



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

Incremental Graph Queries in the Cloud

SCIENTIFIC STUDENTS' ASSOCIATIONS REPORT

Author

Gábor Szárnyas

Supervisors

Dr. István Ráth

Benedek Izsó

Dr. Dániel Varró

October 21, 2013



Contents

1	Introduction	6
1.1	Context	6
1.2	Problem statement, requirements	7
1.3	Objectives	7
1.4	Contribution, added value	7
1.5	Structure of the report	7
2	Background technologies	8
2.1	Big Data and the NoSQL movement	8
2.2	Graph models and storage methods	9
2.2.1	TinkerPop	10
2.2.2	Triplestores	11
2.3	NoSQL technologies	11
2.3.1	Cassandra	12
2.3.2	Hadoop	13
2.3.3	Neo4j	15
2.4	Graph technologies	15
2.4.1	Titan	15
2.4.2	4store	17
2.5	Asynchronous messaging	18
2.5.1	Akka	18
2.6	Eclipse-based technologies	19
2.6.1	EMF	19
2.6.2	EMF-INCQUERY	20
3	Overview of the approach	22
3.1	Case study: TrainBenchmark	22
3.1.1	Domain	23
3.1.2	Queries	23
3.2	Architecture overview	25

3.2.1	Rete in general	25
3.2.2	INCQUERY-D architecture	26
3.3	Initialization and indexing	26
3.3.1	Indexing	26
3.3.2	Graph-like data manipulation	27
3.3.3	Notification mechanisms	27
3.4	Incremental queries and change propagation	28
3.4.1	Distributed storage layer	28
3.4.2	Distributed Rete operation and scalability challenges	28
3.5	Deployment and configuration	30
3.5.1	Tooling	30
3.5.2	Degrees of freedom	31
3.5.3	Workflow	31
3.6	Elaboration of the example	33
3.6.1	Workflow	33
4	Evaluation of scalability and performance	37
4.1	Goals	37
4.1.1	Phases	37
4.1.2	Measure the response time and the scalability	38
4.1.3	Workload profile's difference from standard benchmarks	38
4.2	Benchmark scenario	38
4.3	Generation of models	38
4.4	Benchmark environment	39
4.4.1	Benchmark setup	39
4.4.2	Hardware, software ecosystem	40
4.4.3	Benchmark methodology	40
4.4.4	Data collection	40
4.4.5	Data processing tools	40
4.5	Results	41
4.5.1	Results using the Neo4j graph database	41
5	Related work	45
5.1	Eclipse-based tools	45
5.2	Semantic web and NoSQL	45
5.3	Rete implementations	46
6	Conclusions	47
6.1	Summary of contributions	47
6.1.1	Own work	47

6.2	Limitations	48
6.3	Future work	48
	Bibliography	50
A	Graph formats	55
A.1	Property graph formats	55
A.1.1	GraphML	55
A.1.2	Blueprints GraphSON	56
A.1.3	Faunus GraphSON	57
A.2	Semantic graph formats	57
A.2.1	RDF	57

Chapter 1

Introduction

Nowadays, model-driven software engineering (MDE) plays an important role in the development processes of critical embedded systems. 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. With the dramatic increase in complexity that 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.

1.1 Context

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 [28]) or (ii) applying impact analysis techniques and re-evaluating queries only in contexts that are affected by a change (e.g. the Eclipse OCL Impact Analyzer [35]). 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 [53]). Since single-computer heaps cannot grow arbitrarily (as response times degrade drastically due to garbage collection problems), memory consumption is the most significant scalability limitation.

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 [47, 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 [36]) has crucial conceptual and technological differences to traditional meta-modeling languages such as Ecore [53]. Additionally, while there are initial efforts to overcome the mapping issues between the MDE and Linked Data worlds [39], even the most sophisticated NoSQL storage technologies lack efficient and mature support for executing expressive queries incrementally.

1.2 Problem statement, requirements

[50]

1.3 Objectives

We aim to address these challenges by adapting incremental graph search techniques from EMF-INCQUERY to the cloud infrastructure. We introduce INCQUERY-D, a prototype system based on a distributed Rete network [33] that can scale up from a single-workstation tool to a cluster to handle very large models and complex queries efficiently (??). We carry out an initial performance evaluation in the context of on-the-fly well-formedness validation of software design models (??), discuss related work in ?? and conclude the thesis in ??.

1.4 Contribution, added value

1.5 Structure of the report

Chapter 2

Background technologies

Implementing a scalable graph pattern matcher requires a wide range of technologies. Careful selection of the technologies is critical to the project's success. For INCQUERY-D, we looked for technologies that can form the building blocks of a distributed, scalable model repository and pattern matcher. These technologies must be designed with scalability in mind and deployed at scale successfully. To avoid licensing issues, our search criteria included that, if possible, the technologies should be free and open-source solutions.



During the early phase of the 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 systems' support for asynchronous processing, notification and messaging mechanisms.

While we did not disqualify closed-source operating systems, it quickly became apparent that most tools are best tailored for Unix-like operating systems, e.g. systems based on the Linux kernel.

2.1 Big Data and the NoSQL movement

Since the 1980s, database management systems based on the relational data model [30] 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 [42, 51]. In practice, this renders 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 scien-

tific 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 [15]. During the development of INCQUERY-D's propotype, we experimented with numerous NoSQL databases.

2.2 Graph models and storage methods

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.

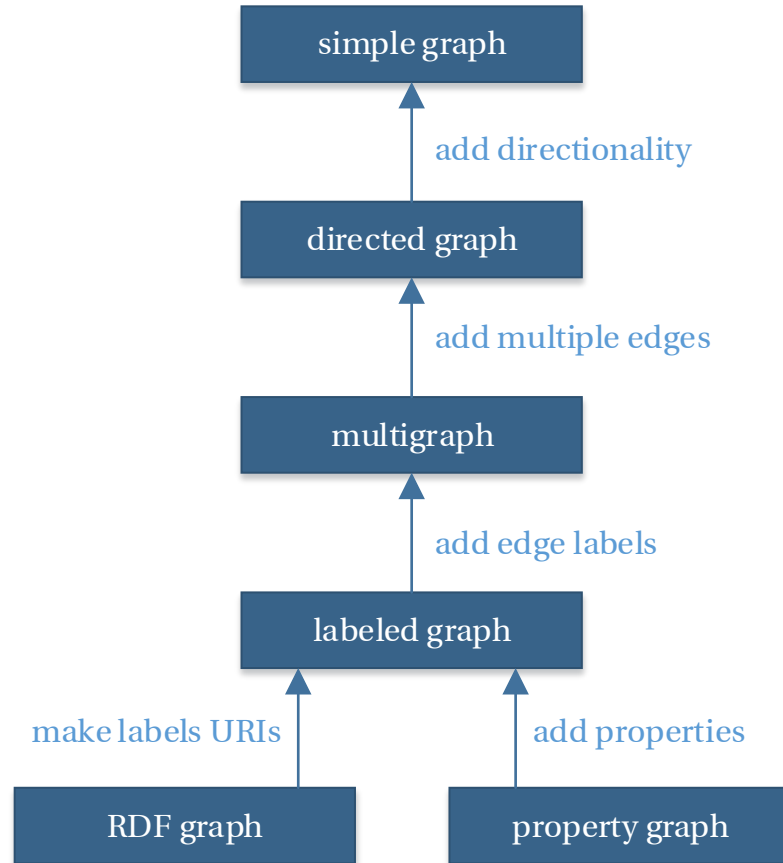


Figure 2.1: Different graph data models (based on [49])

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

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

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*), allow loops and multiple edges (*multi-graph*). To allow different connections, we may label the edges (*labeled graph*). *RDF graphs* use URIs (Uniform Resource Identifiers) instead of labels, otherwise they have similar expressive power as labeled graphs. *Property 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`. Property graphs are powerful enough to describe Java objects or EMF instance models (subsection 2.6.1).

2.2.1 TinkerPop

The *TinkerPop* framework is an open-source software stack for graph storage and processing [23]. TinkerPop includes *Blueprints*, a property graph model interface. Blueprints fulfills 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 (subsection 2.3.3), Titan (subsection 2.4.1), DEX [20], InfiniteGraph [16] and OrientDB [17].

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. 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

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 on Figure 2.2.

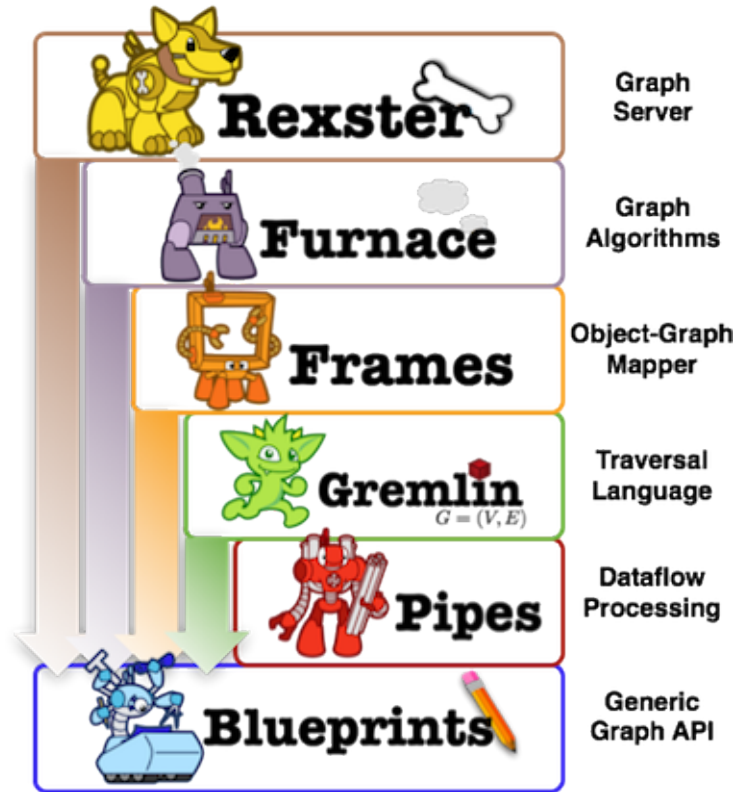


Figure 2.2: The TinkerPop software stack [10]

2.2.2 Triplestores

Triplestores are tailored to store and process triples efficiently. A triple is a data entity composed of a subject, a predicate and an object, e.g. "John instanceof Person", "John is 34". Triplestores are mostly used in semantic technology projects. Also, some triplestores are capable of *reasoning*, i.e. inferring logical consequences from a set of facts or axioms.

Triplestores use the RDF (Resource Description Framework) data model. Although the RDF 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 a RDF. Triplestores are usually queried via the RDF format's query language, SPARQL (recursive acronym for SPARQL Protocol and RDF Query Language).

2.3 NoSQL technologies



In the following section we summarize the core technologies used in INCQUERY-D's prototype implementation. We briefly introduce the goals of each technology, with particular emphasis on the scalability aspects.

2.3.1 Cassandra

Cassandra is one of the most widely used NoSQL databases [7]. Originally developed by Facebook [43], Cassandra is now an open-source Apache project. The whole project is written in Java.

Cassandra's data model is called *column family*. 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.

Cassandra has sophisticated fault-tolerance mechanisms. It allows the application to balance between availability and consistency by allowing it to tune the 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 [18].

Consistent hashing

To distribute the data across the cluster, Cassandra uses a partitioner mechanism. The basic partitioners distribute the rows evenly based on their key's hash value. In a cluster with n nodes, the naïve method for determining the location for a row with key x is computing $h(x) \bmod n$, where $h(x)$ is the hash function. Currently, Cassandra provides partitioners based on the MD5 and the Murmur3 hash function. However, this approach has a serious limitation: if we remove or add nodes to the cluster, we have to recompute the hash values and possibly relocate almost all rows in the cluster. To avoid this, Cassandra uses a special kind of hashing called *consistent hashing*.

The ring is divided into ranges equal to the number of nodes, with each node being responsible for one or more ranges of the data. Before a node can join the ring, it must be assigned a *token*. The token value determines the node's position in the ring and its range of data. The ring is walked clockwise until it locates the node with a token value greater than that of the row key. With the nodes sorted in token order, the last node is considered the predecessor of the first node; hence the ring representation² [5].

For example, consider a simple four-node cluster (Figure 2.3), where all of the row keys managed by the cluster are in the range of 0 to 100. Each node is assigned a token that represents a point in this range. In this example, the token values are 0, 25, 50 and 75. The first node (with token 0), is responsible for the wrapping range (76+), the

²Note "consistent" here is different from both the idea of consistency in data consistency and in the ACID (atomicity, consistency, isolation, durability) properties guaranteed by transactions. It refers to the fact that tries to map the same rows to the same machine, even if the number of machines (n) changes over time slightly.

second node (with token 25) is responsible for the data range 1 – 25, and so on [5].

