

Process Query Language: Design, Implementation and Evaluation

ARTEM POLYVYANY, ARTHUR H.M. TER HOFSTEDÉ, MARCELLO LA ROSA, and
CHUN OUYANG, Queensland University of Technology

Organisations can derive significant benefits from the use of practices, techniques, and tools from the area of business process management. Through the concomitant focus on processes they may acquire a substantial number of process models and such large collections of models require management, including support for versioning, merging, conformance checking, and retrieval. Model retrieval thus far has mostly focussed on querying models using syntactic properties of their control flow rather than on exploiting semantic properties capturing aspects of model execution. While the latter is (much) more challenging, it is also more effective, especially in a context where process models serve as a basis for subsequent automation. The focus of this paper is to overcome the challenges associated with semantic querying of process model collections and thus unlocking its benefits. The first challenge concerns determining decidability of the building blocks of the semantic query language, i.e. semantic relations between tasks. This is important as not all such relations can necessarily be computed, especially if one allows them to be based on unrestricted temporal logic statements. The second challenge concerns useability, specifically, selecting those semantic relations among those that are computable that are also considered relevant and intuitive by domain experts. The third and final challenge is concerned with achieving acceptable performance of query evaluation. The evaluation of a query may require expensive checks on all process models (of which there may be thousands) and to deal with this suitable index structures should be put in place. The effectiveness of these structures should be demonstrated empirically.

Categories and Subject Descriptors: D.1.3 [Programming Techniques]: Concurrent Programming—*Distributed programming; Parallel programming*; D.2.0 [Software Engineering]: General; D.3.1 [Programming Languages]: Formal Definitions and Theory—*Semantics; Syntax*; F.3.2 [Logics and Meanings of Programs]: Semantics of Programming Languages—*Process models*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Query formulation; Relevance feedback; Search process*

General Terms: Design, Languages, Management, Theory

Additional Key Words and Phrases: Process, process model, process instance, process model collection, process model repository, querying, process querying, searching, retrieving

ACM Reference Format:

Artem Polyvyany, Arthur H.M. ter Hofstede, Marcello La Rosa, and Chun Ouyang. 2014–2015. Process Query Language: Design, Implementation and Evaluation. *ACM Trans. Embedd. Comput. Syst.* 9, 4, Article 39 (March 2010), 13 pages.
DOI : <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

Through the application of methods and techniques from the field of business process management, organisations can identify, model, analyse, deploy, and diagnose their business processes. This process-oriented thinking provides great benefits as making processes explicit through formal high-level representations, i.e. process models, allows them to subject these processes to various forms of analysis, to use them as the basis for automated support, and to adapt them more easily as well as more rapidly to continual changes imposed by the organisation's environment, both internal and external. As a consequence, some organisations have collected large numbers of process models and their maintenance poses significant challenges.

Authors' address: A. Polyvyany, A.H.M. ter Hofstede, M. La Rosa, and C. Ouyang, Business Process Management Discipline, Information Systems School, Science & Engineering Faculty, Queensland University of Technology, Brisbane, Australia.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2010 ACM 1539-9087/2010/03-ART39 \$15.00

DOI : <http://dx.doi.org/10.1145/0000000.0000000>

Process models tend to evolve over time, new models may merge, which may need to be informed by existing models, and models may need to be merged. To support these activities, it should be possible to query a potentially large process model repository to retrieve models with certain characteristics. Current process model query languages, of which perhaps the most prominent proponent is BPMN-Q, predominantly focus on *syntactic* aspects of process models. Hence queries are based on paths, which are formed by direct succession relations between model elements.

More powerful than a syntactic approach to querying would be an approach based on *semantic* relations between tasks, e.g. relations that capture that certain tasks need to be executed in order or in parallel or can never be executed as part of the same process instance. Semantic relations are essential when one needs to explore whether certain process behaviours are possible or not.

The added retrieval power of a semantic query language comes at a price. Semantic relations cover a broad spectrum of inter-task dependencies and may be expressed through temporal logic. Temporal logic is a form of modal logic that includes temporal operators and which can be used to reason over the behaviour of systems over time. Temporal logic is powerful enough to be able to express properties that are undecidable [Esparza and Nielsen 1994]. Hence a semantic query language needs to be careful in terms of the semantic inter-task dependencies that it supports.

Decidability of semantic inter-task relations is a factor in choosing which relations to support in a query language, but it is not the one. While some relations are decidable, their use may not be very intuitive to stakeholders, in this case experts that are likely to end up formulating queries over process model collections. It is important that a relation occurs frequently enough in queries to warrant support and that its formal meaning is close to its perceived meaning. Another consideration is that query evaluations are performed in a “reasonable” amount of time as it is anticipated that stakeholders may wish to see the answers to their queries in (almost) real-time.

In this paper a specific process model query language is proposed that is based on a selection of semantic inter-task relations. The Process Query Language (PQL) is a special-purpose programming language for managing process models based on information about process instances that they describe. PQL programs are also called queries. The first version of the PQL language, proposed in this paper, allows formulating search intents for retrieving process models from collections or repositories thereof based on information held in process instances.

The selected semantic inter-task relations supported by PQL, termed *predicates* in PQL, are shown to be decidable through the application of model checking, an automated technique which, given a finite-state model of a system (e.g., a process model) and a formal property (e.g., a temporal logic formula), systematically checks whether this property holds for (a given state in) that model [Baier and Katoen 2008]. In addition, these predicates are validated with domain experts in terms of their perceived usefulness and intuitiveness. To further facilitate query formulation, PQL is not only provided with an abstract syntax but also a concrete one, which is inspired by SQL. The runtime environment for PQL makes use of indexes to enhance the performance of query evaluation. Indexes are special data structures that improve the speed of computations of behavioral operators trading off time for their construction with space for their storage. Performance of query evaluation is demonstrated through a set of experiments with a real-life process model collection.

