



Budapest University of Technology and Economics
Faculty of Electrical Engineering and Informatics
Department of Measurement and Information Systems

Superscalable Modeling

MASTER'S THESIS

Author

Gábor Szárnyas

Supervisors

Dr. István Ráth

Benedek Izsó

Dr. Dániel Varró

March 4, 2014

Kivonat

A lekérdezések központi szerepet játszanak az adatvezérelt alkalmazásokban. A modellvezérelt szoftvertervezés (model-driven engineering, MDE) eszközei és transzformációi erősen támaszkodnak a modell-lekérdezések hatékony kiértékelésére. A szoftvermodellek mérete és komplexitása intenzíven nő, ezért a jelenlegi MDE eszközökkel gyakran komoly skálázhatósági problémák merülnek fel, amelyek csökkentik a fejlesztés hatékonyságát és növelik annak költségeit.

A skálázhatósági kérdések központi témája az adatbázis-kezelés területén végzett kutatásoknak. A NoSQL rendszerek részben megoldást kínálnak több kapcsolódó problémára, de cserébe le kell mondanunk az SQL rendszerek által biztosított deklaratív ad-hoc lekérdezések erejéről. A NoSQL rendszerek modellvezérelt alkalmazásokban történő közvetlen alkalmazása jelenleg is nyitott kutatási kérdés, az ezekben futtatott lekérdezések ugyanis jelentősen bonyolultabbak, mint az általános adatbázis-kezelő alkalmazásokban használtak.

Diplomatervem célja, hogy az EMF-INCQUERY-ben alkalmazott inkrementális gráfmintaillesztő technikákat elosztott, felhőalapú infrastruktúrára adaptáljam. Bemutatok egy olyan újszerű architektúrát, amely elosztott, skálázható módon alkalmas lekérdezések inkrementális kiértékelésére. Az architektúra prototípusa, az INCQUERY-D rendszer képes egyetlen számítógéptől egy többgépes fürtig skálázódni, így képes nagy modelleken komplex lekérdezések hatékony futtatására. Az INCQUERY-D további előnye, hogy a lekérdezőmotor független a mögöttes adatbázis adatmodelljétől.

Az elképzelés működőképességét mérési eredményekkel igazoltam egy RDF- és egy gráfolapú adatbázis rendszerrel. Az eredmények bizonyítják, hogy az inkrementális lekérdezési technikák képesek hatékonyan működni elosztott környezetben is.

Abstract

Queries are the foundations of data intensive applications. In model-driven software engineering (MDE), model queries are core technologies of tools and transformations. As software models are rapidly increasing in size and complexity, traditional MDE tools frequently exhibit scalability issues that decrease productivity and increase costs.

While such scalability challenges are a constantly hot topic in the database community and recent efforts of the NoSQL movement have partially addressed many shortcomings, this happened at the cost of sacrificing the powerful declarative ad-hoc query capabilities of SQL. Unfortunately, this is a critical problem for MDE applications, as their queries can be significantly more complex than in general database applications. The applicability of NoSQL databases in MDE applications is subject for future research.

In my thesis, I aim to address this challenge by adapting incremental graph search techniques, known from the EMF-INCQUERY framework, to a distributed cloud infrastructure. I present a novel architecture for distributed, scalable incremental query evaluation. INCQUERY-D, the prototype system can scale up from a single node to a cluster of nodes that can handle very large models and complex queries efficiently. INCQUERY-D is a backend-agnostic system, meaning that its query engine is independent from the data model of the underlying database.

The feasibility of the approach is supported by early experimental results with both an RDF and a graph database backend. The results prove that incremental query evaluation techniques can work efficiently in a distributed environment as well.

Contents

1	Introduction	10
1.1	Context	10
1.2	Problem Statement and Requirements	10
1.3	Objectives and Contributions	11
1.4	Structure of the Thesis	12
2	Background Technologies	13
2.1	Big Data and the NoSQL Movement	13
2.2	Concepts	14
2.2.1	Graph Data Models	14
2.2.2	Sharding	18
2.2.3	Query Languages and Evaluation Strategies	19
2.3	Graph Storage Technologies	20
2.3.1	EMF Technologies	20
2.3.2	Neo4j	22
2.3.3	Titan	23
2.3.4	4store	26
2.3.5	Overview and Evaluation of Graph Storage Technologies . .	27
2.4	Building Scalable Asynchronous Distributed Systems: Akka	27

3	Overview of the Approach	29
3.1	Incremental Query Evaluation	29
3.2	The Rete Algorithm	30
3.2.1	Incremental Pattern Matching Algorithms	30
3.2.2	Overview of the Rete Algorithm	30
3.2.3	Concepts in the Rete Algorithm	32
3.2.4	Alpha nodes	33
3.2.5	Beta nodes	35
3.3	Incremental Pattern Matching on a Single Workstation: EMF-INCQUERY	38
3.3.1	Architecture	38
3.3.2	Indexing and Initialization	39
3.3.3	Data Representation and Storage	39
3.3.4	Notification Mechanisms	39
3.3.5	Termination Protocol	40
3.3.6	Configuration and Performance Optimization	40
3.4	Extensions for Distributed Scalability: INCQUERY-D	40
3.4.1	Architecture	40
3.4.2	Indexing and Initialization	41
3.4.3	Data Representation and Storage	41
3.4.4	Notification Mechanisms	42
3.4.5	Termination Protocol	42
3.4.6	Configuration and Performance Optimization	43
3.5	Workflow	44
3.5.1	Workflow of EMF-INCQUERY	44
3.5.2	Workflow of INCQUERY-D	45
3.6	Tooling for INCQUERY-D	45
3.7	Elaboration of the Example	46

3.7.1	Case Study: Railroad System Design	47
3.7.2	Local Search-based Query Evaluation	49
3.7.3	Workflow of the Example	50
4	Evaluation of Performance and Scalability	54
4.1	Dimensions of Scalability	54
4.2	Foundations: the Train Benchmark	55
4.2.1	Benchmark Goals	55
4.2.2	Generating Instance Models	56
4.2.3	Original Results for Non-distributed Tools	56
4.3	Distributed Train Benchmark	57
4.3.1	Distributed Architecture	57
4.3.2	Benchmark Limitations	58
4.3.3	Generating Instance Models	58
4.4	Benchmark Environment	58
4.4.1	Benchmark Setup	59
4.4.2	Hardware and Software Ecosystem	59
4.4.3	Benchmark Methodology and Data Processing	60
4.5	Benchmark Results with Neo4j	60
4.6	Benchmark Results with 4store and Titan	61
4.7	Result Analysis	64
4.8	Threats to Validity	65
4.9	Summary	66
5	Related Work	67
5.1	Eclipse-based Tools	67
5.2	Rete Implementations	67
5.3	Benchmarks	68
5.3.1	RDF Benchmarks	68
5.3.2	Model Transformation and Graph Transformation Benchmarks	69

6	Conclusions	70
6.1	Summary of Contributions	70
6.1.1	Scientific Contributions	70
6.1.2	Practical Accomplishments	71
6.2	Limitations and Future Work	72
	Bibliography	76
A	Graph Formats	84
A.1	Property Graph Formats	84
A.1.1	GraphML	84
A.1.2	Blueprints GraphSON	85
A.1.3	Faunus GraphSON	86
A.2	Semantic Graph Formats	87
A.2.1	RDF/XML	87
A.3	Mapping Ecore to Property Graphs	88
A.4	Train Benchmark Queries	90
A.4.1	Relational Schemas	90
A.4.2	Graph Patterns	91
A.4.3	PosLength	91
A.4.4	RouteSensor	92
A.4.5	SignalNeighbor	92
A.4.6	SwitchSensor	94

HALLGATÓI NYILATKOZAT

Alulírott *Szárnyas Gábor*, szigorló hallgató kijelentem, hogy ezt a diplomatervet meg nem engedett segítség nélkül, saját magam készítettem, csak a megadott forrásokat (szakirodalom, eszközök stb.) használtam fel. Minden olyan részt, melyet szó szerint, vagy azonos értelemben, de átfogalmazva más forrásból átvettem, egyértelműen, a forrás megadásával megjelöltem.

Hozzájárulok, hogy a jelen munkám alapadatait (szerző(k), cím, angol és magyar nyelvű tartalmi kivonat, készítés éve, konzulens(ek) neve) a BME VIK nyilvánosan hozzáférhető elektronikus formában, a munka teljes szövegét pedig az egyetem belső hálózatán keresztül (vagy autentikált felhasználók számára) közzétegye. Kijelentem, hogy a benyújtott munka és annak elektronikus verziója megegyezik. Dékáni engedéllyel titkosított diplomatervek esetén a dolgozat szövege csak 3 év eltelte után válik hozzáférhetővé.

Budapest, 2014. március 4.

Szárnyas Gábor
hallgató

Chapter 1

Introduction

1.1 Context

Model-driven software engineering (MDE) plays an important role in the development processes of critical embedded systems¹. Advanced modeling tools provide support for a wide range of development tasks such as requirements and traceability management, system modeling, early design validation, automated code generation, model-based testing and other validation and verification tasks.

Models representing sensor data, reverse engineered software models (e.g. abstract syntax trees of existing source code) and geospatial models can contain well over 10^9 modeling elements [72]. The dramatic increase in complexity is also affecting critical embedded systems in recent years. Modeling toolchains are facing scalability challenges as the size of design models constantly increases, and automated tool features become more sophisticated.

1.2 Problem Statement and Requirements

Many scalability issues can be addressed by improving query performance. *Incremental evaluation* of model queries aims to reduce query response time by limiting the impact of model modifications to query result calculation. Such algorithms work by either (i) building a cache of interim query results and keeping it up-to-date as models change (e.g. EMF-INCQUERY [36]) or (ii) applying impact analysis techniques and reevaluating queries only in contexts that are affected by

¹Section 1.1 and Section 1.2 are based on our earlier publication [55]. This thesis builds on the foundations of [55] and significantly extends them in both theory and practice, as detailed at the end of Section 6.1.

a change (e.g. the Eclipse OCL Impact Analyzer [48]). This technique has been proven to improve performance dramatically in several scenarios (e.g. on-the-fly well-formedness validation or model synchronization), at the cost of increasing memory consumption. Unfortunately, this overhead is combined with the increase in model sizes due to in-memory representation (found in state-of-the-art frameworks such as EMF [76]).

In practice, these scalability problems imply that the users are forced to either split the model into smaller models and handle them separately or cope with the slow query and transformation performance. Both reduce the productivity advantages offered by model-driven design.

A trivial solution would be to increase the amount of available memory. However, the Garbage Collector (GC) used in most modern programming languages, including Java, cannot handle heap sizes larger than 10 GB efficiently, thus introducing long pauses in the application [26]. Of course, this problem is well-known in the Java community. There are alternative Java Virtual Machines (JVMs) with specialized Garbage Collectors, like Azul Systems' JVM. However, the Azul JVM is a proprietary product and has specific hardware requirements. Also, this does not solve the scaling problem entirely – the model size is still limited by the total amount of memory in a single computer.

An alternative approach to tackling MDE scalability issues is to make use of advances in persistence technology. As the majority of model-based tools uses a graph-oriented data model, recent results of the NoSQL and Linked Data movement [65, 1, 2] are straightforward candidates for adaptation to MDE purposes. Unfortunately, this idea poses difficult conceptual and technological challenges: (i) property graph databases lack strong metamodeling support and their query features are simplistic compared to MDE needs, and (ii) the underlying data representation format of semantic databases (RDF [49]) has crucial conceptual and technological differences to traditional metamodeling languages such as Ecore [76]. Additionally, while there are initial efforts to overcome the mapping issues between the MDE and Linked Data worlds [53], even the most sophisticated NoSQL storage technologies lack efficient and mature support for executing expressive queries *incrementally*.

1.3 Objectives and Contributions

We aimed to address the scalability challenge of MDE by adapting incremental graph search techniques from EMF-INCQUERY to the cloud infrastructure—

instead of *vertical scaling* (putting more resources in the same workstation), we decided to opt for *horizontal scaling* (using multiple computers).

On the theoretical side, we adopted EMF-INCQUERY’s existing incremental pattern matching algorithm, the Rete algorithm. We extended the algorithm to work in a distributed environment and designed a *novel architecture*, which is capable of loading, transforming and incrementally querying models, while utilizing the *total amount of memory* in the cluster.

To build a scalable incremental query engine, we needed a distributed software stack. This included a distributed database management system and a messaging framework. We defined the evaluation criteria for these systems and evaluated them accordingly. Based on the architecture and the pattern matcher algorithm, we built a system prototype and compared its performance to existing tools.

We extended an existing *benchmark environment* to evaluate the scalability characteristics of the system and conducted benchmarks with different storage backends and query engines.

1.4 Structure of the Thesis

The thesis is structured as follows. Chapter 2 introduces the background technologies and the motivation for building a distributed, incremental graph pattern matcher. Chapter 3 provides an overview of a single-node incremental pattern matcher, EMF-INCQUERY, and details the architecture of a distributed approach, INCQUERY-D. Chapter 4 shows an initial performance evaluation in the context of on-the-fly well-formedness validation of software design models. Chapter 5 discusses the related work. Chapter 6 concludes the thesis and presents our future plans.

Chapter 2

Background Technologies

Developing a scalable graph pattern matcher requires a wide range of technologies. Careful selection of the technologies is critical to the success of the project. For INCQUERY-D, we looked for technologies that can form the building blocks of a distributed, scalable model repository and pattern matcher. The potential technologies are the ones that were designed with scalability in mind and have been deployed in large-scale distributed systems successfully.

Usually, instance models are graph-like data structures. Therefore, we looked for scalable graph databases. In this context, scalability requires distributed storage and querying capabilities.

During the early phase of our research, we studied the architecture and limitations of the candidate systems. For databases, we inspected the data sharding strategies, consistency guarantees and transaction capabilities, along with the API and query methods. We also checked the support for asynchronous processing, notification and messaging mechanisms.

In this chapter, we introduce the concepts and technologies that can form the basis of a scalable, distributed, asynchronous system.

2.1 Big Data and the NoSQL Movement

Since the 1980s, database management systems based on the relational data model [38] dominated the database market. Relational databases have a number of important advantages: precise mathematical background, understandability, mature tooling and so on. However, due to their rich feature set and the strongly connected nature of their data model, relational databases often have scalability

issues [59, 75]. They are typically optimized for transaction processing, instead of data analysis (see *data warehouses* for an exception). In practice, these render them impractical for a number of use cases, e.g. running complex queries on large data sets.

In the last decade, large organizations struggled to store and process the huge amounts of data they produced. This problem introduces a diverse palette of scientific and engineering challenges, called *Big Data* challenges.

Big Data challenges spawned dozens of new database management systems. Typically, these systems broke with the strictness of the relational data model and utilized simpler, more scalable data models. These systems dropped support for the SQL query language used in relational databases and hence were called *NoSQL databases*¹ [19]. Because relational databases are not suitable for large-scale model-driven applications, we experimented with numerous NoSQL databases.

2.2 Concepts

This section introduces the most important concepts used in this thesis.

2.2.1 Graph Data Models

Along the well-known and widely used relational data model, there are many other data models. NoSQL databases are often categorized based on their data model (e.g. key-value stores, document stores, column families). In this thesis, we focus on *graph data models*.

The graph is a well-known mathematical concept widely used in computer science. For our work, it is important to distinguish between different graph data models.

The most basic graph model is the *simple graph*, formally defined as $G = (V, E)$, where V is the set of vertices and $E \subseteq V \times V$ is the set of edges. Simple graphs are sometimes referred as textbook-style graphs because they are an integral part of academic literature. Simple graphs are useful for modeling homogeneous systems and have plenty of algorithms for processing.

Simple graphs can be extended in several different ways (Figure 2.1). To describe the connections in more detail, we may add directionality to edges (*directed graph*). To allow different connections, we may label the edges (*labeled graph*).

¹The community now mostly interprets NoSQL as “not only SQL”.

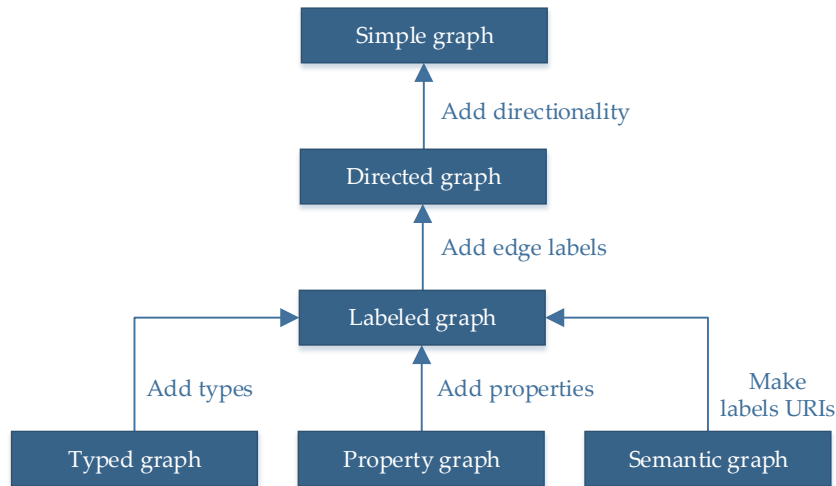


Figure 2.1. *Different graph data models (based on [71])*

Typed graphs introduce types for vertices. *Property graphs* (sometimes called *attributed graphs*) add even more possibilities by introducing properties. Each graph element, both vertices and edges can be described with a collection of properties. The properties are key–value pairs, e.g. `type = 'Person', name = 'John', age = 34`. *Semantic graphs* use URIs (Uniform Resource Identifiers) instead of labels, otherwise they have similar expressive power as labeled graphs.

Graph models are present in many languages and environments. In the following, we will present the ones most important for this thesis: the Ecore metamodeling language, the TinkerPop framework and the Resource Description Framework (RDF).

Metamodeling

Metamodeling is a methodology for the definition of modeling languages. A metamodel specifies the abstract syntax (structure) of a modeling language. Metamodels are expressed using a metamodeling language that itself is a modeling language. The metamodel can also be interpreted as the object-oriented data model of the language under design. Metamodeling can be viewed as the grammar for a *typed property graph*, so the created models are both *typed graphs* and *property graphs*.

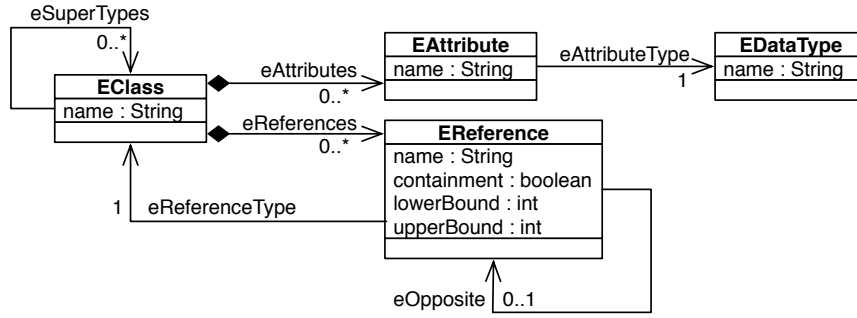


Figure 2.2. The Ecore kernel, a simplified subset of the Ecore metamodel

Ecore

Ecore is the metamodeling language used by EMF. It has been developed in order to provide an approach for metamodel definition that supports the direct implementation of models using a programming language. The main rationale in introducing Ecore separately is that it is the *de facto* standard metamodeling environment of the industry, and several domain-specific languages are defined using this formalism.

Figure 2.2 illustrates the core elements of the Ecore approach. The full metamodel can be found in the EMF documentation [40]. The most important elements are the following.

- **EClass** models classes (or concepts). **EClasses** are identified by name and can have several attributes and references. To support inheritance, a class can refer to a number of *supertype* classes.
- **EAttribute** models attributes, that contain data elements of a class. They are identified by name and have a *data type*.
- **EDataType** is used to represent simple data types that are treated as atomic (their internal structure is not modeled). Data types are also identified by their name.
- **EReference** represents a unidirectional association between **EClasses** and is identified by a name. Lower and upper multiplicities can be specified. It is also possible to mark a reference as a *containment* that represents composition relation between elements. If a bidirectional association is needed, it should be modeled as two **EReference** instances that are mutually connected via their *opposite* references.

The rest of the details of Ecore has been omitted for the sake of clarity, see [40] for further reference.