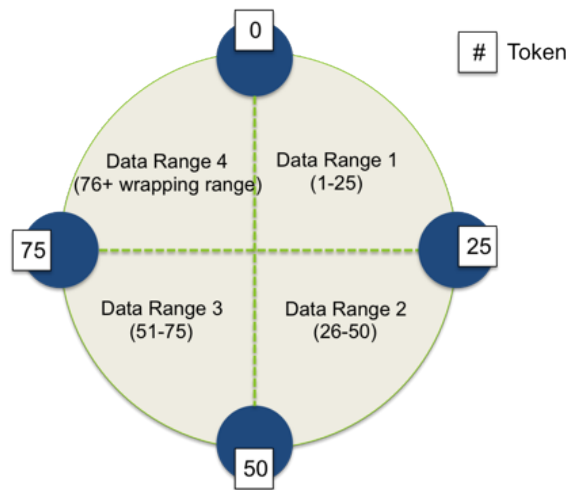


Figure 2.3: *Cassandra's ring for data partitioning [5]*

2.3.2 Hadoop

Hadoop is an open-source, distributed data processing framework inspired by Google's publications about MapReduce [31] and the Google File System [34]. Originally developed at Yahoo!, Hadoop is now an Apache project [8]. 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). All software in the Hadoop framework is written in Java.

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.

Hadoop is often used for sorting, filtering and aggregating data sets. It is also used for fault-tolerant, distributed task execution.

A typical small 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 2.4).

HDFS

The Hadoop Distributed File System (HDFS) is an open-source, distributed file system, inspired by the Google File System and written specifically for Hadoop [8]. Unlike other distributed file systems (e.g. Lustre [13]), which require expensive hardware components, HDFS was designed to run on commodity hardware.

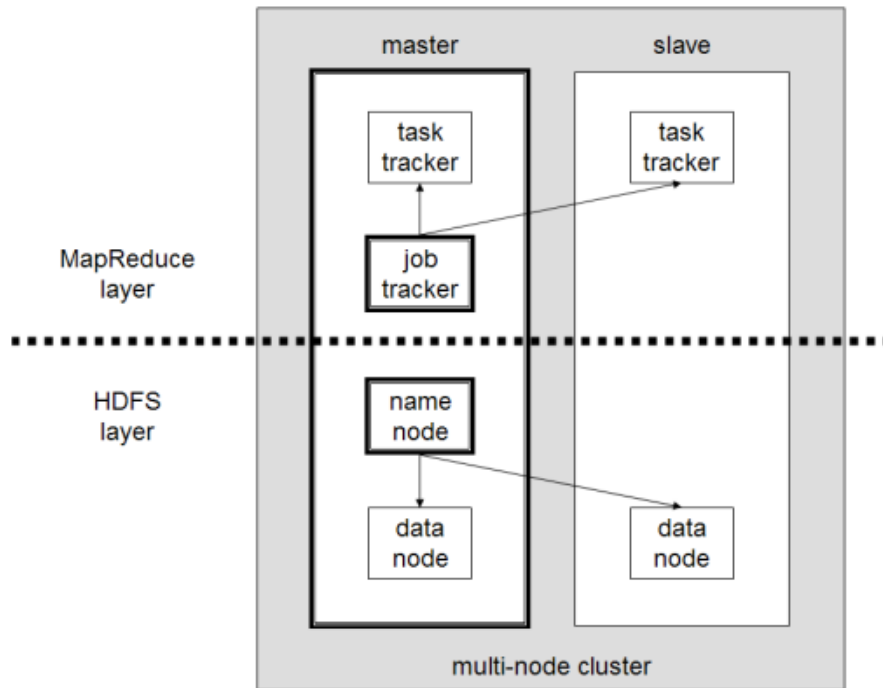


Figure 2.4: Hadoop's architecture [44]

HDFS tightly integrates with Hadoop's architecture (Figure 2.5). The *NameNode* is responsible for storing the metadata of the files and the location of the replicas. The data is stored by the *DataNodes*.

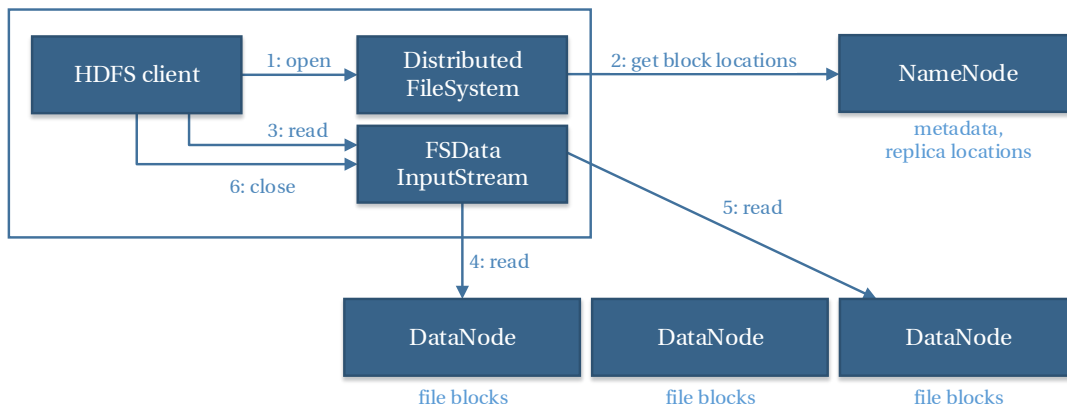


Figure 2.5: HDFS' architecture

HBase

HBase is an open-source, 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 run in Hadoop.

2.3.3 Neo4j

Neo4j, developed by Neo Technology, is the most popular NoSQL graph database. Neo4j implements TinkerPop's Blueprints property graph data model along with Gremlin. It also provides Cypher, a declarative query language for graph pattern matching.

Neo4j is one of the most mature NoSQL databases. It is well documented and provides ample tooling. However, its scalability features are limited: instead of sharding, it only supports replication of data to create a highly available cluster. Of course, the scalability limitations are a hot topic in Neo4j's development. Neo4j's developers make serious efforts to improve Neo4j's scalability in an ongoing project called Rassilon [3].

Neo4j is capable of loading graphs from GraphML [21] and Blueprints GraphSON [12] formats (see section A.1 for examples). Neo4j graphs can be visualized in Neoclipse, an Eclipse RCP application [14]. A part of a TrainBenchmark (section 3.1) instance model is shown on Figure 2.6.

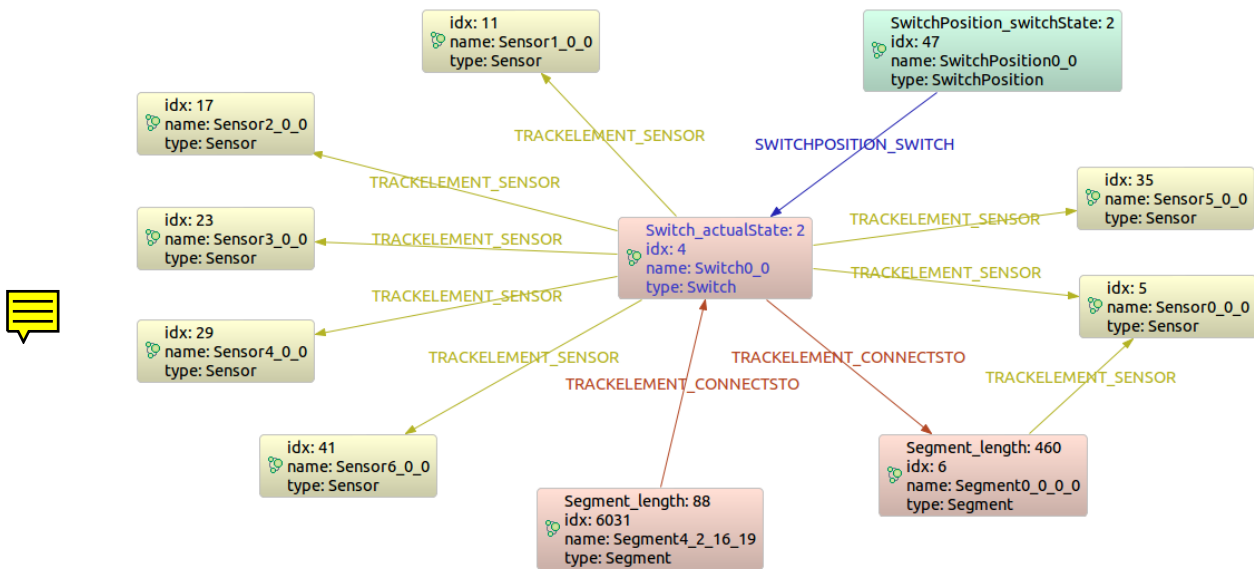


Figure 2.6: TrainBenchmark subgraph visualized in Neoclipse

2.4 Graph technologies

2.4.1 Titan

Titan is an open-source, distributed, scalable graph database from Aurelius, the creators of the TinkerPop framework. Unlike Neo4j, 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.

Mapping and sharding

To store the graph, Titan maps each vertex to a row of a column family (Figure 2.7). The row stores the identifier and the properties of the vertex, along both the incoming and outgoing edges' identifiers, labels and properties.

Titan uses the storage backend's partitioner (e.g. Cassandra's 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 October 2013, but not yet available.

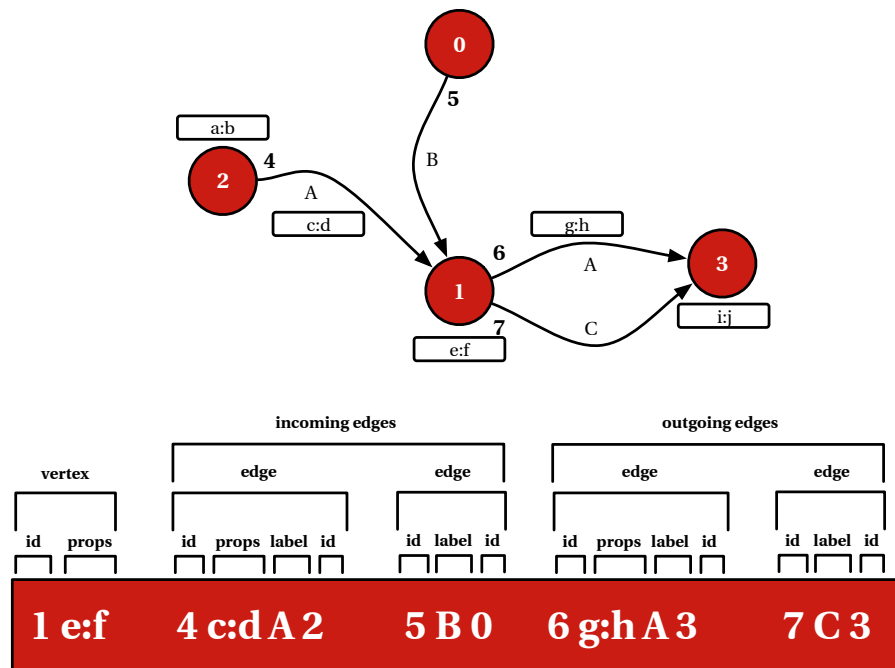


Figure 2.7: Titan graph vertex stored in Cassandra as a row

Deployment

Titan can be deployed in different ways according to the needs of the application. For INCQUERY-D's prototype, we used Titan in *remote server mode* (Figure 2.8). In this setup, Titan runs in the same Java Virtual Machine as the application and communicates with the Cassandra cluster on a low-level protocol (e.g. Thrift).

Faunus

Although Titan was designed with scalability in mind, its query engine does not work in a parallel fashion. Also, it 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.

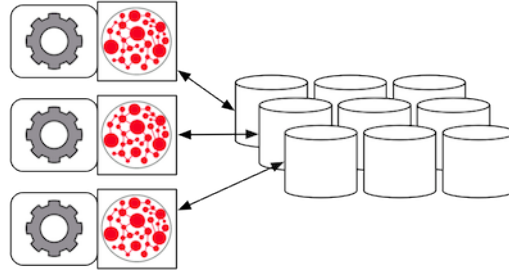


Figure 2.8: Using Titan with Cassandra in remote server mode

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 subsection A.1.3 for an example.

It is important to note that Faunus always traverses the whole graph and does not use its indices. This makes it slow for retrieving nodes or edges by type (see our typical workload in subsection 3.3.1).

