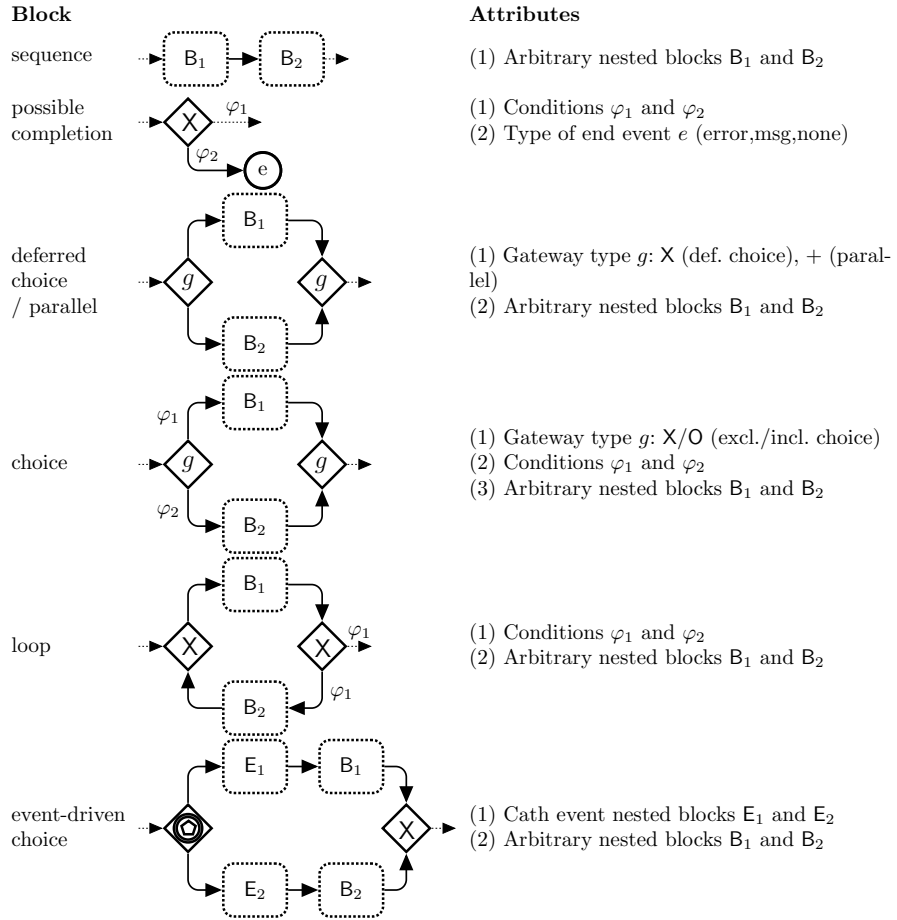
A DAB Blocks

A.1 Basic Blocks

Block		Attributes
empty		
task		(1) Atomic/non-atomic (2) update specification.
catch event		(1) Type of event e (msg, timer, none) (2) update specification.
process block		(1) Type of start event e_s (msg, timer, none) (2) Update specification of e_s (3) Type of end event e_t (msg, none) (4) Update specification of e_t (5) Arbitrary nested block B
subprocess		(1) Inner process block

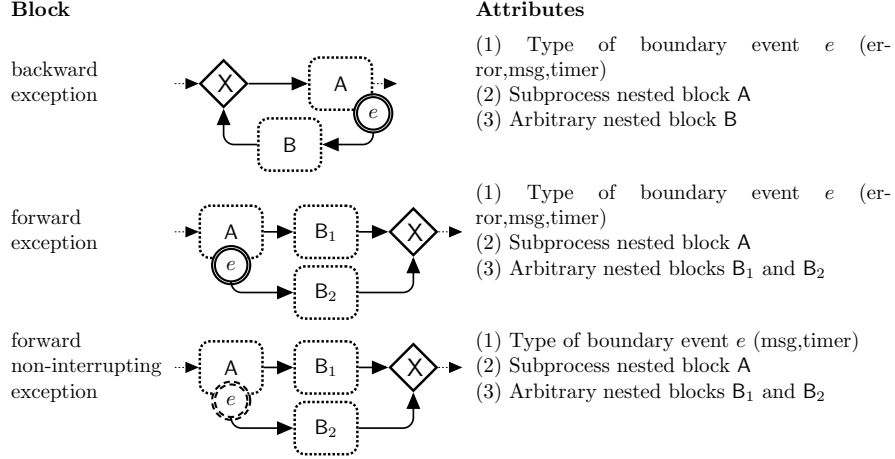
A.2 Flow Blocks

For simplicity, we consider only two nested blocks, but multiple nested blocks can be seamlessly handled.



A.3 Exception Handling Blocks

For simplicity, we show a single boundary event, but multiple boundary events and their corresponding handlers can be seamlessly handled.



B Preliminaries

In this section we give a review of the preliminaries needed to present RASs.

We adopt the usual first-order syntactic notions of signature, term (denoted with t_1, t_2, \dots), atom, (ground) formula, and so on. We use \underline{u} to represent a tuple $\langle u_1, \dots, u_n \rangle$. Our signatures Σ are multi-sorted and include equality for every sort, which implies that variables are sorted as well. Depending on the context, we keep the sort of a variable implicit, or we indicate explicitly in a formula that variable x has sort S by employing notation $x : S$. The notation $t(\underline{x})$, $\phi(\underline{x})$ means that the term t , the formula ϕ has free variables included in the tuple \underline{x} . We are concerned only with constants and function symbols f , each of which has *sources* \underline{S} and a *target* S' , denoted as $f : \underline{S} \rightarrow S'$. We assume that terms and formulae are well-typed, in the sense that the sorts of variables, constants, and function sources/targets match. A formula is said to be *universal* (resp., *existential*) if it has the form $\forall \underline{x}(\phi(\underline{x}))$ (resp., $\exists \underline{x}(\phi(\underline{x}))$), where ϕ is a quantifier-free formula. Formulae with no free variables are called *sentences*.

From the semantic side, we use the standard notions of a Σ -*structure* \mathcal{M} and of *truth* of a formula in a Σ -structure under an assignment to the free variables. A Σ -*theory* T is a set of Σ -sentences; a *model* of T is a Σ -structure \mathcal{M} where all sentences in T are true. We use the standard notation $T \models \phi$ to say that ϕ is true in all models of T for every assignment to the free variables of ϕ . We say that ϕ is *T-satisfiable* iff there is a model \mathcal{M} of T and an assignment to the free variables of ϕ that make ϕ true in \mathcal{M} .

In the following, we specify transitions of an artifact-centric system using first-order formulae. To obtain a more compact representation, we make use

there of definable extensions as a means for introducing so-called *case-defined functions*. We fix a signature Σ and a Σ -theory T ; a T -partition is a finite set $\kappa_1(\underline{x}), \dots, \kappa_n(\underline{x})$ of quantifier-free formulae such that $T \models \forall \underline{x} \bigvee_{i=1}^n \kappa_i(\underline{x})$ and $T \models \bigwedge_{i \neq j} \forall \underline{x} \neg(\kappa_i(\underline{x}) \wedge \kappa_j(\underline{x}))$. Given such a T -partition $\kappa_1(\underline{x}), \dots, \kappa_n(\underline{x})$ together with Σ -terms $t_1(\underline{x}), \dots, t_n(\underline{x})$ (all of the same target sort), a *case-definable extension* is the Σ' -theory T' , where $\Sigma' = \Sigma \cup \{F\}$, with F a “fresh” function symbol (i.e., $F \notin \Sigma$)³, and $T' = T \cup \bigcup_{i=1}^n \{\forall \underline{x} (\kappa_i(\underline{x}) \rightarrow F(\underline{x}) = t_i(\underline{x}))\}$. Intuitively, F represents a case-defined function, which can be reformulated using nested if-then-else expressions and can be written as $F(\underline{x}) := \text{case of } \{\kappa_1(\underline{x}) : t_1; \dots; \kappa_n(\underline{x}) : t_n\}$. By abuse of notation, we identify T with any of its case-definable extensions T' . In fact, it is easy to produce from a Σ' -formula ϕ' a Σ -formula ϕ equivalent to ϕ' in all models of T' : just remove (in the appropriate order) every occurrence $F(\underline{v})$ of the new symbol F in an atomic formula A , by replacing A with $\bigvee_{i=1}^n (\kappa_i(\underline{v}) \wedge A(t_i(\underline{v})))$. We also exploit λ -abstractions (see, e.g., formula (2) below) for a more compact (still first-order) representation of some complex expressions, and always use them in atoms like $b = \lambda y. F(y, \underline{z})$ as abbreviations of $\forall y. b(y) = F(y, \underline{z})$ (where, typically, F is a symbol introduced in a case-defined extension as above).

C Array-based Model

In this section we recall the definitions of the formal setting presented in [4,6] that is exploited in this paper as target model of our translation. This setting relies on array-based systems and provides a general framework where introducing safety verification problems for artifact-centric models called Relational Artifact Systems (RASs). Those models are verbatim of the ones presented [4,6], but we prefer presenting them in detail here for safe of self-containedness.

In the following, we formally define RASs. Since they are array-based systems, we start by recalling the intuition behind them.

In general terms, an array-based system is described using a multi-sorted theory that contains two types of sorts, one accounting for the indexes of arrays, and the other for the elements stored therein. Since the content of an array changes over time, it is referred to using a *function* variable, whose interpretation in a state is that of a total function mapping indexes to elements (so that applying the function to an index denotes the classical *read* operation for arrays). The definition of an array-based system with array state variable a always requires: a formula $I(a)$ describing the *initial configuration* of the array a , and a formula $\tau(a, a')$ describing a *transition* that transforms the content of the array from a to a' . In such a setting, verifying whether the system can reach unsafe configurations described by a formula $K(a)$ amounts to check whether the formula $I(a_0) \wedge \tau(a_0, a_1) \wedge \dots \wedge \tau(a_{n-1}, a_n) \wedge K(a_n)$ is satisfiable for some n .

Following the tradition of artifact-centric systems [14,10,2,15], an array-based Relational Artifact Systems (RAS) consists of a read-only DB, a read-write work-

³ Arity and source/target sorts for F can be deduced from the context (considering that everything is well-typed).

ing memory for artifacts (which are used in our translation for formalizing the set of case variables for every process instance and the shared evolving relations), and a finite set of actions (also called services) that inspect the relational database and the working memory, and determine the new configuration of the working memory.