TinkerPop framework

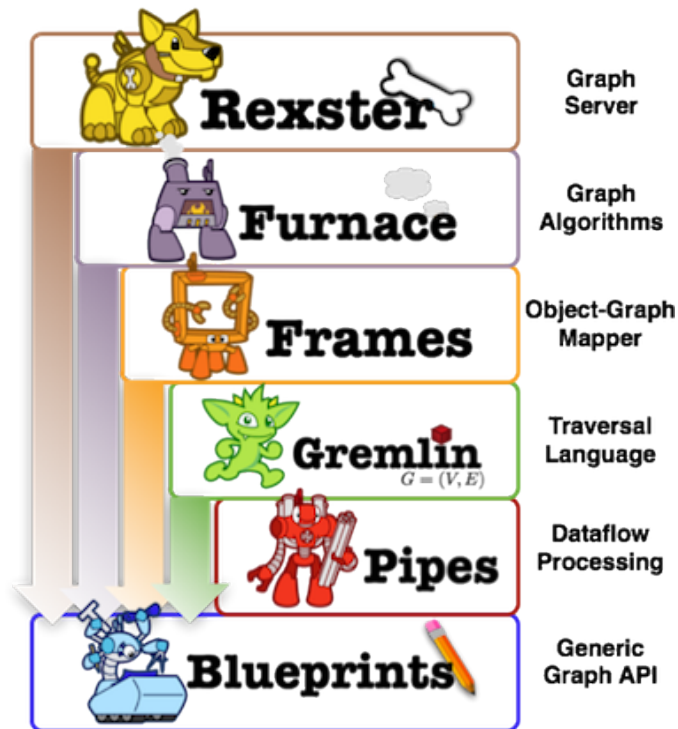


Figure 2.3. The TinkerPop software stack [11]

The *TinkerPop* framework is an open-source software stack for graph storage and processing [28]. TinkerPop includes *Blueprints*, a property graph model interface. *Blueprints* intends to fulfill the same role for graph databases as JDBC does for relational databases. Most NoSQL graph databases implement the property graph interface provided by *Blueprints*, including Neo4j (Section 2.3.2), Titan (Section 2.3.3), DEX [24], InfiniteGraph [20] and OrientDB [21].

TinkerPop also introduces a graph query language, *Gremlin*. *Gremlin* is a domain-specific language based on Groovy, a Java-like dynamic language which runs on the Java Virtual Machine. Unlike most query languages, *Gremlin* is an imperative language with a strong focus on graph traversals.

Gremlin is based on *Pipes*, TinkerPop's dataflow processing framework. Besides traversing, *Gremlin* is capable of analyzing and manipulating the graph as well.

TinkerPop also provides a graph server (*Rexster*), a set of graph algorithms tailored for property graphs (*Furnace*) and an object-graph mapper (*Frames*). The TinkerPop software stack is shown in Figure 2.3.

Resource Description Framework

The Resource Description Framework (RDF) is a family of W3C (World Wide Web Consortium) specifications originally designed as a *metadata data model*.

The RDF data model is based on the idea of making statements about *resources* in the form of triples. A triple is a data entity composed of a *subject*, a *predicate* and an *object*, e.g. "John instanceof Person", "John is 34".

Triples are typically stored in *triplestores*, specialized databases tailored to store and process triples efficiently. Also, some triplestores are capable of *reasoning*, i.e. inferring logical consequences from a set of facts or axioms. Triplestores are mostly used in semantic technology projects.

Triplestores are usually queried via the RDF format's query language, SPARQL (recursive acronym for SPARQL Protocol and RDF Query Language).

The RDF data model is capable of expressing *semantic graphs*. Although the semantic graph data model has less expressive power than the property graph data model, by introducing additional resources for each property, a property graph can be easily mapped to RDF.

Mapping Ecore to Other Data Models

Our intention to reuse EMF-INCQUERY for building INCQUERY-D required us to map EMF's metamodel, Ecore to the domain of property graphs and RDF models.

Ecore concept	Property graph concept	RDF concept
EClass instance	nodes' type property	rdfs:Resource
EAttribute instance	nodes' property names	rdf:Property
EReference instance	edge label	rdf:Property
EDataType instance	Java primitive types	rdfs:Datatype

Table 2.1. Mapping Ecore to property graphs and RDF

Table 2.1 shows a mapping from the concepts defined in Ecore to other modelling domains.

2.2.2 Sharding

To provide scalable persistence and processing for large amounts of data, the data has to be split between multiple computers. This process is known as *data*

sharding. *Graph sharding* is a particularly difficult problem due to the strongly connected and mutable nature of graphs. Efficient sharding of graphs is still an open research area [60].

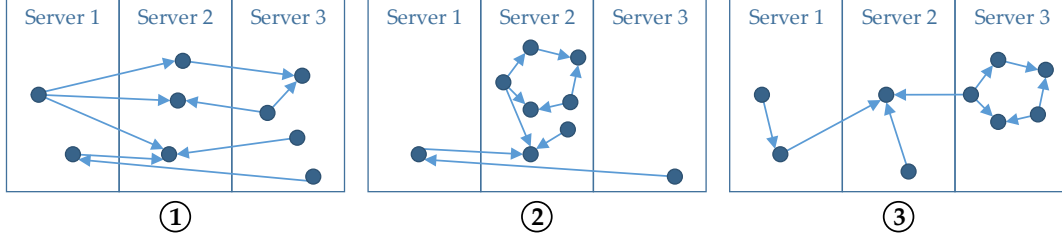


Figure 2.4. *Different partitionings of the same graph*

To illustrate the problem, Figure 2.4 shows different partitionings of the same graph in a three-node cluster. In case ① most edges run between the servers and are therefore expensive to traverse. In case ②, *Server 2* is overloaded, taking more than three quarters of the total load. Case ③ presents a more balanced sharding of the graph. Unfortunately, for large graphs, balanced sharding is hard to achieve in practice.

Most graph partition problems are NP-hard and practical solutions to these problems could be derived using heuristics and approximation algorithms [44]. Unfortunately, open-source database implementations lack support for such algorithms.

2.2.3 Query Languages and Evaluation Strategies

In the context of this thesis, a query defines a *graph pattern*. The result of the query is a *set of subgraphs* of the original graph. Graph patterns are useful for identifying patterns in a set of connected data elements. They are especially widely used in the context of model-driven engineering for formulating well-formedness validation constraints (usually by defining patterns that violate these constraints) and graph transformations.

Query Language

Queries can be defined in both imperative and declarative languages. The theoretical basis for most declarative query languages is first order logic. Both tuple relational calculus and relational algebra (widely used in query processing) are offshoots of first-order logic.

As mentioned earlier, unlike most query languages, Gremlin is an imperative query language describing graph traversals. For example, if John's father is Jack and Jack's father is Scott, we may run the traversals shown on Listing 2.1.

```
1 gremlin> g.V('name', 'John').out('father')
2 ==>Jack
3 gremlin> g.V('name', 'John').out('father').out('father')
4 ==>Scott
```

Listing 2.1. *Simple Gremlin queries*

Query Evaluation Strategies

Query engines can be divided into two core categories: *search-based* and *incremental* engines. The main difference between these approaches is the way they reevaluate queries. While search-based engines process the whole data set (i.e. not just the data elements affected by the change), incremental engines utilize some data structures to be able to reevaluate the query based on the change set.

2.3 Graph Storage Technologies

In this section, we compare different graph storage technologies by systematically discussing their *architecture* and *data model*. We inspect their *sharding* strategies for distributed storage. We also present their *query languages* and *evaluation strategies*, with particular emphasis on the support of distributed operations.

2.3.1 EMF Technologies

Eclipse is a free, open-source software development environment and a platform for plug-in development. Eclipse comes with its own modeling technologies called EMF (Eclipse Modeling Framework). EMF's primary goals are application design and code generation.

Architecture

EMF models can be persisted as XMI (XML Metadata Interchange) documents. By design, EMF models cannot be fragmented, i.e. they can only be used if they fit to a computer's main memory. There are different model repositories and persistence frameworks which can handle large EMF models [73].

- CDO (Connected Data Objects), a distributed shared model framework for EMF models and metamodels [42]. CDO provides an object-relational mapping from Ecore to databases.
- Morsa [66] is a distributed model repository based on MongoDB [17], a popular NoSQL database management system.
- Neo4EMF [30] provides a persistence framework based on Neo4j. It implements an on-demand loading approach which separates EObjects and their data fields, only instantiating the latter if they are accessed by a query.

Data Model

EMF uses the Ecore data model, discussed in Section 2.2.1.

Sharding

Due to the nature of XML documents, EMF models serialized to a single XMI document cannot be sharded. CDO and Neo4EMF do not support automatic sharding, however Morsa does so by using MongoDB’s sharding mechanism.

Query Language and Evaluation

OCL OCL (Object Constraint Language) is a declarative query language to describe well-formedness constraints on UML models. These expressions typically specify invariant conditions that must hold for the system being modeled or queries over objects described in a model.

EMF-INCQUERY EMF-INCQUERY [35] is an Eclipse project developed by the Fault Tolerant Systems Research Group in the Budapest University of Technology and Economics. It provides IQPL (INCQUERY Pattern Language), a declarative language to express queries over EMF models in the form of graph patterns. With the language the user can express combined queries, negative patterns, property conditions, simple calculations, calculate disjunctions and transitive closures, etc. on top of the models. The goal of EMF-INCQUERY is to provide *incremental query evaluation*.

Our research work builds on EMF-INCQUERY, both in theory and practice. We used the Rete algorithm (Section 3.2) which allowed us to reuse some of the existing code base. We also utilized the methodology and environment, originally used to benchmark EMF-INCQUERY (Section 4.2).

2.3.2 Neo4j

Neo4j, developed by Neo Technology, is the most popular NoSQL graph database. Neo4j is one of the most mature NoSQL databases. It is well documented and provides ample tooling, including an Eclipse-based visualization application, Neoclipse [18].

Architecture

Neo4j can be deployed in two scenarios. In *embedded mode*, it runs in the same JVM (Java Virtual Machine), as the client application. In this setup, the database cannot be accessed by other applications. In *server mode*, the database can serve requests from multiple clients over a REST (Representational State Transfer) interface.

Data Model

Neo4j implements the TinkerPop framework's Blueprints property graph data model. Neo4j is capable of loading graphs from GraphML [25] and Blueprints GraphSON [13] formats (see Section A.1 for examples).

Sharding

Instead of sharding, Neo4j only supports replication of data to create a highly available cluster. This implies serious scalability limitations to the system. As of December 2013, the development team of Neo4j is currently working on improving the scalability of the database in an ongoing project called Rassilon [3].

Query Language and Evaluation

Neo4j can be queried in various ways. When deployed in embedded mode, the application can use its Java-based core API. In both embedded and server mode,

Neo4j provides two query languages. The first is the TinkerPop framework's imperative language, Gremlin, primarily targeted for graph traversals. The second is Neo4j's own declarative query language for graph pattern matching, Cypher.

2.3.3 Titan

Titan is a distributed, scalable graph database from Aurelius, the creators of the TinkerPop framework. To understand Titan's complex architecture, we present two additional concepts: the *MapReduce* paradigm and the *column family* data model.

Asynchronous Parallel Processing with MapReduce The *MapReduce paradigm* defines a parallel, asynchronous way of processing the data. As the name implies, MapReduce consists of two phases: the *map* function processes each item of a list. The resulted list is then aggregated by the *reduce* function. MapReduce is often used for sorting, filtering and aggregating data sets. It is also used for fault-tolerant, distributed task execution.

The Column Family Data Model A column family is similar to a table of a relational database: it consists of rows and columns. However, unlike in a relational database's table, the rows do not have to have the same fixed set of columns. Instead, each row can have a different set of columns. This makes the data structure more dynamic and avoids the problems associated with NULL values.

Architecture

Titan is not a standalone database, instead, it builds on top of existing NoSQL database technologies and leverages Hadoop's MapReduce capabilities. Titan supports various storage backends, including Cassandra and HBase. In the following, we shortly cover the technologies Titan builds upon. Both Titan and its dependencies are open-source software, written in Java.

Hadoop Hadoop is a distributed data processing framework inspired by Google's publications about MapReduce [39] and the Google File System [47]. Originally developed at Yahoo!, Hadoop is now an Apache project [7]. Like Google's systems, Hadoop is designed to run on commodity hardware, i.e. server

clusters built from commercial off-the-shelf products. Hadoop provides a distributed file system (HDFS) and a column family database (HBase). A typical Hadoop cluster consists of a single master node which is responsible for the coordination of the cluster and worker nodes which deal with the data processing. The MapReduce job is coordinated by the master's *job tracker* and processed by the slave nodes' *task tracker* modules (Figure 2.5).

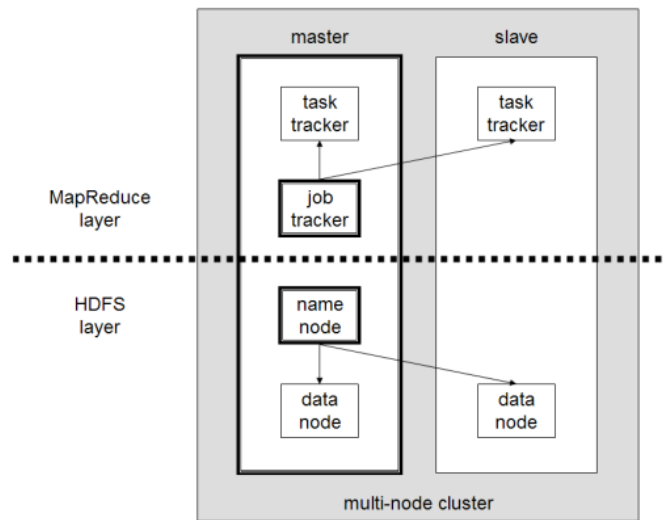


Figure 2.5. *Hadoop's architecture* [62]

HDFS The Hadoop Distributed File System (HDFS) is a distributed file system, inspired by the Google File System and written specifically for Hadoop [7]. Unlike other distributed file systems (e.g. Lustre [16]), which require expensive hardware components, HDFS was designed to run on commodity hardware. HDFS tightly integrates with Hadoop's architecture (Figure 2.5). The *name node* is responsible for storing the metadata of the files and the location of the replicas. The data is stored by the *data nodes*.

HBase HBase [8] is a distributed column family database. It is developed as part of the Hadoop project and runs on top of HDFS. The tables in an HBase database can serve as the input and the output for MapReduce jobs running in Hadoop.

Cassandra Cassandra is one of the most widely used NoSQL databases [6]. Originally developed by Facebook [61], Cassandra is now an Apache project. Cassandra is a column family database with advanced fault-tolerance mechanisms. It allows the application to balance between availability and consistency

by providing tunable consistency constraints. Cassandra is used mainly by Web 2.0 companies, including Digg, Netflix, Reddit, SoundCloud and Twitter. It is also used for research purposes at CERN and NASA [22].

Data Model

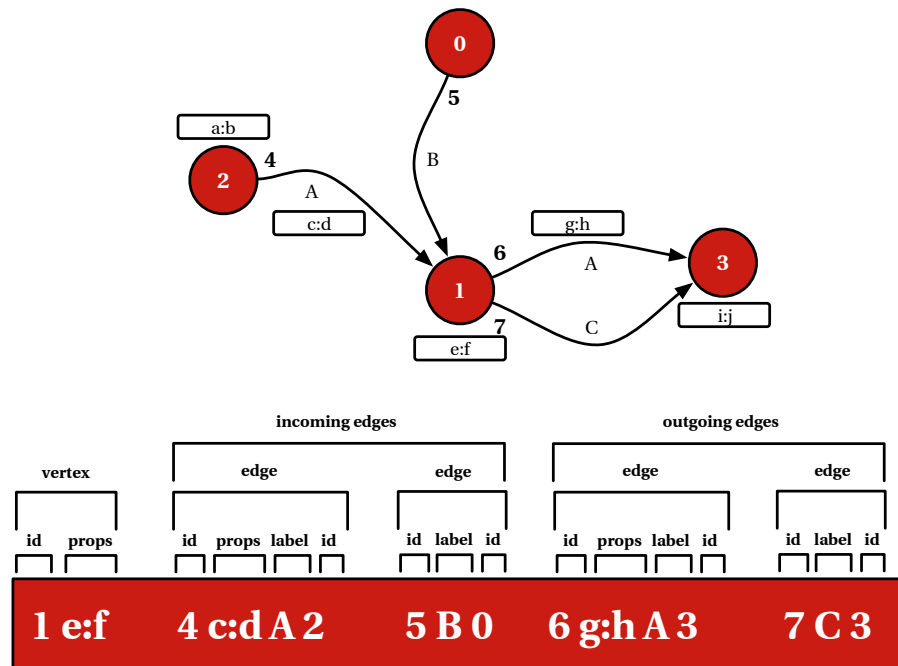


Figure 2.6. Graph vertex mapped by Titan to a row in a Cassandra database

To store the graph, Titan maps each vertex to a row of a column family (Figure 2.6). The row stores the identifier and the properties of the vertex, along both the incoming and outgoing edges' identifiers, labels and properties.

Sharding

Titan uses the storage backend's partitioner, e.g. Cassandra's hash-based RandomPartitioner to shard the data. A more sophisticated partitioning system that will allow for partitioning based on the graph's static and dynamic properties (its domain and connectivity, respectively) is under implementation as of December 2013, but not yet available.

Query Language and Evaluation

Titan supports the TinkerPop framework's Gremlin query language. This implementation, Gremlin/Pipes, utilizes a depth-first search algorithm.

Faunus Although Titan was designed with scalability in mind, its query engine is unable to cope with queries resulting in millions of graph elements. To address this shortcoming, Aurelius developed a Hadoop-based graph analytics engine, Faunus. Faunus has its own format called Faunus GraphSON. The Faunus GraphSON format is vertex-centric: each row represents a vertex of the graph. This way, Hadoop is able to efficiently split the input file and parallelize the load process. See Section A.1.3 for an example of the Faunus GraphSON format. Unlike the Gremlin implementation in Neo4j and Titan, the implementation in Faunus, Gremlin/Faunus, is based on breadth-first search. It is important to note that Faunus always traverses the whole graph and does not use its indices. This makes retrieving nodes or edges by type very slow (see our typical workload in Section 3.4.2).

2.3.4 4store

4store is an open-source, distributed triplestore created by Garlik [4]. Unlike the other tools discussed earlier, 4store is written in C. 4store is primarily applied for semantic web projects.

Architecture

4store was designed to work in a cluster with high-speed networks. 4store server instances are capable of discovering each other using the Avahi configuration protocol [10]. 4store offers a command-line and an HTTP server interface.

Data Model

4store's data model is an RDF graph. It supports the RDF/XML input format, which is processed using the Raptor RDF Syntax Library [68].

Sharding

Similarly to Titan's partitioning, the *segmenting* mechanism in 4store distributes the RDF resources evenly across the cluster. 4store also supports replication by *mirroring* tuples across the cluster.

Query Language and Evaluation

4store uses the Rasqal RDF Query Library [69] to supports SPARQL queries.

2.3.5 Overview and Evaluation of Graph Storage Technologies

Tech-nology	Data model	Distributed operation	Sharding	Queries	Identifier generation
EMF	Ecore	Differs	Differs	OCL, IQPL	Automatic
4store	RDF	Manual	Automatic	SPARQL	Manual
Neo4j	Property graph	Manual	Manual	Cypher	Manual
Titan	Property graph	Automatic	Automatic	Gremlin	Automatic

Table 2.2. *Overview of database technologies*

Table 2.2 summarizes the relevant characteristics of the aforementioned database management systems. These characteristics are crucial for building a distributed pattern matcher. According to these, Titan provides the most complete feature set. 4store and Neo4j lack important features like automatic identifier generation, which has to be implemented in the client application. Neo4j also misses automatic sharding, which seriously hinders its scalability potential. EMF’s distributed operation and sharding capabilities depend on the actual model repository and database backend being used.

2.4 Building Scalable Asynchronous Distributed Systems: Akka

Most distributed, concurrent systems use a messaging framework or message queue service. The INCQUERY-D system also requires a distributed, asynchronous messaging framework. For this purpose, we used the Akka framework.

Akka is an open-source, fault-tolerant, distributed, asynchronous messaging framework developed by Typesafe [5]. Akka is implemented in Scala, a functional and object-oriented programming language which runs on the Java Virtual Machine. Akka provides language bindings for both Java and Scala.

Akka is based on the actor model [52] and provides built-in support for remoting. Unlike traditional remoting solutions, e.g. Java RMI (Remote Method Invocation) and CORBA (Common Object Request Broker Architecture), the remote and local interface is the same for each actor. Actors have both a logical and a physical path

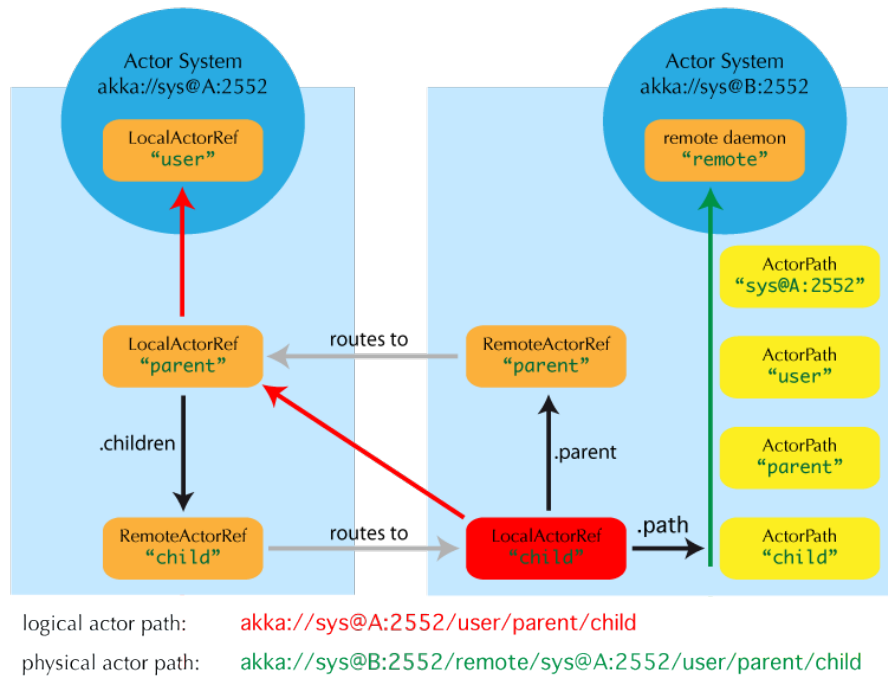


Figure 2.7. Deploying a remote actor in Akka [5]

(Figure 2.7). This way, they can be transparently moved between machines on the network.