2.4.2 4store

4store is an open-source, distributed triplestore created by Garlik [4]. Unlike the other tools discussed earlier, 4store is written in C. While 4store is primarily applied for semantic web projects, its maturity and scalability made it an appropriate storage back-end for INCQUERY-D's prototype.

Similar to Titan's partitioning, 4store's sharding mechanism (called *segmenting*) distributes the RDF resources evenly across the cluster. However, unlike Titan, 4store's data model is an RDF graph. Hence, 4store's input format is RDF/OWL.

Table 2.1: Overview of database technologies

Technology	Data model	Sharding	Distributed operation	DML facility	Identifier generation
4store	RDF	Automatic	Manual	SPARQL	Manual
Neo4j	Property graph	Manual	Manual	Cypher	Manual
Titan	Property graph	Automatic	Automatic	Gremlin	Automatic

Table 2.1 summarizes the relevant characteristics of the aforementioned database management systems. 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.

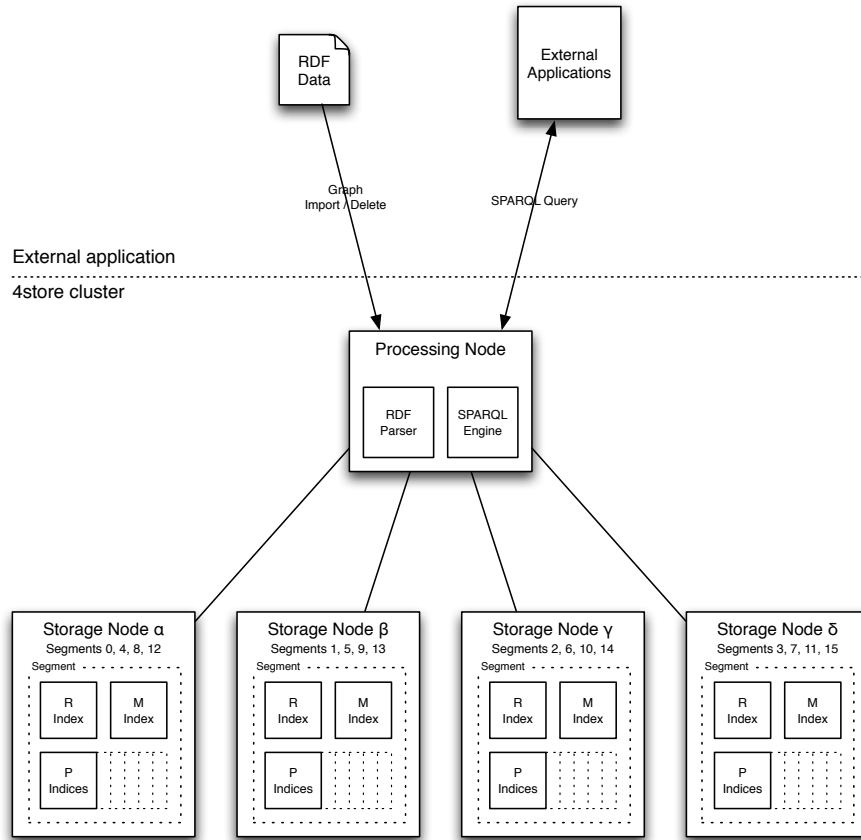


Figure 2.9: 4store's distributed architecture [37]

2.5 Asynchronous messaging

Most distributed, concurrent systems use a messaging framework or message queue service. Because of the nature of the Rete algorithm (subsection 3.2.1), INCQUERY-D requires a distributed, asynchronous messaging framework.

2.5.1 Akka

Akka is an open-source, fault-tolerant, distributed, asynchronous messaging framework developed by Typesafe [6]. 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 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 2.10). This way, they can be transparently moved between machines on the network.

As of October 2013, the latest version (Akka 2.2) also features *pluggable transport support* to use various transports to communicate with remote systems [6]. For se-

realizing the messages, Akka supports different frameworks, including Java's built-in serialization, Google Protobuf [19] and Thrift [9].

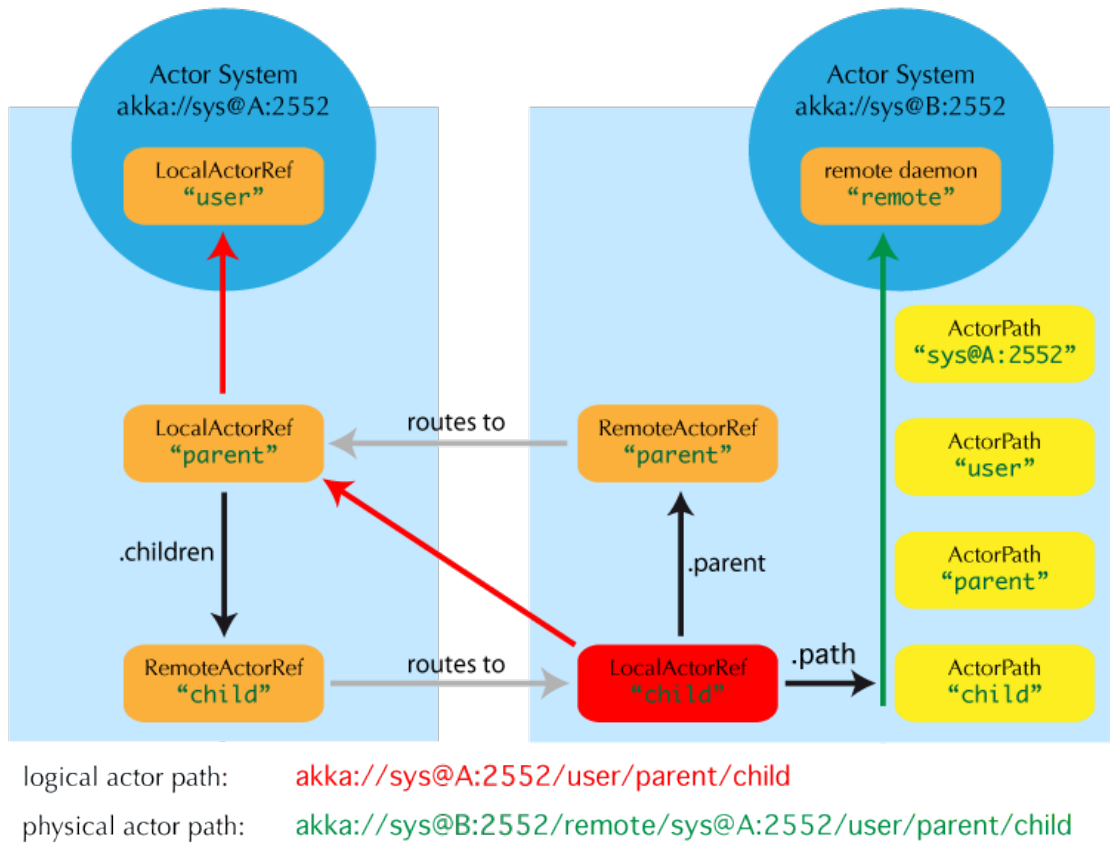


Figure 2.10: Deploying a remote actor in Akka



2.6 Eclipse-based technologies

Eclipse is a free, open-source software development environment and a platform for plug-in development. Members of the Eclipse Foundation include industry giants like IBM, Intel, Google and SAP.

INCQUERY-D's single workstation predecessor, EMF-INCQUERY is built around Eclipse-based technologies. To reap the benefits of a mutual code base, we designed INCQUERY-D to use as much of EMF-INCQUERY's components as possible. In the following section, we introduce the Eclipse-based technologies most important for our work.

2.6.1 EMF

Eclipse comes with its own modeling technologies called EMF (Eclipse Modeling Framework). EMF provides a metamodel (Ecore) for designing applications and a code generation facility to produce the Java classes for the model.

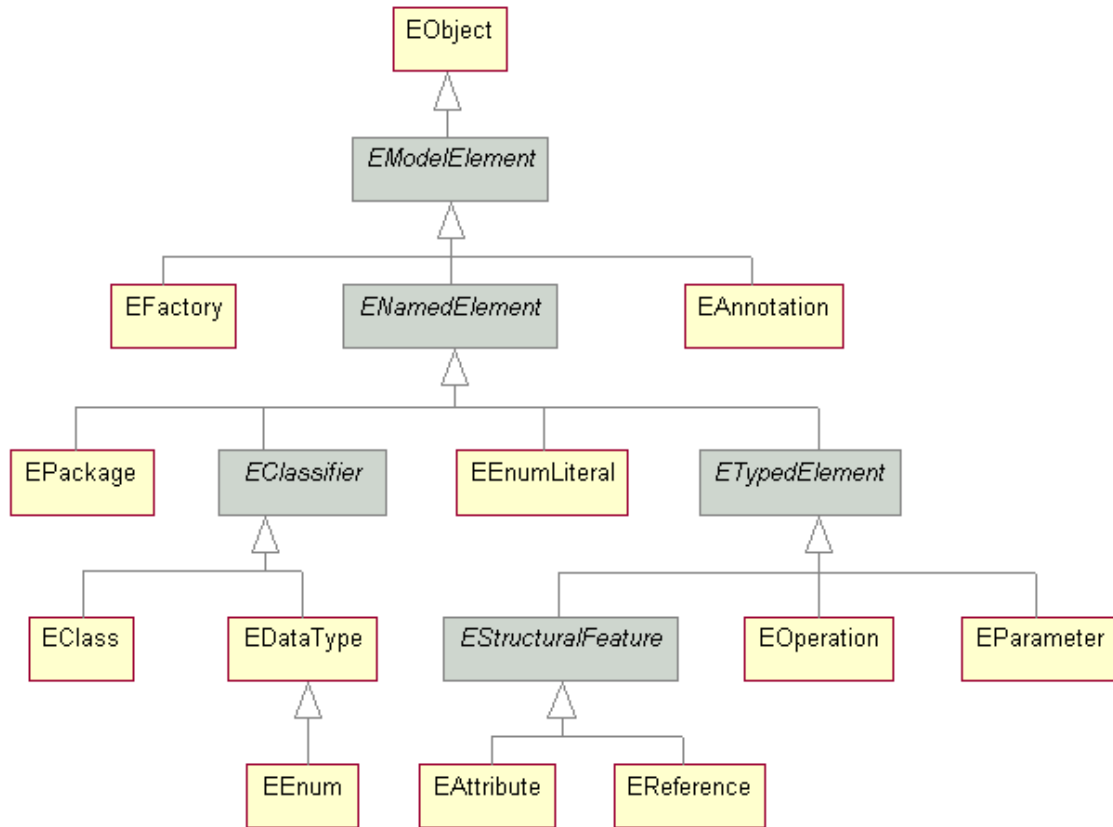


Figure 2.11: *The metamodel of Ecore*

2.6.2 EMF-INCQUERY

EMF-INCQUERY is developed by the Fault Tolerant Systems Research Group (FTSRG) in the Budapest University of Technology and Economics. EMF-INCQUERY is an open-source Eclipse project which provides incremental query evaluation on EMF models. The queries (graph patterns) are defined in INCQUERY Pattern Language (IQPL), a domain-specific language implemented in Xtext and Xtend. EMF-INCQUERY's architecture is shown on Figure 2.12.