C.1 Read-only DB schemata

We now provide a formal definition of (read-only) DB-schemas by relying on an algebraic, functional characterization.

Definition 14. A DB schema is a pair $\langle \Sigma, T \rangle$, where: (i) Σ is a DB signature, that is, a finite multi-sorted signature whose only symbols are equality, unary functions, and constants; (ii) T is a DB theory, that is, a set of universal Σ -sentences. \triangleleft

Next, we refer to a DB schema simply through its (DB) signature Σ and (DB) theory T , and denote by Σ the set of sorts and by Σ_{fun} the set of functions in Σ .

Remark 1. If desired, we can freely extend DB schemas by adding arbitrary n -ary relation symbols to the signature Σ . For this purpose, we give the following definition.

Definition 15. A DB extended-schema is a pair $\langle \Sigma, T \rangle$, where: (i) Σ is a DB extended-signature, that is, a finite multi-sorted signature whose only symbols are equality, n -ary relations, unary functions, and constants; (ii) T is a DB extended-theory, that is, a set of universal Σ -sentences. \triangleleft

Since for our application we are only interested in relations with primary and foreign key dependencies (even if our implementation takes into account also the case of “free” relations, i.e. without key dependencies), we restrict our focus on DB schemas, which are sufficient to capture those constraints (as explained in the following subsection). We notice that, in case the assumptions over DB schemas that discussed below hold for DB extended-theories, all the results presented in Section D (and Theorem 7) still hold even considering DB extended-schemas instead of DB schemas. \triangleleft

We associate to a DB signature Σ a characteristic graph $G(\Sigma)$ capturing the dependencies induced by functions over sorts. Specifically, $G(\Sigma)$ is an edge-labeled graph whose set of nodes is Σ , and with a labeled edge $S \xrightarrow{f} S'$ for each $f : S \rightarrow S'$ in Σ_{fun} . We say that Σ is *acyclic* if $G(\Sigma)$ is so. The *leaves* of Σ are the nodes of $G(\Sigma)$ without outgoing edges. These terminal sorts are divided in two subsets, respectively representing *unary relations* and *value sorts*. Non-value sorts (i.e., unary relations and non-leaf sorts) are called *id sorts*, and are conceptually used to represent (identifiers of) different kinds of objects. Value sorts, instead, represent datatypes such as strings, numbers, clock values, etc. We denote the set of id sorts in Σ by Σ_{ids} , and that of value sorts by Σ_{val} , hence $\Sigma = \Sigma_{ids} \uplus \Sigma_{val}$.

We now consider extensional data.

Definition 16. A DB instance of DB schema $\langle \Sigma, T \rangle$ is a Σ -structure \mathcal{M} that is a model of T and such that every id sort of Σ is interpreted in \mathcal{M} on a finite set. \triangleleft

What may appear as not customary in Definition 16 is the fact that value sorts can be interpreted on infinite sets. This allows us, at once, to reconstruct the classical notion of DB instance as a finite model (since only finitely many values can be pointed from id sorts using functions), at the same time supplying a potentially infinite set of fresh values to be dynamically introduced in the working memory during the evolution of RASs.

We respectively denote by $S^{\mathcal{M}}$, $f^{\mathcal{M}}$, and $c^{\mathcal{M}}$ the interpretation in \mathcal{M} of the sort S (this is a set), of the function symbol f (this is a set-theoretic function), and of the constant c (this is an element of the interpretation of the corresponding sort). Obviously, $f^{\mathcal{M}}$ and $c^{\mathcal{M}}$ must match the sorts in Σ . E.g., if f has source S and target U , then $f^{\mathcal{M}}$ has domain $S^{\mathcal{M}}$ and range $U^{\mathcal{M}}$.

We close the formalization of DB schemas by discussing DB theories, whose role is to encode background axioms. We illustrate a typical background axiom, required to handle the possible presence of *undefined identifiers/values* in the different sorts. This axiom is essential to capture artifact systems whose working memory is initially undefined, in the style of [15,20]. To specify an undefined value we add to every sort S of Σ a constant \mathbf{undef}_S (written from now on, by abuse of notation, just as \mathbf{undef} , used also to indicate a tuple). Then, for each function symbol f of Σ , we add the following axiom to the DB theory:

$$\forall x (x = \mathbf{undef} \leftrightarrow f(x) = \mathbf{undef}) \quad (1)$$

This axiom states that the application of f to the undefined value produces an undefined value, and it is the only situation for which f is undefined.

As discussed in [4], the theory T from Definition 14 must satisfy few crucial model-theoretic requirements for our approach to work: these requirements are *finite model property*, *decidable constraint satisfiability* and the *existence of T^** , i.e. the *model completion of T* . Specifically, the backward reachability procedure requires the availability of suitable quantifier elimination algorithms. However, a DB theory T does not necessarily have quantifier elimination; nevertheless, it is often possible to strengthen T in a conservative way (with respect to constraint satisfiability) and get quantifier elimination. In order to do that, in [4] we consider the theory T^* (when it exists, it is unique), and we show that model completion turns out to be quite effective to attack the verification of dynamic systems operating over relational databases. In all this paper we assume that DB theories T have finite model property, decidable constraint satisfiability and that admit a model completion T^* . Specifically, from now on we assume that T consists of only Axioms (1): in this case, all the assumptions hold.

C.2 Working Memory

In array-based RASs, the working memory consist of *function* variables. These variables (usually called *arrays*) are supposed to model evolving relations, so-called *artifact relations* in [15,20]. The idea is to treat artifact relations in a

uniform way as we did for the read-only DB, where we used function symbols for representing relations with key dependencies: for the working memory, we need extra sort symbols (as explained in the translation section, each sort symbol corresponds to a database relation symbol) and extra unary function symbols. variables.

Given a DB schema Σ , an *artifact extension* of Σ is a signature Σ_{ext} obtained from Σ by adding to it some extra sort symbols⁴. These new sorts (usually indicated with letters E, F, \dots) are called *artifact sorts* (or *artifact relations* by some abuse of terminology), while the old sorts from Σ are called *basic sorts*. In array-based BPMN models, artifacts and basic sorts correspond, respectively, to the index and the elements sorts mentioned in the literature on array-based systems. Below, given $\langle \Sigma, T \rangle$ and an artifact extension Σ_{ext} of Σ , when we speak of a Σ_{ext} -model of T , a DB instance of $\langle \Sigma_{ext}, T \rangle$, or a Σ_{ext} -model of T^* , we mean a Σ_{ext} -structure \mathcal{M} whose reduct to Σ respectively is a model of T , a DB instance of $\langle \Sigma, T \rangle$, or a model of T^* .

An *artifact setting* over Σ_{ext} is a pair $(\underline{x}, \underline{a})$ given by a finite set \underline{x} of individual variables and a finite set \underline{a} of unary function variables: *the latter are required to have an artifact sort as source sort and a basic sort as target sort*. Variables in \underline{x} are called *artifact variables*, and variables in \underline{a} *artifact components*. Given a DB instance \mathcal{M} of Σ_{ext} , an *assignment* to an artifact setting $(\underline{x}, \underline{a})$ over Σ_{ext} is a map α assigning to every artifact variable $x_i \in \underline{x}$ of sort S_i an element $x^\alpha \in S_i^\mathcal{M}$ and to every artifact component $a_j : E_j \rightarrow U_j$ (with $a_j \in \underline{a}$) a set-theoretic function $a_j^\alpha : E_j^\mathcal{M} \rightarrow U_j^\mathcal{M}$. In our array-based RASs, artifact components and artifact variables correspond, respectively, to *arrays* and *constant arrays* (i.e., arrays with all equal elements) mentioned in the literature on array-based systems. Intuitively, an artifact setting represents the “working” memory of an array-based RAS.

We can view an assignment to an artifact setting $(\underline{x}, \underline{a})$ as a DB instance *extending* the DB instance \mathcal{M} as follows. Let all the artifact components in $(\underline{x}, \underline{a})$ having source E be $a_{i_1} : E \rightarrow S_1, \dots, a_{i_n} : E \rightarrow S_n$. Viewed as a relation in the artifact assignment (\mathcal{M}, α) , the artifact relation E “consists” of the set of tuples $\{\langle e, a_{i_1}^\alpha[e], \dots, a_{i_n}^\alpha[e] \rangle \mid e \in E^\mathcal{M}\}$. Thus each element of E is formed by an “entry” $e \in E^\mathcal{M}$ (uniquely identifying the tuple, and called “internal identifier” of the tuple $(e, a_{i_1}^\alpha[e], \dots, a_{i_n}^\alpha[e])$) and by “data” $\underline{a}_i^\alpha(e)$ taken from the read-only database \mathcal{M} . When the system evolves, the set $E^\mathcal{M}$ of entries remains fixed, whereas the components $\underline{a}_i^\alpha(e)$ may change: typically, we initially have $\underline{a}_i^\alpha(e) = \mathbf{undef}$, but these values are changed when some defined values are inserted into the relation modeled by E ; the values are then repeatedly modified (and possibly also reset to \mathbf{undef} , if the tuple is removed and e is re-set to point to undefined values)⁵.

⁴ By ‘signature’ we always mean ‘signature with equality’, so as soon as new sorts are added, the corresponding equality predicates are added too.

⁵ In accordance with MCMT conventions, we denote the application of an artifact component a to a term (i.e., constant or variable) v as $a[v]$ (standard notation for arrays), instead of $a(v)$.

C.3 RASs

In order to introduce verification problems in the symbolic setting of array-based systems, one first has to specify which formulae are used to represent

- sets of states,
- the system initializations, and
- system evolution.

To introduce relational dynamic systems we discuss the kind of formulae we use. In such formulae, we use notations like $\phi(\underline{z}, \underline{a})$ to mean that ϕ is a formula whose free individual variables are among the \underline{z} and whose free unary function variables are among the \underline{a} . Let $(\underline{x}, \underline{a})$ be an artifact setting over Σ_{ext} , where $\underline{x} = x_1, \dots, x_n$ are the artifact variables and $\underline{a} = a_1, \dots, a_m$ are the artifact components (their source and target sorts are left implicit).

An *initial formula* is a formula $\iota(\underline{x})$ of the form⁶ $(\bigwedge_{i=1}^n x_i = c_i) \wedge (\bigwedge_{j=1}^m a_j = \lambda y.d_j)$, where c_i, d_j are constants from Σ (typically, c_i and d_j are **undef**).