As of December 2013, the latest version (Akka 2.2) also supports *pluggable transport* to use various transports to communicate with remote systems [5]. For serializing the messages, Akka supports different frameworks, including Java’s built-in serialization, Google Protobuf [23] and Apache Thrift [9].

Chapter 3

Overview of the Approach

The primary goal of INCQUERY-D is to provide a scalable architecture for executing incremental queries over large models. Our approach is based on the following foundations: (i) a distributed model storage system that (ii) supports a graph-oriented data representation format, and (iii) a graph query language adapted from the EMF-INCQUERY framework. The novel contribution of this thesis is an architecture that consists of a (i) distributed model management middleware, and a (ii) distributed and stateful pattern matcher network based on the Rete algorithm.

INCQUERY-D provides incremental query execution by *indexing model contents* and *capturing model manipulation operations* in the middleware layer, and *propagating change tokens* along the pattern matcher network to *produce query results* and *query result changes* (corresponding to model manipulation transactions) efficiently. As the primary sources of memory consumption, i.e. both the indexing and intermediate Rete nodes can be distributed in a cloud infrastructure, the system is expected to scale well beyond the limitations of the traditional single workstation setup.

3.1 Incremental Query Evaluation

Some queries, e.g. well-formedness constraints in MDE are evaluated many times, while the data set they are evaluated on only changes to a small degree. In these cases, the idea of incremental query evaluation arises naturally: to speed up queries, we should not start the evaluation all over again. Instead, we should rely on the (partial) results derived during the previous executions of the query and process only the changes that occurred.

In practice, incremental query evaluation algorithms typically use data structures for caching the interim results. This means that they usually consume more memory, in other words, they trade memory consumption for execution speed. This approach, called *space–time tradeoff*, is well-known and widely used in computer science.

3.2 The Rete Algorithm

In the following, we provide an overview of the *Rete algorithm*, which forms the theoretical basis of EMF-INCQUERY and INCQUERY-D.

3.2.1 Incremental Pattern Matching Algorithms

Numerous algorithms were invented for the purpose of incremental pattern matching. Mostly, these algorithms originate from the field of rule-based expert systems.

One of the most well-known is the *Rete algorithm*¹, which creates a propagation network. The network stores the partial matches found in the graph. The *TREAT algorithm* [63] aims at minimizing memory usage, while having the same algorithmic complexity as Rete. It stores only the input facts and the conflict sets, and does not store partial pattern matches. Another candidate is the *LEAPS algorithm* [31], which is claimed to provide better space–time complexity. However, we found that LEAPS is difficult to understand and implement even on a single workstation, not to mention the distributed case.

Rete has many improved versions (e.g. Rete II, Rete III, Rete-NT), however, unlike the original algorithm, these are not publicly available. Because the original Rete algorithm is well-understood by the EMF-INCQUERY team, we decided to build INCQUERY-D on the same foundation. Experimenting with improved versions or alternative approaches is subject to future work.

3.2.2 Overview of the Rete Algorithm

The Rete algorithm was originally created by Charles Forgy [45] for rule-based expert systems. Gábor Bergmann adapted the algorithm for EMF models and added many tweaks and improvements to it [32].

¹*Rete* is Latin for *net*.

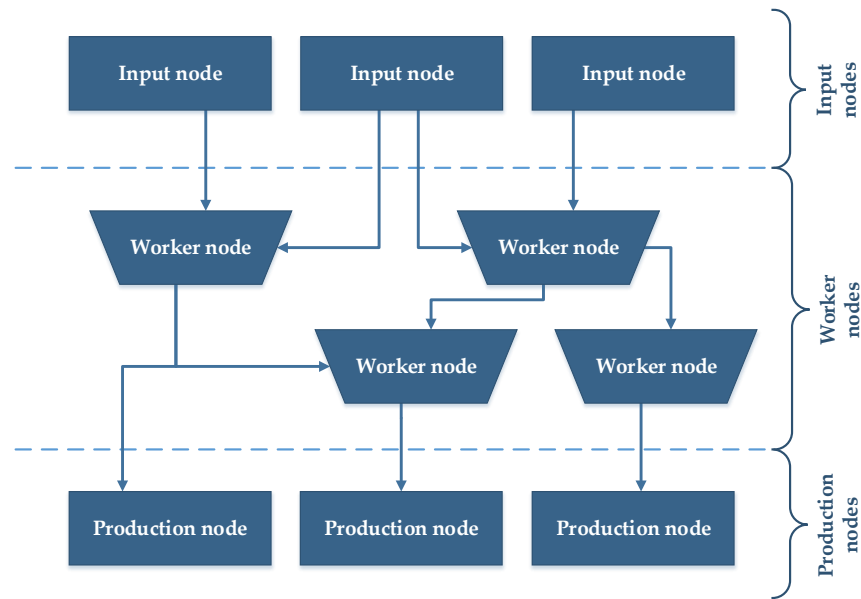


Figure 3.1. The structure of the Rete propagation network

The Rete algorithm defines an asynchronous network of communicating nodes (Figure 3.1). This is essentially a dataflow network. First, the network computes the set of pattern matches in the graph. The main feature of the algorithm is that it is capable of incrementally maintaining the match set by propagating *update messages* (also known as *deltas*, *notification objects* or *tokens*). Creating new graph elements (vertices or edges) results in *positive update messages*, while removing graph elements results in *negative update messages*.

The network consists of three types of nodes:

- *Input nodes* are responsible for indexing the model by type, i.e. they store the appropriate tuples for the vertices and edges. They are also responsible for spawning the update messages and propagating them to the *worker nodes*.
- *Worker nodes* perform a transformation on the output of their parent node(s) and propagate the results. The worker nodes store partial query results in their own memory (the details of this may vary depending on the implementation).
- *Production nodes* are terminators that provide an interface for fetching the queries' results (the match sets) and also the queries' change sets (the changes introduced by the latest update message).

3.2.3 Concepts in the Rete Algorithm

Here, we define the most important concepts in the Rete algorithm.

Data Representation

A *tuple* is an ordered list of elements, e.g. $\langle 15, 23, 81, 2 \rangle$, $\langle 7, 'red', true \rangle$ are tuples. The Rete algorithm defined in [32] uses tuples to represent the vertices, edges and matches of the (sub)patterns in the graph.

The items in a tuple are referenced by their index. Note that in our notation, the first element has an index of 0. A *pattern mask* is an array of indices, $\mu = (\mu_0, \mu_1, \dots, \mu_{n-1})$, which can be used to select certain elements in a tuple. For example, extracting the items defined by the pattern mask $\mu = (1, 3)$ from tuple $\langle 15, 23, 81, 2 \rangle$ results in the tuple $\langle 23, 2 \rangle$.

To use tuples for graph pattern matching, the vertices and edges in the graph have to be mapped to tuples. In the following, we present these mappings, using the graph in Figure 3.2 as an example. We presume that each vertex in the graph has a unique identifier.

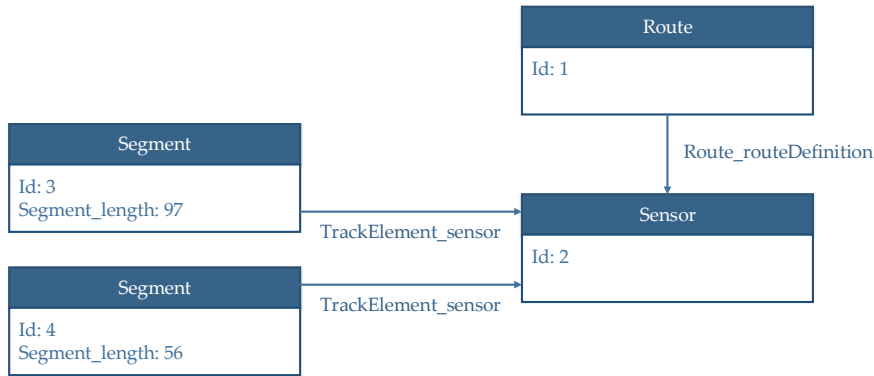


Figure 3.2. Example graph

Mapping Vertices to Tuples

A possible mapping of vertices is creating $\langle id, property\ value_0, property\ value_1, \dots \rangle$ tuples. It is a common practice to store the tuples in different *relations*, each containing a specific type:

- $Route = \{\langle 1 \rangle\}$

- $Sensor = \{\langle 2 \rangle\}$
- $Segment = \{\langle 3, 97 \rangle, \langle 3, 56 \rangle\}$

Here, the *Segment* relation consists of $\langle Id, Segment_length \rangle$ tuples. Note that the semantics of the property values have to be stored separately from the tuples. This is closely related to the concept of *relational schemas* in relational database design.

Mapping Edges to Tuples

Edges are mapped in a straightforward way: each (directed) edge is represented by a $\langle source\ vertex\ id, target\ vertex\ id \rangle$ pair. The edges on the example graph constitute the following tuples:

- $Route_routeDefinition = \{\langle 1, 2 \rangle\}$
- $TrackElement_sensor = \{\langle 3, 2 \rangle, \langle 4, 2 \rangle\}$

In the following, we present the most common Rete nodes. We define each node's operation in relational algebra and also formulate the operations which are performed upon receiving a positive or a negative update.

We use the conventional notations for the operations in relational algebra. These can be found in any computer science textbook discussing relational database management systems, e.g. [78].

3.2.4 Alpha nodes

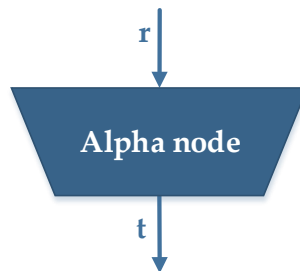


Figure 3.3. Graphical notation of an alpha node

Alpha nodes (Figure 3.3) have one input slot. They filter the content of the parent node according to some criteria. In the following, the relation representing the input tuples is denoted with r , the relation representing the output tuples is denoted with t . In general, an alpha node performs an operation on the input relation:

$$t = \alpha(r).$$

If the alpha node receives an update Δr , it performs the operation and computes the change set. For *positive updates*, the result (t') and the changeset (Δt) are:

$$t' = \alpha(r \cup \Delta r) = \alpha(r) \cup \alpha(\Delta r) = t \cup \alpha(\Delta r) \Rightarrow \Delta t = \alpha(\Delta r).$$

Similarly, for *negative updates*:

$$t' = \alpha(r \setminus \Delta r) = \alpha(r) \setminus \alpha(\Delta r) = t \setminus \alpha(\Delta r) \Rightarrow \Delta t = \alpha(\Delta r).$$

Trimmer Node

The *trimmer node* is similar to the projection operation defined in relational algebra. Specifically, if the pattern mask is denoted with λ , the result relation is

$$t = \pi_{R.\lambda_0, R.\lambda_1, \dots, R.\lambda_{n-1}}(r).$$

Equality Node

The *equality node* checks whether certain elements in the tuple, selected by a pattern mask μ are *all equal*:

$$t = \sigma_{(R.\mu_0=R.\mu_1) \wedge (R.\mu_0=R.\mu_2) \wedge \dots \wedge (R.\mu_0=R.\mu_{n-1})}(r).$$

Inequality Node

The *inequality node* works similarly to the equality node. However, it checks for tuples where the selected elements (specified by μ_1, μ_2, \dots) are *all different* from the subject element (specified by μ_0):

$$t = \sigma_{(R.\mu_0 \neq R.\mu_1) \wedge (R.\mu_0 \neq R.\mu_2) \wedge \dots \wedge (R.\mu_0 \neq R.\mu_{n-1})}(r).$$

Predicate Evaluator Node

The result of the *predicate evaluator node* is defined as

$$t = \sigma_{\theta}(r),$$

where the predicate θ is a propositional formula. The formula consist of atoms: binary operations ($=, \neq, \leq, <, >, \geq$) between attributes and constants. The atoms are connected by logical operators (\wedge, \vee and \neg). This selects all tuples in r for which θ holds.

3.2.5 Beta nodes

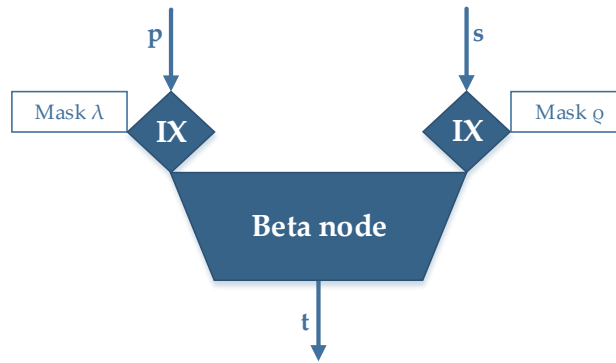


Figure 3.4. Graphical notation of a beta node

Beta nodes (Figure 3.4) have two input slots: the *primary* (p) and the *secondary* (s). Each input slot has its own pattern mask: λ for the primary mask and ρ for the secondary mask. Beta node implementations typically store the input relations in *indexers* (denoted with IX in the figure). The relation representing the result tuples is denoted with t .

Join Node

The *join node* implements the equijoin operation. This operation is often denoted similarly to the natural join operation.

$$t = p \bowtie s$$

The equijoin is determined by creating the Cartesian product of the relations and filtering those tuples which are equal on the attributes defined by the pattern masks.

$$\sigma_{(P.\lambda_0=S.\rho_0) \wedge (P.\lambda_1=S.\rho_1) \wedge \dots (P.\lambda_{n-1}=S.\rho_{n-1})} (p \times s)$$

The combined tuples are projected: from the attributes selected by the pattern mask, we only keep the ones in p and drop the ones in s . If the join node is used to implement a *natural join* operation, i.e. it joins the relations based on common attributes, the join is defined as

$$p \bowtie s = \pi_{P \cup S} \left(\sigma_{(P.\lambda_0=S.\rho_0) \wedge (P.\lambda_1=S.\rho_1) \wedge \dots (P.\lambda_{n-1}=S.\rho_{n-1})} (p \times s) \right),$$

where $\pi_{P \cup S}$ denotes a projection operation which only keeps one of each common attribute of the schemas over relations p and s .

Note that if the pattern masks are empty, the join operation is equivalent to creating the Cartesian product of the relations. This may cause large (interim) results sets and therefore, if possible, should be avoided. In query engine implementations, this is usually assisted by query optimization techniques. Also, the user is sometimes forced to reformulate her queries for better performance.

If the join node receives a *positive update* Δp on the *primary* input slot, the result (t') and the change set (Δt) is determined as:

$$t' = (p \cup \Delta p) \bowtie s = (p \bowtie s) \cup (\Delta p \bowtie s) = t \cup (\Delta p \bowtie s) \quad \Rightarrow \quad \Delta t = \Delta p \bowtie s$$

If a *positive update* Δs is received on the *secondary* input slot, the result and the change set are:

$$t' = p \bowtie (s \cup \Delta s) = (p \bowtie s) \cup (p \bowtie \Delta s) = t \cup (p \bowtie \Delta s) \quad \Rightarrow \quad \Delta t = p \bowtie \Delta s$$

For *negative updates*, the changeset is the same, but it is propagated as a *negative update*. The result is $t' = t \setminus (\Delta p \bowtie s)$ and $t' = t \setminus (p \bowtie \Delta s)$, for updates messages on the primary and the secondary input slots, respectively.

Antijoin Node

The *antijoin node* implements the antijoin operation which collects the tuples from the primary relation p which have no matching pair in the secondary relation s .

The antijoin operation may be expressed with the equijoin operation as follows:

$$t = p \triangleright s = p \setminus \pi_P (p \bowtie s),$$

where π_P denotes a projection operation which only keeps the attributes of the schema over relation p .

As the antijoin operation is not commutative, handling the update messages requires us to distinguish between the following cases:

- **Update on the primary slot.**

- **Positive update:** send a *positive update* for the incoming tuples for which there are no matches on the secondary indexer.

$$t' = (p \cup \Delta p) \triangleright s = (p \triangleright s) \cup (\Delta p \triangleright s) = t \cup (\Delta p \triangleright s) \Rightarrow \Delta t = \Delta p \triangleright s$$

- **Negative update:** send a *negative update* for the incoming tuples for which there are no matches on the secondary indexer.

$$t' = (p \setminus \Delta p) \triangleright s = (p \triangleright s) \setminus (\Delta p \triangleright s) = t \setminus (\Delta p \triangleright s) \Rightarrow \Delta t = \Delta p \triangleright s$$

- **Update on the secondary slot.**

- **Positive update:** send a *negative update* for the tuples in the primary indexer which match the incoming tuples. Using the definition of antijoin:

$$\begin{aligned} t' &= p \triangleright (s \cup \Delta s) \\ &= p \setminus \pi_P (p \bowtie (s \cup \Delta s)) \\ &= p \setminus \pi_P ((p \bowtie s) \cup (p \bowtie \Delta s)) \\ &= p \setminus \pi_P (p \bowtie s) \setminus \pi_P (p \bowtie \Delta s) \\ &= (p \triangleright s) \setminus \pi_P (p \bowtie \Delta s) \\ &= t \setminus \pi_P (p \bowtie \Delta s) \end{aligned}$$

So $\Delta t = \pi_P (p \bowtie \Delta s)$. We may optimize this further by using t instead of p : $\Delta t = \pi_P (t \bowtie \Delta s)$.

- **Negative update:** send a *positive update* with the tuples ...

$$\begin{aligned}
t' &= p \triangleright (s \setminus \Delta s) \\
&= p \setminus \pi_P (p \bowtie (s \setminus \Delta s)) \\
&= p \setminus \pi_P ((p \bowtie s) \setminus (p \bowtie \Delta s)) \\
&= p \setminus \pi_P (p \bowtie s) \cup \pi_P (p \bowtie \Delta s) \\
&= (p \triangleright s) \cup \pi_P (p \bowtie \Delta s)
\end{aligned}$$

3.3 Incremental Pattern Matching on a Single Workstation: EMF-INCQUERY

In the following, we will overview the architecture of a *single-node* incremental pattern matcher, specifically EMF-INCQUERY.

3.3.1 Architecture

The Rete algorithm forms the foundation of EMF-INCQUERY's query engine. Figure 3.5 shows the architecture of EMF-INCQUERY and the role of the Rete network in the system.

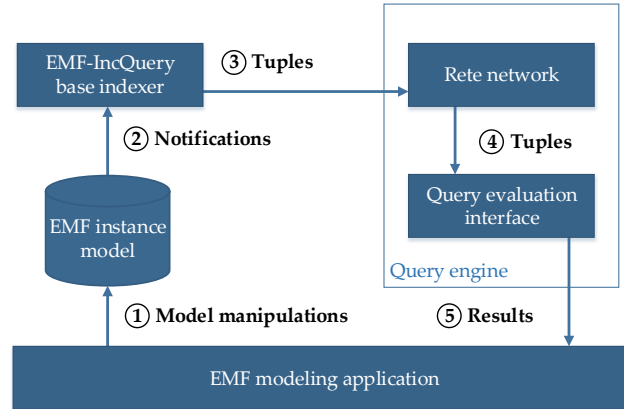


Figure 3.5. EMF-INCQUERY's architecture

A typical model transformation sequence is the following. The modeling application manipulates the EMF instance model ①. Upon modification, the model sends notifications to EMF-INCQUERY's base indexer ②. The indexer propagates the modified tuples to the Rete network as update messages ③, which processes the updates and sends the resulting tuples to the query evaluation interface ④. The modeling application can retrieve the results from the interface ⑤.

3.3.2 Indexing and Initialization

Indexing is a common technique for decreasing the execution time of database queries. In MDE, *model indexing* has a key role to high performance model queries. As MDE primarily uses a metamodeling infrastructure, all queries utilize some type attribute. Typical elementary queries are listed below.

- Retrieving all node instances of a given type (e.g. get all nodes with the type `Person`).
- Retrieving all edge instances of a given label (e.g. get all edges with the label `child`).
- Retrieving a given node's all incoming and/or outgoing edges of a given type (e.g. get all outgoing `child` edges of a given node).
- Reverse navigation: retrieving the node on the other end of an edge (e.g. the `child` relation is identical to the inverse of the `parent` relation).

EMF-INCQUERY uses the EMF API to run these queries efficiently.

3.3.3 Data Representation and Storage

EMF-INCQUERY works on in-memory EMF models. The Rete network represents the data in *tuples*. Basically, the network's tuples can contain two sorts of values: (i) pointers to an EMF model, (ii) Ecore scalar values (`EString`, `EInt`, etc. instances). This data representation principle intends to keep the Rete network's size as small as possible, while allowing efficient processing. Because of the tuple representation, various operations, e.g. projection (π) and join (\bowtie), can be simply defined using pattern masks [32].

3.3.4 Notification Mechanisms

Model change notifications are required by incremental query evaluation, thus model changes are captured and their effects are propagated in the form of *notification objects* (NOs). The notifications generate *tokens* that keep the Rete network's state consistent with the model.

3.3.5 Termination Protocol

As the Rete algorithm’s change propagation is asynchronous, the system must also implement a *termination protocol* to ensure that the query results can be retrieved consistently with the model state after a given transaction (i.e. by signaling when the update propagation has been terminated). The correctness of the protocol is proved in [32].