2. PQL SYNTAX

2.1. Abstract Syntax

This section discusses the syntax of the PQL language in the form of an *abstract syntax*, which is also often referred to as an (*abstract*) *grammar*. The grammar of the PQL language is defined using the notation introduced in [Meyer 1990]. In this notation, the abstract grammar of a programming language consists of a finite set of names of *constructs* and a finite set of *productions*, each associated with a construct. Each construct describes the structure of a set of objects, also called *specimens* of the language, using productions of three types; these are *aggregate*, *choice*, and *list* productions.

The top construct of the abstract grammar of the PQL language is Query. It captures the core structure of all PQL programs, i.e., queries.

$$\text{Query} \triangleq \text{vars}:\text{Variables}; \text{atts}:\text{Attributes}; \text{locs}:\text{Locations}; \text{pred}:\text{Predicate}$$

The Query construct (see above, on the left hand side) is defined as an *aggregate* production composed of four components (see above, on the right hand side); in general, an aggregate production defines a construct that is made of a fixed number of components. The components are separated by semicolons, each preceded by a *tag* indicating its *role* within the construct. Thus, every PQL query is composed of variables, attributes, locations, and a predicate, which are distinguished via tags *vars*, *atts*, *locs*, and *pred*, respectively. Intuitively, a PQL query specifies a search intent to discover the attributes of *all* process models in the collection of models identified by the locations that satisfy the predicate, where the evaluation of the predicate relies on information stored in the variables. The order in which the various specimens are listed in aggregate productions is irrelevant for the sake of the abstract grammar specification. This order is important in the context of the next section, in which the concrete syntax of the PQL language is proposed.

The Query construct seen above defines a class of PQL queries. One can specify an instance of this class using abstract syntactic expressions. For example, the statement $q \triangleq \text{Query}(\text{vars} : vs; \text{atts} : as; \text{locs} : ls; \text{pred} : p)$ defines a query having *vs*, *as*, *ls*, and *p*, as variables, attributes, locations, and a predicate, respectively (assuming that all the specimens, i.e., *vs*, *as*, *ls*, and *p*, are provided).

In PQL, variables, attributes, and locations are defined as *list* productions, where a list production defines a sequence of zero, one, or more specimens of another construct.

$$\begin{aligned} \text{Variables} &\triangleq \text{Variable}^* \\ \text{Attributes} &\triangleq \text{Attribute}^+ \\ \text{Locations} &\triangleq \text{Location}^+ \end{aligned}$$

Therefore, a PQL query defines a sequence of zero, one, or more variables, denoted by Variable^* ; the asterisk symbol stands for the *Kleene star*—its standard language theory meaning. Every sequence of attributes must contain at least one attribute, denoted by Attribute^+ ; note that the asterisk symbol is replaced by a plus sign to signify that the list of locations cannot be empty. Similarly, every sequence of locations must contain at least one location specimen.

PQL introduces a dedicated construct, denoted by Variable , to define *variables*.

$$\text{Variable} \triangleq \text{name}:\text{VariableName}; \text{tasks}:\text{SetOfTasks}$$

A PQL variable associates a symbolic name with a set of *tasks*, or to be more precise, with a collection of abstract concepts that represent tasks. Tasks are introduced in the PQL language to refer to atomic units of observable behavior that are captured in process models. Each variable is an aggregate of two constructs: a variable name ($\text{name}:\text{VariableName}$), and a collection of tasks ($\text{tasks}:\text{SetOfTasks}$). Such a separation of the variable name from its associated content allows the name to be used independently of the exact information it represents. Thus, a variable name can be bound to a set of tasks during run time, and the content of the set may change during evaluation of the query. When a predicate of some PQL query gets evaluated, every variable name that is mentioned in the predicate is replaced by the corresponding set of tasks.

The PQL language introduces the Attribute construct to allow specifying those process model properties that must be retrieved in a response to a successful query matching exercise.

$$\text{Attribute} \triangleq \text{Universe} \mid \text{AttributeID} \mid \text{AttributeName} \mid \text{AttributeModel}$$

A PQL attribute can be anything that identifies a single property or a collection of properties of a process model. In the first version of the PQL language, the Attribute construct is associated with a *choice* production that allows for four alternatives; in general, a choice production defines a construct as a set of alternatives. The alternatives are separated by vertical bar symbols. Every attribute is either the *universe* attribute, denoted by Universe , the *identifier* property, denoted by AttributeID , the *name* property, denoted by AttributeName , or a formal *specification* property,

denoted by `AttributeModel`. The universe attribute has a special meaning. It refers to a collection of all properties that are associated with process models in the process model repository.

In the PQL language, locations are used to address process models that are of interest to the search intents of queries, i.e., those models that should be matched against the queries.

$$\text{Location} \triangleq \text{Universe} \mid \text{LocationID} \mid \text{LocationDirectory}$$

A location can generally be anything that identifies a single process model or a collection of process models. It is defined as a *choice* production that allows for three alternatives. A location is either the *universe* location, denoted by `Universe`, an *identifier* location, denoted by `LocationID`, or a *directory* location, denoted by `LocationDirectory`. The universe location is designed to address all process model in the scope of the query (usually, all process models in the repository). Identifier locations are introduced in the PQL language to allow fine-grained targeting of models based on their unique identifiers. We assume that repositories do indeed tag models with unique identifiers, e.g., universally unique identifiers (UUIDs) or integer identifiers. Finally, a directory location allows addressing models that are stored in a particular directory of the repository. For example, a directory location can be a name of a directory, specified as a character string, an URI [URI Planning Interest Group 2001], or an XPath expression [W3C XSL/XML Query Working Groups 2007].