A *state formula* has the form $\exists \underline{e} \phi(\underline{e}, \underline{x}, \underline{a})$, where ϕ is quantifier-free and the \underline{e} are individual variables of artifact sorts.

A *transition formula* $\hat{\tau}$ has the form

$$\exists \underline{e} (\gamma(\underline{e}, \underline{x}, \underline{a}) \wedge \bigwedge_i x'_i = F_i(\underline{e}, \underline{x}, \underline{a}) \wedge \bigwedge_j a'_j = \lambda y.G_j(y, \underline{e}, \underline{x}, \underline{a})) \quad (2)$$

where the \underline{e} are individual variables (of *both* basic and artifact sorts), γ (the ‘guard’) is quantifier-free, $\underline{x}', \underline{a}'$ are renamed copies of $\underline{x}, \underline{a}$, and the F_i, G_j (the ‘updates’) are case-defined functions.ed

Definition 17. An array-based RAS is

$$\mathcal{S} = \langle \Sigma, T, \Sigma_{ext}, \underline{x}, \underline{a}, \iota(\underline{x}, \underline{a}), \tau(\underline{x}, \underline{a}, \underline{x}', \underline{a}') \rangle$$

where: (i) $\mathcal{DB} := \langle \Sigma, T \rangle$ is a (read-only) DB schema, (ii) Σ_{ext} is an artifact extension of Σ , (iii) $(\underline{x}, \underline{a})$ is an artifact setting over Σ_{ext} , (iv) ι is an initial formula, and (v) τ is a disjunction of transition formulae $\hat{\tau}$. \triangleleft

D Parameterized Safety via Backward Reachability in RAS.

All the result presented in this section come from [4], where all the proofs and the details are provided.

Given a RAS \mathcal{S} , we say that a *safety* formula for \mathcal{S} is a state formula $v(\underline{x})$ describing undesired states of \mathcal{S} . As usual in array-based systems, we say that \mathcal{S} is *safe with respect to* v if intuitively the system has no finite run leading from ι to v . Formally, there is no DB-instance \mathcal{M} of $\langle \Sigma_{ext}, T \rangle$, no $k \geq 0$, and no assignment in \mathcal{M} to the variables $\underline{x}^0, \underline{a}^0, \dots, \underline{x}^k, \underline{a}^k$ such that the formula

$$\iota(\underline{x}^0, \underline{a}^0) \wedge \tau(\underline{x}^0, \underline{a}^0, \underline{x}^1, \underline{a}^1) \wedge \dots \wedge \tau(\underline{x}^{k-1}, \underline{a}^{k-1}, \underline{x}^k, \underline{a}^k) \wedge v(\underline{x}^k, \underline{a}^k) \quad (3)$$

is true in \mathcal{M} (here $\underline{x}^i, \underline{a}^i$ are renamed copies of $\underline{x}, \underline{a}$). The *safety problem* for \mathcal{S} is the following: *given a safety formula* v *decide whether* \mathcal{S} *is safe with respect to* v .

⁶ Recall that $a_j = \lambda y.d_j$ abbreviates $\forall y a_j(y) = d_j$.

In order to assess safety of Data-aware BPMN models, we run the backward reachability procedures on RASs, by exploiting the translation of Data-Aware BPMN models into the array-based relational setting presented in the previous sections.

We shall introduce an algorithm that semi-decides safety problems for \mathcal{S} and then we shall examine some interesting cases where the algorithm terminates and gives a decision procedure. Algorithm 1

describes the *backward reachability algorithm* for handling the safety problem for array-based systems \mathcal{S} . An integral part of the algorithm is to compute preimages. For that purpose, for any $\phi_1(\underline{x}, \underline{x}')$ and $\phi_2(\underline{x})$, we define $Pre(\phi_1, \phi_2)$ to be the formula $\exists \underline{x}'(\phi_1(\underline{x}, \underline{x}') \wedge \phi_2(\underline{x}'))$. The *preimage* of the set of states described by a state formula $\phi(\underline{x})$ is the set of states described by $Pre(\tau, \phi)$.

The subprocedure $QE(T^*, \phi)$ in Line 6 applies the quantifier elimination algorithm of T^* (the model completion of T) to the existential formula ϕ . Algorithm 1 computes iterated preimages of v and applies to them quantifier elimination, until a fixpoint is reached or until a set intersecting the initial states (i.e., satisfying ι) is found. *Inclusion* (Line 2) and *disjointness* (Line 3) tests produce proof obligations that can be discharged via proof obligations to be handled by SMT solvers. The fixpoint is reached when the test in Line 2 returns *unsat*, which means that the preimage of the set of the current states is included in the set of states reached by the backward search so far.

We obtain the following Theorem (to understand the statement of the theorem, notice that by *partial correctness* we mean that, when the algorithm terminates, it gives a correct answer and by *effectiveness* we means that all sub-procedures in the algorithm can be effectively executed):

Theorem 6. *Backward search (cf. Algorithm 1) is effective and partially correct for solving safety problems for RASs. Specifically, it is sound and complete for detecting unsafety.* \triangleleft

Theorem 6 shows that backward search is a semi-decision procedure: if the system is unsafe, backward search always terminates and discovers it; if the system is safe, the procedure can diverge (but it is still correct). Notice that the role of quantifier elimination (Line 6 of Algorithm 1) is twofold: (i) It allows to discharge the fixpoint test of Line 2; (ii) it ensures termination in significant cases, namely those where (*strongly*) *local formulae*, introduced in the next section, are involved.

An interesting class of RASs is the one where the working memory consists *only* of artifact variables (without artifact relations): we call SASs such systems. For SASs, the following termination result holds:

Algorithm 1: Backward reachability algorithm

Function BReach(v)

- 1 $\phi \leftarrow v; B \leftarrow \perp;$
- 2 **while** $\phi \wedge \neg B$ *is T-satisfiable*
- 3 **do**
- 4 **if** $\iota \wedge \phi$ *is T-satisfiable.*
- 5 **then**
- 6 \perp **return unsafe**
- $B \leftarrow \phi \vee B;$
- $\phi \leftarrow Pre(\tau, \phi);$
- $\phi \leftarrow QE(T^*, \phi);$
- return (safe, B);**

Theorem 7. *Let $\langle \Sigma, T \rangle$ be a DB schema with Σ acyclic. Then, for every SAS $\mathcal{S} = \langle \Sigma, T, \underline{x}, \iota, \tau \rangle$, backward search terminates and decides safety problems for \mathcal{S} in PSPACE in the combined size of \underline{x} , ι , and τ . \triangleleft*

We remark that acyclicity of Σ is a strong condition and it is not needed in general, and that Theorem 7 holds also for DB extended-schemas (so, even adding “free relations” to the DB signatures). In fact, analyzing the proof of Theorem 7, it is clear that the decidability of the safety problems is guaranteed when in T there are only finitely many quantifier-free formulae in which \underline{x} occur: this happens, for instance, in case T has a purely relational signature or, more generally, even when T is a generic first-order theory (and not just a DB (extended)-schema) that is *locally finite*⁷.

E Termination with local transitions

We briefly recall the notion of *locality* and *strong locality* of transitions as presented in [4]. All the following notions (and the following theorem) are presented in [4].

Consider an acyclic signature Σ , a DB theory T and an artifact setting $(\underline{x}, \underline{a})$ over an artifact extension Σ_{ext} of Σ . We call a state formula *local* if it is a disjunction of the formulae

$$\exists e_1 \cdots \exists e_k (\delta(e_1, \dots, e_k) \wedge \bigwedge_{i=1}^k \phi_i(e_i, \underline{x}, \underline{a})), \quad (4)$$

and *strongly local* if it is a disjunction of the formulae

$$\exists e_1 \cdots \exists e_k (\delta(e_1, \dots, e_k) \wedge \psi(\underline{x}) \wedge \bigwedge_{i=1}^k \phi_i(e_i, \underline{a})). \quad (5)$$

In (4) and (5), δ is a conjunction of variable equalities and inequalities, ϕ_i, ψ are quantifier-free, and e_1, \dots, e_k are individual variables varying over artifact sorts. The key limitation of local state formulae is that they cannot compare entries from different tuples of artifact relations: each ϕ_i in (4) and (5) can contain only the existentially quantified variable e_i .

A transition formula $\hat{\tau}$ is *local* (resp., *strongly local*) if whenever a formula ϕ is local (resp., strongly local), so is $Pre(\hat{\tau}, \phi)$ (modulo the axioms of T^*). Examples of (strongly) local $\hat{\tau}$ are discussed in Appendix F in [4].

Theorem 8. *If Σ is acyclic, backward search (cf. Algorithm 1) terminates when applied to a local safety formula in a RAS whose τ is a disjunction of local transition formulae. \triangleleft*

F Translation of Data-aware BPMN models into array-based Systems

In this section, we define the translation of a Data-aware BPMN model into array-based RASs.

⁷ We say that T is locally finite iff for every finite tuple of variables \underline{x} there are only finitely many non T -equivalent atoms $A(\underline{x})$ involving only the variables \underline{x} .

F.1 Translation of the Data Schema

Catalogue. We now clarify how the classical, relational database Cat can be actually translated to the algebraic, functional characterization of symbolic DB schemata and instances. To technically explain the correspondence, we adopt the *named perspective*, where each relation schema is defined by a signature containing a *relation name* and a set of *typed attribute names*.

First of all, take the set \mathcal{S} of sorts of a DAB as set of basic sorts for the translated DB schema that we want to define. Let Cat be a catalogue as defined in Section 2.1. For every $n + 1$ -ary relation R in Cat , every attribute A of R is defined over a corresponding basic sort S_A in \mathcal{S} . Since each relation $R(id, \vec{A})$ in Cat must have a unary primary key as its first attribute $R.id$, we define a mechanism to correctly reference other attributes \vec{A} in R by means of unary functions in the DB signature $f_{R,A_i} : S_{R.id} \rightarrow S_{A_i}$ (where $i = 1, \dots, n$ and $\vec{A} = (A_1, \dots, A_n)$): $S_{R.id}$ is set as the id sort of R and, in the corresponding DB instance \mathcal{M} , $f_{R,A_i}^{\mathcal{M}}$ maps, for every tuple \underline{x} in $R^{\mathcal{M}}$, its identifier element (i.e., the first component of \underline{x}) to the unique element in the tuple corresponding to the attribute A_i . If A_i is an id sort of some other relation R' , f_{R,A_i} represents the foreign key referencing to R' .