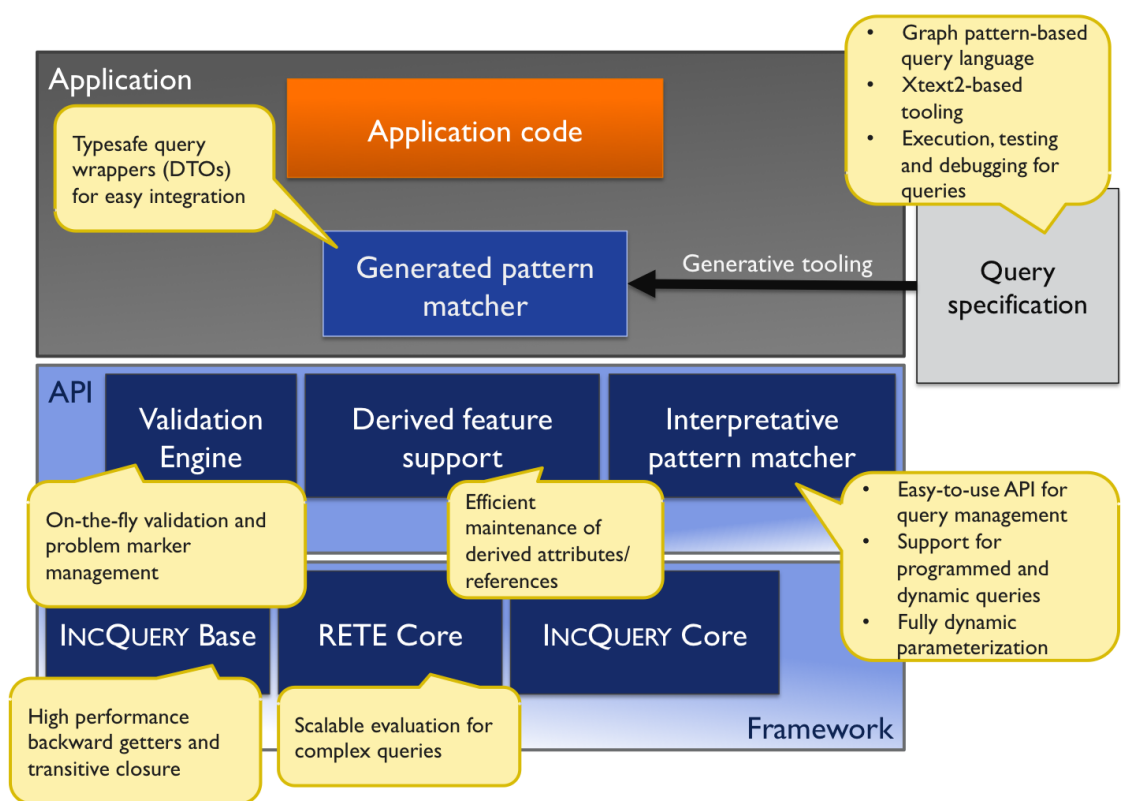


Figure 2.12: EMF-INCQUERY's architecture

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 report is an architecture ~~and an implementation~~ 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 Case study: TrainBenchmark

Due to both confidentiality and technical reasons, it is difficult to obtain real-world industrial models and queries. Also, using confidential data sets hampers the reproducibility of the conducted benchmarks. Therefore, we used an artificial data set which mimics real-world models.

In the following section we present the *TrainBenchmark*, a benchmark scheme and environment. The TrainBenchmark was designed and implemented by Benedek Izsó, István Ráth and Zoltán Szatmári [41]. The original goal of the TrainBenchmark was to compare various incremental and non-incremental query engines' performance to EMF-INCQUERY's. For INCQUERY-D, we used an improved version of the TrainBench-

mark, which also works in a distributed environment.

3.1.1 Domain

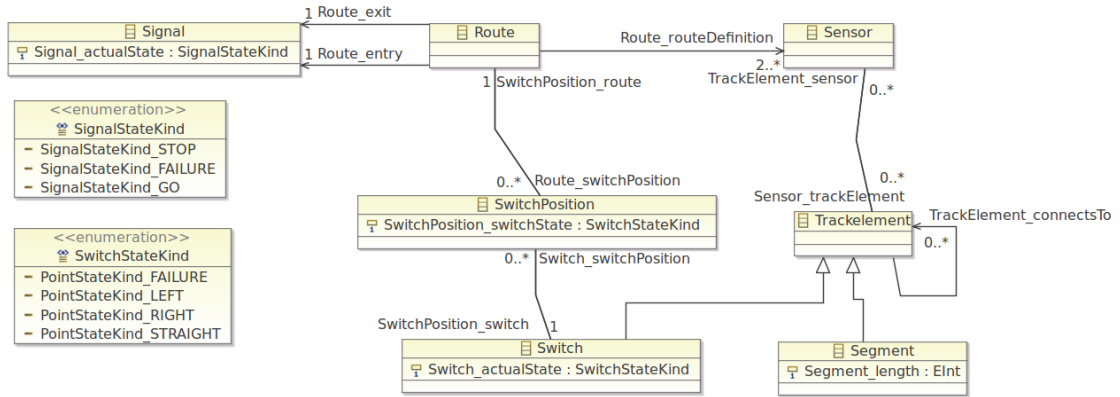


Figure 3.1: The EMF metamodel of the TrainBenchmark

The TrainBenchmark models is an ~~imaginary~~ railroad network. The network is composed of typical railroad items, including signals, segments, switches and sensors. The complete EMF metamodel of the TrainBenchmark is shown on Figure 3.1.

The **generator** project of TrainBenchmark is capable of generating railroad instance models of different sizes. It is capable of generating models in different formats, including EMF, OWL, RDF and SQL.

For Neo4j (subsection 2.3.3) and Titan (subsection 2.4.1), we expanded the generator with a module that can generate property graphs based on the TrainBenchmark's meta-model. It supports the GraphML [21], the Blueprints GraphSON [12] and the Faunus GraphSON [11] output formats.

3.1.2 Queries

The TrainBenchmark consists of queries that resemble a typical MDE application's workload. In general, MDE queries are more complex than those used in traditional databases. They often define large patterns with multiple join operations. The TrainBenchmark's queries look for violations of well-formedness constraints in the model. Although the TrainBenchmark defines a total of four queries, in this report we only discuss the *RouteSensor* query in detail.

RouteSensor

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. The graphical repre-

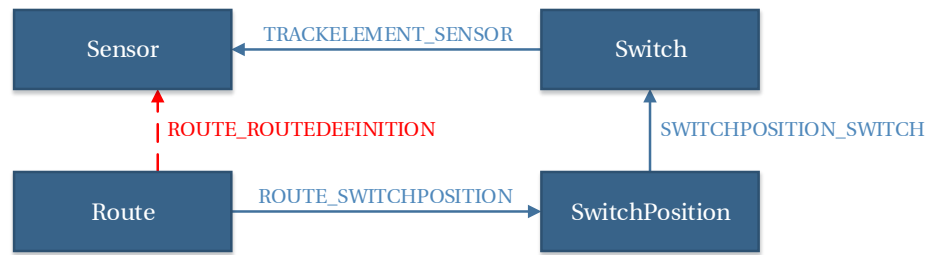


Figure 3.2: Graphical representation of the RouteSensor query's pattern. The dashed red arrow defines a negative condition.

sensation of the RouteSensor query is shown on Figure 3.2. The **IQPL** query (subsection 2.6.2) is shown on Listing 3.1, while the SPARQL query is shown on Listing 3.2.¹

```

1 package hu.bme.mit.train.constraintcheck.inquiry
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

```

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 .
15   ?xSwitch base:TrackElement_sensor ?xSensor .

```

¹Note 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.