PQL provides several alternatives for specifying a set of tasks. A set of tasks can be defined as an enumeration of tasks, a result of standard operations on sets of tasks, information stored in a variable, a construction macro, or a dynamically-valued constant. These various possibilities are captured in the `SetOfTasks` construct of the PQL grammar.

$$\begin{aligned} \text{SetOfTasks} \triangleq & \text{VariableName} \mid \text{Universe} \mid \text{SetOfTasksLiteral} \mid \text{SetOfTasksConstruction} \\ & \mid \text{Union} \mid \text{Intersection} \mid \text{Difference} \end{aligned}$$

The `SetOfTasks` construct is defined as a choice production. One can use the `VariableName` construct to refer to the set of tasks associated with a name of some variable. Alternatively, one can specify a set of tasks using the `Universe` construct. The `Universe` construct, when used in the context of a reference to a set of tasks, constitutes a dynamically-valued constant that refers to the set of all tasks of the process model currently being matched to the query. The content of this set should be created at initialization time, freshly for every new process model that gets matched to the query.

The PQL language proposes a notation to specify set of tasks literals, i.e., a notation for representing sets of tasks as fixed values. Set of tasks literals can be defined using the `SetOfTasksLiteral` construct, which is specified as a list production of zero, one, or more tasks.

$$\text{SetOfTasksLiteral} \triangleq \text{Task}^*$$

As mentioned above, tasks are abstract representations of atomic units of observable behavior.

$$\text{Task} \triangleq \text{label}:\text{Label}; \text{sim}:\text{Similarity}$$

A PQL task is defined as an aggregate of two components: a *label*, denoted by `label:Label`, and a *similarity* degree threshold, denoted by `sim:Similarity`. The idea is that given a label of a task one may be interested in all the tasks of which the label has (at least) a certain degree of similarity to the given label.

Another way to specify a set of tasks is to construct it. For this purpose, one can rely on the `SetOfTasksConstruction` construct, which is defined as a choice production below.

$$\text{SetOfTasksConstruction} \triangleq \text{UnaryPredicateConstruction} \mid \text{BinaryPredicateConstruction}$$

Given a set of tasks and a unary behavioral primitive, the `UnaryPredicateConstruction` construct can be used to compose a set of tasks that contains every task from the given set and for which the given behavioral primitive holds. The given behavioral primitive must be evaluated in the context of the process model that is being matched to the query. Similarly, the `BinaryPredicateConstruction` construct is introduced in the PQL language to allow selecting those tasks from a given set of tasks for which certain binary behavioral primitive holds, either with at least one or with all tasks taken from another given set of tasks. The choice of a quantifier type,

either the existential or universal, to be used during the above described selections is implemented via the `AnyAll` construct.

$$\begin{aligned} \text{UnaryPredicateConstruction} &\triangleq \text{name} : \text{UnaryPredicateName}; \text{tasks} : \text{SetOfTasks} \\ \text{BinaryPredicateConstruction} &\triangleq \text{name} : \text{BinaryPredicateName}; \text{tasks}_1 : \text{SetOfTasks}; \\ &\quad \text{tasks}_2 : \text{SetOfTasks}; q : \text{AnyAll} \\ \text{AnyAll} &\triangleq \text{Any} \mid \text{All} \end{aligned}$$

Both, the `UnaryPredicateConstruction` construct and the `BinaryPredicateConstruction` construct, are associated with aggregate productions. The `AnyAll` construct is specified as a choice between the `Any` qualifier versus the `All` qualifier, where `Any` and `All` stand for the existential quantifier type and the universal quantifier type, respectively. The PQL language uses the `UnaryPredicateName` construct and the `BinaryPredicateName` construct to refer to unary behavioral primitives and binary behavioral primitives, respectively. The first edition of the PQL language supports two unary and six binary behavioral primitives. These are the `CanOccur` and `AlwaysOccurs` unary behavioral primitives, and the `CanConflict`, `CanCooccur`, `Conflict`, `Cooccur`, `TotalCausal`, and `TotalConcurrent` binary behavioral primitives.

$$\begin{aligned} \text{UnaryPredicateName} &\triangleq \text{CanOccur} \mid \text{AlwaysOccurs} \\ \text{BinaryPredicateName} &\triangleq \text{CanConflict} \mid \text{CanCooccur} \mid \text{Conflict} \\ &\quad \mid \text{Cooccur} \mid \text{TotalCausal} \mid \text{TotalConcurrent} \end{aligned}$$

Finally, a set of tasks can be constructed from other sets of tasks via the application of the fundamental set operations of *union*, *intersection*, and *difference*, denoted by the `Union`, `Intersection`, and `Difference` construct, respectively.

The PQL language proposes several ways to specify predicates; all the options are captured in the choice production that is associated with the `Predicate` construct.

$$\begin{aligned} \text{Predicate} &\triangleq \text{UnaryPredicate} \mid \text{BinaryPredicate} \mid \text{UnaryPredicateMacro} \\ &\quad \mid \text{BinaryPredicateMacro} \mid \text{SetPredicate} \mid \text{TruthValue} \mid \text{Negation} \\ &\quad \mid \text{Conjunction} \mid \text{Disjunction} \mid \text{LogicalTest} \end{aligned}$$

For instance, predicates can be captured as specimens of `UnaryPredicate` or `BinaryPredicate`.

$$\begin{aligned} \text{UnaryPredicate} &\triangleq \text{name} : \text{UnaryPredicateName}; \text{task} : \text{Task} \\ \text{BinaryPredicate} &\triangleq \text{name} : \text{BinaryPredicateName}; \text{task}_1 : \text{Task}; \text{task}_2 : \text{Task} \end{aligned}$$

The `UnaryPredicate` construct and the `BinaryPredicate` construct are introduced in the PQL language to allow checking the unary behavioral primitives and the binary behavioral primitives, respectively. Both these constructs are aggregations of a name (specified by the `UnaryPredicateName` construct or the `BinaryPredicateName` construct) and a respective number of `Task` constructs; one for the `UnaryPredicate` construct and two for the `BinaryPredicate` construct.