3.3.6 Configuration and Performance Optimization

For a given model, the system’s performance for a query is mainly determined by the layout of the generated Rete network. Similarly to relational query optimization, we can also optimize the Rete network’s layout. Currently, EMF-INCQUERY supports basic optimizations. It utilizes node sharing, i.e. it detects if two Rete nodes would store the same partial matches and merges them to a single node. More details are available in [33].

3.4 Extensions for Distributed Scalability: INCQUERY-D

Developing a distributed, scalable, incremental pattern matcher introduces numerous challenges. In the following, we will cover the INCQUERY-D’s architecture and our main extensions to EMF-INCQUERY.

3.4.1 Architecture

The INCQUERY-D architecture in an example configuration is shown in Figure 3.6. INCQUERY-D’s architecture consists of three layers: the storage layer, the middleware and the production network. The *storage layer* is a distributed database which is responsible for persisting the model (Section 3.4.3). The client application communicates with the *middleware* ①. The middleware provides a unified API for accessing the database ②. It also sends change notifications ③ (Section 3.3.4) to the production network and retrieves the query results from the production network ④. The *production network* is implemented with a distributed Rete network which provides incremental query evaluation (Section 3.2).

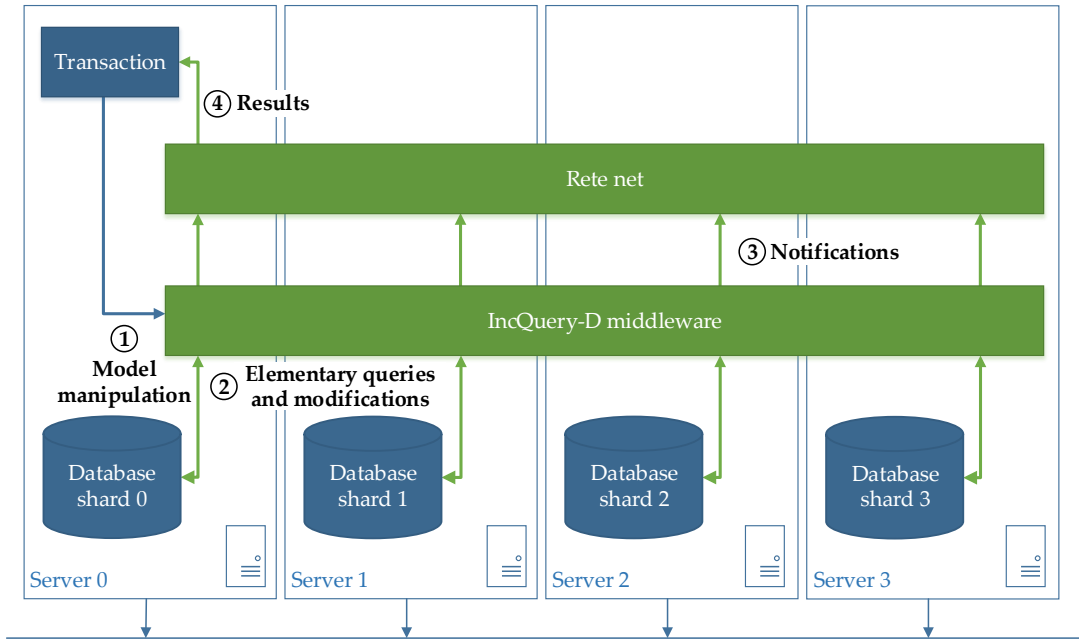


Figure 3.6. INCQUERY-D's architecture on a four-node cluster

3.4.2 Indexing and Initialization

Efficient query processing requires some sort of indexing mechanism. Hence, the INCQUERY-D middleware maintains type-instance indexes so that all instances of a given type (both edges and graph nodes) can be enumerated quickly. These indexers form the bottom layer of the Rete production network. During initialization, these indexers are filled from the database backend (Figure 3.6 ②). In order to reduce the initialization time, the underlying storage layer must be able to process these queries efficiently.

3.4.3 Data Representation and Storage

Conceptually, the architecture of INCQUERY-D allows the usage of a wide scale of model representation formats. Our prototype has been evaluated in the context of the *property graph* and the *RDF* data model, but other mainstream metamodeling and knowledge representation languages such as relational databases' SQL dumps and Ecore instance models (Section 2.2.1) could be supported, as long as they can be mapped to an efficient and distributed storage backend.

For the storage layer, the most important issue from an incremental query evaluation perspective is that the indexers of the middleware should be filled as quickly as possible. This favors technologies where model sharding can be performed

appropriately (i.e. with balanced shards in terms of type-instance relationships), and elementary queries can be executed efficiently.

INCQUERY-D's middleware exposes an API that provides methods for manipulating the graph. By allowing graph-like data manipulation we allow the user to focus on the domain-specific challenges, thus increasing her productivity. The middleware translates the user's operation to the backend's query language and forwards it to the underlying data storage.

In order to allow the Rete algorithm to work, each model element has to have a unique identifier. Issuing a unique identifier in a distributed system is a non-trivial task and is subject to future work (Section 6.2).

To support different data models, we only have to supply the appropriate connector class to INCQUERY-D's middleware. The current prototype supports 4store, Neo4j and Titan.

3.4.4 Notification Mechanisms

While relational databases usually provide *triggers* for generating notifications, most triplestores and graph databases lack this feature. Among our primary database backends, 4store provides no triggers at all. Titan and Neo4j incorporate Blueprints, which provides an `EventGraph` class capable of generating notification events, but the events are only propagated in a single JVM (Java Virtual Machine). Implementing distributed notifications would require us to extend the `EventGraph` class and use a messaging framework. This is subject to future work (Section 6.2).

Because of the lack of support for distributed notifications, in INCQUERY-D's prototype, notifications are controlled by the middleware by providing a facade for all model manipulation operations (Figure 3.6 ③). The notification messages are propagated through the Rete network using the Akka messaging framework.

3.4.5 Termination Protocol

INCQUERY-D's termination protocol works by adding a stack to the message. The stack registers each Rete node the message passes through. After the message reaches the production node, the termination protocol starts. Based on the content of the stack, acknowledgement messages are propagated back on the network. When all relevant input nodes (where the original notification token(s)

started from) receive the acknowledge messages, the termination protocol finishes.

3.4.6 Configuration and Performance Optimization

The Rete algorithm (Section 3.2) utilizes both indexing and caching to provide fast incremental query evaluation. INCQUERY-D's horizontal scalability is supported by the distribution of the pattern matcher's Rete network. To enable this, the system must be able to allocate the Rete nodes to different hosts in a cloud computing infrastructure.

The deployment and configuration of a distributed pattern matcher involves many degrees of freedom and design decisions. The overall performance of the system is influenced by a number of factors.

- For the storage layer, we may choose different database implementations due to INCQUERY-D's backend-agnostic nature. In this thesis, we used property graph databases (Neo4j, Titan) and triplestores (4store).
- We may use different database sharding strategies (e.g. random partitioners or more sophisticated sharding methods based on domain-specific knowledge).
- Using query optimization methods, we can derive *Rete networks with different layouts* for the same query. The most efficient layout can be chosen based on both query and instance model characteristics, e.g. to keep the resource requirement of intermediate join operations to a minimum. [33] discusses the possible optimizations in detail.
- We may choose different strategies to *allocate the Rete nodes* in the distributed system. The optimization strategy may choose to optimize local resource usage, or to minimize the amount of remote network communication. Note that in theory, this is *orthogonal* to the database's sharding strategy, i.e. these are two distinct levels of distribution that do not directly depend upon each other. However, we expect that keeping the Rete network's type indexer nodes and the instances of the given type on the same server would improve the speed of the initialization and modification tasks significantly.
- We may implement *dynamic adaptability* to changing conditions. For example, when the model size and thus query result size grows rapidly, the Rete

network may require *dynamic reallocation* or *node sharding* due to local resource limitations.

3.5 Workflow

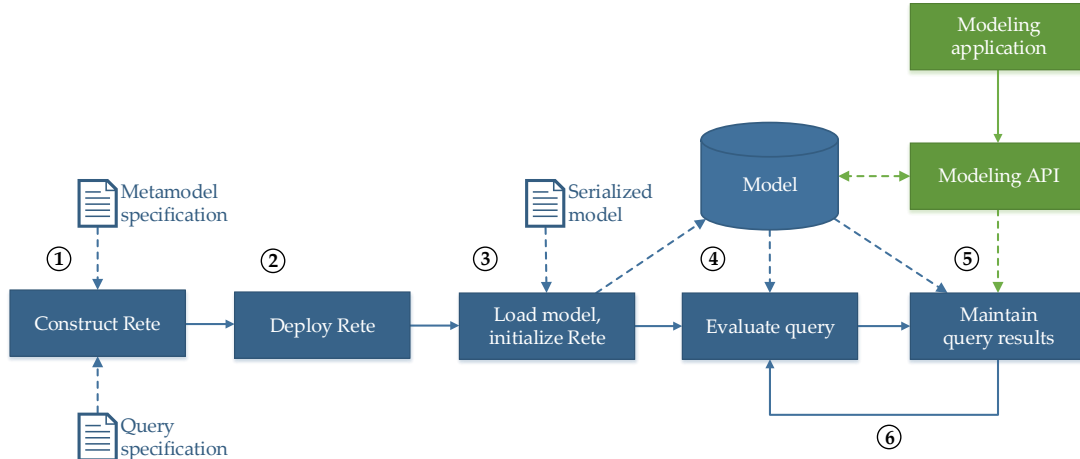


Figure 3.7. The general workflow of incremental pattern matching with the Rete algorithm

In the following, we will describe the workflow behind the pattern matching process. Starting from a metamodel, an instance model and a graph pattern, we will cover the problem pieces that need to be solved for setting up an incremental, distributed pattern matcher. The workflow is shown in Figure 3.7. First, we describe the workflow of EMF-INCQUERY and then compare it to the workflow of INCQUERY-D.

3.5.1 Workflow of EMF-INCQUERY

Based on the *metamodel* and the *query specification*, EMF-INCQUERY first constructs a Rete network ① and deploys it ②. It loads the model (from the persistent storage) to an *in-memory storage* ③ and traverses it to initialize the Rete network’s indexers. The Rete network evaluates the query by processing the incoming tuples ④. If the modeling application modifies the model (through the EMF API), the modifications are propagated through the Rete network, hence keeping it in a consistent state ⑤. The query results can be retrieved from the Rete network ⑥. The modeling application may modify the model and reevaluate the query again.

3.5.2 Workflow of INCQUERY-D

By design, the workflow steps of INCQUERY-D are similar to EMF-INCQUERY's, discussed in Section 3.5.1. However, due to the system's distributed nature, they are more difficult to design and implement.

The main differences are the following. In INCQUERY-D, deploying the Rete network ② requires the deployment of remote actors (Section 2.4) on the servers. Both the Rete indexers and the database are distributed across the cluster. Hence, loading the model and initializing the Rete network needs network communication ③. The Rete network works using Akka's remote messaging feature. The query results can be retrieved from the Rete network (this may also require network communication) ④. The database shards can only be accessed through the middleware, which is responsible for sending notifications to the Rete network's appropriate indexers. After the notifications are processed and the distributed termination algorithm finishes, the Rete network is in a consistent state ⑤. The results can be retrieved by the client which may modify the model and reevaluate the query again ⑥.

3.6 Tooling for INCQUERY-D

As mentioned earlier, we aimed to build INCQUERY-D on top of EMF-INCQUERY's pattern language (IQPL) and its Rete network generator. Because EMF-INCQUERY has an Eclipse-based user interface for defining and executing queries, we plan to provide the same tooling environment. Also, for the allocation of Rete nodes, we created an Eclipse-based editor and viewer.

To aid the system's dynamic capabilities, we plan to develop a *runtime model-based dashboard* to monitor the state of INCQUERY-D's nodes. Currently, the INCQUERY-D tooling generates an architecture definition file (`arch`), which is used for deploying the distributed pattern matcher.

This file contains the Rete network's layout and its allocation in the cloud (as of now, the latter is defined manually). INCQUERY-D uses the architecture description for instantiating the Rete network and initializing the middleware (Figure 3.8).

Using an architecture definition file has a number of advantages. It allows us decouple the client-side dashboard and the server-side runtime. It also simplifies the experiments with different Rete network layouts and mappings.

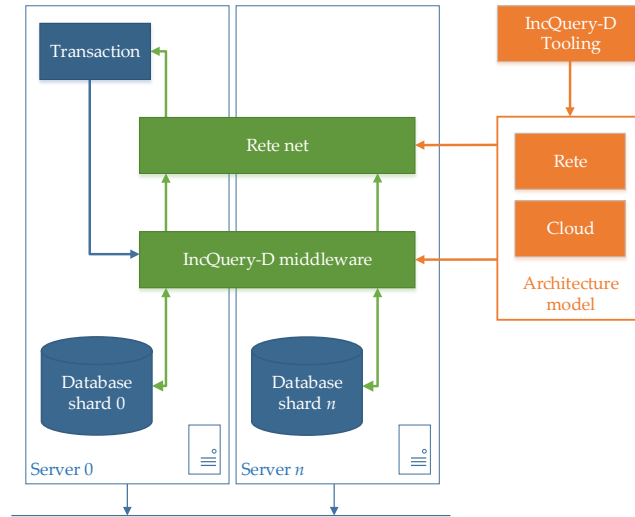


Figure 3.8. Architecture of INCQUERY-D with a runtime dashboard

To provide live feedback, we will adopt a *live* architecture model. The live model will provide real-time details about the system’s current state, including the local resources on each server, the Rete nodes’ memory consumption and so on.

3.7 Elaboration of the Example

To demonstrate INCQUERY-D’s approach, we elaborate an example in detail. We introduce a case study, then formulate a query and show the workflow that executes the distributed, incremental evaluation of the pattern defined by query.

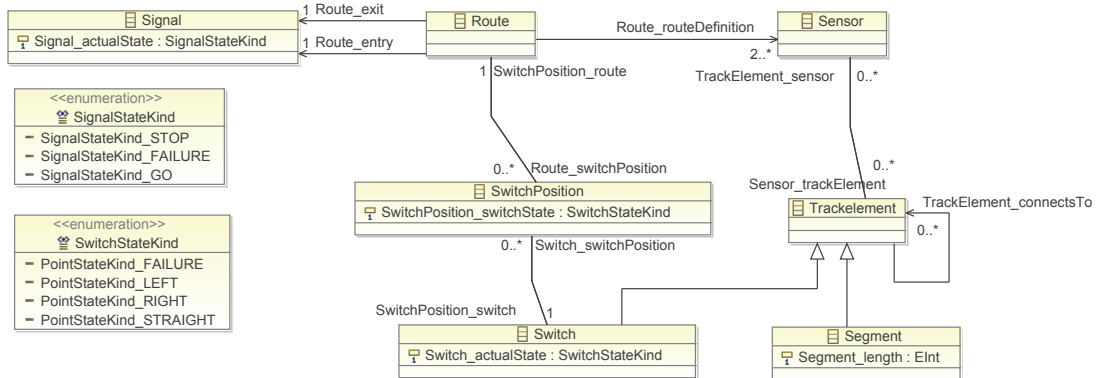


Figure 3.9. The EMF metamodel of the railroad system

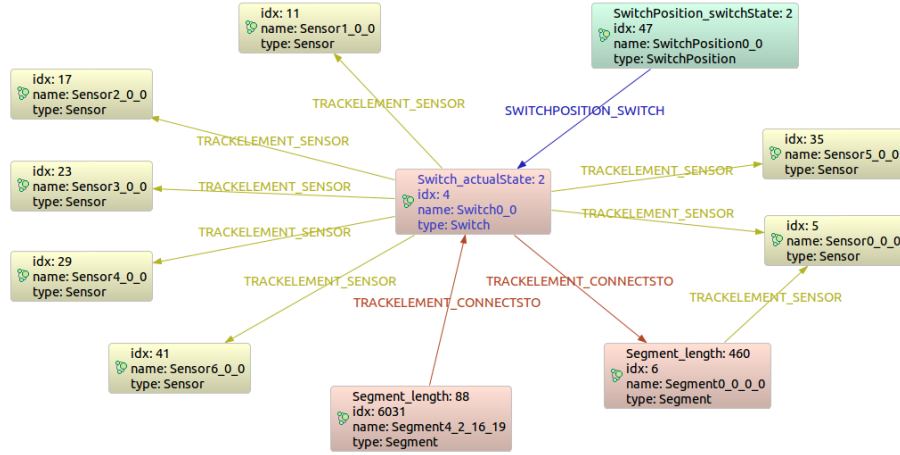


Figure 3.10. A subgraph of a railroad system visualized

3.7.1 Case Study: Railroad System Design

The example is built around an imaginary railroad system defined in the MOGENTES EU FP7 [77] project. The system's network is composed of typical railroad items, including signals, segments, switches and sensors. The complete EMF metamodel is shown in Figure 3.9. A subgraph of an instance model is shown in Figure 3.10.

We used a query that resembles a typical MDE application's workload. The query, called *RouteSensor*, looks for violations of a *well-formedness constraint* in the model.

RouteSensor

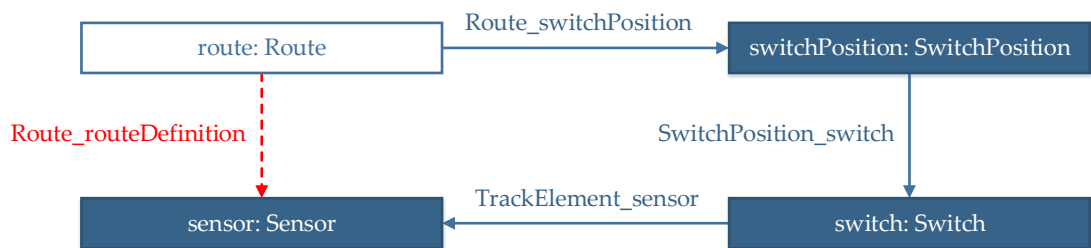


Figure 3.11. Graphical representation of the *RouteSensor* query's pattern. The dashed red arrow defines a negative application condition.

The *RouteSensor* query looks for Sensors that are connected to a Switch, but the Sensor and the Switch are *not* connected to the same Route. In other words, all sensors that are associated with a switch that belongs to a route must also be associated directly with the same route.

The graphical representation of the query is shown in Figure 3.11. Basically, the query binds the type of the vertices, defines three edges and one negative edge, called NAC (negative application condition).

```

1 package hu.bme.mit.train.constraintcheck.incquery
2
3 import "http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl"
4
5 pattern routeSensor(Sen, Sw, Sp, R) = {
6   Route(R);
7   SwitchPosition(Sp);
8   Switch(Sw);
9   Sensor(Sen);
10
11   Route.Route_switchPosition(R, Sp);
12   SwitchPosition.SwitchPosition_switch(Sp, Sw);
13   Trackelement.TrackElement_sensor(Sw, Sen);
14
15   neg find head(Sen, R);
16 }
17
18 pattern head(Sen, R) = {
19   Route.Route_routeDefinition(R, Sen);
20 }

```

Listing 3.1. *The RouteSensor query in IQPL*

The RouteSensor query in IQPL (INCQUERY Pattern Language) is shown on Listing 3.1. This query binds the variables (Sen, Sw, Sp, R) to the appropriate type. It defines the three edges as relationships between the variables and defines the negative application condition as a negative pattern (neg find).

For comparison, we also present the RouteSensor query in SPARQL (RDF's query language) on Listing 3.2. Here, the types are defined with the `rdf:type` predicate, while the edges are defined with base predicates. The negative application condition is defined with the `FILTER NOT EXISTS` construction².

```

1 PREFIX base: <http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
5
6 SELECT DISTINCT ?xSensor
7 WHERE
8 {
9   ?xRoute rdf:type base:Route .
10   ?xSwitchPosition rdf:type base:SwitchPosition .
11   ?xSwitch rdf:type base:Switch .
12   ?xSensor rdf:type base:Sensor .
13   ?xRoute base:Route_switchPosition ?xSwitchPosition .
14   ?xSwitchPosition base:SwitchPosition_switch ?xSwitch .

```

²Note that the two queries are slightly different: the SPARQL query returns only a set of Sensors, while the IQPL query returns a set of (Sensor, Switch, SwitchPosition, Route) tuples.


```

15   ?xSwitch base:TrackElement_sensor ?xSensor .
16
17   FILTER NOT EXISTS {
18       ?xRoute ?Route_routeDefinition ?xSensor .
19   } .
20 }

```

Listing 3.2. *The RouteSensor query in SPARQL*

Given the mapping defined in Section 3.2.3, the *RouteSensor* query can be formalized in relational algebra as:

$$\pi_{Route} \left(\begin{array}{l}
 Route_switchPosition \bowtie SwitchPosition_switch \bowtie \\
 TrackElement_sensor \triangleright Route_routeDefinition
 \end{array} \right)$$

3.7.2 Local Search-based Query Evaluation

Query evaluation is often implemented using local search-based algorithms [54]. Search-based approaches derive a *search plan* from the query specification and execute the plan on the graph. A search plan is a totally ordered list of *search operations*, where search operations represent the atomic units of pattern matching (a single step in the matching process). A unit is either an *extend* operation which extends the matching by a new element (e.g. match the target node along an edge), or a *check* operation used for checking constraints between pattern elements (e.g. whether an edge runs between two nodes) [54].