Conversely, starting from a symbolic DB schema, we see how Definition 14 naturally corresponds to the definition of relational database schema equipped with single-attribute *primary keys* and *foreign keys*

Let $\langle \Sigma, T \rangle$ be a DB schema. Each id sort $S \in \Sigma_{ids}$ corresponds to a dedicated relation R_S with the following attributes: (i) one identifier attribute id_S with type S ; (ii) one dedicated attribute a_f with type S' for every function symbol $f \in \Sigma_{fun}$ of the form $f : S \rightarrow S'$.

The fact that R_S is built starting from functions in Σ naturally induces different database dependencies in R_S . In particular, for each non-id attribute a_f of R_S , we get a *functional dependency* from id_S to a_f ; altogether, such dependencies in turn witness that id_S is the (*primary*) *key* of R_S . In addition, for each non-id attribute a_f of R_S whose corresponding function symbol f has id sort S' as image, we get an *inclusion dependency* from a_f to the id attribute $id_{S'}$ of $R_{S'}$; this captures that a_f is a *foreign key* referencing $R_{S'}$.

Given a DB instance \mathcal{M} of $\langle \Sigma, T \rangle$, its corresponding *relational instance* \mathcal{I} is the minimal set satisfying the following property: for every id sort $S \in \Sigma_{ids}$, let f_1, \dots, f_n be all functions in Σ with domain S ; then, for every identifier $\circ \in S^{\mathcal{M}}$, \mathcal{I} contains a *labeled fact* of the form $R_S(id_S : \circ^{\mathcal{M}}, a_{f_1} : f_1^{\mathcal{M}}(\circ^{\mathcal{M}}), \dots, a_{f_n} : f_n^{\mathcal{M}}(\circ^{\mathcal{M}}))$. With this interpretation, the *active domain* of \mathcal{I} is the set

$$\bigcup_{S \in \Sigma_{ids}} (S^{\mathcal{M}} \setminus \{\mathbf{undef}^{\mathcal{M}}\}) \cup \left\{ v \in \bigcup_{V \in \Sigma_{val}} V^{\mathcal{M}} \mid v \neq \mathbf{undef}^{\mathcal{M}} \text{ and there exist } f \in \Sigma_{fun} \right. \\ \left. \text{and } \circ \in \text{dom}(f^{\mathcal{M}}) \text{ s.t. } f^{\mathcal{M}}(\circ) = v \right\}$$

consisting of all (proper) identifiers assigned by \mathcal{M} to id sorts, as well as all values obtained in \mathcal{M} via the application of some function. Since such values are necessarily *finitely many*, one may wonder why in Definition 16 we allow for interpreting value sorts over infinite sets. The reason is that, in our framework, an

evolving artifact system may use such infinite provision to inject and manipulate new values into the working memory. From the definition of active domain above, exploiting Axioms (1) we get that the membership of a tuple (x_0, \dots, x_n) to a generic $n + 1$ -ary relation R_S with key dependencies (corresponding to an id sort S) can be expressed in our setting by using just unary function symbols and equality:

$$R_S(x_0, \dots, x_n) \text{ iff } x_0 \neq \mathbf{undef} \wedge x_1 = f_1(x_0) \wedge \dots \wedge x_n = f_n(x_0) \quad (6)$$

Hence, the representation of negated atoms is the one that directly follows from negating (6):

$$\neg R_S(x_0, \dots, x_n) \text{ iff } x_0 = \mathbf{undef} \vee x_1 \neq f_1(x_0) \vee \dots \vee x_n \neq f_n(x_0) \quad (7)$$

Formula (6) exactly summarizes and explicitly shows the equivalence between symbolic DB schemata and relational Catalogues. Thus, in the following we will make use of relational Catalogues or DB schemata interchangeably.

This relational interpretation of DB schemas exactly reconstructs the requirements posed by [15,20] on the schema of the *read-only* database: (i) each relation schema has a single-attribute primary key; (ii) attributes are typed; (iii) attributes may be foreign keys referencing other relation schemas; (iv) the primary keys of different relation schemas are pairwise disjoint.

We stress that all such requirements are natively captured in our functional definition of a DB signature, and do not need to be formulated as axioms in the DB theory. The DB theory is used to express additional constraints, like that in Axioms (1) Notice that, in order to translate Data-aware BPMN models into the array-based setting, we just need to consider the DB theory with Axioms (1) only: this is because the empty DB theory itself is able to capture the requirements of Section 2.1, and Axioms (1) are needed to handle the undefined values in a correct way, that is every function symbol f maps the undefined value of one sort to the undefined value of another one, and it is the only case when f is undefined.

Repository and Case Variables. In the following, we denote an artifact assignment and a DB instance by using α and \mathcal{M} respectively.

Consider the set of process instances \mathcal{PI} of a data-aware BPMN model. We associate to \mathcal{PI} a fresh sort symbol PI_{index} . We call an artifact component with source PI_{index} “*case variables* artifact component” and we say that all the case variables artifact components form “the *case variables* artifact relation”. Intuitively, every tuple with first component $i \in PI_{index}$ of this artifact relation is used to formalize the values of the case variables for the specific process instance represented by the element i that is the *internal identifier* of the tuple. We associate to every case variable $v \in V_C$ with sort S an artifact component a_v with PI_{index} as its source sort and S as its target sort. Intuitively, we are associating to every case variable an array whose locations are indexed by the process instances from PI_{index} and whose components contain a value from the

interpretation of the sort S in Cat . For every process instance $i \in PI_{index}^{\mathcal{M}}$, the tuple $(i, a_{v_1}^\alpha[i], \dots, a_{v_n}^\alpha[i])$ denotes the set of case variables for the process instance i . All the “the *case variables* artifact relation” are usually initialized with undefined values.

Then, we associate to every relation $R := \langle R_1, \dots, R_m \rangle$ in $Repo$ an artifact sort symbol R_{index} different from PI_{index} and m artifact components a_{R_1}, \dots, a_{R_m} with R_{index} as their common source and respectively the sorts of R_1, \dots, R_m as their target: sometimes, we denote the tuple $(a_{R_1}, \dots, a_{R_m})$ by writing a_R . Intuitively, given a tuple $(z_1^{\mathcal{M}}, \dots, z_m^{\mathcal{M}}) \in R^{\mathcal{M}}$, we associate to it an element $e \in R_{index}^{\mathcal{M}}$ and the tuple $(e, a_{R_1}^\alpha[e], \dots, a_{R_m}^\alpha[e])$, where $a_{R_1}^\alpha[e] = z_1^{\mathcal{M}}, \dots, a_{R_m}^\alpha[e] = z_m^{\mathcal{M}}$: the element e is the “internal identifier” of that tuple in $R^{\mathcal{M}}$. Artifact relations that have a sort E different from PI_{index} as their artifact sort are called “shared (or repository) artifact relations”. All the “shared artifact relations” are usually initialized with undefined values.

F.2 Translation of Update Logic and Process Schema

We inductively translate into the array-based setting the blocks from our data-aware BPMN models that are used to construct workflows. We associate to the lifecycle status of every block B a *case variable artifact component* “*lifecyleStateB*” with PI_{index} as their artifact sort. These function variables constitute the *control variables* of the translated workflows. For every DAB \mathcal{S} , we define $\mathcal{B} := \{lifecyleStateB \mid B \in Blocks\}$, where $Blocks$ is the set of all the blocks that form the Process Component of \mathcal{S} .

The first block that should be translated is the block Task. Since it involves preconditions and update, we preliminarily discuss how to translate them into the array-based setting.

Given a guard $q(\vec{x}) \leftarrow Q$ over \mathcal{D} as defined in Section 2.1, for the purpose of our translation we can assume that Q , instead of a disjunction of conjunctive queries with filters, only consists of a *conjunction of atoms or negated atoms* (i.e., *cubes*), where every atom is like in Section 2.1. In fact, we can first put Q in disjunctive normal form: so, Q is equivalent to a disjunction of cubes. Then, it can be easily seen that, since existential quantifiers commute with disjunctions, it is always possible to preprocess a precondition of a task that is a disjunction of cubes so as to split that task into new tasks with preconditions that are cubes (every disjunct is the precondition of one of the new tasks). The resulting DAB is equivalent to the original one.

We say that an *extended guard* is a cube Q whose variables v from $\mathcal{D}.cvars \cup \mathcal{B}$ are substituted with terms $v[I]$, where v are function variables that keep the same name of v and I is a process instance in \mathcal{PI} . Analogously, we define *extended update for a task* a formula of type “update” as defined in Definition 2.2 whose variables v from $\mathcal{D}.cvars \cup \mathcal{B}$ are substituted with terms $v[I]$, where v are function variables that keep the same name of v and I is a process instance in \mathcal{PI} . In general, given a formulae ϕ , we call $\phi_{ext(I)}$ the formula obtained from ϕ by substituting every variable v from $\mathcal{D}.cvars \cup \mathcal{B}$ in ϕ with $v[I]$, where $I \in \mathcal{PI}$.

Given an extended formulae, the translation of the query language works as follows: a variable $v \in \mathcal{D}.cvars \cup \mathcal{B}$ is associated to a *case variable artifact component* v (that keeps the same name) in the array-based setting, and every term $v[I]$ (with $I \in \mathcal{PI}$) of v is associated to the corresponding function application (read-operation) $v[i]$ (with $i \in PI_{index}$) in the array-based setting.

We give the formal translation of the semantics of data-aware BPMN models by introducing *rule-based transitions* that fit the format of transition formulae in array-based RASs. In fact, for the purpose of this paper, we can simplify the form of transition formulae from RASs by focusing on rule-based formulae of the following form:

```

rule Transition =
    if Guard
    then Update

```

when “Guard” is an *extended guard* and “Update” is an *extended update for a task*.