The PQL language utilizes a well-known mechanism of macros for combining results of several `UnaryPredicate` or `BinaryPredicate` checks into a result of a single statement.

$$\begin{aligned} \text{UnaryPredicateMacro} &\triangleq \text{name} : \text{UnaryPredicateName}; \text{tasks} : \text{SetOfTasks}; q : \text{AnyAll} \\ \text{BinaryPredicateMacro} &\triangleq \text{BinaryPredicateMacroTaskSet} \mid \text{BinaryPredicateMacroSetSet} \end{aligned}$$

The aggregate production associated with the `UnaryPredicateMacro` construct is composed of a reference to a unary behavioral primitive ($\text{name} : \text{UnaryPredicateName}$), a set of tasks ($\text{tasks} : \text{SetOfTasks}$), and a quantifier ($q : \text{AnyAll}$). Intuitively, a single macro statement $p \triangleq \text{UnaryPredicateMacro}(\text{name} : n; \text{tasks} : ts; q : x)$ is equivalent to a complex check of whether it holds that for at least one (if x is set to `Any`) or for every (if x is set to `All`) task t in set of tasks ts statement `UnaryPredicate(p.name; task : t)` evaluates to *true*. Similarly, one can rely on the

BinaryPredicateMacro construct to combine results of multiple BinaryPredicate checks.

$$\begin{aligned} \text{BinaryPredicateMacroTaskSet} &\triangleq \text{name} : \text{BinaryPredicateName}; \text{task} : \text{Task}; \\ &\quad \text{tasks} : \text{SetOfTasks}; \text{q} : \text{AnyAll} \\ \text{BinaryPredicateMacroSetSet} &\triangleq \text{name} : \text{BinaryPredicateName}; \text{tasks}_1 : \text{SetOfTasks}; \\ &\quad \text{tasks}_2 : \text{SetOfTasks}; \text{q} : \text{AnyEachAll}; \\ \text{AnyEachAll} &\triangleq \text{Any} \mid \text{Each} \mid \text{All} \end{aligned}$$

The BinaryPredicateMacroTaskSet construct is designed to allow checking whether a certain binary behavioral primitive (*name* : BinaryPredicateName) holds between a given task (*task* : Task) and either at least one (if the AnyAll construct is instantiated with the Any specimen) or every (if the AnyAll construct is instantiated with the All specimen) task in a given set of tasks (*tasks* : SetOfTasks). Similarly, the BinaryPredicateMacroSetSet construct can be used to check whether a binary behavioral primitive of interest evaluates to *true* for certain pairs of tasks in the Cartesian product of two given sets of tasks. Note for the option to use the Each qualifier as a specimen of the AnyEachAll construct in the respective production above. When employed, this option induces a check of whether for every task in one given set of tasks the specified behavioral relation holds with some task from the other given set of tasks.

The PQL language supports checks of basic binary relations between sets of tasks. These are captured by the choice production associated with the SetPredicate construct.

$$\begin{aligned} \text{SetPredicate} &\triangleq \text{TaskInSetOfTasks} \mid \text{SetComparison} \\ \text{TaskInSetOfTasks} &\triangleq \text{task} : \text{Task}; \text{tasks} : \text{SetOfTasks} \\ \text{SetComparison} &\triangleq \text{tasks}_1 : \text{SetOfTasks}; \text{oper} : \text{SetComparisonOperator}; \text{tasks}_2 : \text{SetOfTasks} \\ \text{SetComparisonOperator} &\triangleq \text{Identical} \mid \text{Different} \mid \text{OverlapsWith} \\ &\quad \mid \text{SubsetOf} \mid \text{ProperSubsetOf} \end{aligned}$$

The PQL language allows checking if a task is a member of a given set of tasks. This can be accomplished using the TaskInSetOfTasks construct, which is specified as an aggregation of a task (*task* : Task) and a set of tasks (*tasks* : SetOfTasks). Moreover, the PQL language allows checking several binary relations between sets of tasks using the SetComparison construct. The SetComparison construct is composed of two sets of tasks (*tasks*₁ : SetOfTasks and *tasks*₂ : SetOfTasks) and a reference to a comparison operator (*oper* : SetComparisonOperator). The PQL language supports five comparison operations. They specify checks of whether two sets of tasks are identical (Identical), different (Different), overlap (OverlapsWith), or whether one set of tasks is a subset (SubsetOf) or a proper subset (ProperSubsetOf) of the other set of tasks.

As PQL is designed to utilize three-valued reasoning, it operates with three truth values: *true*, *false*, and *unknown*. This is reflected in the three literals of the choice production associated with the TruthValue construct, which is proposed below.

$$\text{TruthValue} \triangleq \text{True} \mid \text{False} \mid \text{Unknown}$$

To allow complex logical statements on atomic propositions, PQL supports standard logical operations. These are negation (Negation), conjunction (Conjunction), and disjunction (Disjunction).

To permit for a three-valued logic that is used with PQL to be functionally complete, the language includes a test of whether a given three-valued logic value is *unknown*. This check is reflected in the IsUnknown option of the LogicalTest construct proposed below, which for the sake of completeness allows for the total of six different tests of whether a three-valued logic value is or is not equal to a certain truth value.

$$\text{LogicalTest} \triangleq \text{IsTrue} \mid \text{IsNotTrue} \mid \text{IsFalse} \mid \text{IsNotFalse} \mid \text{IsUnknown} \mid \text{IsNotUnknown}$$

For a grammar of a language to be complete, all its constructs must be specified in terms of well-defined components, called the *terminal* constructs. The following constructs are the terminal constructs of the PQL grammar: Any, All, Each, Universe, AttributeID, AttributeName, AttributeModel, Identical, Different, OverlapsWith, SubsetOf, ProperSubsetOf, True,

False, Unknown, as well as all the constructs that are parts of the choice productions associated with the UnaryPredicateName and BinaryPredicateName constructs. All the above mentioned constructs do not have an internal structure and, thus, are atomic constructs of the PQL language.