```

16
17     FILTER NOT EXISTS {
18         ?xRoute ?Route_routeDefinition ?xSensor .
19     } .
20 }

```

Listing 3.2: *The RouteSensor query in SPARQL*

3.2 Architecture overview

In the following section, we provide an overview of the Rete algorithm, which forms the theoretical basis of EMF-INCQUERY and INCQUERY-D. We also describe INCQUERY-D's architecture.

3.2.1 Rete in general

INCQUERY-D is based on the Rete algorithm, which provides incremental graph pattern matching. The algorithm was originally created by Charles Forgy [33] for expert systems. Gábor Bergmann adapted it for EMF models and added many ~~tweaks and~~ improvements to the algorithm [27].

The Rete algorithm defines an asynchronous network of communicating nodes. This is essentially a dataflow network, with two types of nodes. Change notification objects (*tokens*) are propagated to intermediate *worker nodes* that perform operations, like filtering tokens based on constant expressions and performing join or antijoin operations based on their contents. The worker nodes store partial (interim) query results in their own memory. In contrast, *production nodes* are terminators that provide an interface for fetching query results and also their changes (*deltas*). Connections between nodes can be *local* (within one host) or *remote* (when two Rete nodes are allocated to different hosts).

It is important to emphasize that the database shards and Rete nodes are two distinct levels of distribution that do not directly depend upon each other.

Similar algorithms

Along the original Rete algorithm, many algorithms were developed for incremental pattern matching.

TODO

- variations of Rete: RETE II, RETE*, ...
- TREAT [45]
- LEAPS [26]

3.2.2 INCQUERY-D architecture

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 graph. The client application communicates with the *middleware*. The middleware provides a unified API for accessing the database. It also sends change notifications to the production network and retrieves the query results from the production network. The *production network* is implemented with a distributed Rete net which provides incremental query evaluation.

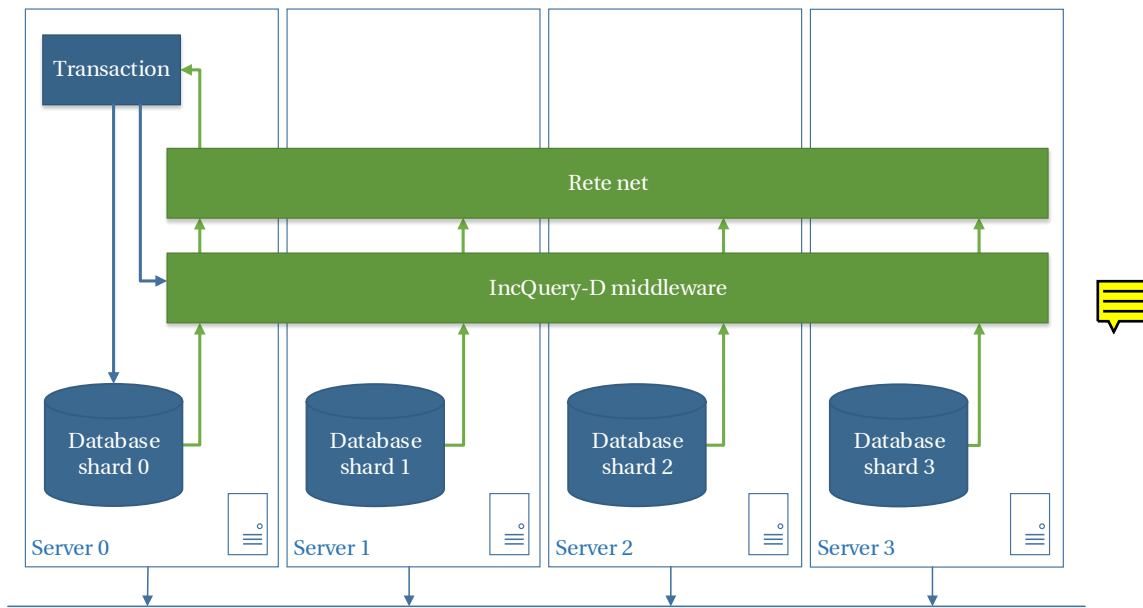


Figure 3.3: INCQUERY-D's architecture on a four-node cluster

The INCQUERY-D architecture in a four-node cluster configuration is shown in Figure 3.3.

3.3 Initialization and indexing

In the following section we will cover the challenges that arise around the indexing and initialization of INCQUERY-D.

3.3.1 Indexing

Indexing is a common technique for decreasing the execution time of database queries. In MDE, *model indexing* is the key to high performance model queries. As MDE primarily uses a metamodeling infrastructure, 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.

3.3.2 Graph-like data manipulation

INCQUERY-D's middleware exposes an API that provides methods to manipulate 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 and forwards it to the underlying data storage (e.g. SPARQL queries for 4store and Gremlin queries for Titan).

Data representation

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 [53] could be supported, as long as they can be mapped to an efficient and distributed storage backend (e.g. triplestores, key-value stores or column-family databases).

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

3.3.3 Notification mechanisms

Model change notifications are required by incremental query evaluation, thus model changes are captured and their effects propagated in the form of *notification objects* (NOs). The notifications generate *tokens* that keep the Rete network's state consistent with the model. INCQUERY-D's middleware layer facilitates notifications by providing a facade for model manipulation operations.



Current database management systems

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 (see section 6.3).

Because the lack of support for distributed notifications, in INCQUERY-D's current implementation, notifications are controlled by the middleware. The notification messages are propagated through the Rete network via the Akka messaging framework.

3.4 Incremental queries and change propagation

Distributed incremental query evaluation introduces a number of challenges. In the following section, we will describe these and present INCQUERY-D's solutions for them.

3.4.1 Distributed storage layer

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

3.4.2 Distributed Rete operation and scalability challenges

The Rete algorithm (subsection 3.2.1) 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 net. To enable this, the system must be able to allocate the Rete nodes to different hosts in a cloud computing infrastructure. 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 overall performance of the system is influenced by a number of factors.

- The *layout of the Rete network*. This can be optimized depending on both query and instance model characteristics, e.g. to keep the resource requirement of intermediate join operations to a minimum.
- The *allocation* of Rete nodes to host computers. The optimization strategy may choose to optimize local resource usage, or to minimize the amount of remote network communication.
- The *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.

Practice

In the prototype of INCQUERY-D, the distributed middleware and Rete network were implemented using Akka (subsection 2.5.1). The communication protocol was built on top of Akka's built-in serialization support.



Rete operation