Formally, it can be easily seen that the previous rule-based formulae can be translated into transition formulae of RASs as follows. First of all, given an extended guard G and an extended update U , we rewrite G and U into the array-based setting as follows (when we say “add to G (or to U) the sub-formula ϕ ”, it means that we conjunct the sub-formula ϕ with G (or U)):

1. If U is the constructor **INSERT \vec{u} INTO R AND SET $\vec{x} = \vec{v}$** for $R \in Repo$, and the required semantics is the set-theoretic one, substitute it in U with the sub-formula:

$$\underline{a}'_R := \forall j (\text{if } j = e_{ins} \text{ then } \vec{u} \text{ else } (\text{if } \underline{a}_R[j] = \vec{u} \text{ then } \mathbf{undef} \text{ else } \underline{a}_R[j])) \wedge \vec{x}' := \vec{v}$$
 where $\underline{a}_R := (a_{R_1}, \dots, a_{R_n})$ are the artifact components of the shared artifact relation R_{index} and e_{ins} is in R_{index} , and add to G the sub-formula $(a_{R_1}[e_{ins}] = \mathbf{undef} \wedge \dots \wedge a_{R_n}[e_{ins}] = \mathbf{undef})$;
2. If U is the constructor **INSERT \vec{u} INTO R AND SET $\vec{x} = \vec{v}$** for $R \in Repo$, and the required semantics is the multiset-theoretic one, substitute it in U with the sub-formula $\underline{a}'_R[e_{ins}] := \vec{u} \wedge \vec{x}' := \vec{v}$, where $\underline{a}_R := (a_{R_1}, \dots, a_{R_n})$ are the artifact components of the shared artifact relation R_{index} and e_{ins} is in R_{index} , and add to G the sub-formula $(a_{R_1}[e_{ins}] = \mathbf{undef} \wedge \dots \wedge a_{R_n}[e_{ins}] = \mathbf{undef})$;
3. If U is the constructor **DEL \vec{u} FROM R AND SET $\vec{x} = \vec{v}$** for $R \in Repo$, substitute it in U with the sub-formula $\underline{a}'_R[e_{del}] := \mathbf{undef} \wedge \vec{x}' := \vec{v}$, where $\underline{a}_R := (a_{R_1}, \dots, a_{R_n})$ are artifact components of the shared artifact relation R_{index} and e_{del} is in R_{index} , and add to G the subformula $\underline{a}_R[e_{del}] = \vec{u}$.
4. If U is the constructor **UPDATE $R(\vec{v})$ IF $\psi(\vec{u}, \vec{v})$ THEN $R(\vec{u}')$ ELSE $R(\vec{u}'')$** for $R \in Rep$, substitute it in U with $\underline{a}'_R := \forall j (\text{if } \psi(\vec{u}, \underline{a}_R[j]) \text{ then } \vec{u}_1 \text{ else } \vec{u}_2)$, where $\underline{a}_R = (a_{R_1}, \dots, a_{R_n})$ are artifact components of the shared artifact relation R_{index} , j is in R_{index} and u_1, u_2 are u', u'' where every occurrence of variable v_k from \vec{v} has been substitute with $a_{R_k}[j]$.

5. We substitute every term of the kind $v[I]$ in G or in U , where $v \in \mathcal{D}.cvars$, with the term $a_v[i]$, where a_v is the case variables artifact component associated to v and $i \in PI_{index}$;
6. We substitute every atom $R(t_1, \dots, t_n)$ in G , where $R \in Repo$, with the sub-formula $(a_{R_1}[e] = t_1 \wedge \dots \wedge a_{R_n}[e] = t_n)$, where a_{R_1}, \dots, a_{R_n} are artifact components of the shared artifact relation R_{index} and $e \in R_{index}$ is a fresh variable.
7. We substitute every atom $R(t_0, \dots, t_n)$ in G and U , where $R \in Cat$, with the sub-formula $(t_0 \neq \mathbf{undef} \wedge t_1 = f_1(t_0) \wedge \dots \wedge t_n = f_n(t_0))$, where each f_k ($k = 1, \dots, n$) is the unary function f_{R, A_k} associated to R and its k -th attribute A_k (here, we employ Formula (6)).

In Step (4) above, we translated only the *flat* “if-then-else” constructor, since the “nested” one has an analogous translation: in fact, after the keywords “then” or “else”, instead of a term there will be the iterated translation of the same “if-then-else” constructor.

After this rewriting phase of G and U , we obtain the following formula:

$$\exists \underline{e}_{ins} \exists \underline{e}_{del} \exists \underline{e} \exists i \exists \underline{y} \left(G(v_1[i], \dots, v_m[i], \underline{y}_1, \underline{e}_{ins}, \underline{e}_{del}, \underline{e}) \wedge U(v_1[i], \dots, v_m[i], \underline{y}, \underline{e}_{ins}, \underline{e}_{del}, \underline{e}) \right) \quad (8)$$

where $\underline{e}_{ins}, \underline{e}_{del}, \underline{e}$ contains all the variables of artifact sorts that have been called e_{ins}, e_{del}, e respectively during the rewriting phase, v_k are case variables artifact components associated to case variables in $\mathcal{D}.cvars$ and i is the only variable of artifact sort PI_{index} introduced during the rewriting phase: notice that i is different from all the $\underline{e}_{ins}, \underline{e}_{del}, \underline{e}$ variables. Notice that we add the existential quantifier $\exists i$ in front of the transition formula and, if some variable $y \in V \setminus V_C$ occurs free in G or in U , we also add (avoiding redundancies) the existential quantifiers $\exists \underline{y}$ in front of the transition formula. Then, we eliminate the quantifiers of the form $\exists y$ that bind variables of type y that allow at least one definition like $y := a_R[e]$.

Formula (8) fits the format of Formula (2). Thus, from now on it is sufficient to show that, in a rule-based transition, the formulae corresponding to “Guard” and “Update” are respectively an extended guard and an extended update. As already noticed, it is straightforward to see that preconditions and updates as presented in Section 2.1 can be transformed into extended guards and extended updates respectively.

Every block B has the control variable $lifecycleStateB[i]$ for every process instance i that can take at least three distinct values: “idle”, “enabled”, “completed”. Blocks that have a boundary event can also have $lifecycleStateB[i] := \text{“error”}$. We now provide the formal translation of every blocks, by exploiting rule-base transitions.⁸

⁸ In the following, the constant “error” refers to finitely many different labels of type error, and each of them is linked to an exception handler: hence, every error handler is labeled using its unique error label and every constant “error” refers uniquely to it.

Base case: an atomic task T . In case of an atomic task, the variable $lifecycleStateT[i]$ (for every process instance i) takes three distinct values: “idle”, “enabled”, “completed”. An atomic task T is made up of two transitions:

- when T is “enabled” for a process instance i , preconditions over data are evaluated and, in case they are true, T can non-deterministically update the working memory (data + control variables) and become “completed” for i .

We can express formally the lifecycle of an atomic task T as follows:

```

rule  $T_1$  =
  if    $lifecycleStateT[I] = \text{enabled}$ 
        PRECONDITION ON DATA
  then   $lifecycleStateT'[I] = \text{completed}$ 
        UPDATES OVER THE WORKING MEMORY

```

Since “PRECONDITION ON DATA” and “UPDATES OVER THE WORKING MEMORY” are the corresponding extended versions of a guard and an update of a data-aware BPMN model respectively, and since they are conjuncted with formulae of the kind $lifecycleStateT[I] = \text{constant}$, the previous transitions fit the format of (8).

In the following paragraphs, it can be easily seen that all the transitions fit the format of (8): specifically, the formulae ϕ, ϕ_1, ϕ_2 that appear in some Guards, are *conditions* in the sense of Section 2.1, hence the claim is true: all the other cases are straightforward.

Remark 2. A task T could be also formalized in a non-atomic way: in this case, the variable $lifecycleStateT[i]$ for every process instance i takes four distinct values: “idle”, “enabled”, “active”, “completed”. An atomic task T is made up of two transitions:

- when T is “enabled” for a process instance i , preconditions over data are evaluated and, in case they are true, T can non-deterministically become “active” for i ;
- when T is “active” for a process instance i , T becomes “completed” for i and the updates over the working memory (data + control variables) are performed. ◁

Formally, we have:

```

rule  $T_1$  =
  if    $lifecycleStateT[I] = \text{enabled}$ 
        PRECONDITION
  then   $lifecycleStateT'[I] = \text{active}$ 

rule  $T_2$  =
  if    $lifecycleStateT[I] = \text{active}$ 
  then   $lifecycleStateT'[I] = \text{completed}$ 
        UPDATES OVER THE WORKING MEMORY

```

Base case: an event E . For every event E , the variable $lifecycleStateT[i]$ for every process instance i takes three distinct values: “idle”, “enabled”, “completed”. An event E is made up of one transition: when T is “enabled” for a process instance i , T can non-deterministically become “completed” for i and the updates over the working memory (data + control variables) are performed.

We can express formally the lifecycle of an event E as follows:

```

rule  $T_1 =$ 
  if  $lifecycleStateE[I] = \text{enabled}$ 
  then  $lifecycleStateE'[I] = \text{completed}$ 
        UPDATES OVER THE WORKING MEMORY

```

By using an argument similar to the previous one, we conclude that the transitions translating the behavior of an event E fit the format of (8).

In the following, we give the translation of the DABs blocks. Whenever a block B has some sub-components B_i (that are still blocks), we assume that they are defined by inductive hypothesis. For blocks, we use the label “waiting” for denoting that it is “active”.