Several of the proposed PQL constructs can be defined in terms of well-known sets. For instance, in the discussion above, we hint at the fact that LocationID can be specified as an integer. This can be captured rigorously in the production $\text{LocationID} \triangleq \text{value} : \mathbb{Z}$, where \mathbb{Z} is the symbol often used in mathematics to denote the set of all integers. Similarly, we specify LocationDirectory, VariableName, and Similarity, as $\text{LocationDirectory} \triangleq \text{value} : \mathbb{S}$, $\text{VariableName} \triangleq \text{id} : \mathbb{V}$, and $\text{Similarity} \triangleq \text{value} : [0..1]$, respectively; here, \mathbb{S} and \mathbb{V} are the set of all character strings and the set of all legal variable names, respectively. Note that set \mathbb{V} is defined in the next section.

Some of the PQL constructs are still not defined in terms of terminal constructs. The PQL language trivially defines the Negation construct and all the six options associated with the LogicalTest construct in terms of a single Predicate component, e.g., $\text{Negation} \triangleq \text{pred} : \text{Predicate}$, $\text{IsTrue} \triangleq \text{pred} : \text{Predicate}$, etc. Finally, for the sake of space considerations, at this stage we omit rigorous definitions of five PQL constructs: Conjunction, Disjunction, Union, Intersection, and Difference. Intuitively, Conjunction and Disjunction can be defined as sets of predicates, whereas Union, Intersection, and Difference can be specified as collections of sets of tasks. However, any definition of priorities for the operations that the above stated constructs represent in terms of grammar rules is rather lengthy and is driven by semantic, rather than syntactic, rules. In the next section, we discuss priorities of various operations that are supported in PQL, whereas missing rigorous specifications of the five mentioned constructs can be found in Appendix A.

2.2. Concrete Syntax

The abstract syntax of PQL is independent of any particular representation. This section proposes a mapping from the abstract syntax of PQL to its specific encoding. This encoding constitutes one possible concrete syntax of the PQL language, i.e., its machine- and human-readable representation.

The first concrete syntax of the PQL language proposed in this section is inspired by SQL—a programming language for managing data stored in a relational database management system (DBMS) [Date and Darwen 1997]. Being inspired by SQL, we intent to keep the core structure of concrete PQL queries as similar as possible to that of SQL queries and to reuse SQL keywords in PQL, given that the contexts are similar. The reason for this is threefold:

- Despite addressing different domains, i.e., dynamic processes versus static data, both languages serve the same purpose—the purpose of querying for information. Note that SQL was originally proposed to *retrieve* data stored in quasi-relational DBMS [Chamberlin and Boyce 1974].
- SQL is a widely used standard that is supported by just about every DBMS on the market. As a result, its syntax is well-recognized by technical specialists and analysts. By closely following the concrete syntax of SQL, PQL becomes readily usable by a wide range of stakeholders.
- As suggested by several interviewees, it would be beneficial for the syntax of the envisioned query language to resemble that one of SQL. For example, one interviewee commented: “From an overall strategic point of view it’ll bring a lot of benefits because different parts of the organization will be able to work together by using some kind of a structured query language (SQL)”.

Given a construct of an abstract grammar, one can specify its concrete syntax as a function that yields all its specific forms. In this section, PQL is defined as a textual language. Hence, for each PQL construct, its concrete syntax is given as a function that takes a specimen of the respective abstract construct as input and returns a collection of character strings that are accepted as its concrete encodings. We shall denote such a function by the name of the respective construct with subscript c .

For example, the concrete syntax of a specimen of the Query construct is defined as follows.

$$\begin{aligned} \text{Query}_c(q : \text{Query}) \triangleq & \text{Variables}_c(q.\text{vars}) \\ & \text{'SELECT' } \text{Attributes}_c(q.\text{atts}) \\ & \text{'FROM' } \text{Locations}_c(q.\text{locs}) \\ & (\text{'WHERE' } \text{Predicate}_c(q.\text{pred}))? \text{' ;' } \end{aligned}$$

We use regular expressions [Aho and Ullman 1992] to define the concrete syntax of PQL specimens. Hence, as per the above definition, a PQL query is a character string that starts with a specification of variables, followed by the **SELECT** keyword, followed by a specification of attributes, followed by the **FROM** keyword, followed by a specification of locations, followed by the **WHERE** keyword, followed by a specification of a predicate, followed by the semicolon mark, i.e., ‘;’. There can be an arbitrary number of whitespace characters between any two subsequent components of a query string. The order of various components is fixed. Note that the presence of the **WHERE** clause in a PQL query is optional, i.e., the **WHERE** keyword and the specification of the predicate can be skipped.

The reader might have already noticed that the core structure of a PQL query is similar to that one of an SQL query that is signified with the declarative **SELECT** statement and is used to formulate an intent for retrieving data from one or more database tables or expressions.