We present the Rete algorithm's operation on an actual instance model. Suppose we have the graph shown on the top of Figure 3.4 loaded to the Rete net and we decide to delete the ROUTE_ROUTEDEFINITION edge between vertices 2 and 1.

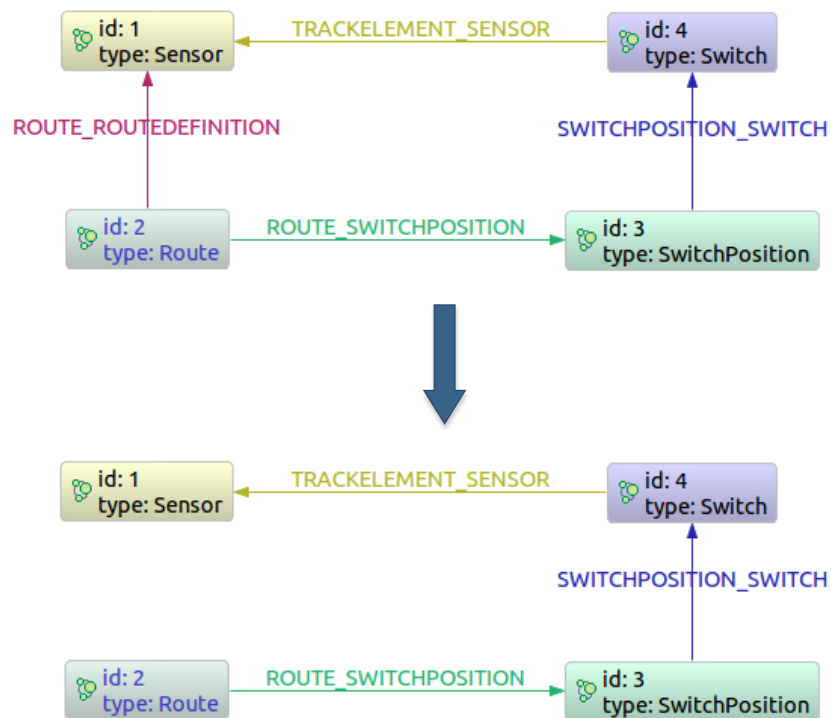


Figure 3.4: A modification on a TrainBenchmark instance model

Figure 3.5 shows the Rete net containing the partial matches of the original graph. When we delete the edge between vertices 2 and 1, the ROUTE_ROUTEDEFINITION type indexer receives a notification from the middleware and sends a *negative update* ① with the tuple (2, 1). The antijoin node processes the negative update and propagates a negative update ② with the tuple (3, 4, 2, 1). **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 "deltas", i.e. only the the tuples that have been added or deleted since the last modification.



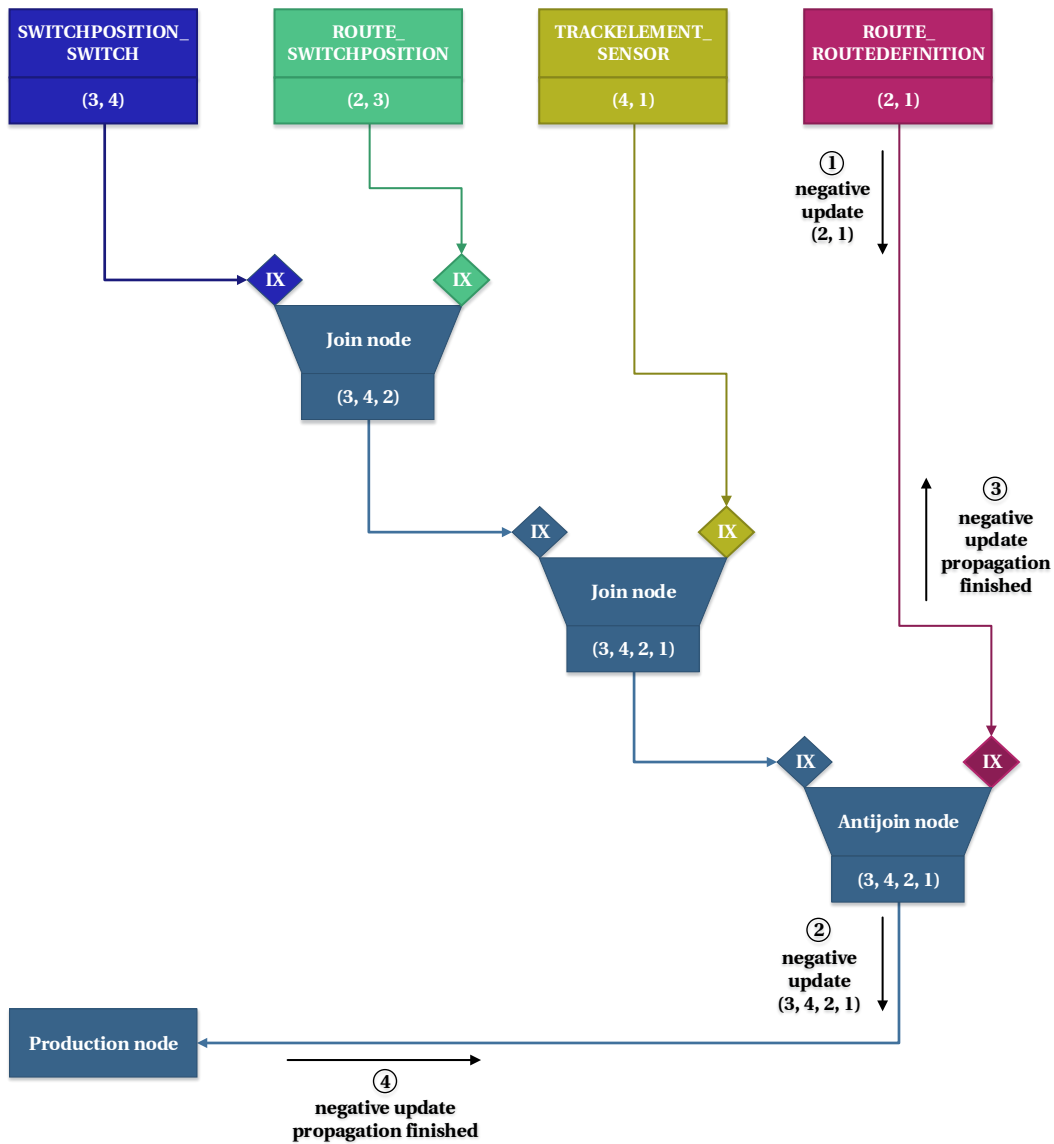


Figure 3.5: The Rete net and the partial matches stored in its nodes

3.5 Deployment and configuration

Deploying, configuring and operating a distributed pattern matcher is a complex task. In the following, we will break down this task to smaller steps and present our tools for solving them.

3.5.1 Tooling

We aimed to build INCQUERY-D on top of EMF-INCQUERY's pattern language (IQPL) and its Rete net generator. For the allocation of Rete nodes, we created an Eclipse-based editor and viewer.

3.5.2 Degrees of freedom

The deployment and configuration of a distributed pattern matcher involves many degrees of freedom:

- We may choose different database implementations due to the INCQUERY-D's backend-agnostic nature. Until now, we experimented with 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).
- We may choose different strategies to allocate the Rete nodes. Note that in theory, this is orthogonal to the database's sharding strategy. 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.

3.5.3 Workflow

In the following part, 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 on Figure 3.6.

Analyze the metamodel and the query

Task. First, we determine the constraints defined by the query pattern. The matches satisfying these constraints will define the results of the query.

Implementation. The pattern is defined in an IQPL (INCQUERY Pattern Language) text file. Using Xtext [24], an Eclipse-based framework for creating domain-specific languages, the textual representation of the pattern is parsed to an EMF model. Based on the EMF model's pattern and the metamodel, a constraint network called *PSystem* (short for *Pattern System*) is generated.

Build a Rete layout

Task. To allow incremental query evaluation, we create a Rete net based on the constraints derived from the query.

Implementation. As we mentioned earlier, we aim to reuse as much of EMF-INCQUERY's existing code base as possible. As part of this attempt, we introduced the concept of *Rete recipes* which define the layout of a Rete network.

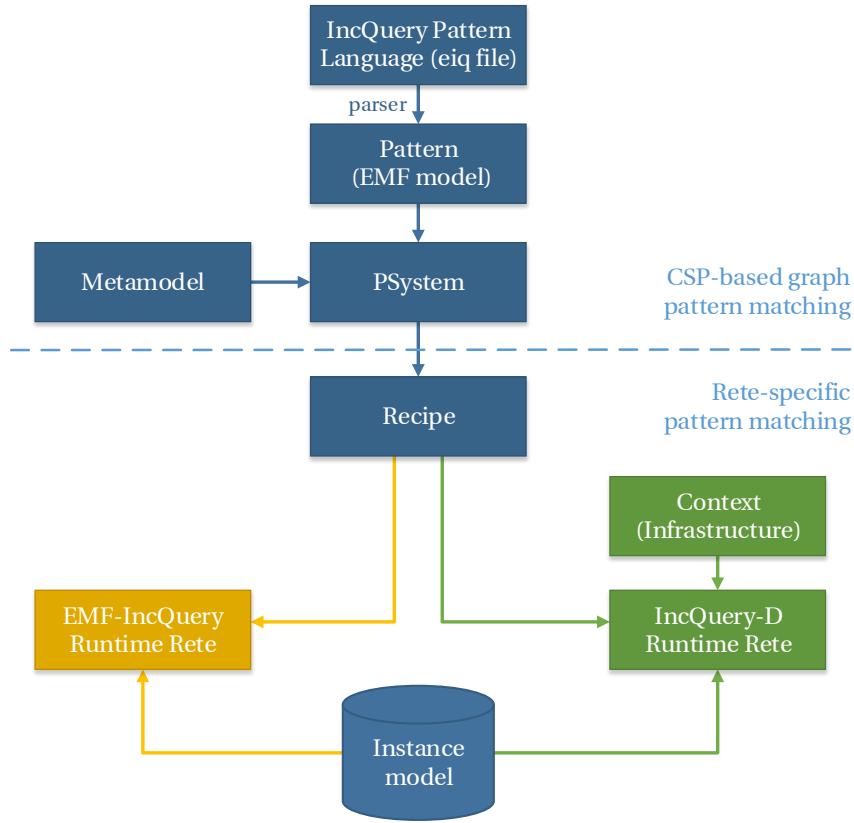


Figure 3.6: The workflow of EMF-INCQUERY (blue) and INCQUERY-D (green)

Allocate the Rete network in the cloud's nodes

Task. Because of its single workstation nature, EMF-INCQUERY simply unfolds the Rete net based on the derived Rete recipe. At the same time, INCQUERY-D operates in a distributed environment where local resource exhaustion, network latency and throughput are critical aspects.

Implementation. Currently, the allocation of the Rete nodes is done manually. To address this limitation, we plan to utilize CSP (Constraint Satisfaction Problem) solvers, or dynamic techniques like DSE (Design Space Exploration) [38].

Bootstrap the system

Task. Based on the Rete network's allocation, we have to deploy the Rete nodes in the distributed systems. After the successful deployment, the Rete network has to be initialized. Due to the Rete algorithm's asynchronous nature, it uses a termination protocol to signal when the data processing is finished.

Implementation. In INCQUERY-D's prototype, both the bootstrapping and the Rete network's operation is carried out automatically. The Akka actors representing the Rete nodes are deployed and initiated using Akka's *programmatic remote deployment* feature. For signalling the end of data processing, an asynchronous termination protocol was implemented.



3.6 Elaboration of the example

We use the *RouteSensor* query as our example. The query is shown as a graph pattern definition on Listing 3.1 and visualized on Figure 3.2. Queries like this are typical in MDE applications (such as well-formedness validation or complex model transformations).

3.6.1 Workflow

Following the workflow defined in subsection 3.5.3, we will cover the actual steps for deploying and operating a distributed pattern matcher for the *RouteSensor* query.

Analyze the metamodel and the query

The metamodel is shown on Figure 3.1. Using EMF-INCQUERY's tooling, the textual representation (`routeSensor.arch`, see Listing 3.1) is analyzed and parsed to an EMF model (Figure 3.7).

Build a Rete layout

Based on the query's EMF model, EMF-INCQUERY's tooling builds PSystem and creates a Rete layout, that guarantees the satisfaction of the constraints. The Rete layout is shown on Figure 3.5.

Allocate the Rete network in the cloud's nodes

In INCQUERY-D's current implementation, the Rete recipe's nodes are allocated manually on the cloud servers (called *Machines*). The allocation is currently defined in an architecture file (e.g. `routeSensor.arch`). The Rete nodes are associated with the machines with *infrastructure mapping* relationships.