Table 3.1 presents a possible search plan for the *RouteSensor* pattern. A possible implementation of the NAC (negative application condition) is the definition of a negative subpattern (a *Route_routeDefinition* edge from *Route* nodes to *Sensor* nodes).

Compared to incremental approaches, search-based algorithms do not cache the partial matches and are therefore not able to maintain the match set without reevaluating the whole query. Search-based approaches are often assisted by advanced search plan generation techniques which use metamodel analysis and other heuristics.

	operation	type
1	route instance of Route	check
2	source of e1 is route	extend
3	e1 instance of Route_switchPosition	check
4	target of e1 is switchPosition	extend
5	switchPosition instance of SwitchPosition	check
6	source of e2 is switchPosition	extend
7	e2 instance of SwitchPosition_switch	check
8	target of e2 is switch	extend
9	switch instance of Switch	check
10	source of e3 is switch	extend
11	e3 instance of TrackElement_sensor	check
12	target of e3 is sensor	extend
13	sensor instance of Sensor	check
14	NAC for the Route_routeDefinition edge	check

Table 3.1. Search plan for the evaluation of the RouteSensor query's pattern

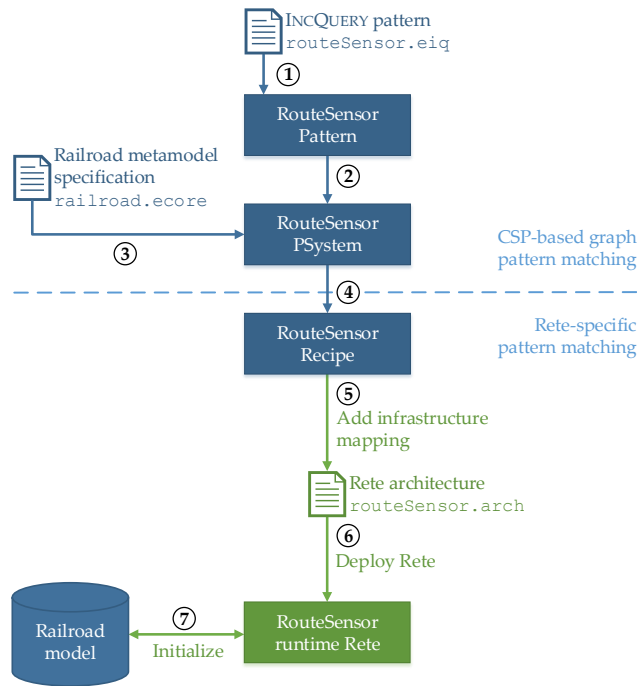


Figure 3.12. INCQUERY-D's workflow

3.7.3 Workflow of the Example

Following the workflow defined in Section 3.5, we will cover the steps for deploying and operating a distributed pattern matcher. The actual workflow for the *RouteSensor* query is shown in Figure 3.12.

Constructing a Rete Network

First, using EMF-INCQUERY's tooling, the query (`routeSensor.iqpl`, see Listing 3.1) is analyzed and parsed to an EMF model ①. The metamodel (`railroad.ecore`) is shown in Figure 3.9. Based on the query ② and the metamodel ③ EMF-INCQUERY builds a *pattern system* (PSystem). The PSystem is translated to a Rete recipe, which defines a Rete layout ④, that guarantees the satisfaction of the constraints. The Rete layout is shown in Figure 3.13.

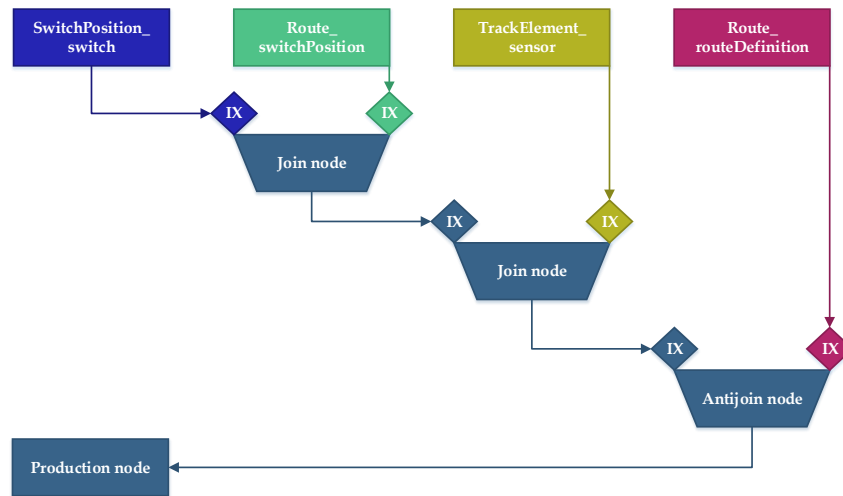


Figure 3.13. The RouteSensor query's layout

Deploying the Rete Network

The Rete nodes are allocated to the cluster's servers by providing the infrastructure mapping ⑤. In INCQUERY-D's prototype, the Rete nodes defined in the recipe are allocated manually on the cloud servers (called *Machines*). The Rete nodes are associated with the machines by *infrastructure mapping* relationships. INCQUERY-D's tooling currently provides an Eclipse-based tree editor to define machines and the infrastructure mapping relationships.

The tooling is capable of visualizing the Rete network and its mapping to the machines (see Figure 3.14). The Rete network is deployed to the Akka instances running on the servers ⑥.

Evaluating the Query

The query is evaluated by initializing the Rete network ⑦ and reading the results from its production node.

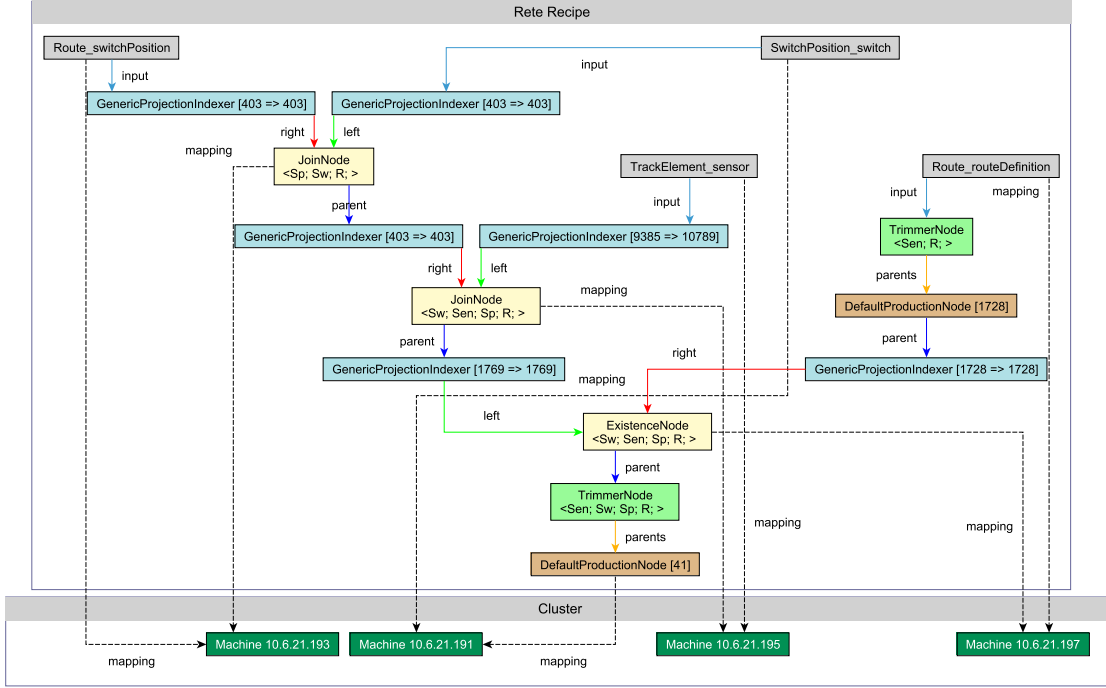


Figure 3.14. The yFiles viewer in INCQUERY-D's tooling

Maintaining the Query Results

In order to provide query results that are consistent with the model, we need maintain the Rete network's state. Suppose we have the graph shown on Figure 3.15 and we decide to delete the `Route_routeDefinition` edge between vertices 2 and 1.

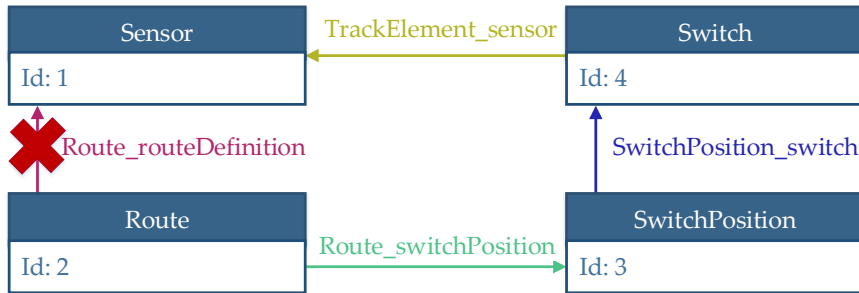


Figure 3.15. A modification on a Train Benchmark instance model

Figure 3.16 shows the distributed Rete network containing the partial matches of the original graph. When we delete the edge between vertices 2 and 1, the `Route_routeDefinition` type indexer (an input node) receives a notification from the middleware and sends a *negative update* ① with the tuple $\langle 2, 1 \rangle$. The anti-join node processes the negative update and propagates a negative update ②

with the tuple $\langle 3, 4, 2, 1 \rangle$. This is received by the production node, which initiates the *termination protocol* ③, ④. After the termination protocol finishes, the indexer signals the client about the successful update. The client can now retrieve the results from the production node. The client may choose to retrieve only the *change set*, i.e. only the tuples that have been added or deleted since the last modification.

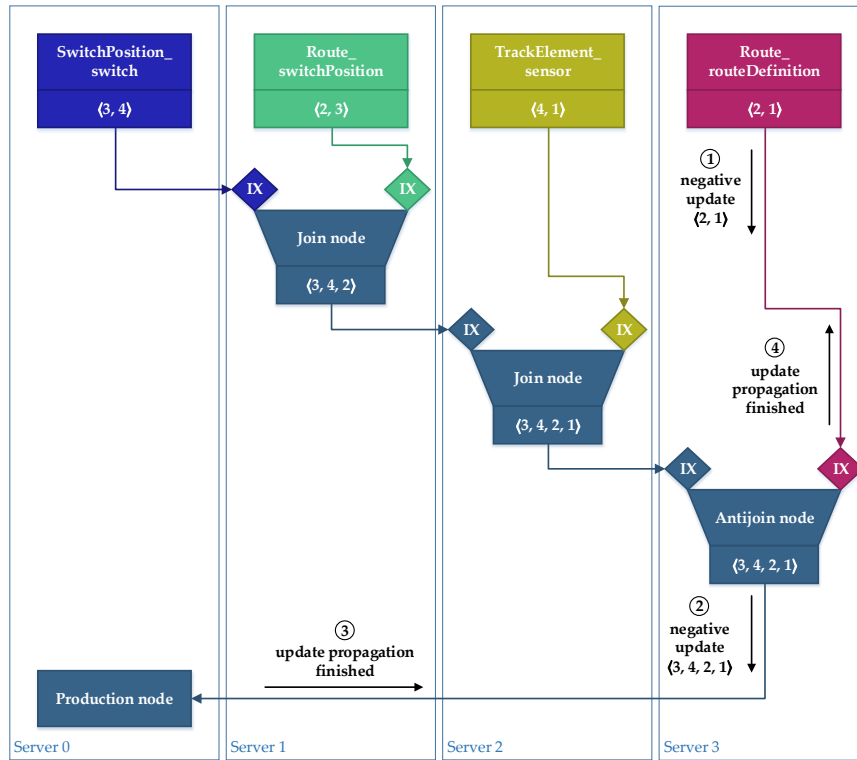


Figure 3.16. Operation sequence on a distributed Rete network

Chapter 4

Evaluation of Performance and Scalability

We developed a prototype of INCQUERY-D to evaluate the feasibility of the approach. In the following chapter, based on the work for EMF-INCQUERY, we introduce a distributed performance benchmark. We present the benchmark environment and analyze the results, with particular emphasis on the scalability of our approach.

The prototype of INCQUERY-D is based on the architecture presented in Chapter 3. A working prototype is beneficial for a number of reasons. First, it serves as a proof concept by demonstrating that a distributed, incremental pattern matcher is feasible with the technologies currently available. On the other hand, it gives us the opportunity to define and run benchmarks, so that we can evaluate the scalability aspects of the system.

4.1 Dimensions of Scalability

A distributed system's *scalability* has multiple dimensions. Usually, when aiming for *horizontal scalability*, the most emphasized dimension is the *number of processing nodes* (computers) in the system. However, there are other important aspects that include *local resources* of the servers, *network communication overhead*, etc. The main goal of our benchmark was to measure the scalability of INCQUERY-D with respect to the model size and compare it to other non-incremental query technologies.

4.2 Foundations: the Train Benchmark

The Train Benchmark was designed at the Fault Tolerant Systems Research Group [58, 56] to measure the efficiency of model queries and manipulation operations in different tools. The Train Benchmark is primarily targeted for typical MDE workloads, more specifically for well-formedness validations.

4.2.1 Benchmark Goals

The Train Benchmark measures the response time of the system under load. The benchmark models a “real-world” MDE workload by simulating a user’s interaction with the model. In this sequence, the user loads the model and validates it against a set queries (defining well-formedness constraints). The user edits the model in small steps. The user’s work is more productive and less error-prone if she receives the results of the validation instantly after each edit. Therefore, the user would like to reevaluate well-formedness queries quickly.

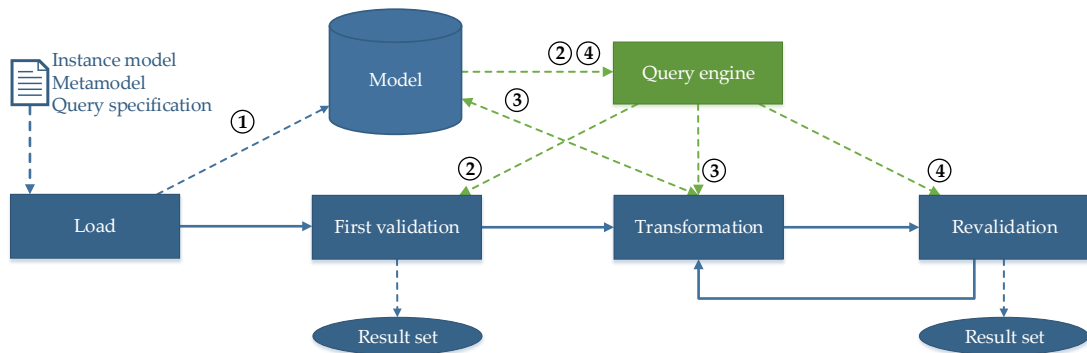


Figure 4.1. Execution sequence of the Train Benchmark

The benchmark defines four distinct phases, also shown in Figure 4.1.

1. *Load*: load the serialized instance model to the database ①.
2. *First validation*: execute the well-formedness query on the model ②.
3. *Transformation*: modify the model ③.
4. *Revalidation*: execute the well-formedness query again ④.

To assess the scalability of the tools, the benchmark uses instance models of growing sizes, each model containing about twice as many model elements as the

previous one (Section 4.2.2). Running the same validation sequence on different model sizes highlighted the limitations of the tested query engines.

Scalability is also evaluated against the complexity of the queries. The benchmark defines different queries, each testing different aspects of the query engine (filtering, join and antijoin operations, etc.). To achieve a successful run, the tested tool is expected to evaluate the query and return the *identifiers* of the model elements in the result set.

4.2.2 Generating Instance Models

Due to both confidentiality and technical reasons, it is difficult to obtain real-world industrial models and queries. Also, using confidential data sets hinders the reproducibility of the conducted benchmarks. Therefore, a generator was developed which creates instance models which mimic real-world models.

We used the *railway system metamodel*, defined in Section 3.7. The instance models are generated pseudorandomly, with pre-defined structural constraints and a regular fan-out structure (i.e. nodes of a given type have similar indegree and outdegree) [56]. The generator is capable of generating models of different sizes and formats, including EMF, OWL, RDF and SQL. We also developed a generator for the property graph data model. In Section A.3, we provide some examples about mapping the EMF metamodel to the framework of property graphs.

4.2.3 Original Results for Non-distributed Tools

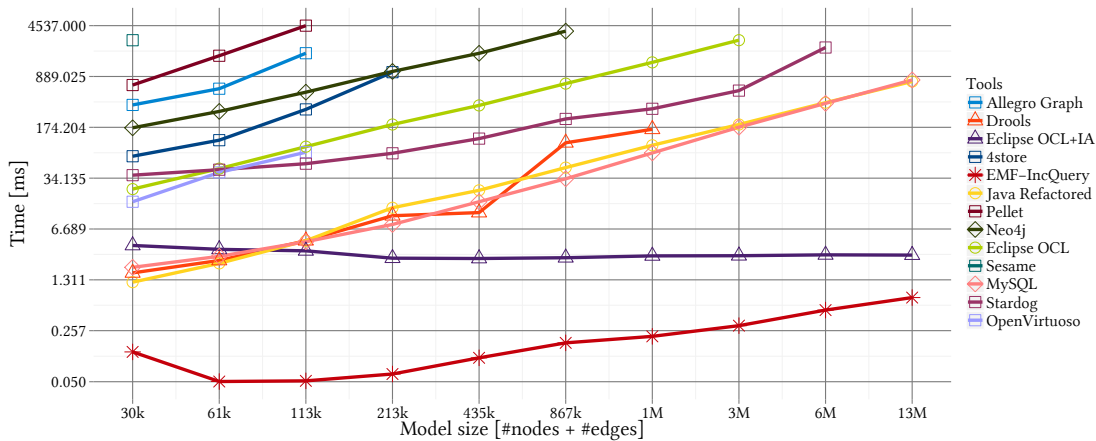


Figure 4.2. Train Benchmark: response times for incremental query evaluation, measured on a single node [58]

The Train Benchmark was designed to work with different tools originating from various technological spaces, e.g. EMF-based tools (EMF-INCQUERY, Eclipse OCL), semantic web technologies (AllegroGraph, Sesame, 4store), NoSQL databases (Neo4j), etc.

Figure 4.2 shows the incremental transformation and validation time for the *RouteSensor* query, discussed in Section 3.7.1. The results clearly show the advantage of incremental query engines. Both Eclipse OCL Impact Analyzer and EMF-INCQUERY scale very well (their characteristic is almost constant to the model size and linear to the size of the result set), while non-incremental tools scale linearly at best, which renders them inefficient for large models.

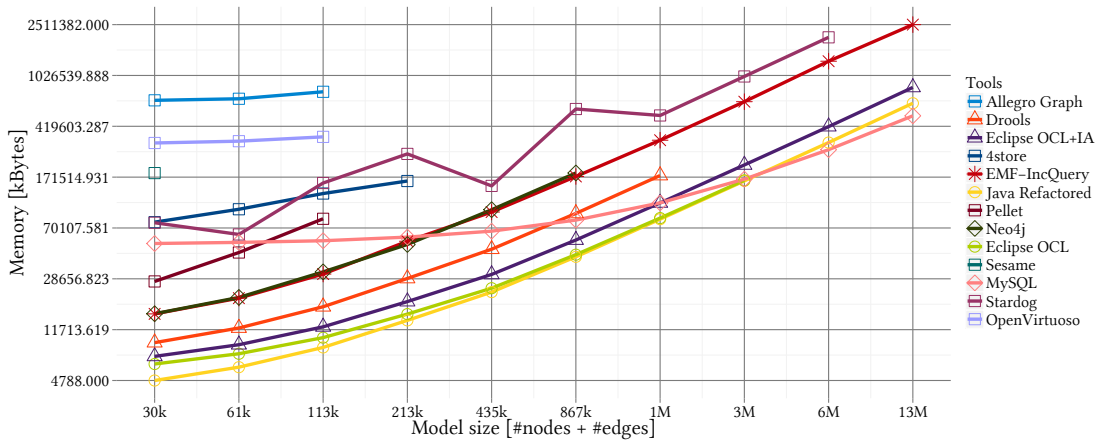


Figure 4.3. Train Benchmark: memory consumption of the tools [58]

Figure 4.3 shows the memory consumption of the different tools. It is apparent that the incremental tools' space-time tradeoff causes them to consume more memory.

4.3 Distributed Train Benchmark

Based on the Train Benchmark, discussed in Section 4.2, we created an extended version for distributed systems. The main goal of the distributed Train Benchmark is the same as the original's: measure the response time and inspect the scalability of different tools. Specifically, the main goal was to compare the performance of INCQUERY-D to distributed, non-incremental query technologies.

4.3.1 Distributed Architecture

The distributed benchmark defines the same phases as the original Train Benchmark (Figure 4.1). The benchmark is controlled by a distinguished node of the

system, called the *coordinator*. The coordinator delegates the operations (e.g. loading the graph) to the distributed system. The queries and the model manipulation operations are handled by the database management system which runs them distributedly and waits for the distributed operation to finish (effectively creating a synchronization point after each operation).