SEQF: sequence flow.

```

rule  $T_1 =$ 
  if  $lifecycleStateB[I] = \text{enabled}$ 
  then  $lifecycleStateB'_1[I] = \text{enabled}$ 
         $lifecycleStateB'[I] = \text{waiting}$ 

```

```

rule  $T_2 =$ 
  if  $lifecycleStateB_1[I] = \text{completed}$ 
  then  $lifecycleStateB'_1[I] = \text{idle}$ 
         $lifecycleStateB_2[I] = \text{enabled}$ 

```

```

rule  $T_3 =$ 
  if  $lifecycleStateB_2[I] = \text{completed}$ 
  then  $lifecycleStateB'_2[I] = \text{idle}$ 
         $lifecycleStateB'[I] = \text{completed}$ 

```

PAR: parallel block B .

```

rule  $T_1 =$ 
  if  $lifecycleStateB[I] = \text{enabled}$ 
  then  $lifecycleStateB'_1[I] = \text{enabled}$ 
         $lifecycleStateB'_2[I] = \text{enabled}$ 
         $lifecycleStateB'[I] = \text{waiting}$ 

```

```

rule  $T_2 =$ 
  if  $lifecycleStateB_1[I] = \text{completed}$ 
         $lifecycleStateB_2[I] = \text{completed}$ 
  then  $lifecycleStateB'_1[I] = \text{idle}$ 

```

$$\begin{aligned} \text{lifecycleStateB}'_2[I] &= \text{idle} \\ \text{lifecycleStateB}'[I] &= \text{completed} \end{aligned}$$

OR: conditional inclusive block B .

rule T_1 =
if $\text{lifecycleStateB}[I] = \text{enabled}$
 $\phi_1 \wedge (\neg\phi_2)$
then $\text{lifecycleStateB}'_1[I] = \text{enabled}$
 $\text{lifecycleStateB}'[I] = \text{waiting1}$

rule T_2 =
if $\text{lifecycleStateB}[I] = \text{enabled}$
 $\phi_2 \wedge (\neg\phi_1)$
then $\text{lifecycleStateB}'_2[I] = \text{enabled}$
 $\text{lifecycleStateB}'[I] = \text{waiting1}$

rule T_3 =
if $\text{lifecycleStateB}[I] = \text{enabled}$
 $\phi_1 \wedge \phi_2$
then $\text{lifecycleStateB}'_1[I] = \text{enabled}$
 $\text{lifecycleStateB}'_2[I] = \text{enabled}$
 $\text{lifecycleStateB}'[I] = \text{waiting2}$

rule T_4 =
if $\text{lifecycleStateB}_1[I] = \text{completed}$
 $\text{lifecycleStateB}[I] = \text{waiting1}$
then $\text{lifecycleStateB}'_1[I] = \text{idle}$
 $\text{lifecycleStateB}'[I] = \text{completed}$

rule T_5 =
if $\text{lifecycleStateB}[I] = \text{waiting1}$
 $\text{lifecycleStateB}_2[I] = \text{completed}$
then $\text{lifecycleStateB}'_2[I] = \text{idle}$
 $\text{lifecycleStateB}'[I] = \text{completed}$

rule T_6 =
if $\text{lifecycleStateB}_1[I] = \text{completed}$
 $\text{lifecycleStateB}_2[I] = \text{completed}$
then $\text{lifecycleStateB}'_1[I] = \text{idle}$
 $\text{lifecycleStateB}'_2[I] = \text{idle}$
 $\text{lifecycleStateB}'[I] = \text{completed}$

CHOICE: conditional exclusive block with choice B .

rule $T_1 =$
if $lifecycleStateB[I] = \text{enabled}$
 ϕ
then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
if $lifecycleStateB[I] = \text{enabled}$
 $\neg\phi$
then $lifecycleStateB'_2[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_3 =$
if $lifecycleStateB_1[I] = \text{completed}$
then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule $T_4 =$
if $lifecycleStateB_2[I] = \text{completed}$
then $lifecycleStateB'_2[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

DEF-CHOICE: conditional exclusive block with deferred choice B .

rule $T_1 =$
if $lifecycleStateB[I] = \text{enabled}$
then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
if $lifecycleStateB[I] = \text{enabled}$
then $lifecycleStateB'_2[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_3 =$
if $lifecycleStateB_1[I] = \text{completed}$
then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule $T_4 =$
if $lifecycleStateB_2[I] = \text{completed}$
then $lifecycleStateB'_2[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

LOOP: a loop block B .

rule $T_1 =$
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
 if $lifecycleStateB_1[I] = \text{completed} \wedge \phi$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB_2[I] = \text{enabled}$

rule $T_3 =$
 if $lifecycleStateB_2[I] = \text{completed}$
 then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB_2[I] = \text{idle}$

rule $T_4 =$
 if $lifecycleStateB_1[I] = \text{completed} \wedge (\neg\phi)$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

PROC: a process block B .

rule $T_1 =$
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
 if $lifecycleStateB_1[I] = \text{completed}$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

BLOCK-ERR: a block with error B .

rule $T_1 =$
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
 if $lifecycleStateB_1[I] = \text{completed} \wedge \phi$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule $T_3 =$
 if $lifecycleStateB_1[I] = \text{completed} \wedge (\neg\phi)$

then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateH'[I] = \text{idle}$
 $lifecycleStateB'_{\text{bound-hand}}[I] = \text{error}$

where the H s are all the sub-blocks of the block $B_{\text{bound-hand}}$ whose boundary is directly connected to the handler block for “error”.

Alternatively, T_3 could also be

rule $T_{3,alt} =$
 if $lifecycleStateB_1[I] = \text{completed} \wedge (\neg\phi)$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateH'[I] = \text{idle}$
 $errorB'_{\text{bound-hand}}[I] = \text{true}$

where $errorB_{\text{bound-hand}}$ is a boolean variable linked to $B_{\text{bound-hand}}$.

BACK-EXCP: a backward exception handling block B .

Here, A is a subprocess.

rule $T_1 =$
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateA'[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule $T_2 =$
 if $lifecycleStateA[I] = \text{completed}$
 then $lifecycleStateA'[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule $T_3 =$
 if $lifecycleStateA[I] = \text{error}$
 then $lifecycleStateA'[I] = \text{idle}$
 $lifecycleStateB'_1[I] = \text{enabled}$

rule $T_{err1} =$
 if $lifecycleStateA[I] = \text{enabled}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

rule $T_{err2} =$
 if $lifecycleStateA[I] = \text{waiting}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

rule $T_{err3} =$
 if $lifecycleStateA[I] = \text{active}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

where the H s are all the sub-blocks of the block A .

FOR-EXCP: a forward exception handling block B .

Here, A is a subprocess.

rule T_1 =
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateA'[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule T_2 =
 if $lifecycleStateA[I] = \text{completed}$
 then $lifecycleStateA'[I] = \text{idle}$
 $lifecycleStateB'_1[I] = \text{enabled}$

rule T_3 =
 if $lifecycleStateB_1[I] = \text{completed}$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule T_4 =
 if $lifecycleStateA[I] = \text{error}$
 then $lifecycleStateA'[I] = \text{idle}$
 $lifecycleStateB'_2[I] = \text{enabled}$

rule T_5 =
 if $lifecycleStateB_2[I] = \text{completed}$
 then $lifecycleStateB'_2[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

rule T_{err1} =
 if $lifecycleStateA[I] = \text{enabled}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

rule T_{err2} =
 if $lifecycleStateA[I] = \text{waiting}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

rule T_{err3} =
 if $lifecycleStateA[I] = \text{active}$
 then $lifecycleStateA'[I] = \text{error}$
 $lifecycleStateH'[I] = \text{idle}$

where the H s are all the sub-blocks of the block A .

NON-INTERR: a non-interrupting exception handling block B .

Here, we also need a boolean variable $errorA$ that changes its value when the error occurs. A is a subprocess.

```

rule  $T_1 =$ 
  if  $lifecycleStateB[I] = \text{enabled}$ 
  then    $lifecycleStateA'[I] = \text{enabled}$ 
           $lifecycleStateB'[I] = \text{waiting}$ 

rule  $T_2 =$ 
  if  $lifecycleStateA[I] = \text{completed}$ 
  then    $lifecycleStateA'[I] = \text{idle}$ 
           $lifecycleStateB'_1[I] = \text{enabled}$ 

rule  $T_3 =$ 
  if  $lifecycleStateB_1[I] = \text{completed}$ 
       $errorA = \text{false}$ 
  then    $lifecycleStateB'_1[I] = \text{idle}$ 
           $lifecycleStateB'[I] = \text{completed}$ 

rule  $T_4 =$ 
  if  $errorA = \text{true}$ 
  then    $lifecycleStateB'_2[I] = \text{enabled}$ 
           $errorA' = \text{true}$ 

rule  $T_5 =$ 
  if    $lifecycleStateB_1[I] = \text{completed}$ 
         $lifecycleStateB_2[I] = \text{completed}$ 
  then    $lifecycleStateB'_1[I] = \text{idle}$ 
           $lifecycleStateB'_2[I] = \text{idle}$ 
           $lifecycleStateB'[I] = \text{completed}$ 

rule  $T_{err1} =$ 
  if  $lifecycleStateA[I] = \text{enabled}$ 
  then    $lifecycleStateA'[I] = \text{error}$ 
           $lifecycleStateH'[I] = \text{idle}$ 

rule  $T_{err2} =$ 
  if  $lifecycleStateA[I] = \text{waiting}$ 
  then    $lifecycleStateA'[I] = \text{error}$ 
           $lifecycleStateH'[I] = \text{idle}$ 

rule  $T_{err3} =$ 
  if  $lifecycleStateA[I] = \text{active}$ 
  then    $lifecycleStateA'[I] = \text{error}$ 
           $lifecycleStateH'[I] = \text{idle}$ 

```

where the H s are all the sub-blocks of the block A .

F.3 Alternative translation of specific blocks

For sake of simplicity (and in order to make the translation more efficient in performance), sometimes we can freely employ the following alternative translation for specific blocks that are useful in practice.

***n*-SEQ: *n*-iterated sequence flow block *B*.**

In case a sequence flow block *B* is formed of $n > 2$ sub-blocks B_k , we can make use of the following rule-based transitions:

rule T_1 =
 if $lifecycleStateB[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{enabled}$
 $lifecycleStateB'[I] = \text{waiting}$

rule T_2 =
 if $lifecycleStateB_1[I] = \text{completed}$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateB'_2[I] = \text{enabled}$

rule T_3 =
 if $lifecycleStateB_k[I] = \text{completed}$
 then $lifecycleStateB_k[I] = \text{idle}$
 $lifecycleStateB'_{k+1}[I] = \text{enabled}$

where $2 \leq k < n$

rule T_4 =
 if $lifecycleStateB_n[I] = \text{completed}$
 then $lifecycleStateB'_n[I] = \text{idle}$
 $lifecycleStateB'[I] = \text{completed}$

ERR&EVENT: “either error or event” block.

In case of a event-based fork block, where the first branch has an event B_1 and the second one is an error event B_2 , we have the following translation:

rule T_1 =
 if $lifecycleStateB_1[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{completed}$
 UPDATES OVER THE WORKING MEMORY

rule T_2 =
 if $lifecycleStateB_1[I] = \text{enabled}$
 then $lifecycleStateB'_1[I] = \text{idle}$
 $lifecycleStateH'[I] = \text{idle}$
 $lifecycleStateB'_{\text{bound-hand}}[I] = \text{error}$
 UPDATES OVER THE WORKIN MEMORY PERFORMED BY B_2

where the *H*s are all the sub-blocks of the block $B_{\text{bound-hand}}$ whose boundary is directly connected to the handler block for “error”.