INCQUERY-D's tooling currently provides an Eclipse-based tree editor to define machines and the infrastructure mapping edges (Figure 3.8).

The tooling is capable of visualizing the Rete network and its mapping to the machines (see Figure 3.9)

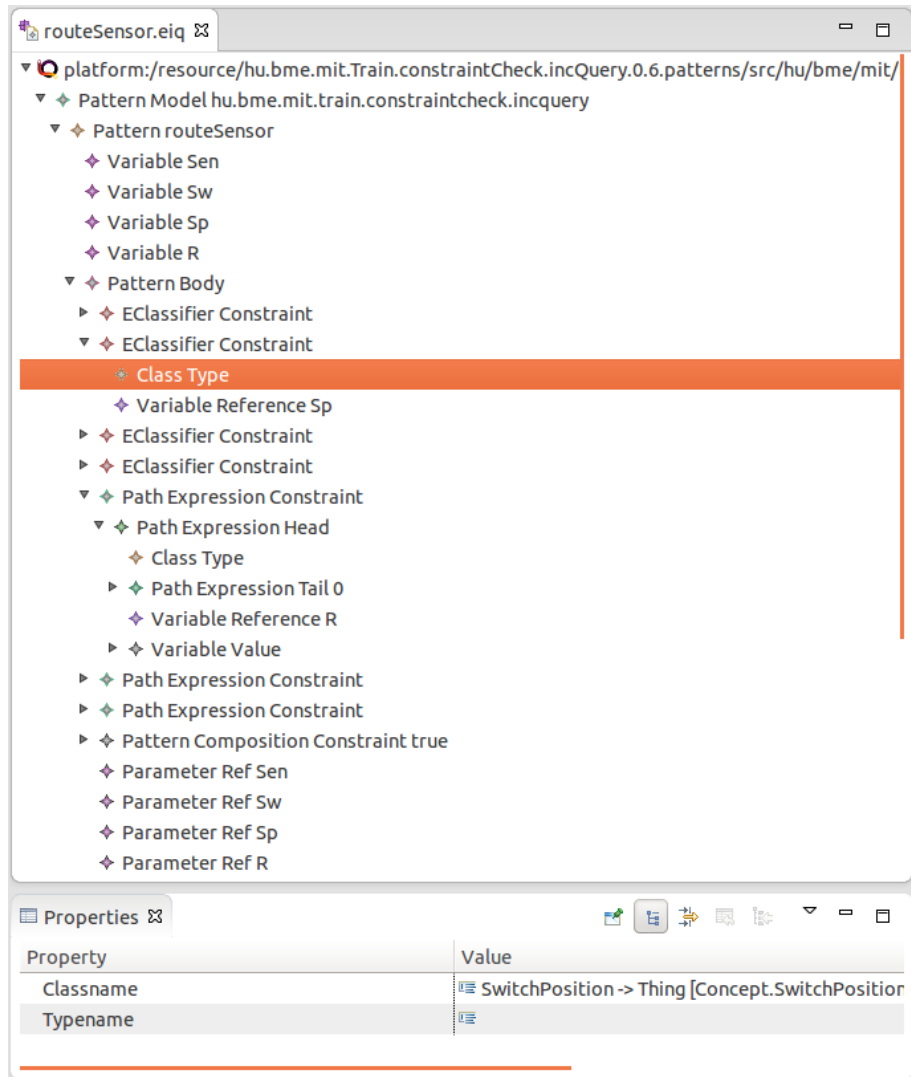


Figure 3.7: The EMF model generated from the pattern

Bootstrap the system

INCQUERY-D's current implementation, the distributed system is initiated with a Bash script which launches the Akka microkernel on the appropriate nodes. The Akka actors representing the Rete network's nodes are deployed automatically by the INCQUERY-D *Coordinator* node.

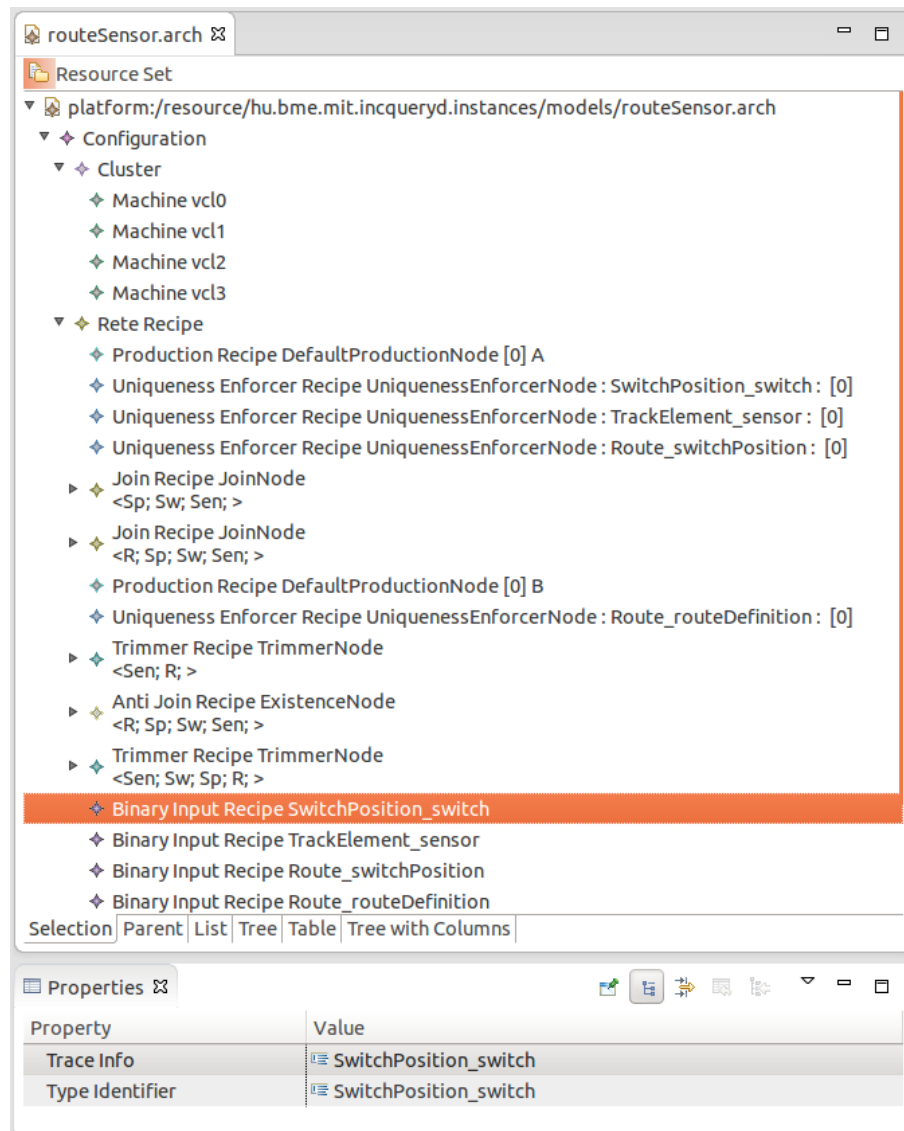


Figure 3.8: The tree editor in INCQUERY-D's tooling

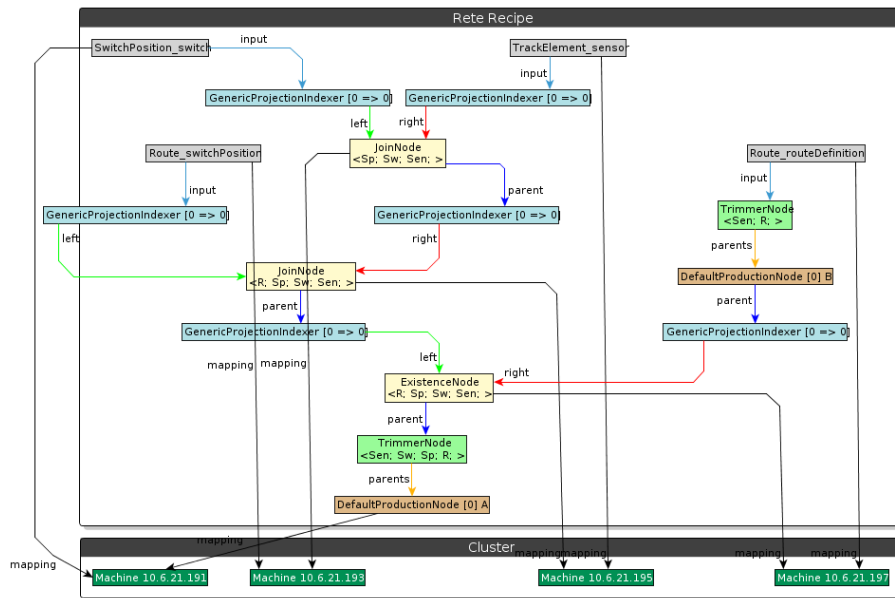


Figure 3.9: The yFiles viewer in INCQUERY-D's tooling

Chapter 4

Evaluation of scalability and performance

In this chapter, ...

4.1 Goals

4.1.1 Phases

The TrainBenchmark consists of the following phases:

1. *read*: loading the model,
2. *check₀*: running the queries,
3. *edit_i*: editing the model,
4. *check_i*: running the queries again.

In a "real-world" model editing sequence, the user typically edits the model in small steps (*edit_i* phases). The user's work is much more productive if she receives an instant feedback, hence we would like to run re-evaluate well-formedness queries quickly (preferably in sub-second time). This creates the need for an incremental pattern matcher tools.

4.1.2 Measure the response time and the scalability

4.1.3 Workload profile's difference from standard benchmarks

4.2 Benchmark scenario

In order to measure the efficiency of model queries and manipulation operations over the distributed architecture, we designed a benchmark to measure tool response times in a well-formedness validation use case. The benchmark transaction sequence consists of four phases: (i) during the *load* phase, the serialization of the model is loaded into the database; (ii) a test query (??) is executed (*check₀*); finally, in a cycle consisting of several repetitions, some elements are programmatically modified (*edit_i*) and the test query is re-evaluated (*check_i*). We ran the benchmark on pseudo-randomly generated instance models of growing size, each model containing about twice as many elements (vertices and edges) as the previous one and having a regular fan-out structure. As the current version of Neo4j does not have built-in support for graph sharding, the benchmark uses a manually sharded strategy where each shard contains a disjoint partition of the model.

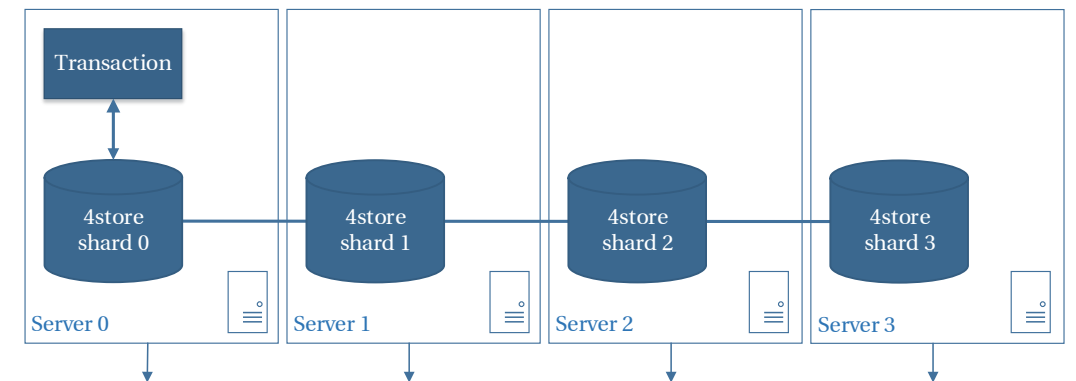


Figure 4.1: The non-incremental baseline benchmark's setup

4.3 Generation of models

We created a property graph generator project based on the previous TrainBenchmark generators. The generator creates a graph in a Neo4j database and uses the Blueprints library's `GraphMLWriter` and `GraphSONWriter` classes to serialize the graph to GraphML and Blueprints GraphSON formats. It is also capable of serializing the graph to Faunus GraphSON format.

4.4 Benchmark environment

describe the benchmark environment

Big MDE article [40]

4.4.1 Benchmark setup

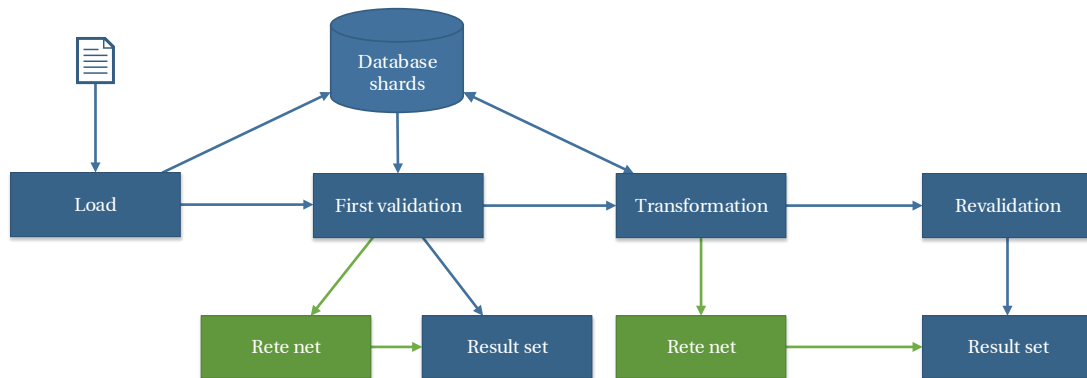


Figure 4.2: *The benchmark scenario*

distributed the Rete net as shown on Figure 4.3

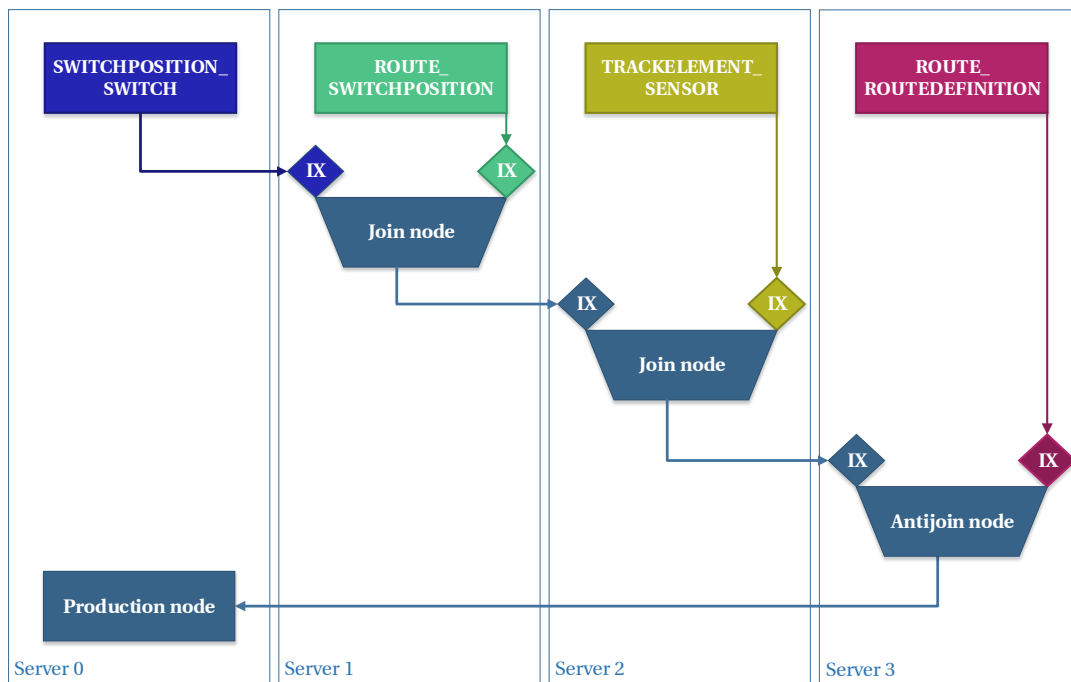


Figure 4.3: *The layout of the distributed Rete net*

4.4.2 Hardware, software ecosystem

As the testbed, we deployed our system to a private cloud consisting of 4 virtual machines on separate host computers.

Hardware

For Titan and 4store, each virtual machine used dual 2.5 GHz Intel Xeon L5420 CPUs with 8 GBs of RAM, running on Ubuntu 12.10 64-bit.

For Neo4j, each virtual machine had the same parameters but twice as much, 16 GBs of RAM.

Software

We used the following technologies:

- Ubuntu 12.10 64-bit
- Oracle Java 7 64-bit
- 4store 1.1.5
- Titan 0.3.2
- Faunus 0.3.2
- Hadoop 1.1.2
- Cassandra 1.2.2
- Neo4j 1.8
- Blueprints 2.3.0
- Akka 2.1.2
- Eclipse 4.3 (Kepler)

4.4.3 Benchmark methodology

five runs

some load on servers, so used the *minimum* of the results

4.4.4 Data collection

TrainBenchmark's utility classes

R script [22] developed by Benedek Izsó

4.4.5 Data processing tools

To compare the performance characteristics of INCQUERY-D to a traditional case, we defined two scenarios.

The *batch* scenario uses only 4store to manage models and evaluate the queries (depicted as ① in ??). This serves as a baseline for the *incremental* scenario, which uses INCQUERY-D (shown as ② in ??). For these initial experiments, the layout and allocation of the Rete network was determined manually.

4.5 Results

The measurement results of our experiments are shown in ?? (aggregated from several complete sets to filter transient effects). As expected, the *load* phase take about the same time for both scenarios, and INCQUERY-D is about half an order of magnitude slower when evaluating the query at first (*check₀* phase) due to the Rete construction overhead. However, INCQUERY-D is several orders of magnitude faster during the *edit_i – check_i* cycles, making on-the-fly query (re)evaluation feasible even for models larger than 50 million elements. Once initialized, INCQUERY-D scales linearly, since query response times for growing models can be kept low by adding additional computers for hosting Rete nodes.

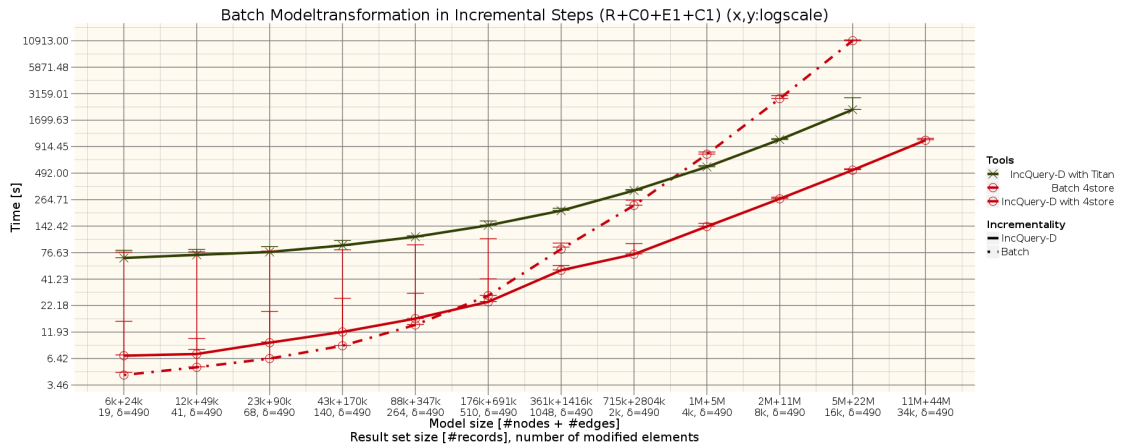


Figure 4.4: Batch transformation

Figure 4.4

Figure 4.5

Figure 4.6

Figure 4.7