Specimens of PQL constructs that are associated with list productions must be encoded as string concatenations of concrete forms of their components and whitespace characters. Often, we inject special symbols between every two subsequent components and/or at the beginning and end of the respective encodings. For example, the concrete syntax of a list of variables is defined as follows.

$$\text{Variables}_c(vs: \text{Variables}) \triangleq \text{isEmpty}(vs) ? '': \text{Variable}_c(vs.FIRST) \text{Variables}_c(vs.TAIL)$$

That is, the encoding of the empty list of variables is the empty string. However, if a list of variables contains at least one element, its encoding is constructed as a concatenation of a concrete form of its first element, denoted by $vs.FIRST$, and an encoding of the list of its all other elements, denoted by $vs.TAIL$. The concrete syntax of a PQL variable is defined below.

$$\text{Variable}_c(v: \text{Variable}) \triangleq \text{VariableName}_c(v.name) \text{'='} \text{SetOfTasks}_c(v.tasks) \text{';'}$$

The concrete syntax of all other specimens of PQL constructs that are associated with list productions is defined similar to that one of the **Variables** construct seen above. However, all these encodings expect to include special symbols between every two subsequent components. This is the comma symbol, i.e., ‘,’, for specimens of **Attributes**, **Locations**, and **SetOfTasksLiteral**, and the PQL keywords **UNION**, **INTERSECT**, **EXCEPT**, **AND**, and **OR**, for specimens of **Union**, **Intersection**, **Difference**, **Conjunction**, and **Disjunction**, respectively. Additionally, every encoding of a specimen of the **SetOfTasksLiteral** construct must begin with the opening curly bracket, i.e., ‘{’, and end with the closing curly bracket, i.e., ‘}’. For instance, the character string ‘{“Buy item”, “Purchase product”}’ is a valid encoding of a specimen of the **SetOfTasksLiteral** construct that contains two elements; here, we follow the standard notation for specifying fixed sets. In the example above, strings “Buy item” and “Purchase product” are valid encodings of tasks. In general, the concrete encoding of a PQL task is defined as follows.

$$\begin{aligned} \text{Task}_c(t: \text{Task}) \triangleq & \text{'~'} \text{'"} \text{Label}_c(t.label) \text{'"} \\ & | \text{'"} \text{Label}_c(t.label) \text{'"} \text{'(['} \text{Similarity}_c(t.sim) \text{' ']} \text{'?} \end{aligned}$$

Labels of PQL tasks are always enclosed in double quotes. A label can be preceded by the tilde symbol, i.e., ‘~’, or succeeded by an encoding of a similarity degree threshold enclosed in square brackets. The tilde symbol denotes the fact that one is interested in all the tasks of which the label has a degree of similarity to the specified label that is equal or larger than some preconfigured value. A degree of similarity must be specified as a decimal representation of a real number greater or equal to zero and less than or equal to one, e.g., 0.5 or .95.

Every specimen of a construct that is associated with a choice production is a specimen of one of the constructs from the list of alternatives of the choice production. Hence, a concrete syntax of an abstract grammar can (and often does) omit special encodings to signify choice productions, which is the case for the concrete syntax of PQL that is being proposed here. Thus, in the sequel, we only propose concrete encodings for the remaining aggregate productions of the PQL grammar.

A specimen of the **Attribute** construct is a specimen of one out of four terminal constructs. These are **Universe**, **AttributeID**, **AttributeName**, and **AttributeModel**. In PQL, they are denoted by character strings ‘*’, ‘id’, ‘name’, and ‘model’, respectively. Similarly, a specimen of the **Location** construct is a specimen of either **Universe**, or **LocationID**, or **LocationDirectory**.

A specimen of `LocationID` is an integer, whereas a specimen of `LocationDirectory` is a character string; PQL character strings are sequences of characters enclosed in double quotes.

A name $x \in \mathbb{V}$ of a PQL variable, i.e., a concrete encoding of a specimen of the `VariableName` construct, may contain lower case letters from the English alphabet, digits, and the underscore symbol, i.e., ‘_’. It is necessary to use a letter or the underscore symbol at the start of a variable name; digit at start are not allowed. Subsequent characters may be letters, digits, or underscore symbols.

Next, we propose the concrete syntax for specimens of the `UnaryPredicateConstruction` construct and the `BinaryPredicateConstruction` construct.

$$\text{UnaryPredicateConstruction}_c(\text{upc} : \text{UnaryPredicateConstruction}) \triangleq \\ \text{‘GetTasks’UnaryPredicateName}_c(\text{upc.name}) \text{ ‘(’ SetOfTasks}_c(\text{upc.tasks}) \text{ ‘)’}$$

$$\text{BinaryPredicateConstruction}_c(\text{bpc} : \text{BinaryPredicateConstruction}) \triangleq \\ \text{‘GetTasks’BinaryPredicateName}_c(\text{bpc.name}) \\ \text{‘(’ SetOfTasks}_c(\text{bpc.tasks}_1) \text{ ‘,’ SetOfTasks}_c(\text{bpc.tasks}_2) \text{ ‘,’ AnyAll}_c(\text{bpc.q}) \text{ ‘)’}$$

The concrete encodings of specimens of these constructs follow the syntax for specifying function calls that is used in many well-celebrated programming languages, i.e., a name of a function to be called is followed by a comma-separated list of parameters which is enclosed in parentheses. Here, the names of functions are obtained by prefixing ‘GetTasks’ to names of unary and binary predicates. The remaining components are specified as parameters of the respective function calls.

PQL exercises similar principles when specifying the concrete syntax of predicates and macros, both for unary and binary cases. The concrete syntax for specifying predicates proceeds as follows.

$$\text{UnaryPredicate}_c(\text{up} : \text{UnaryPredicate}) \triangleq \\ \text{UnaryPredicateName}_c(\text{up.name}) \text{ ‘(’ Task}_c(\text{up.task}) \text{ ‘)’}$$

$$\text{BinaryPredicate}_c(\text{bp} : \text{BinaryPredicate}) \triangleq \\ \text{BinaryPredicateName}_c(\text{bp.name}) \text{ ‘(’ Task}_c(\text{bp.task}_1) \text{ ‘,’ Task}_c(\text{bp.task}_2) \text{ ‘)’}$$

The only difference is that the names of these imitated function calls are not prefixed, but are solely composed of the concrete encodings of the respective predicate names.