F.4 Translation of Reachability Queries

The translation of a reachability query is totally analogous to the translation presented in Subsection F.2 for guards and updates, with the proviso that case variables associated to different process instances must be translated into case variables artifact component applied to different indexes from PI_{index} . More formally, given a reachability query $Q := \bigwedge_{i \in I} G_i[i]$, where $G_i := \bigvee_k G_k$ is a guard, we associate to every $i \in I$ and index $e_i \in PI_{index}$, we substitute every case variable v in $G_i[i]$ with the term $v[i]$ (read-operation of a function variable x) and we apply the same rewriting policy of Subsection F.2 to every its conjunct G_k : thus, we obtain a quantifier-free formula of the kind:

$$\exists \underline{e} \phi(\underline{e}, \underline{x}, \underline{a})$$

, where ϕ is quantifier-free and the \underline{e} are individual variables of artifact sorts (where each e is either in PI_{index} or in some “repository” artifact relation), i.e. we get a state formula of RASs, as required.

G Soundness and Completeness Results

In this section we sketch the proof of the soundness and completeness results.

Remark 3. Notice that, in our translation to array-based systems, case variables of case-bounded DABs can be treated as proper *artifact variables* of RASs, instead of arrays. This is trivial since we do not need (undounded) indexes (taken from some artifact sort) in case of 1-case DABs, and, analogously, in case of k -bounded DABs, it is sufficient to associate every case variable to k corresponding artifact variables. \triangleleft

Theorem 1 and Theorem 2 clearly follow from Theorem 6 and from the translation into RASs described in the previous section.

H Termination Results

Notice that clearly Cat is acyclic iff its corresponding DB schema in RASs is acyclic. Hence, a DAB is acyclic iff its translated RAS has an acyclic DB schema.

Thanks to Remark 3, Theorem 3 simply follows from Theorem 7 and from our translation in RASs.

It can be easily seen, by exploiting our translation into RASs, that the conditions posed over updates in Theorem 4 and Theorem 5 exactly correspond to suitable requirements over the translated updates in RASs that guarantee strong locality. Hence, it is possible to apply Theorem 8 in order to get termination: in fact, it is clear that the translation of a separated reachability query is a strongly local formula.

We devote the following subsection to sketch the proof of strong locality of the transitions translating the restricted updates of Theorem 4.

H.1 Translated Updates that are strongly local

We give a sketch of the proof of strong locality of transitions that are the translations into the array-based setting of “Insert&set rule” and ”Delete&set rule” with the restriction of Section 3.5: the other updates presented in Section 3.5 can be proved in a similar way to have a strongly local translation. All the proofs sketched in this section concerning the fact that the format of those transitions fits the definition of strong locality are similar to the ones in Appendix F of [4], where all the details are deeply analyzed and *all* the restricted updates corresponding to the ones presented in Section 3.5 are proved to be strongly local transitions. Specifically, by adopting the conventions of [4], we notice that the proofs of strong locality of the translations of “Insert&set rule”, “Set rule”, “Delete&set rule” and “Conditional update rule” correspond and are analogous (with trivial adaptations) to the proofs of strong locality (provided in Appendix F in [4]) of “Insertion Updates”, “Propagation Updates”, “Deletion Updates” and “Bulk Updates”. In the following, when we say that a translated array-based formula is “repository-free”, “over the Catalogue” or etc., we mean, by abuse of notation, that this formula is the translation of a corresponding “repository-free”, “over the Catalogue” or etc. query of a DAB.

Delete&set rule. We want to remove a tuple $\underline{t} := (t_1, \dots, t_m)$ from an m -ary relation R of the Repository and assign the values t_1, \dots, t_m to some of the case variables (let $\underline{x} := \underline{x}_1, \underline{x}_2$, where $\underline{x}_1 := (x_{i_1}, \dots, x_{i_m})$ are the variables where we want to transfer the tuple \underline{t}). This operation has to be applied only if the current case variables \underline{x} satisfy the repository-free pre-condition $\pi(\underline{x}_1, \underline{x}_2)$ and additional variables $\underline{y} := \underline{y}_1, \underline{y}_2$ (where the \underline{y}_1 are elements from the tuple \underline{t}) satisfy the post-condition $\psi(\underline{y}_1, \underline{y}_2)$ over the Catalogue. The variables \underline{x}_2 are not propagated, i.e. they are reassigned (possibly with the same values that they have before). Let $\underline{r} := r_1, \dots, r_m$ be the artifact components of R in the translated array-based setting. Such an update can be formalized in the translated array-based formalism as follows:

$$\exists \underline{y}_1 \exists \underline{y}_2 \exists i \exists e \exists \underline{i} \left(\begin{array}{l} \pi(\underline{x}_1[i], \underline{x}_2[i]) \wedge \underline{r}[e] = \underline{y}_1 \wedge \psi(\underline{y}_1, \underline{y}_2) \wedge r_1[e] \neq \text{undef} \wedge \dots \\ \wedge r_n[e] \neq \text{undef} \wedge (\underline{x}'_1[i] := \underline{r}[e] \wedge \underline{x}'_2[i] := \underline{y}_2 \wedge \underline{s}' := \underline{s} \wedge \\ \wedge \underline{r}' := \lambda j. (\text{if } j = e \text{ then undef else } \underline{r}[j])) \end{array} \right) \quad (9)$$

where \underline{s} are the artifact components of the relations from the Repository different from R , and π and ψ are free-repository conjunctive queries. Notice that the $\underline{y}_1, \underline{y}_2$ are non deterministically produced values for the updated \underline{x}'_2 . In the terminology of [20], notice that no case variable variable is “propagated” in a deletion update.

We sketch the proof of the fact that the preimage along (9) of a strongly local formula is strongly local. Consider a strongly local formula

$$K := \psi'(\underline{x}[i]) \wedge \exists \underline{e} \left(\text{Diff}(\underline{e}) \wedge \bigwedge_{e_r \in \underline{e}} \phi_{e_r}(\underline{r}[e_r]) \wedge \Theta \right)$$

where Θ is a formula involving the artifact components \underline{s} (which are not updated) such that no e_r occurs in it.

Computing the preimage $Pre(9, K)$, we get (with a computation analogous to the one done in Appendix F in [4]) the disjunction of the formulae:

$$\begin{aligned} & - \exists e, \underline{e} (\text{Diff}(\underline{e}, e) \wedge \pi(\underline{x}_1[i], \underline{x}_2[i]) \wedge \theta(\underline{r}[e]) \wedge \bigwedge_{e_r \in \underline{e}} \phi_{e_r}(\underline{r}[e_r]) \wedge \Theta) \\ & - \exists \underline{d} \exists \underline{e} \left(\text{Diff}(\underline{e}) \wedge \pi(\underline{x}_1[i], \underline{x}_2[i]) \wedge \theta(\underline{r}[e_j]) \wedge \right. \\ & \quad \left. \wedge \bigwedge_{e_r \in \underline{e}, e_r \neq e_j} \phi_{e_r}(\underline{r}[e_r]) \wedge \phi_{e_j}(\mathbf{undef}) \wedge \Theta \right) \end{aligned}$$

which is strongly local, where θ is a quantifier-free Σ -formula (Σ is the DB signature of the read-only DB that is translation of the Catalogue).

Insert&set rule. We want to insert a tuple of values $\underline{t} := (t_1, \dots, t_m)$ from the case variables $\underline{x}_1 := (x_{i_1}, \dots, x_{i_m})$ (let $\underline{x} := \underline{x}_1, \underline{x}_2$ as above) into an m -ary relation R of the Repository. This operation has to be applied only if the current case variables \underline{x} and additional variables $\underline{y} := y_1, y_2$ satisfy the repository-free pre-condition $\pi(\underline{x}_1, \underline{x}_2, \underline{y})$. Let $\underline{r} := r_1, \dots, r_m$ be the artifact components of R . Such an update can be formalized in the translated array-based formalism as follows:

$$\exists \underline{d}_1, \underline{d}_2 \exists e \left(\begin{array}{l} \pi(\underline{x}_1[i], \underline{x}_2[i], \underline{y}) \wedge \underline{r}[e] = \mathbf{undef} \\ \wedge (\underline{x}'[i] := \underline{y} \wedge \underline{s}' := \underline{s} \wedge \\ \wedge \underline{r}' := \lambda j. (\mathbf{if } j = e \mathbf{ then } \underline{x}_1[i] \mathbf{ else } \underline{r}[j])) \end{array} \right) \quad (10)$$

where \underline{s} are the artifact components of the relations from the Repository different from R . Notice that \underline{y} are used to produce values for the updated case variables \underline{x}' . In the terminology of [20], notice that no artifact variable is propagated in a insertion update. Notice that it is allowed that some case variables are propagated (i.e., that some, or all, $x'_k := x_k$)

Notice also that the following arguments remain the same even if $\underline{r}[e] = \mathbf{undef}$ is replaced with a conjunction of *some* literals of the form $r_j[e] = \mathbf{undef}$, for some $j = 1, \dots, m$, or even if $\underline{r}[e] = \mathbf{undef}$ is replaced with a generic constraint $\chi(\underline{r}[e])$.

We sketch the proof of the fact that the preimage along (10) of a strongly local formula is strongly local. Consider a strongly local formula

$$K := \psi'(\underline{x}[i]) \wedge \exists \underline{e} \left(\text{Diff}(\underline{e}) \wedge \bigwedge_{e_r \in \underline{e}} \phi_{e_r}(\underline{r}[e_r]) \wedge \Theta \right)$$

where Θ is a formula involving the translation of the relations \underline{s} (which are not updated) from the Repository such that no e_r occurs in it.