4.3.2 Benchmark Limitations

It is possible that the incoming data sets lack a globally unique identifier. In this case, we need to automatically generate unique identifiers. While some systems (e.g. Titan) support this, other systems (e.g. 4store) do not have such feature. For these systems, the INCQUERY-D middleware should be able to generate unique identifiers. This feature is subject to future work (Section 6.2). In the current benchmark, we worked around this by enforcing the generator to create models with numeric unique identifiers¹.

A common reason for designing and implementing distributed systems is that they are capable of handling a large number of concurrent requests. This way, more users can use the system at the same time. In the distributed Train Benchmark, the system is only used by a single user. Simulating multiple users and issuing concurrent requests is also subject to future work (Section 6.2).

4.3.3 Generating Instance Models

For Neo4j, we expanded the generator with a *property graph generator* module. The generator creates a graph in a Neo4j database and uses the Blueprints library's `GraphMLWriter` and `GraphSONWriter` classes to serialize it to GraphML (Section A.1.1) and Blueprints GraphSON (Section A.1.2) formats.

Titan's Faunus framework requires a third format called Faunus GraphSON (Section A.1.3). To use Faunus, we extended the property graph generator to generate Faunus GraphSON files as well.

4.4 Benchmark Environment

We used the distributed Train Benchmark (Section 4.3) to evaluate INCQUERY-D's performance and compare it to non-incremental solutions. In the following sec-

¹Unlike for property graphs, numeric unique identifiers are not required by the RDF data model.

tion, we will discuss the benchmark setup and the environment in detail.

4.4.1 Benchmark Setup

We tested INCQUERY-D with three storage backends: first with Neo4j (Section 2.3.2), then with 4store (Section 2.3.4) and Titan (Section 2.3.3). In both cases, the system was deployed on a four-node cluster.

As a *non-incremental baseline*, we used Neo4j’s and 4store’s own query engines. While we also planned to use Titan’s query engine, our experiments showed that even for medium-sized graphs, the system was unable to run even the elementary queries (e.g. retrieving vertices by type), not to mention the more complex ones.

The benchmark follows the phases defined in the distributed Train Benchmark. Note that the main difference between the batch and incremental scenarios is that the latter maintain a distributed Rete network, which allows efficient query (re)evaluation.

4.4.2 Hardware and Software Ecosystem

As the testbed, we deployed our system to a private cloud. The cloud is managed by Apache VCL (Virtual Computing Lab) and is also used for educational purposes. Therefore, during the benchmark, the network and the host machines could be under load from other users as well. We consider the effect of these in Section 4.8.

The detailed configuration of the servers are provided below.

Hardware

Each virtual machine used two cores of an Intel Xeon L5420 CPU running at 2.5 GHz and had 8 GBs of RAM. The host machines were connected with gigabit Ethernet network connection.

Software

For the benchmarks, we used the following software stack. The technologies are discussed in Chapter 2.

- Ubuntu 12.10 64-bit
- Oracle Java 7 64-bit
- Neo4j 1.8
- 4store 1.1.5
- Titan 0.3.2
- Faunus 0.3.2
- Hadoop 1.1.2
- Cassandra 1.2.2
- Akka 2.1.2

4.4.3 Benchmark Methodology and Data Processing

Both during the development and in runtime we ensured the *functional equivalence* of the measured tools. *During the development*, we followed the Train Benchmark’s well-defined specification [56]. This precisely defines the steps for each phase, e.g. the number of elements to modify in each transformation and the amount of transformation–validation cycles. *In runtime*, we checked the result set for correctness against the reference implementation.

The benchmark coordinator software used the Train Benchmark’s framework to collect data about the results of the benchmark. We measured the execution time of the predefined phases. The execution time includes the time required for the coordinator’s operation, the computation and IO operations of the cluster’s computers and the network communication (to both directions). The execution times were determined using Java’s `System.nanoTime()` method.

The results were processed by an R script [27] capable of aggregating and visualizing the results.

4.5 Benchmark Results with Neo4j

During the earlier phases of the research, we conducted measurement using only Neo4j. These results were published in [57]. The benchmark’s setup was slightly different, with the main difference being that due to the lack of sharding in Neo4j, we *sharded the graph manually*. This had some important implications.

- The *batch* queries ran on all shards separately and their results were aggregated by the coordinator. The transformations also ran separately.

- The *incremental* queries were evaluated with a distributed Rete network. The elementary model queries (for filling the indexers) were ran on all shards separately and aggregated by the indexers. The transformations also ran separately.

Because the graph was sharded to disjoint partitions with no edges between them, this can be viewed as an ideal case of graph sharding. Therefore, we can use the results to inspect an “ideal” sharding strategy’s impact on the performance. We present the most important results of the benchmark.

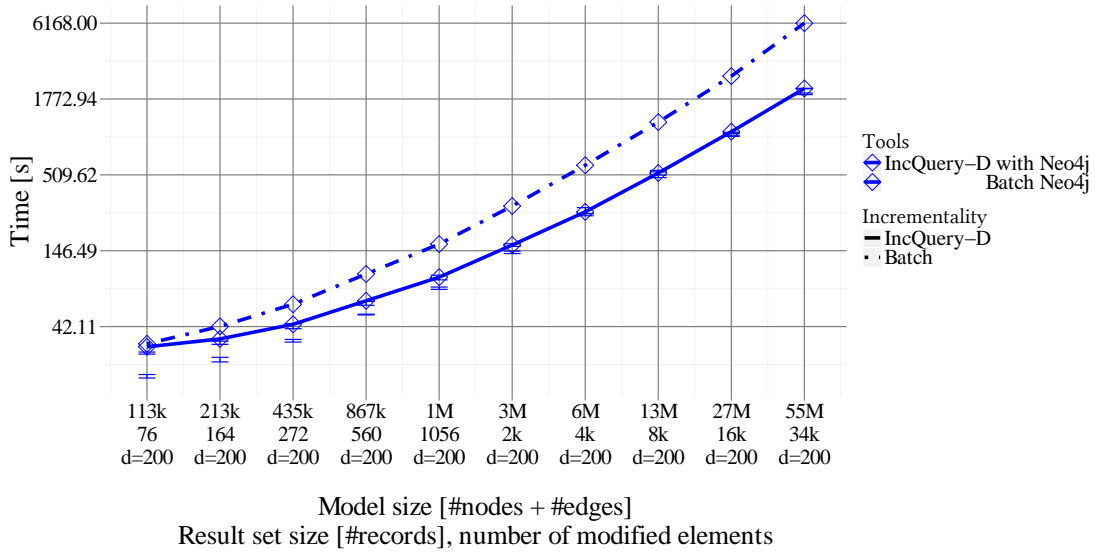


Figure 4.4. Total execution times for 50 validations

Figure 4.4 shows that for validation, INCQUERY-D with Neo4j consistently outperforms Neo4j’s query engine.

Figure 4.5 shows that for transformation and revalidation, INCQUERY-D with Neo4j is about two orders of magnitude faster than Neo4j’s query engine.

4.6 Benchmark Results with 4store and Titan

This section presents the benchmark results with 4store and Titan. Unlike the benchmark with Neo4j (Section 4.5), this benchmark used truly distributed storage backends. Unfortunately, we found that the query engine in Titan is not capable of executing queries with large (interim) result sets, hence the *batch* Titan tool is not present in the results.

The execution times for the *load and first validation* phases are shown in Figure 4.6. As expected, due to the overhead of the Rete network’s construction, the *batch*

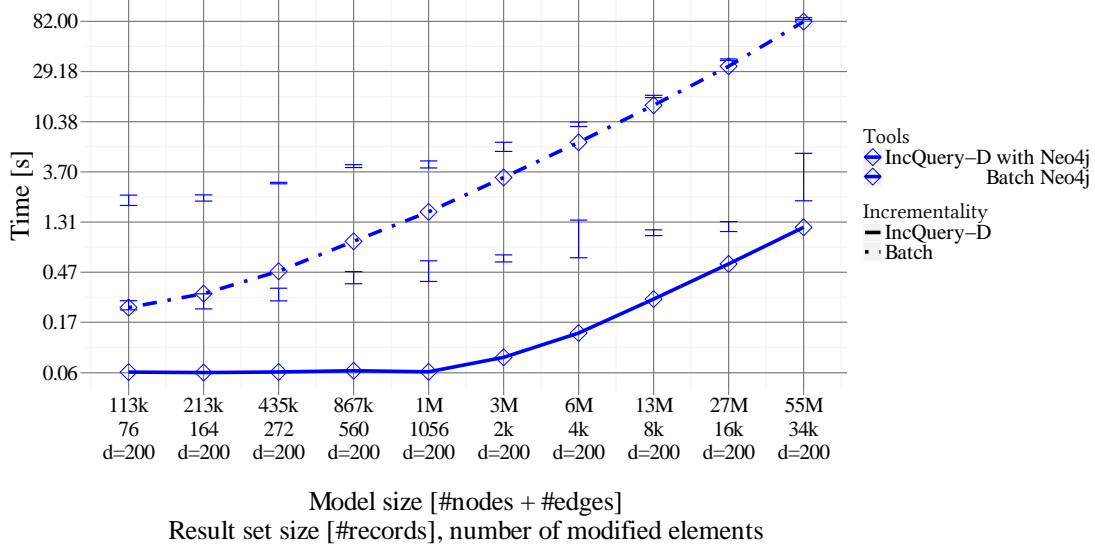


Figure 4.5. Execution times for transformation and revalidation

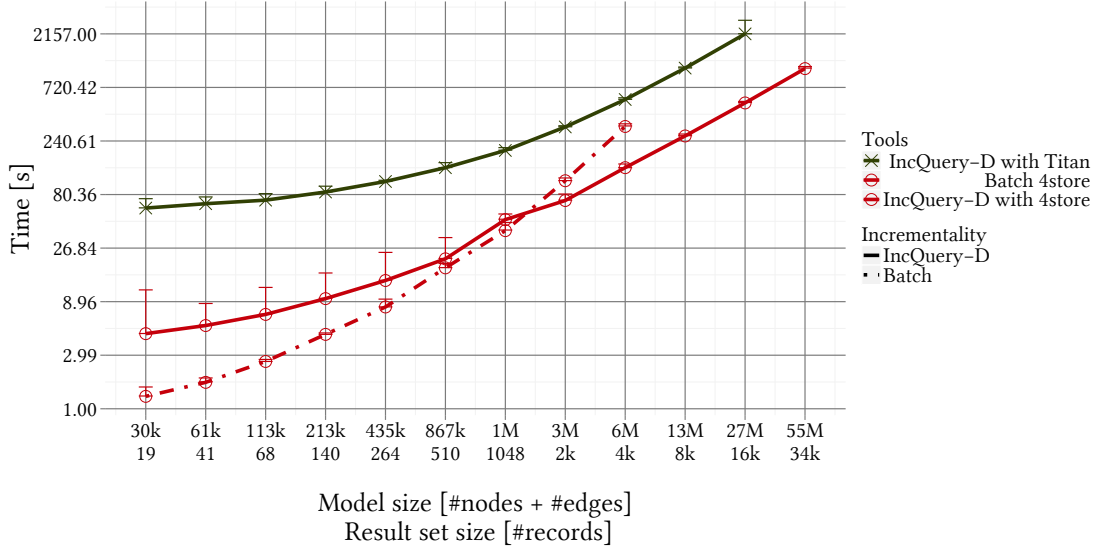


Figure 4.6. Execution times for load and first validation

tool is faster for small models. However, it is important to observe that even for medium-sized models (with a couple of million elements), the INCQUERY-D tools start to edge ahead. This shows that the Rete network's construction overhead already pays off for the first validation.

The execution times for the *transformation* phase are shown in Figure 4.7. The incremental tools provide faster transformation times due to the fact that instead of querying the database, the modeling application can rely on the query layer's indexers. Even for medium-sized models, the INCQUERY-D tools are more than two orders of magnitude faster than the batch tool.

The incremental tools have an even greater advantage for *revalidation* times,

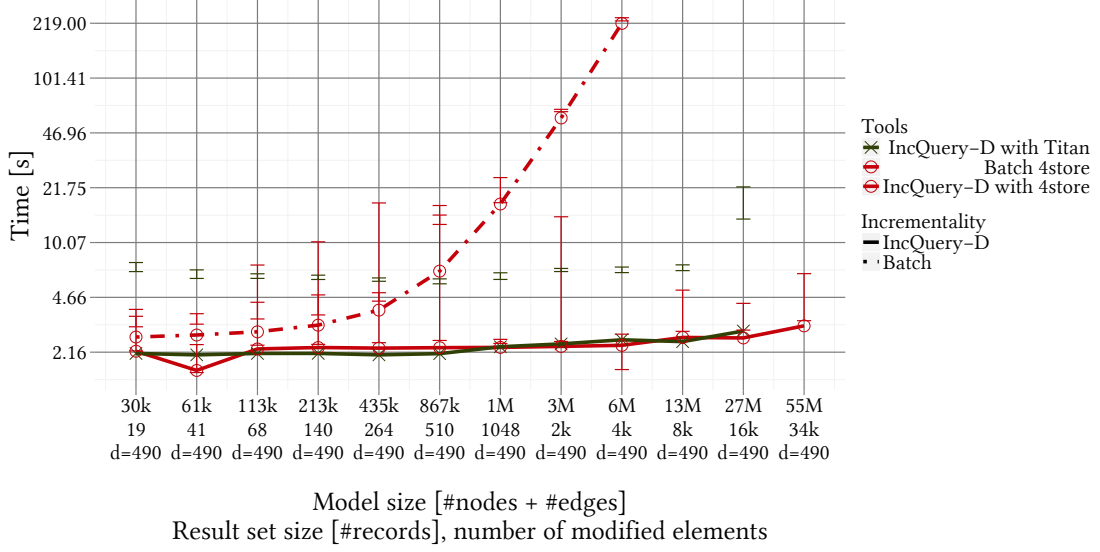


Figure 4.7. Execution times for transformation

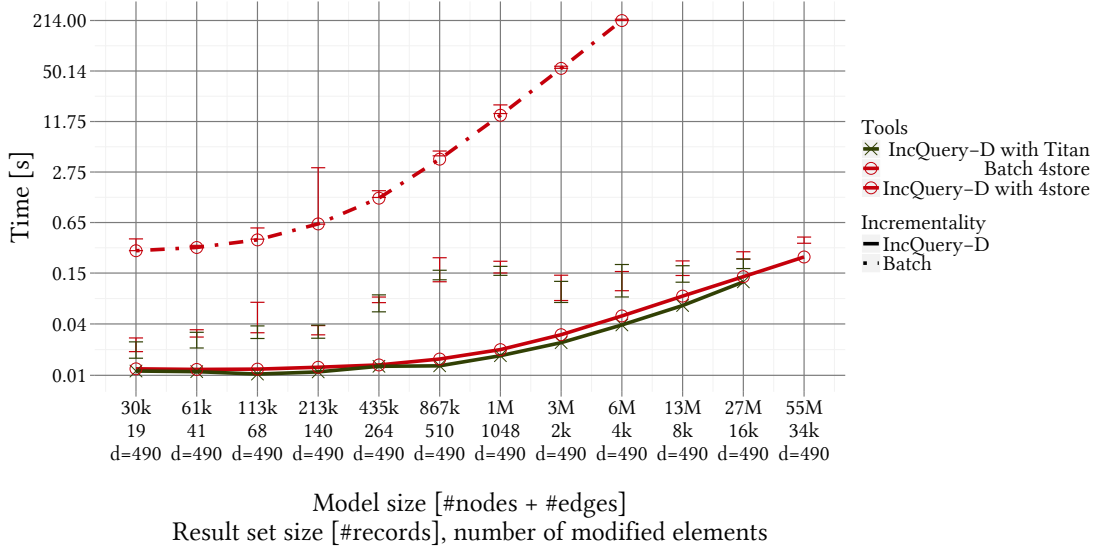


Figure 4.8. Execution times of the revalidation

shown in Figure 4.8. For medium-sized models, they are more than three orders of magnitude faster than the batch tool.

This shows that INCQUERY-D is not just capable of processing models with tens of millions of elements (well beyond the capabilities of single-node tools), but also, it provides sub-second revalidation times.

Figure 4.9 shows the total execution time for a sequence: loading the model, then running transformations and revalidations 50 times. Due to the large number of transformations and revalidations, incremental tools are significantly faster. For example, for a model with 6 million elements, the batch tool took almost 6 hours, while the 4store-based incremental tool took less than 5 minutes.

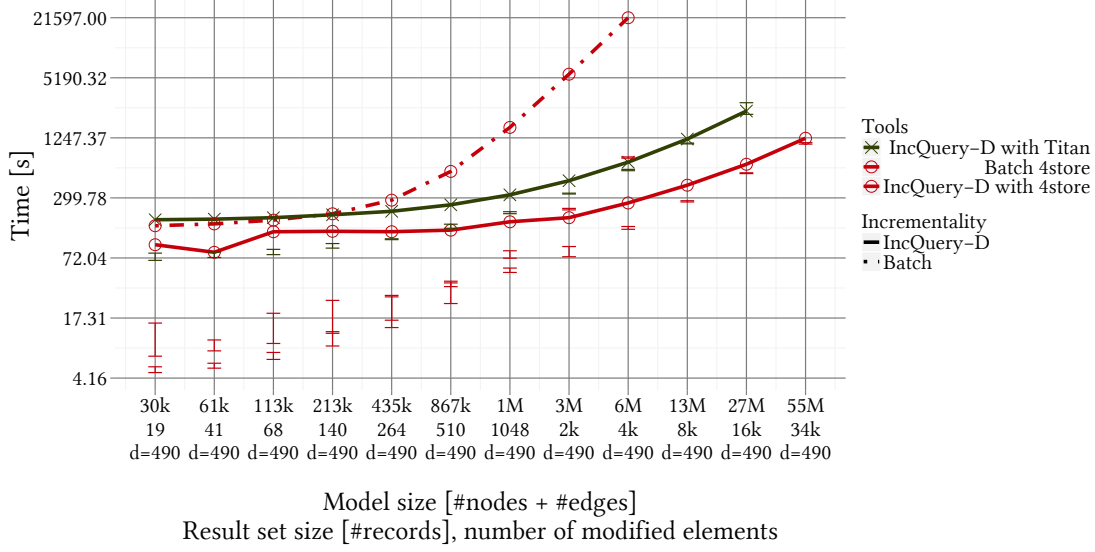


Figure 4.9. Total execution times for 50 validations

4.7 Result Analysis

The results clearly show that the initialization of the Rete network adds some overhead during the *load and first validation phases*. However, even for medium-sized models, this is easily outweighed by the high query performance of the Rete network.

The almost constant characteristic of the execution times of the INCQUERY-D tools' *transformation* and *validation* phases confirm that a distributed, scalable, incremental pattern matcher is feasible with current technologies. Based on the results, we can conclude that while network latency is present, the distributed Rete network still allows sub-second on-the-fly model validation operations. It is also important to observe the similar characteristic of INCQUERY-D's and EMF-INCQUERY's transformation and validation times (Figure 4.2, Figure 4.7 and Figure 4.8).

Another important observation is that for INCQUERY-D tools, the execution time is approximately proportional to *the size of the change*. For batch tools, it is proportional to the *size of the model*.

Note that these results and scalability characteristics do not apply for every workload profile. For example, if the user modifies large chunks of the model and issues queries infrequently, batch query evaluation methods may result in faster execution.

The high memory consumption of the Rete algorithm was one of our main motivations to build a distributed system. For very large models (beyond 10^8 model

elements), we ran into cases where the Java Virtual Machine ran out of or had just enough memory. This resulted in `OutOfMemoryError: Java heap space` and `OutOfMemoryError: GC overhead limit exceeded` exceptions, respectively. Introducing a *Rete node sharding* or other fault-tolerance mechanisms for these cases is subject to future work (Section 6.2).

The results show that the 4store-based INCQUERY-D prototype is consistently faster in the *load* phase than the Titan-based one. This is due to 4store’s simpler architecture and different data model, which is better suited to the INCQUERY-D middleware’s elementary model queries.

In accordance with the original Train Benchmark’s results, the distributed Train Benchmark proved that incremental tools have an advantage for transformation and well-formedness validation sequences. Compared to the Train Benchmark, we managed to work with significantly larger models with more than 50 million model elements. Based on the results, we expect INCQUERY-D to also perform well on different data sets and queries.

4.8 Threats to Validity

To guarantee the correctness of our benchmarks, we laid out some rules to ensure the precision of the results.

First, to start each benchmark sequence independently, we turned the operating system’s caching mechanisms off. The execution time of the *validation and transformation phases* were determined by running them 50 times and taking the *median* values (we decided to take the median instead of the mean value, because the former is less sensitive to transient effects). This way, we could measure the Java Virtual Machine’s warmup effect, which would also occur in a real-world model query engine running for several hours or even longer.

As discussed in Section 4.4.2, our servers could be influenced from the workload caused by other users of the same cloud. To minimize the effect of this and other transient loads, we ran the benchmark five times and took the *minimum* value for each phase. We also disabled file caching in the operating system, so that the serialized model always must be read from the disk.

Despite our efforts, transient effects could still be present in the results. However, their effect is only a threat for smaller model sizes, where the measured execution times are low. For larger models, the main targets of our work, due to longer

execution times, the transient effects do not threaten the validity of the benchmark results.