It is proposed that the concrete syntax for denoting the PQL macros overloads the syntax for specifying function calls which encode the PQL predicates, i.e., names of functions and types of outputs are the same, both for a predicate and the respective macro, but types of inputs differ.

$$\text{UnaryPredicateMacro}_c(\text{upm} : \text{UnaryPredicateMacro}) \triangleq \\ \text{UnaryPredicateName}_c(\text{upm.name}) \text{ ‘(’ SetOfTask}_c(\text{upm.tasks}) \text{ ‘,’ AnyAll}_c(\text{upm.q}) \text{ ‘)’}$$

$$\text{BinaryPredicateMacroTaskSet}_c(\text{bpm} : \text{BinaryPredicateMacroTaskSet}) \triangleq \\ \text{BinaryPredicateName}_c(\text{bpm.name}) \\ \text{‘(’ Task}_c(\text{bpm.task}) \text{ ‘,’ SetOfTask}_c(\text{bpm.tasks}) \text{ ‘,’ AnyAll}_c(\text{bpm.q}) \text{ ‘)’}$$

$$\text{BinaryPredicateMacroSetSet}_c(\text{bpm} : \text{BinaryPredicateMacroSetSet}) \triangleq \\ \text{BinaryPredicateName}_c(\text{bpm.name}) \\ \text{‘(’ SetOfTask}_c(\text{bpm.tasks}_1) \text{ ‘,’ SetOfTask}_c(\text{bpm.tasks}_2) \text{ ‘,’ AnyEachAll}_c(\text{bpm.q}) \text{ ‘)’}$$

The above proposed syntax rules rely on the concrete encodings of `AnyAll` and `AnyEachAll`, which when instantiated are specified as specimens of `Any`, `Each`, or `All`. These are denoted by the PQL keywords `ANY`, `EACH`, and `ALL`, respectively.

A specimen of the `TaskInSetOfTasks` construct can be typed in as follows.

$$\text{TaskInSetOfTasks}_c(\text{in} : \text{TaskInSetOfTasks}) \triangleq \text{Task}_c(\text{in.task}) \text{ ‘IN’ SetOfTask}_c(\text{in.tasks})$$

That is, every character string that starts with an encoding of a task, which is followed by the PQL keyword `IN`, and ends with an encoding of a set of tasks, specifies a specimen of `TaskInSetOfTasks`.

A specimen of the `SetComparison` construct can be specified in this concrete syntax of PQL as two encodings of sets of tasks with a representation of a comparison operator in between.

$$\text{SetComparison}_c(\text{comp} : \text{SetComparison}) \triangleq \\ \text{SetOfTask}_c(\text{comp.tasks}_1) \text{SetComparisonOperator}_c(\text{comp.oper}) \text{SetOfTask}_c(\text{comp.tasks}_2)$$

A set comparison operator is instantiated in PQL via a choice between specimens of terminal constructs `Identical`, `Different`, `OverlapsWith`, `SubsetOf`, and `ProperSubsetOf`, which are defined in this concrete syntax of PQL via the keywords `EQUALS`, `NOT EQUALS`, `OVERLAPS WITH`, `IS SUBSET OF`, and, `IS PROPER SUBSET OF`, respectively. The specimens of the terminal constructs `True`, `False`, and `Unknown` get encoded as the keywords `TRUE`, `FALSE`, and `UNKNOWN`, respectively.

Predicate names in PQL are encoded as names of the respective terminal constructs, e.g., the unary predicate construct `CanOccur` and the binary predicate construct `TotalCausal` have concrete encodings `CanOccur` and `TotalCausal`, respectively.

Finally, the concrete encodings of the specimens of the `Negation` and the six logical test constructs are given below.

$$\begin{aligned} \text{Negation}_c(\text{not} : \text{Negation}) &\triangleq \text{'NOT'} \text{Predicate}_c(\text{not.pred}) \\ \text{IsTrue}_c(\text{test} : \text{IsTrue}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'TRUE'} \\ \text{IsNotTrue}_c(\text{test} : \text{IsNotTrue}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'NOT'} \text{'TRUE'} \\ \text{IsFalse}_c(\text{test} : \text{IsFalse}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'FALSE'} \\ \text{IsNotFalse}_c(\text{test} : \text{IsNotFalse}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'NOT'} \text{'FALSE'} \\ \text{IsUnknown}_c(\text{test} : \text{IsUnknown}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'UNKNOWN'} \\ \text{IsNotUnknown}_c(\text{test} : \text{IsNotUnknown}) &\triangleq \text{Predicate}_c(\text{test.pred}) \text{'IS'} \text{'NOT'} \text{'UNKNOWN'} \end{aligned}$$

These encodings rely on the use of the PQL keywords `NOT`, `IS`, `TRUE`, `FALSE`, `UNKNOWN`, and the concrete encoding of the `Predicate` construct, which is defined by concrete encodings of all the alternatives associated with the corresponding choice production.

We envision introduction of other specific encodings of the abstract syntax of the PQL language, e.g., a visual encoding of PQL queries. We believe that availability of different concrete encodings will make the PQL language accessible to a wider audience.