Computing the preimage $Pre(10, K)$, we get (with a computation analogous to the one done in Appendix F in [4]) the disjunction of the formulae:

$$\begin{aligned} & - \exists e, \underline{e} (\text{Diff}(\underline{e}, e) \wedge \theta(\underline{x}_1[i], \underline{x}_2[i]) \wedge \underline{r}[e] = \mathbf{undef} \wedge \bigwedge_{e_r \in \underline{e}} \phi_{e_r}(\underline{r}[e_r]) \wedge \Theta) \\ & - \exists \underline{e} \left(\text{Diff}(\underline{e}) \wedge \theta(\underline{x}_1[i], \underline{x}_2[i]) \wedge \phi_{e_j}(\underline{x}_1[i]) \wedge \underline{r}[e] = \mathbf{undef} \wedge \bigwedge_{e_r \in \underline{e}, e_r \neq e_j} \phi_{e_r}(\underline{r}[e_r]) \wedge \Theta \right) \end{aligned}$$

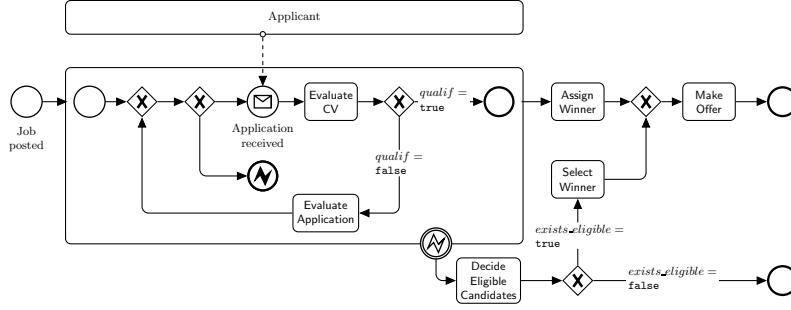


Fig. 2. Job hiring process

which is a strongly local formula, where θ is a quantifier-free Σ -formula (Σ is the DB signature of the read-only DB that is translation of the Catalogue).

We remark that, in a “Insert&set” update, the insertion of the same content in correspondence to different entries is allowed. *If we want to avoid this kind of multiple insertions*, the update r' must be modified as follows:

$$r' := \lambda j. \left(\begin{array}{l} \text{if } j = e \text{ then } \underline{x}_1 \text{ else} \\ \text{(if } r[j] = \underline{x}_1 \text{ then undef else } r[j]) \end{array} \right)$$

which is not strongly local.

Example 8. We consider a very slight variant of the example of a job hiring process in a company presented in the paper. The human resource (HR) branch of the company stores in its catalog database information relevant to the process. Specifically, the company’s catalog Cat is composed of the following relations:

- $JobCategory(Jcid : \text{jobcatID})$ is used to access different types of jobs that are available in the company ;
- $User(Uid : \text{userID}, Name : \text{StringName}, Age : \text{NumAge})$ stores data about users registered to the company website, who might be potentially interested in job positions offered by the company. \triangleleft

To manage information about submitted applications, including data on users, the score they receive after having been interviewed and their eligibility, the company employs repository Rep that consists of one relation $Application(Jcid : \text{jobcatID}, Uid : \text{userID}, Name : \text{StringName}, Age : \text{NumAge}, Score : \text{NumScore}, Eligible : \text{BoolString})$ that stores all the information. Notice that NumScore contains 100 values in the range $(0, 101)$, where a score from 1 to 100 indicates the actual one assigned after evaluating the application. For readability, we use the usual predicates $<$, $>$, and $=$ to compare variables of type NumScore : this is syntactic sugar and does not require to introduce rigid predicates in our framework.

Since the job posting is created using a dedicated company’s portal, the information related to this posting does not have to be stored persistently and thus can be maintained just for a given case. To represent it we use a set of case variables V_C , where $jid : \text{jobcatID}$ references a job type, $uid : \text{userID}$ user’s identifier together with her name $name : \text{StringName}$ and age $age : \text{NumAge}$, a user that wins the position $winner : \text{userID}$ together with the check of

her eligibility *result* : BoolString. Moreover, we use three more auxiliary variables: *qualif* : Bool identifies whether a currently selected applicant is qualified or not, *exists_eligible* : Bool indicates successful (or not) termination of the eligible candidate selection process and *tPassed* : StringDate indicates the time passed from the moment when company started receiving applications for the open position. In the following we described data updates issued by tasks and events in our process model. The execution starts by receiving a new job posting that generates a new posting identifier using the following effect: $Eff(\text{Job posted}) = \{\text{SET } jcid' = id_j\}$. As soon the job offer has been published, the company starts a process of receiving and evaluating applications. Such process runs until a qualified candidate is found: nevertheless, if no qualified candidate is found after a non-deterministically assigned deadline has been reached, the process is interrupted anyway. In our case the deadline is modeled using an error event with effect $Eff(\text{ErrorEvent}) = \{\text{SET } tPassed' = \mathbf{1Month}\}$.

Whenever a new application is received, the *Application received* event gets triggered and assigns user data that came together with the application to designated case variables using the following effect:

$$Eff(\text{Application received}) = \left\{ \begin{array}{l} \text{SET } uid' = id_u, \\ \text{SET } name' = n, \\ \text{SET } age' = a \end{array} \right\}$$

Next, a CV attached to the application undergoes a preliminary evaluation with the sole purpose of detecting a candidate that may be a perfect fit for the position, and thus should immediately win the competition. This is modeled by the Evaluate CV task, s.t. $G(\text{Evaluate CV}) = \mathbf{true}$ ⁹ and $Eff(\text{Evaluate CV}) = \{\text{SET } qualif' := o\}$.

If a candidate is not perfectly apt for the position in absentia, we proceed directly to the thorough evaluation of the application followed by the process immediately recording the interview result using the Evaluate Application task. This task requires a precondition on data updates, i.e.

$$G(\text{Evaluate Application}) = \left\{ \begin{array}{l} (exists_eligible = \mathbf{true} \wedge 0 < s < 101 \wedge y = \mathbf{true}) \vee \\ \vee (exists_eligible = \mathbf{false} \wedge 80 < s < 101 \wedge y = \mathbf{true}) \vee \\ \vee (exists_eligible = \mathbf{false} \wedge 0 < s \leq 80 \wedge y = \mathbf{false}) \end{array} \right\}$$

and its update inserts a new tuple to the process repository with a score *s* such that $0 < s < 101$:

$$Eff(\text{Evaluate Application}) = \left\{ \begin{array}{l} \text{INSERT}(jcid, uid, name, age, s, \mathbf{undef}) \\ \text{INTO } Application \text{ AND SET } exists_eligible \text{ TO } y \end{array} \right\}$$

Notice that the the case variable *exists_eligible* becomes **true** in case at least an applicant is evaluated with score greater than 80.

⁹ Hereinafter we shall avoid putting explicitly trivial (i.e., containing only **true**) preconditions.

In case a candidate is considered to be perfectly qualified for the position, it immediately gets considered as a winner of the selection process. Such a functionality is carried over by the **Assign Winner** task that simply assigns the winners user identifier to the dedicated case variable:

$$Eff(\text{Assign Winner}) = \left\{ \begin{array}{l} \text{SET } winner' = uid, \\ \text{SET } result' = qualified \\ \text{SET } tPassed' = LessThan1Month \end{array} \right\}$$

If the process of the application evaluation has ended up due to the deadline, the process runs a task that decides on the eligible candidates among all those that have sent applications. Here as eligible we consider only those candidates whose interview score is greater than 80, whereas others are regarded as not eligible. This is done using the **Decide Eligible Candidate** task with the following update:

$$Eff(\text{Decide Eligible Candidate}) = \left\{ \begin{array}{l} \text{UPDATE } Application(JCID, UID, N, A, S, Elig) \\ \text{IFS } > 80 \\ \text{THEN } Application(JCID, UID, N, A, S, 'eligible') \\ \text{ELSE } Application(JCID, UID, N, A, S, 'noteligible') \end{array} \right\}$$

Note that the *if-then-else* clause allows us to perform a sort of a bulky update over the repository relation *Application* by changing the eligibility status of its entries.

If case there is at least one eligible candidate, she can be selected as a winner. This is done by the **Select Winner** task that nondeterministically selects one such candidate from *Application* (i.e., $G(\text{Select Winner}) = Application(Jcid, Uid, Name, Age, Score, Eligible) \wedge Eligible = 'eligible'$) and moves her data to the case variables:

$$Eff(\text{Select Winner}) = \left\{ \begin{array}{l} \text{DEL } (Jcid, Uid, Name, Age, Eligible) \text{ FROM } Application \\ \text{TO } (jcid, uid, name, age, result) \\ \text{AND SET } tPassed = 1MonthPlus1Week \\ \wedge qualif' = false \wedge exists_eligible' = undef \end{array} \right\}$$

Here we also take into account that in order to decide on the eligibility of candidates as well as a winning candidate, the HR staff of the company may require some time. This is duly represented by updating the amount of time passed since the beginning of the candidate selection process.

At last, when a winning candidate has been selected, the HR office prepares a official offer that is then sent to the winner. In our model this is represented with the **Make Offer** activity that does not issue any updates on the data component (i.e., $Eff(\text{Make Offer}) = \text{true}$).

We analyze four examples of reachability queries, taken from the MCMT specifications that we tested in our experimental evaluation: we focus on 1-case safety verification.

The first query expresses that the job hiring process has completed, i.e. it has reached its final state: the tool returns **UNSAFE** (since it exists a sequence

of configurations starting from the initial states to the final one), as expected. Formally, we have:

$$\exists i:PI_{index} (lifecycleProcess[i] = \mathbf{completed})$$

The second query formalizes the situation where, after the evaluation of an application (i.e., EvaluateApplication is completed), there exists at least an applicant with score greater than 0: the tool returns UNSAFE, as expected. Formally, we have:

$$\exists i:PI_{index} \left(lifecycleEvaluateApplication[i] = \mathbf{completed} \wedge \right. \\ \left. \wedge Application(Jcid, Uid, Name, Age, Score, Eligible) \wedge Score > 0 \right)$$

The third query represents the configuration of the system in which a winner has been selected after deadline (i.e., SelectWinner is completed), but the case variable *result* witnesses that the winner was a not eligible candidate: the tool returns SAFE (since this configuration is not reachable from the initial states), as expected. Formally, we have

$$\exists i:PI_{index} (lifecycleSelectWinner[i] = \mathbf{completed} \wedge result[i] = \mathbf{noteligible})$$

The final query describes the configuration in which, after the evaluation of an application, there exists an applicant with score greater than 100: the tool returns SAFE, as expected. Formally, we have:

$$\exists i:PI_{index} \left(lifecycleEvaluateApplication[i] = \mathbf{completed} \wedge \right. \\ \left. \wedge Application(Jcid, Uid, Name, Age, Score, Eligible) \wedge Score > 100 \right)$$

All these queries has been checked running MCMT over the running example.