4.9 Summary

Our benchmarks proved that the proposed architecture is capable of providing scalable, incremental query evaluation. INCQUERY-D's scalability characteristics confirmed that despite the additional network latency, it is possible to keep EMF-INCQUERY's almost constant performance characteristics in a distributed environment. The results show the model size barrier, primarily caused by limitations of memory, can be pushed further using a horizontal scaling approach.

It is important to note that our benchmark did not cover all aspects of distributed scalability. For example, simulating multiple users, measuring the exact memory consumption and network traffic of each server is subject to future work (Section 6.2).

Chapter 5

Related Work

A wide range of special languages have been developed to support *graph-based* representation and querying of computer data. This chapter collects the research and development works that are related to INCQUERY-D.

5.1 Eclipse-based Tools

A class-diagram like modeling language is Ecore of the EMF (Eclipse Modeling Framework, discussed in Section 2.3.1), where classes, references between them and attributes of classes describe the domain. Extensive tooling helps the creation and transformation of such domain models. For EMF models, OCL (Object Constraint Language) is a declarative constraint description and query language that can be evaluated with the local-search based Eclipse OCL [41] engine. To address scalability issues, *incremental* impact analysis tools [48] have been developed as extensions or alternatives to Eclipse OCL.

5.2 Rete Implementations

As a very recent development, Rete-based caching approaches have been proposed for the processing of Linked Data (bearing the closest similarity of our approach). INSTANS [70] uses this algorithm to perform complex event processing (formulated in SPARQL) on RDF data, gathered from distributed sensors.

Diamond [64] uses a *distributed Rete network* to evaluate SPARQL queries on Linked Data, but it lacks an indexing middleware layer so their main challenge is efficient data traversal.

The conceptual foundations of our approach are based on EMF-INCQUERY [36], a tool that evaluates graph patterns over EMF models using Rete. Up to our best knowledge, INCQUERY-D is the first approach to promote distributed scalability by *distributed incremental query evaluation* in the context of model-driven engineering. As the architecture of INCQUERY-D separates the data store from the query engine, we believe that the scalable processing of RDF and property graphs can open up interesting applications outside of the MDE world.

Acharya et al. described a Rete network mapping for fine-grained and medium-grained message-passing computers [29]. The medium-grained computer connected processors in a crossbar architecture, while our approach uses computers connected by gigabit Ethernet. The paper published benchmark results of the medium-grained solution, but these are based only on simulations.

5.3 Benchmarks

This section is based on [56]. Benchmarks have been proposed earlier, mainly to track improvements of a query engine, or to compare tool performance for a given use case.

5.3.1 RDF Benchmarks

SP²Bench [74] is a SPARQL benchmark that measures only query throughput. The goal of this benchmark is to measure query evaluation performance of different tools for a single set of SPARQL queries that contain most language elements. The artificially generated data is based on the real world DBLP bibliography; this way instance models of different sizes reflect the structure and complexity of the original real world dataset. However, other model element distributions or queries were not considered, and the complexity of queries were not analyzed.

The Berlin SPARQL Benchmark (BSBM) [37] measures SPARQL query evaluation throughput for an e-commerce case study modeled in RDF. The benchmark uses a single dataset, but recognizes several use cases with their own query mix. The dataset scales in model size (10 million to 150 billion), but does not vary in structure.

SPLODGE [50] is an approach, where SPARQL queries were generated systematically, based on metrics for a predefined dataset. The method supports distributed SPARQL queries (via the *SERVICE* keyword), however the implementation scaled only up to three steps of navigation, due to the resource consumption

of the generator. The paper did not mention instance model complexity, and only the adequacy of the approach was demonstrated with the RDF3X engine, the effect of queries with different metrics combinations to different engines was not tested.

5.3.2 Model Transformation and Graph Transformation Benchmarks

There are numerous graph transformation benchmarks that do not focus specifically on query performance. However [34] aims to design and evaluate graph transformation benchmark cases corresponding to three usage patterns for the purpose of measuring the performance of incremental approaches on different model sizes and workloads. These scenarios are conceptual continuations of the comprehensive graph transformation benchmark library proposed earlier in [79], which gave an overview on typical application scenarios of graph transformation together with their characteristic features. [46] suggested some improvements to the benchmarks described in [79] and reported measurement results for many graph transformation tools. As model validation is an important use case of incremental model queries, several model query and/or validation tools have been measured in incremental constraint validation benchmarks [67, 43, 36].

Chapter 6

Conclusions

This chapter summarizes the contributions presented in the thesis.

6.1 Summary of Contributions

We presented INCQUERY-D, a novel approach to adapt distributed incremental query techniques to large and complex model-driven software engineering scenarios. Our proposal is based on a distributed Rete network that is decoupled from the sharded graph databases by a middleware layer. The feasibility of the approach has been evaluated using a benchmarking scenario of on-the-fly well-formedness validation of software design models. The results are promising as they show nearly instantaneous query re-evaluation as model sizes grow well beyond 50 million elements.

During the research and development of INCQUERY-D so far, I achieved the following results.

6.1.1 Scientific Contributions

I achieved the following scientific contributions:

- I proposed a novel architecture for building a distributed, scalable, incremental graph query engine over different storage backends. The architecture was published in [55].
- I designed and implemented a *distributed, asynchronous version of the Rete algorithm*.

- I extended the termination protocol used EMF-INCQUERY to work in a distributed environment.
- I extended the Train Benchmark to work in a distributed environment.
- I conducted a benchmark to measure INCQUERY-D's *response time and scalability characteristics*. For the benchmark's baseline, I created *distributed non-incremental benchmark scenarios*.

6.1.2 Practical Accomplishments

I achieved the following practical accomplishments:

- Based on the Rete algorithm, I created a *distributed incremental query engine's prototype*, which is not only detached from the data storage backend, but also agnostic to the storage backend's data model. To prove this, the query engine was tested with both property graphs and RDF graphs.
- I created JUnit [15] unit tests for the Rete algorithm's implementation.
- I used the Akka framework's test kit to create test cases for the Akka actors which operate the distributed Rete network.
- I extended the Train Benchmark with a *new instance model generator*, which can produce property graphs and serialize them in various formats: GraphML, Blueprints GraphSON and Faunus GraphSON.
- I developed INCQUERY-D's prototype, including the query layer, the middleware and the integration to different storage technologies. The prototype consists of more than 5000 lines of Java code and approximately 500 lines of configuration and deployment scripts.
- I elaborated automated deployment tools based on EMF-INCQUERY's existing technologies.
- I experimented with modern non-relational database management systems with a focus on NoSQL graph databases and triple stores. For the purpose of benchmarking different tools, I created scripts to install various graph storages.
 - I deployed a manually sharded *Neo4j* cluster. I formulated the appropriate Cypher queries and created the connector class in INCQUERY-D's middleware to access Neo4j.

- I implemented scripts to install the *Titan graph database and its ecosystem* on a cluster. Titan’s ecosystem includes technologies on different maturity levels, including the Apache Cassandra database, the Apache Hadoop MapReduce framework with the HDFS distributed file system, the TinkerPop graph framework and the Faunus graph analytics engine. I formulated the necessary Gremlin queries and created the connector class in INCQUERY-D’s middleware.
- I implemented scripts to install the *4store triplestore* on a cluster. I formulated the necessary the SPARQL queries and created the connector class in INCQUERY-D’s middleware.
- The analysis of the results confirmed the feasibility of the approach and showed its superiority to current open-source graph query engines.
- I implemented scripts for *automating the benchmark* and *operating a cluster of Akka microkernels*.

This thesis significantly extends the contributions of our earlier results [55]. It discusses the theoretical background and practical difficulties, including the detailed presentation of the architecture (Chapter 3). The benchmark evaluation has been extended to include truly distributed graph database backends (Chapter 4) as INCQUERY-D, as presented here, features backends with automatic sharding, automatic deployment of the Rete nodes and an Eclipse-based tooling environment.

6.2 Limitations and Future Work

INCQUERY-D’s current implementation has some limitations, the most important ones are the following.

1. The Rete nodes are allocated manually. The user has to define the mapping between the Rete network and the infrastructure. However, given a mapping, the system is capable of automatically deploying the Rete network.
2. Only a subset of the nodes defined in the Rete algorithm are implemented. For example, the current implementation does not support recursive patterns and transitive closures.
3. The Eclipse-based tooling does not cover the whole workflow. The user is required to do some manual work, e.g. running scripts manually.

For future work, we plan to address the aforementioned limitations.

1. The allocation of the Rete nodes will be supported using techniques like CSP (Constraint Satisfaction Problem) solvers and DSE (Design Space Exploration) [51]. We plan to further explore advanced optimization challenges such as dynamic reconfiguration and fault tolerance.
2. We will complete the implementation of the nodes defined in Rete algorithm.
3. The tooling is under active development with plans for a *live monitoring* feature.

We also plan to extend the distributed Train Benchmark to model different real-world workloads, e.g. simulating multiple users issuing concurrent requests. We will add a *unique identifier generator* to the middleware, which will allow us to use multiple data sources (which may have different element with the same identifiers). INCQUERY-D already has been designed with different metamodeling languages in mind. In the future, we will create the mapping from more modeling languages, thus working towards a query engine on a federated database system.

Another direction is experimenting with programming languages that are better suited to asynchronous algorithms, e.g. Erlang and Scala, a Java-based functional object-oriented programming language. For our storage layer, we plan to test distributed in-memory databases, e.g. Hazelcast [14] and to extend proven solutions, like Titan, with a distributed notification layer. Also, we are constantly looking for alternative scalable persistent graph database technologies.

Acknowledgements

I would like to thank my supervisors Dr. István Ráth, Benedek Izsó and Dr. Dániel Varró for their friendly advice and enthusiasm. I would also like to thank Dr. Gábor Bergmann for answering all my questions about the details of the Rete algorithm, Dr. Ákos Horváth for providing insight into the search-based approach, Áron Tóth for helping my work in the cloud and all other colleagues in the Fault Tolerant Systems Research Group who provided numerous valuable observations and suggestions.

I would like to extend my appreciation to Dr. Sándor Gajdos for introducing me to the world of database theory and database management systems.

Last but not least, I am deeply grateful to my family and friends for their continuous support.

List of Figures

2.1	Different graph data models (based on [71])	15
2.2	The Ecore kernel, a simplified subset of the Ecore metamodel	16
2.3	The TinkerPop software stack [11]	17
2.4	Different partitionings of the same graph	19
2.5	Hadoop’s architecture [62]	24
2.6	Graph vertex mapped by Titan to a row in a Cassandra database . .	25
2.7	Deploying a remote actor in Akka [5]	28
3.1	The structure of the Rete propagation network	31
3.2	Example graph	32
3.3	Graphical notation of an alpha node	33
3.4	Graphical notation of a beta node	35
3.5	EMF-INCQUERY’s architecture	38
3.6	INCQUERY-D’s architecture on a four-node cluster	41
3.7	The general workflow of incremental pattern matching with the Rete algorithm	44
3.8	Architecture of INCQUERY-D with a runtime dashboard	46
3.9	The EMF metamodel of the railroad system	46
3.10	A subgraph of a railroad system visualized	47
3.11	Graphical representation of the RouteSensor query’s pattern. The dashed red arrow defines a negative application condition.	47
3.12	INCQUERY-D’s workflow	50

3.13	The RouteSensor query's layout	51
3.14	The yFiles viewer in INCQUERY-D's tooling	52
3.15	A modification on a Train Benchmark instance model	52
3.16	Operation sequence on a distributed Rete network	53
4.1	Execution sequence of the Train Benchmark	55
4.2	Train Benchmark: reponse times for incremental query evaluation, measured on a single node [58]	56
4.3	Train Benchmark: memory consumption of the tools [58]	57
4.4	Total execution times for 50 validations	61
4.5	Execution times for transformation and revalidation	62
4.6	Execution times for load and first validation	62
4.7	Execution times for transformation	63
4.8	Execution times of the revalidation	63
4.9	Total execution times for 50 validations	64
A.1	An example graph based on the railway system metamodel	84
A.2	The PosLength query's pattern	92
A.3	The SignalNeighbor query's pattern	93
A.4	The SwitchSensor query's pattern	95

Bibliography

- [1] OpenLink Software: Virtuoso Universal Server. <http://virtuoso.openlinksw.com/>.
- [2] Sesame: RDF API and Query Engine. <http://www.openrdf.org/>.
- [3] 2013: What's Coming Next in Neo4j! <http://blog.neo4j.org/2013/01/2013-whats-coming-next-in-neo4j.html>, January 2013.
- [4] 4store. <http://4store.org/>, October 2013.
- [5] Akka. <http://akka.io/>, May 2013.
- [6] Apache Cassandra. <http://cassandra.apache.org/>, May 2013.
- [7] Apache Hadoop. <http://hadoop.apache.org/>, May 2013.
- [8] Apache HBase. <http://hbase.apache.org/>, May 2013.
- [9] Apache Thrift. <http://thrift.apache.org/>, October 2013.
- [10] Avahi. <http://www.avahi.org/>, October 2013.
- [11] Blueprints. <http://blueprints.tinkerpop.com/>, May 2013.
- [12] GraphSON Format. <https://github.com/thinkaurelius/faunus/wiki/GraphSON-Format>, October 2013.
- [13] GraphSON Reader and Writer Library. <https://github.com/tinkerpop/blueprints/wiki/GraphSON-Reader-and-Writer-Library>, October 2013.
- [14] In-Memory Data Grid – Hazelcast. <http://www.hazelcast.com/index.jsp>, October 2013.
- [15] JUnit. <http://junit.org/>, December 2013.
- [16] Lustre. <http://lustre.org/>, October 2013.

- [17] MongoDB. <http://www.mongodb.org/>, October 2013.
- [18] Neoclipse. <https://github.com/neo4j/neoclipse>, May 2013.
- [19] NoSQL Databases. <http://nosql-database.org/>, May 2013.
- [20] Objectivity – InfiniteGraph. <http://www.objectivity.com/infinitegraph>, October 2013.
- [21] OrientDB Graph-Document NoSQL DBMS. <http://www.orientdb.org/>, October 2013.
- [22] Planet Cassandra – Companies. <http://planetcassandra.org/Company/ViewCompany>, October 2013.
- [23] Protocol Buffers – Google’s data interchange format. <https://code.google.com/p/protobuf/>, October 2013.
- [24] Sparsity-technologies: DEX high-performance graph database. <http://www.sparsity-technologies.com/dex>, October 2013.
- [25] The GraphML File Format. <http://graphml.graphdrawing.org/>, October 2013.
- [26] The Innovative Zing JVM. http://www.azulsystems.com/sites/default/files/images/Innovative_Zing_JVM_v2.pdf, August 2013.
- [27] The R Project for Statistical Computing. <http://www.r-project.org/>, October 2013.
- [28] TinkerPop. <http://www.tinkerpop.com/>, May 2013.
- [29] Acharya, A. et al. Implementation of production systems on message-passing computers. *IEEE Trans. Parallel Distr. Syst.*, 3(4):477–487, July 1992.
- [30] Atlanmod research team. NEO4EMF. <http://neo4emf.com/>, October 2013.
- [31] Don Batory. The LEAPS Algorithm. Technical report, Austin, TX, USA, 1994.
- [32] Gábor Bergmann. Incremental graph pattern matching and applications. Master’s thesis, Budapest University of Technology and Economics, http://mit.bme.hu/~rath/pub/theses/diploma_bergmann.pdf, 2008.

- [33] Gábor Bergmann. *Incremental Model Queries in Model-Driven Design*. Ph.D. dissertation, Budapest University of Technology and Economics, Budapest, 10/2013 2013.
- [34] Gábor Bergmann, Ákos Horváth, István Ráth, and Dániel Varró. A benchmark evaluation of incremental pattern matching in graph transformation. In Hartmut Ehrig, Reiko Heckel, Grzegorz Rozenberg, and Gabriele Taentzer, editors, *Proc. 4th International Conference on Graph Transformations, ICGT 2008*, volume 5214 of *Lecture Notes in Computer Science*, pages 396–410. Springer, Springer, 2008. Acceptance rate: 40%.
- [35] Gábor Bergmann, Ákos Horváth, István Ráth, and Dániel Varró. *Incremental Evaluation of Model Queries over EMF Models: A Tutorial on EMF-IncQuery*, volume 6698 of *Lecture Notes in Computer Science*, pages 389–390. Springer Berlin / Heidelberg, 2011. 10.1007/978-3-642-21470-7_32.
- [36] Bergmann, Gábor et al. Incremental evaluation of model queries over EMF models. In *MODELS*, volume 6394 of *LNCS*. Springer, 2010.
- [37] Christian Bizer and Andreas Schultz. The Berlin SPARQL Benchmark. *International Journal On Semantic Web and Information Systems*, 5(2), 2009.
- [38] E. F. Codd. A relational model of data for large shared data banks. *Commun. ACM*, 13(6):377–387, June 1970.
- [39] Jeffrey Dean and Sanjay Ghemawat. MapReduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, January 2008.
- [40] EMF documentation. Package org.eclipse.emf.ecore. <http://download.eclipse.org/modeling/emf/emf/javadoc/2.9.0/org/eclipse/emf/ecore/package-summary.html>, October 2013.
- [41] Eclipse MDT Project. Eclipse OCL, 2011. <http://eclipse.org/modeling/mdt/?project=ocl>.
- [42] Eclipsepedia. CDO. <http://wiki.eclipse.org/CDO>, October 2013.
- [43] Jean-Rémy Falleri, Xavier Blanc, Reda Bendraou, Marcos Aurélio, Almeida da Silva, and Cédric Teyton. Incremental inconsistencies detection with low memory overhead. *Software: Practice and Experience*, 43, 2013.
- [44] Andreas Emil Feldmann and Luca Foschini. Balanced partitions of trees and applications. In Christoph Dürr and Thomas Wilke, editors, *STACS*,

volume 14 of *LIPICs*, pages 100–111. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2012.

- [45] Charles Forgy. Rete: A fast algorithm for the many pattern/many object pattern match problem. *Artificial Intelligences*, 19(1):17–37, 1982.
- [46] Rubino Geiß and Moritz Kroll. On improvements of the Varro benchmark for graph transformation tools. Technical Report 2007-7, Universität Karlsruhe, IPD Goos, 12 2007. ISSN 1432-7864.
- [47] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google file system. In *Proceedings of the nineteenth ACM symposium on Operating systems principles, SOSP '03*, pages 29–43, New York, NY, USA, 2003. ACM.
- [48] Thomas Goldschmidt and Axel Uhl. Efficient OCL impact analysis, 2011.
- [49] RDF Core Working Group. Resource Description Framework (RDF). <http://www.w3.org/RDF/>, 2004.
- [50] Olaf Görlitz, Matthias Thimm, and Steffen Staab. SPLODGE: Systematic generation of SPARQL benchmark queries for Linked Open Data. In Philippe Cudré-Mauroux, Jeff Heflin, Evren Sirin, Tania Tudorache, Jérôme Euzenat, Manfred Hauswirth, JosianeXavier Parreira, Jim Hendler, Guus Schreiber, Abraham Bernstein, and Eva Blomqvist, editors, *The Semantic Web – ISWC 2012*, volume 7649 of *Lecture Notes in Computer Science*, pages 116–132. Springer Berlin Heidelberg, 2012.
- [51] Ábel Hegedüs, Ákos Horváth, István Ráth, and Dániel Varró. A Model-driven Framework for Guided Design Space Exploration. In *26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011)*, Lawrence, Kansas, USA, 11/2011 2011. IEEE Computer Society, IEEE Computer Society. ACM Distinguished Paper Award, Acceptance rate: 15%.
- [52] Carl Hewitt, Peter Bishop, and Richard Steiger. A universal modular actor formalism for artificial intelligence. In *Proceedings of the 3rd international joint conference on Artificial intelligence, IJCAI'73*, pages 235–245, San Francisco, CA, USA, 1973. Morgan Kaufmann Publishers Inc.
- [53] Guillaume Hillairet, Frédéric Bertrand, Jean Yves Lafaye, et al. Bridging EMF applications and RDF data sources. In *Proceedings of the 4th International Workshop on Semantic Web Enabled Software Engineering, SWESE*, 2008.

- [54] Ákos Horváth, Gábor Bergmann, István Ráth, and Dániel Varró. Experimental assessment of combining pattern matching strategies with VIATRA2. *International Journal on Software Tools for Technology Transfer*, 12(3-4):211–230, 2010.
- [55] Benedek Izsó, Gábor Szárnyas, István Ráth, and Dániel Varró. IncQuery-D: incremental graph search in the cloud. In *Proceedings of the Workshop on Scalability in Model Driven Engineering, BigMDE '13*, pages 4:1–4:4, New York, NY, USA, 2013. ACM.
- [56] Benedek Izsó, Zoltán Szatmári, Gábor Bergmann, Ákos Horváth, and István Ráth. Towards Precise Metrics for Predicting Graph Query Performance. In *28th IEEE/ACM International Conference on Automated Software Engineering (ASE 2013)*, 2013. Accepted.
- [57] Benedek Izsó, Zoltán Szatmári, Gábor Bergmann, Ákos Horváth, István Ráth, and Varro Daniel. Ontology driven design of EMF metamodels and well-formedness constraints. In *Proceedings of the 12th Workshop on OCL and Textual Modelling, OCL '12*, page 37–42, New York, NY, USA, 2012. ACM, ACM.
- [58] Benedek Izsó, Zoltán Szatmári, and István Ráth. High performance queries and their novel applications. <http://incquery.net/publications/trainbenchmark/full-results>, 2012. OMG Workshop.
- [59] Adam Jacobs. The pathologies of big data. *Commun. ACM*, 52(8):36–44, August 2009.
- [60] Jim Webber. On Sharding Graph Databases. <http://jim.webber.name/2011/02/on-sharding-graph-databases/>, February 2011.
- [61] Avinash Lakshman and Prashant Malik. Cassandra: a decentralized structured storage system. *SIGOPS Oper. Syst. Rev.*, 44(2):35–40, April 2010.
- [62] Michael G. Noll. Running Hadoop on Ubuntu Linux (Multi-Node Cluster). <http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-multi-node-cluster/>, October 2013.
- [63] D. P. Miranker and B. J. Lofaso. The Organization and Performance of a TREAT-Based Production System Compiler. *IEEE Trans. on Knowl. and Data Eng.*, 3(1):3–10, March 1991.