A. PQL GRAMMAR

This appendix specifies the complete grammar of the PQL language captured in the ANTLR notation. ANTLR (ANOther Tool for Language Recognition) is a parser generator for reading and translating structured text or binary files [Parr and Quong 1995; Parr 2013]. ANTLR can take a grammar of a language as input and generate source code for a parser that can build and walk syntax trees [Meyer 1990]. The language must be specified using a context-free grammar which is expressed using extended Backus-Naur Form [Hopcroft et al. 2006].

```

1 // A PQL v1 grammar for ANTLR v4
2
3 // [The "BSD licence"]
4 // Copyright (c) 2014-2015 Artem Polyvyanyy
5 // All rights reserved.
6
7 // A PQL v1 grammar for ANTLR v4
8
9 grammar PQL;
10
11 query      : variables
12           : SELECT attributes
13           : FROM locations
14           : (WHERE predicate)? EOS ;
15
16 variables : variable* ;
17 variable  : varName ASSIGN
18           : setOfTasks EOS ;
19
20 varName   : VARIABLE_NAME ;
21
22 attributes : attribute (SEP attribute)* ;
23 attribute  : universe
24           : | attributeID
25           : | attributeName
26           : | attributeModel ;
27
28 locations  : location (SEP location)* ;
29 location   : universe
30           : | locationID
31           : | locationDirectory ;
32

```

```

33 universe          : UNIVERSE ;
34 attributeID       : ATTRIBUTE_ID ;
35 attributeName     : ATTRIBUTE_NAME ;
36 attributeModel    : ATTRIBUTE_MODEL ;
37 locationID        : INTEGER ;
38 locationDirectory : STRING ;
39
40 setOfTasks        : tasks
41                   | union
42                   | intersection
43                   | difference ;
44
45 tasks             : ( varName | universe )
46                   | setOfTasksLiteral
47                   | setOfTasksConstruction
48                   | setOfTasksParentheses ;
49
50 setOfTasksLiteral :
51   LB (task (SEP task)*)? RB ;
52
53 task              : approximate label
54                   | label (LSB similarity RSB)? ;
55 approximate: TILDE ;
56 label             : STRING ;
57 similarity        : SIMILARITY ;
58
59 setOfTasksConstruction :
60   unaryPredicateConstruction
61   | binaryPredicateConstruction ;
62
63 unaryPredicateConstruction :
64   (GET_TASKS)unaryPredicateName
65   LP setOfTasks RP ;
66
67 binaryPredicateConstruction :
68   (GET_TASKS)binaryPredicateName
69   LP setOfTasks SEP setOfTasks
70   SEP anyAll RP ;
71
72 anyAll           : ANY | ALL ;
73
74 unaryPredicateName : CAN_OCCUR
75                   | ALWAYS_OCCURS ;
76
77 binaryPredicateName : CAN_CONFLICT
78                   | CAN_COOCCUR
79                   | CONFLICT
80                   | COOCCUR
81                   | TOTAL_CAUSAL
82                   | TOTAL_CONCUR ;
83
84 predicate        : proposition
85                   | conjunction
86                   | disjunction
87                   | logicalTest ;
88
89 proposition: unaryPredicate
90             | binaryPredicate
91             | unaryPredicateMacro
92             | binaryPredicateMacro
93             | setPredicate
94             | truthValue
95             | parentheses
96             | negation ;
97
98 unaryPredicate   : unaryPredicateName
99                 LP task RP ;
100
101 binaryPredicate  : binaryPredicateName
102                 LP task SEP task RP ;
103
104 unaryPredicateMacro : unaryPredicateName
105                   LP setOfTasks SEP anyAll RP ;
106
107 binaryPredicateMacro:
108   binaryPredicateMacroTaskSet
109   | binaryPredicateMacroSetSet ;
110
111 binaryPredicateMacroTaskSet :
112   binaryPredicateName LP task
113   SEP setOfTasks SEP anyAll RP ;
114
115 binaryPredicateMacroSetSet :
116   binaryPredicateName
117   LP setOfTasks SEP setOfTasks
118   SEP anyEachAll RP ;
119
120 anyEachAll      : ANY | EACH | ALL ;
121
122 setPredicate    : taskInSetOfTasks
123                 | setComparison ;
124
125 taskInSetOfTasks : task IN setOfTasks ;
126
127 setComparison   : setOfTasks
128                 setComparisonOperator
129                 setOfTasks ;
130
131 setComparisonOperator : identical
132                       | different
133                       | overlapsWith
134                       | subsetOf
135                       | properSubsetOf ;
136
137 truthValue      : TRUE
138                 | FALSE
139                 | UNKNOWN ;
140
141 logicalTest     : isTrue
142                 | isNotTrue
143                 | isFalse
144                 | isNotFalse
145                 | isUnknown
146                 | isNotUnknown ;
147
148 union           : (tasks | difference |
149                 intersection) UNION (tasks |
150                 difference | intersection)
151                 (UNION (tasks | difference
152                 | intersection))* ;
153
154 intersection   : (tasks | difference)
155                 INTERSECTION
156                 (tasks | difference)
157                 (INTERSECTION (tasks
158                 | difference))* ;

```


- Christel Baier and Joost-Pieter Katoen. 2008. *Principles of model checking*. MIT Press.
- D. D. Chamberlin and R. F. Boyce. 1974. SEQUEL: A Structured English Query Language. In *ACM SIGMOD*. 249–264.
- C.J. Date and H. Darwen. 1997. *A Guide to the SQL Standard: A User's Guide to the Standard Database Language SQL*. Addison-Wesley.
- Javier Esparza and Mogens Nielsen. 1994. Decidability Issues for Petri Nets - A survey. *Bulletin of the EATCS* 52 (1994), 244–262.
- John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. 2006. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- Bertrand Meyer. 1990. *Introduction to the Theory of Programming Languages*. Prentice-Hall.
- Terence John Parr. 2013. *The Definitive ANTLR 4 Reference*. Pragmatic Programmers, LLC.
- Terence John Parr and Russell W. Quong. 1995. ANTLR: A Predicated-*LL(k)* Parser Generator. *Software — Practice and Experience (SPE)* 25, 7 (1995), 789–810.
- URI Planning Interest Group. 2001. *URIs, URLs, and URNs: Clarifications and Recommendations 1.0*. Technical Report. W3C. <http://www.w3.org/TR/uri-clarification/>
- W3C XSL/XML Query Working Groups. 2007. The XPath 2.0 Standard. (2007). <http://www.network-theory.co.uk/w3c/xpath/>

Received February 2007; revised March 2009; accepted June 2009