- [64] Miranker, Daniel P et al. Diamond: A SPARQL query engine, for linked data based on the Rete match. *AIImWD*, 2012.
- [65] Neo Technology. Neo4j. <http://neo4j.org/>, 2013.
- [66] Javier Espinazo Pagán, Jesús Sánchez Cuadrado, and Jesús García Molina. Morsa: a scalable approach for persisting and accessing large models. In *Proceedings of the 14th international conference on Model driven engineering languages and systems*, MODELS'11, pages 77–92, Berlin, Heidelberg, 2011. Springer-Verlag.
- [67] Alexander Reder and Alexander Egyed. Incremental consistency checking for complex design rules and larger model changes. In *MODELS'12*. Springer-Verlag, 2012.
- [68] Redland RDF Libraries. Raptor RDF Syntax Library. <http://librdf.org/raptor/>, October 2013.
- [69] Redland RDF Libraries. Rasqal RDF Query Library. <http://librdf.org/rasqal/>, October 2013.
- [70] Mikko Rinne. SPARQL update for complex event processing. In *ISWC'12*, volume 7650 of *LNCS*. 2012.
- [71] Marko A. Rodriguez and Peter Neubauer. Constructions from dots and lines. *CoRR*, abs/1006.2361, 2010.
- [72] Markus Scheidgen. How Big are Models – An Estimation. Technical report, Department of Computer Science, Humboldt Universität zu Berlin, 2012.
- [73] Markus Scheidgen, Anatolij Zubow, Joachim Fischer, and Thomas H. Kolbe. Automated and transparent model fragmentation for persisting large models. In *Proceedings of the 15th international conference on Model Driven Engineering Languages and Systems*, MODELS'12, pages 102–118, Berlin, Heidelberg, 2012. Springer-Verlag.
- [74] Michael Schmidt, Thomas Hornung, Georg Lausen, and Christoph Pinkel. SP2Bench: A SPARQL performance benchmark. In *Proc. of the 25th International Conference on Data Engineering*, pages 222–233, Shanghai, China, 2009. IEEE.
- [75] Sherif Sakr. Supply cloud-level data scalability with NoSQL databases. <http://www.ibm.com/developerworks/cloud/library/cl-nosqldatabase/index.html>, March 2013.

- [76] The Eclipse Project. Eclipse Modeling Framework. <http://www.eclipse.org/emf>, October 2012.
- [77] The MOGENTES project. Model-Based Generation of Tests for Dependable Embedded Systems. <http://www.mogentes.eu/>.
- [78] Jeffrey D. Ullman. *Principles of Database Systems, 2nd Edition*. Computer Science Press, 1982.
- [79] Gergely Varró, Andy Schürr, and Dániel Varró. Benchmarking for graph transformation. In *Proc. IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC 05)*, pages 79–88, Dallas, Texas, USA, September 2005. IEEE Press.

Appendix A

Graph Formats

In this chapter, we provide examples for the different graph serialization formats, including property graphs and RDF graphs. The examples describe a small instance model based on the railway system metamodel, shown in Figure A.1.

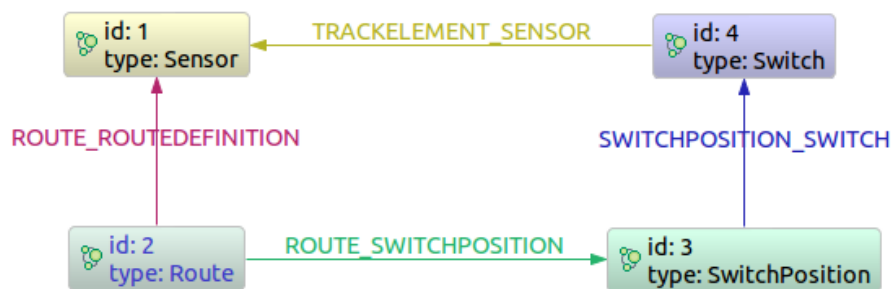


Figure A.1. An example graph based on the railway system metamodel

A.1 Property Graph Formats

A.1.1 GraphML

The GraphML format [25] is the most widely used graph representation format, based on XML (Extensible Markup Language). It has strong tooling support between graph databases and graph visualizing tools.

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <graphml xmlns="http://graphml.graphdrawing.org/xmlns" xmlns:xsi="http://www.w3.org
   /2001/XMLSchema-instance" xsi:schemaLocation="http://graphml.graphdrawing.org/xmlns
   http://graphml.graphdrawing.org/xmlns/1.1/graphml.xsd">
3   <key id="type" for="node" attr.name="type" attr.type="string" />
4   <graph id="G" edgedefault="directed">
5     <node id="1">
6       <data key="type">Sensor</data>
```

```

7   </node>
8   <node id="2">
9     <data key="type">Route</data>
10  </node>
11  <node id="3">
12    <data key="type">SwitchPosition</data>
13  </node>
14  <node id="4">
15    <data key="type">Switch</data>
16  </node>
17  <edge id="0" source="2" target="1" label="ROUTE_ROUTEDEFINITION" />
18  <edge id="1" source="2" target="3" label="ROUTE_SWITCHPOSITION" />
19  <edge id="2" source="3" target="4" label="SWITCHPOSITION_SWITCH" />
20  <edge id="3" source="4" target="1" label="TRACKELEMENT_SENSOR" />
21 </graph>
22 </graphml>

```

Listing A.1. *A graph based on the railway system metamodel stored in GraphML format*

A.1.2 Blueprints GraphSON

Blueprints GraphSON [13] is a JSON-based (JavaScript Object Notation) format. It is not as well supported as the GraphML format (Section A.1.1), but it is less verbose and more readable.

```

1  {
2    "vertices": [
3      {
4        "type": "Sensor",
5        "_id": 1,
6        "_type": "vertex"
7      },
8      {
9        "type": "Route",
10       "_id": 2,
11       "_type": "vertex"
12     },
13     {
14       "type": "SwitchPosition",
15       "_id": 3,
16       "_type": "vertex"
17     },
18     {
19       "type": "Switch",
20       "_id": 4,
21       "_type": "vertex"
22     }
23   ],
24   "edges": [
25     {
26       "_id": 0,
27       "_type": "edge",
28       "_outV": 2,
29       "_inV": 1,
30       "_label": "ROUTE_ROUTEDEFINITION"

```

```

31     },
32     {
33         "_id":1,
34         "_type":"edge",
35         "_outV":2,
36         "_inV":3,
37         "_label":"ROUTE_SWITCHPOSITION"
38     },
39     {
40         "_id":2,
41         "_type":"edge",
42         "_outV":3,
43         "_inV":4,
44         "_label":"SWITCHPOSITION_SWITCH"
45     },
46     {
47         "_id":3,
48         "_type":"edge",
49         "_outV":4,
50         "_inV":1,
51         "_label":"TRACKELEMENT_SENSOR"
52     }
53 ]
54 }

```

Listing A.2. A graph based on the railway system metamodel stored in Blueprints GraphSON format

A.1.3 Faunus GraphSON

In the Faunus GraphSON format [12], each line is a separate JSON (JavaScript Object Notation) document representing a vertex in the graph. This way, the file can be splitted to blocks efficiently and processed on Hadoop nodes in a parallel way.

```

1 {"type":"Sensor","_id":1,"_outE":[],"_inE":[{"_id":0,"_outV":2,"_label":"
   ROUTE_ROUTEDEFINITION"}, {"_id":3,"_outV":4,"_label":"TRACKELEMENT_SENSOR"}]}
2 {"type":"Route","_id":2,"_outE":[{"_id":0,"_inV":1,"_label":"ROUTE_ROUTEDEFINITION"}, {"
   _id":1,"_inV":3,"_label":"ROUTE_SWITCHPOSITION"}], "_inE":[]}
3 {"type":"SwitchPosition","_id":3,"_outE":[{"_id":2,"_inV":4,"_label":"
   SWITCHPOSITION_SWITCH"}], "_inE":[{"_id":1,"_outV":2,"_label":"ROUTE_SWITCHPOSITION"}
   ]}
4 {"type":"Switch","_id":4,"_outE":[{"_id":3,"_inV":1,"_label":"TRACKELEMENT_SENSOR"}, "
   _inE":[{"_id":2,"_outV":3,"_label":"SWITCHPOSITION_SWITCH"}]}

```

Listing A.3. A graph based on the railway system metamodel stored in Faunus GraphSON format

A.2 Semantic Graph Formats

A.2.1 RDF/XML

RDF/XML is an XML-based (Extensible Markup Language) format for serializing RDF triples.

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <rdf:RDF
3   xmlns="http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl#"
4   xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
5   xmlns:swrl="http://www.w3.org/2003/11/swrl#"
6   xmlns:swrlb="http://www.w3.org/2003/11/swrlb#"
7   xmlns:xsd="http://www.w3.org/2001/XMLSchema#"
8   xmlns:owl="http://www.w3.org/2002/07/owl#"
9   xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
10
11 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
12   TrainRequirementOntology.owl">
13   <rdf:type rdf:resource="http://www.w3.org/2002/07/owl#Ontology"/>
14 </rdf:Description>
15
16 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
17   TrainRequirementOntology.owl#Segment">
18   <rdf:type rdf:resource="http://www.w3.org/2002/07/owl#Class"/>
19   <rdfs:subClassOf rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
20     TrainRequirementOntology.owl#Trackelement"/>
21 </rdf:Description>
22
23 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
24   TrainRequirementOntology.owl#Switch">
25   <rdf:type rdf:resource="http://www.w3.org/2002/07/owl#Class"/>
26   <rdfs:subClassOf rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
27     TrainRequirementOntology.owl#Trackelement"/>
28 </rdf:Description>
29
30 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
31   TrainRequirementOntology.owl#1">
32   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
33     TrainRequirementOntology.owl#Sensor"/>
34 </rdf:Description>
35
36 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
37   TrainRequirementOntology.owl#2">
38   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
39     TrainRequirementOntology.owl#Route"/>
40 </rdf:Description>
41
42 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
43   TrainRequirementOntology.owl#3">
44   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
45     TrainRequirementOntology.owl#Switch"/>
46 </rdf:Description>
47
48 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
49   TrainRequirementOntology.owl#4">
```

```

38 <rdf:type rdMapping Ecoref:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#SwitchPosition"/>
39 </rdf:Description>
40
41 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#3">
42 <TrackElement_sensor rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#1"/>
43 </rdf:Description>
44
45 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#4">
46 <SwitchPosition_switch rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#3"/>
47 </rdf:Description>
48
49 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#2">
50 <Route_routeDefinition rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#1"/>
51 <Route_switchPosition rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#4"/>
52 </rdf:Description>
53
54 </rdf:RDF>

```

Listing A.4. *A graph based on the railway system metamodel stored in RDF format*

A.3 Mapping Ecore to Property Graphs

Mapping the Ecore kernel's concepts to property graphs is not a trivial task. We developed the property graph generator module for the Train Benchmark based on the railroad system's Ecore metamodel (Section 3.7.1), which meant the Ecore concepts had to be mapped to property graphs. Following the mapping defined in Section 2.2.1, we created the equivalent instance models for property graphs as well. Below, we provide some examples about the mapping:

- `Segment` is an `EClass` instance. In a property graph, types cannot be represented explicitly. Instead, for each node representing a `Segment` instance, we add a `type` property with the value `Segment`.
- `Segment_length` is an `EAttribute` instance. For each graph node representing a `Segment`, we define a property with the value `Segment_length`.
- `TrackElement_Sensor` is an `EReference` instance. For each edge representing a `TrackElement_Sensor` instance, we add the `TRACKELEMENT_SENSOR` label.

- `EInt` in an `EDataType` instance. Each attribute with this type, e.g. the `Sensor` class' `Segment_length` attribute, is defined with the Java primitive type `int`.

A.4 Train Benchmark Queries

In the following, we present the queries defined in the Train Benchmark. These queries were widely used during the development of INCQUERY-D, both for functional tests and performance benchmarks of different query engines, including INCQUERY-D itself.

We describe the semantics and the goal of each query. We also show the associated graph pattern and relational algebra query. The metamodel of the railroad system is shown in Figure 3.9.

A.4.1 Relational Schemas

For formulating the queries in relational algebra we define the following relational schemas for representing the vertices (objects) in the graph (instance model).

- *Route* (*id*)
- *Sensor* (*id*, *Segment_length*)
- *Signal* (*id*)
- *Switch* (*id*)
- *SwitchPosition* (*id*)
- *TrackElement* (*id*)

The edges (relationships) are represented with the following relational schemas:

- *Route_entry* (*Route*, *Signal*)
- *Route_exit* (*Route*, *Signal*)
- *Route_switchPosition* (*Route*, *SwitchPosition*)
- *Route_routeDefinition* (*Route*, *Sensor*)
- *SwitchPosition_switch* (*SwitchPosition*, *Switch*)
- *TrackElement_sensor* (*Switch*, *Sensor*)
- *TrackElement_connectsTo* (*TrackElement*, *TrackElement*)

A.4.2 Graph Patterns

Blue rectangles and arrows mark simple constraints, while red rectangles and arrows represent negative application conditions. The query returns the nodes in hollow blue rectangles. Additional constraints (e.g. arithmetic comparisons) are shown in the figure in text.

A.4.3 PosLength

Description

The *PosLength* well-formedness constraint requires that a segment must have positive length. Therefore, the query (Figure A.2) checks for segments with a length less than or equal to zero. The SPARQL representation of the query is shown in Listing A.5.

Goal

The query checks whether an object has an attribute. If it does, the value is checked. Checking attributes is a real world use case, although a very simple one. Note that simple string checking is also measured in the Berlin SPARQL Benchmark [37], and it concludes that the string comparison algorithm dominates the query time.

```
1 PREFIX base: <http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
5
6 SELECT DISTINCT ?xSegment
7 WHERE
8 {
9     ?xSegment rdf:type base:Segment .
10    ?xSegment base:Segment_length ?xSegment_length .
11
12    FILTER (?xSegment_length <= 0)
13 }
```

Listing A.5. The *PosLength* query in SPARQL

Relational algebraic form

The *PosLength* query can be formalized in relational algebra as:



Figure A.2. The *PosLength* query's pattern

$$\pi_{Sensor_id} (\sigma_{Segment_length \leq 0} (Sensor))$$

A.4.4 RouteSensor

The *RouteSensor* query is discussed in Section 3.7.1.

A.4.5 SignalNeighbor

Description

The *SignalNeighbor* well-formedness constraint requires that routes that are connected through sensors and track elements have to belong to the same signal. Therefore, the query (Figure A.3) checks for routes which have an exit signal and a sensor connected to another sensor (which is in a definition of another route) by two track elements, but there is no other route that connects the same signal and the other sensor. The SPARQL representation of the query is shown in Listing A.6.

Goal

This pattern checks for the absence of circles, so the efficiency of the join operation is tested. One-way navigable references are also present in the constraint, so the efficient evaluation of these are also tested. Subsumption inference is required, as the two track elements can be switches or segments.

```

1 PREFIX base: <http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl
  #>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
5
6 SELECT DISTINCT ?xRoute1
7 WHERE
8 {
9   ?xRoute1 rdf:type base:Route .
10  ?xSensor1 rdf:type base:Sensor .
  
```

```

11  ?xSensor2 rdf:type base:Sensor .
12  ?xSignal rdf:type base:Signal .
13  ?xTrackElement1 rdf:type base:Trackelement .
14  ?xTrackElement2 rdf:type base:Trackelement .
15
16  ?xRoute1 base:Route_exit ?xSignal .
17  ?xRoute1 base:Route_routeDefinition ?xSensor1 .
18  ?xTrackElement1 base:TrackElement_sensor ?xSensor1 .
19  ?xTrackElement1 base:TrackElement_connectsTo ?xTrackElement2 .
20  ?xTrackElement2 base:TrackElement_sensor ?xSensor2 .
21
22  ?xRoute3 rdf:type base:Route .
23  ?xRoute3 base:Route_routeDefinition ?xSensor2 .
24  FILTER ( ?xRoute3 != ?xRoute1 )
25
26  OPTIONAL {
27      ?xRoute2 rdf:type base:Route .
28      ?xRoute2 base:Route_entry ?xSignal .
29      ?xRoute2 base:Route_routeDefinition ?xSensor2 .
30  } .
31  FILTER (!BOUND(?xRoute2))
32 }

```

Listing A.6. *The SignalNeighbor query in SPARQL*

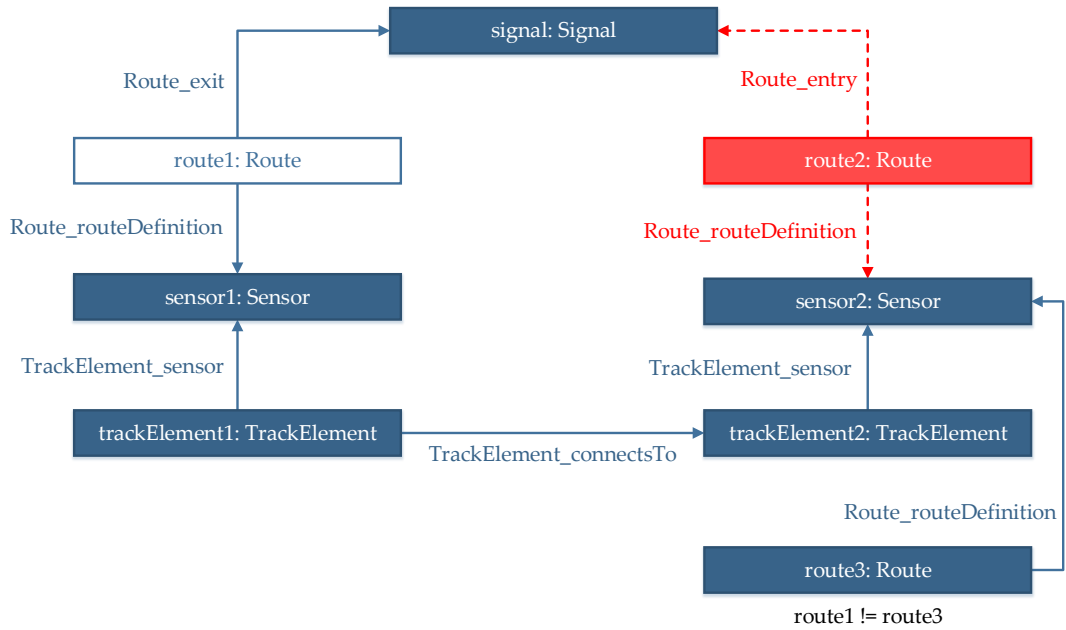


Figure A.3. *The SignalNeighbor query's pattern*

Relational algebraic form

The *SignalNeighbor* query can be formalized in relational algebra as:

$$\pi_{Route_entry.Route} \left(\sigma_{Route_entry.Route \neq Route_routeDefinition_2.Route} \left(\begin{aligned} &Route_entry \bowtie Route_routeDefinition_1 \bowtie TrackElement_sensor_1 \bowtie \\ &TrackElement_connectsTo \bowtie TrackElement_sensor_2 \bowtie Route_routeDefinition_2 \triangleright \\ &(Route_exit \bowtie Route_routeDefinition_3) \end{aligned} \right) \right)$$

A.4.6 SwitchSensor

Description

The *SwitchSensor* well-formedness constraint requires that every switch must have at least one sensor connected to it. Therefore, the query (Figure A.4) checks for switches that have no sensors associated with them. The SPARQL representation of the query is shown in Listing A.7.

Goal

This query checks whether an object is connected to a relation. This pattern is common in more complex queries, e.g. it is used the *RouteSensor* and the *Signal-NearNeighbor* queries.

```

1 PREFIX base: <http://www.semanticweb.org/ontologies/2011/1/TrainRequirementOntology.owl
  #>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
5
6 SELECT DISTINCT ?xSwitch
7 WHERE
8 {
9   ?xSwitch rdf:type base:Switch .
10
11   OPTIONAL {
12     ?xSensor rdf:type base:Sensor .
13     ?xSwitch base:TrackElement_sensor ?xSensor .
14   } .
15   FILTER (!BOUND (?xSensor))
16 }
```

Listing A.7. The *RouteSensor* query in SPARQL

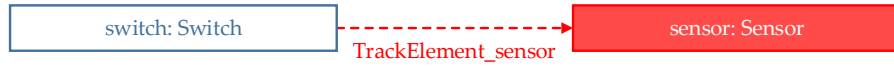


Figure A.4. *The SwitchSensor query's pattern*

Relational algebraic form

The *SwitchSensor* query can be formalized in relational algebra as:

$$Switch \triangleright Sensor$$