Figure 4.8

Figure 4.9

Figure 4.10

4.5.1 Results using the Neo4j graph database

??

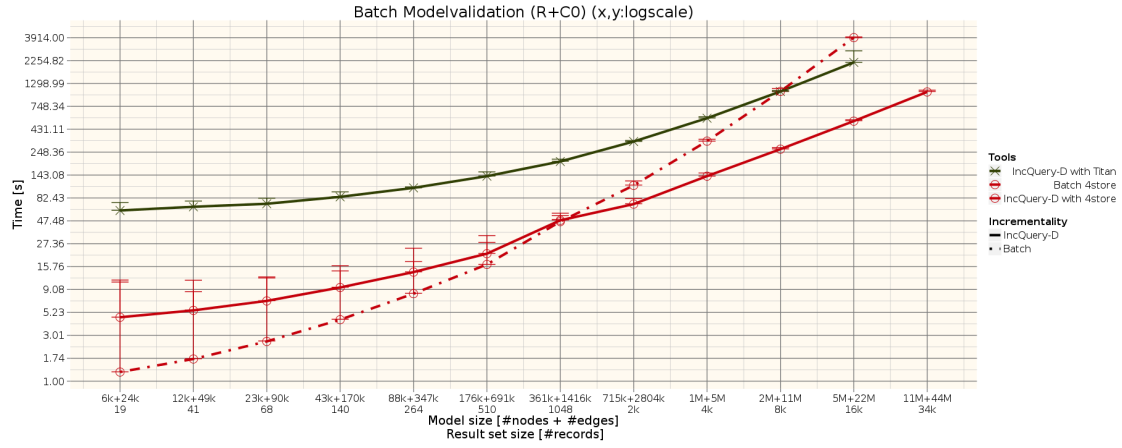


Figure 4.5: Batch validation

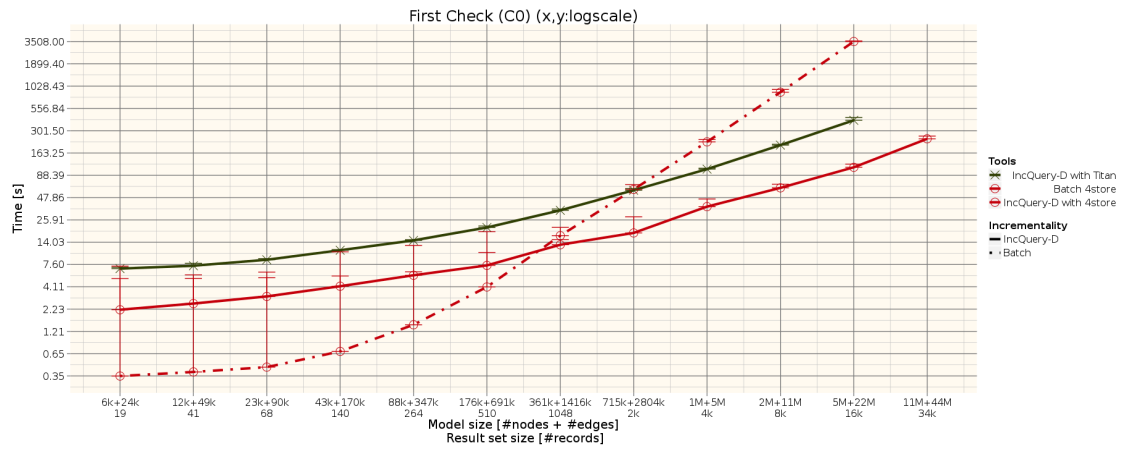


Figure 4.6: check₀ phase

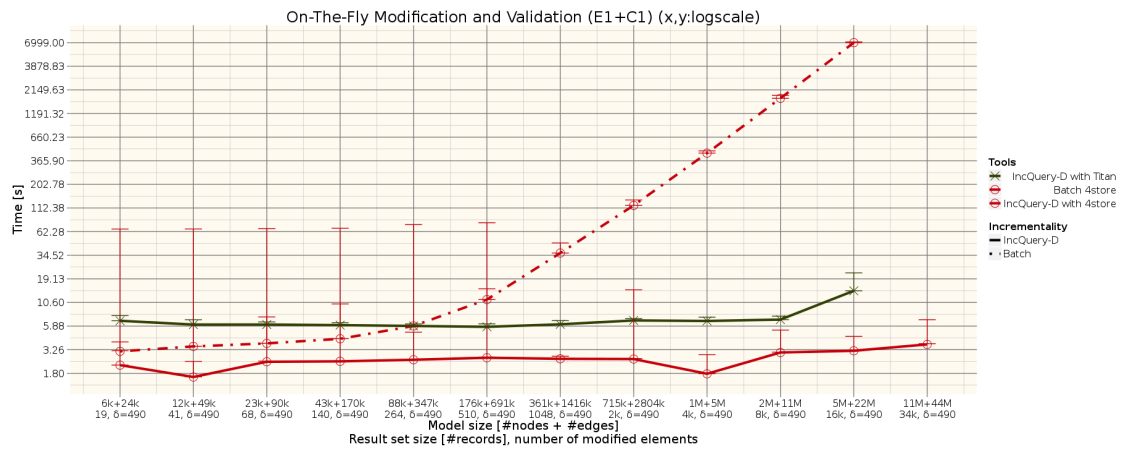


Figure 4.7: On-the-fly revalidation (edit and check₁ phase)

??

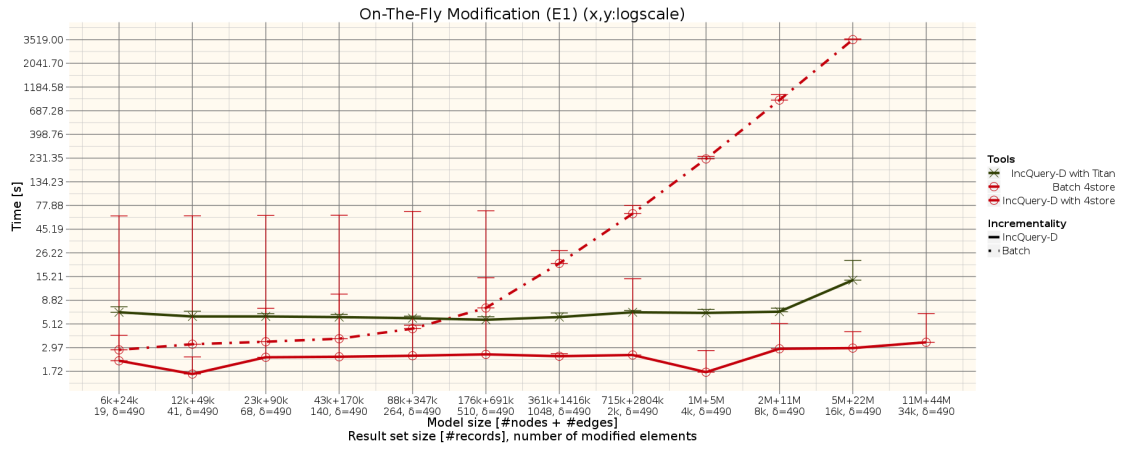


Figure 4.8: On-the-fly revalidation, edit phase

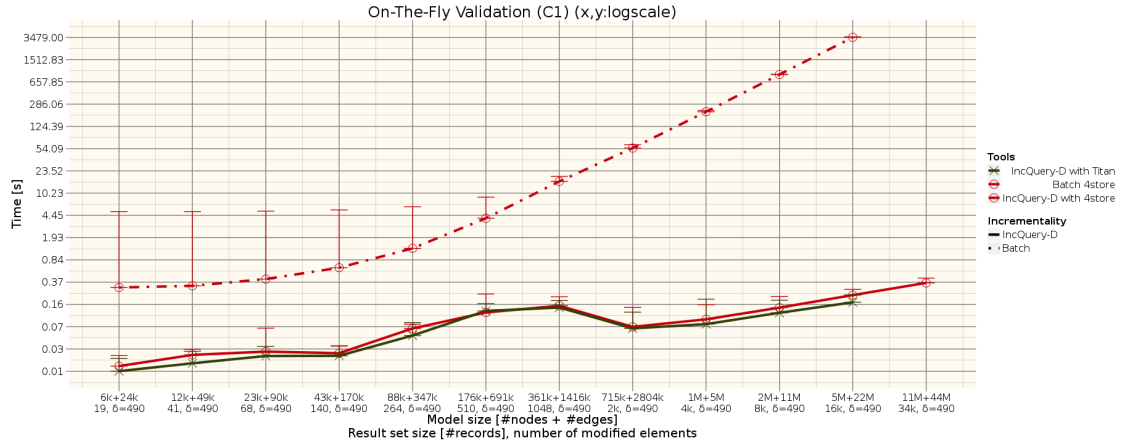


Figure 4.9: On-the-fly revalidation, check₁ phase

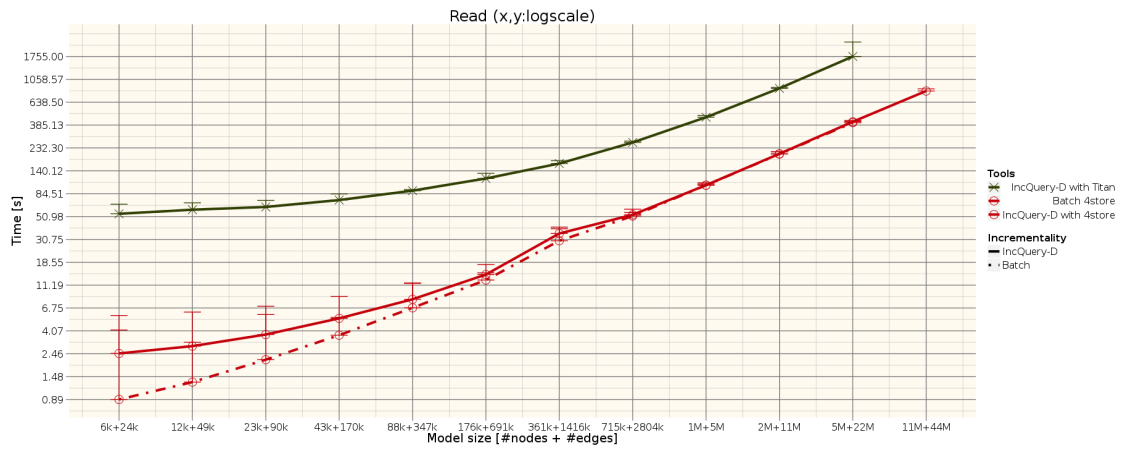


Figure 4.10: Read phase

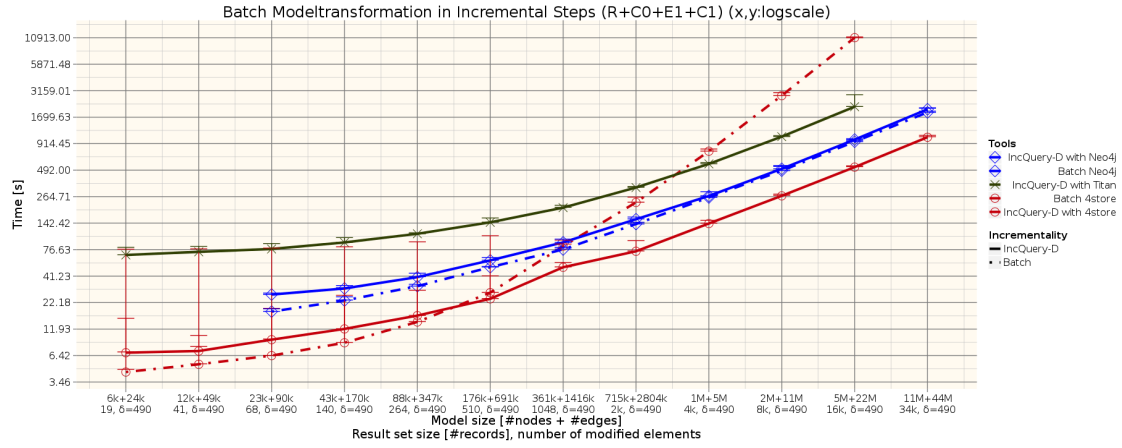


Figure 4.11: Batch transformation

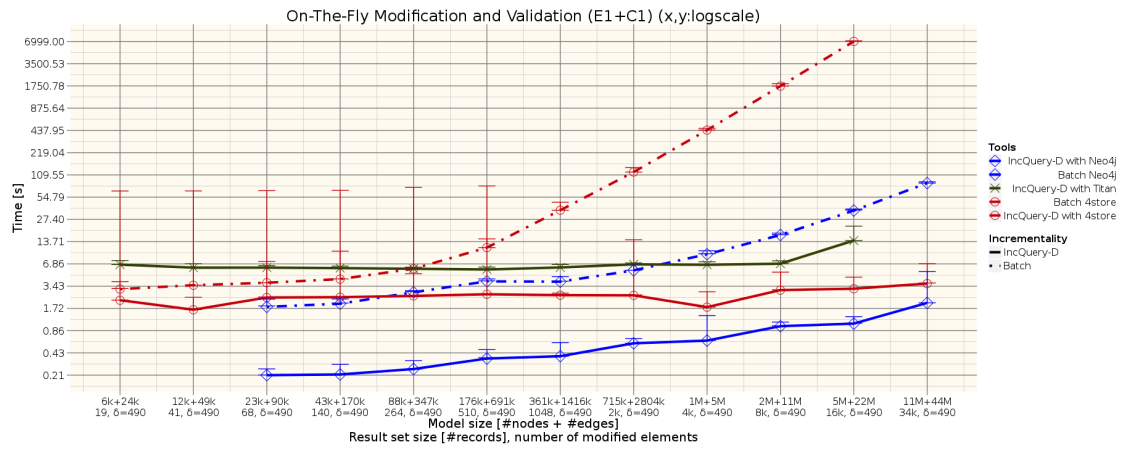


Figure 4.12: On-the-fly revalidation (edit and check₁ phase)

Chapter 5

Related work

This chapter collects the research and development works that are related to INCQUERY-D.

5.1 Eclipse-based tools

A wide range of special languages have been developed to support *graph based* representation and querying of computer data. A class-diagram like modeling language is Ecore of the Eclipse Modeling Framework (EMF [53]), 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 is a declarative constraint description and query language that can be evaluated with the local-search based Eclipse OCL [32] engine. To address scalability issues, impact analysis tools [35] have been developed as extensions or alternatives to Eclipse OCL.

5.2 Semantic web and NoSQL

Outside the Eclipse ecosystem, the Resource Description Framework (RDF [36]) is developed to support the description of instances of the semantic web, assuming sparse, ever-growing, incomplete data. Semantic models are built up from triple statements and they can be queried using the SPARQL [54] graph pattern language with tools like Sesame [2] or Virtuoso [1]. Property graphs [49] provide a more general way to describe graphs by annotating vertices and edges with key-value properties. Such data structures can be stored in graph databases like Neo4j [47] which provides the Cypher [52] query language. Even though big data storage (usually based on MapReduce) provides fast object persistence and retrieval, query engines realized directly on these data structures do not provide dedicated support for incremental query evaluation.

5.3 Rete implementations

In the context of event-based systems, distributed evaluation engines were proposed earlier [25]. However they scaled up in the number of rules [29] rather than in the number of data elements. 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 [48] uses this algorithm to perform complex event processing (formulated in SPARQL) on RDF data, gathered from distributed sensors. Diamond [46] evaluates 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 as based on EMF-INCQUERY [28], 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 MDE context. As the architecture of INCQUERY-D separates the data store from the query engine, we believe that the scalable processing of 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 [25]. The medium-grained computer connected processors in a crossbar architecture, while our approach use computers connected by gigabit Ethernet. The paper published benchmark results of the medium-grained solution, but these are based only on simulations.

Chapter 6

Conclusions

6.1 Summary of contributions

In this chapter, ...

6.1.1 Own work

During the research and development of INCQUERY-D's prototype, I achieved the following results:

- I expanded the TrainBenchmark with a *new instance model generator*, which can produce property graphs and serialize them in various formats: GraphML, Blueprints GraphSON and Faunus GraphSON.
- I implemented a *distributed, asynchronous version of the Rete algorithm* with a termination protocol. 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.
- Based on EMF-INCQUERY and István Ráth's work, I created an *Eclipse-based environment* to provide automated deployment of the Rete network. This allows the user to define complex queries in IQPL, a high-level pattern language.
- I experimented with contemporary non-relational database management systems with a focus on NoSQL graph databases and triple stores.
 - 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 with the Gremlin query language and the Faunus graph analytics engine.
 - I implemented scripts to install the *4store triplestore* on a cluster. I created

the connector class in INCQUERY-D's middleware and formulated the necessary the SPARQL queries.

- I deployed a manually sharded *Neo4j cluster*. I created the connector class in INCQUERY-D's middleware to use Neo4j's REST interface and formulated the appropriate Cypher queries.
- I implemented scripts for *automating the benchmark and operating a cluster of Akka microkernels*.
- 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*, with Neo4j and 4store.
- For INCQUERY-D's implementation, I wrote more than 3000 lines of Java code and approximately 500 lines of configuration and deployment scripts.

6.2 Limitations

- manual allocation of Rete nodes
- only a subset of Rete nodes implemented
- lack of complete tooling

6.3 Future work

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 sharded graph databases by a middleware layer, and its feasibility 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 10^7 elements.

For future work, we plan on providing more sophisticated automation for sharded Ecore models, and further exploring advanced optimization challenges such as dynamic reconfiguration and fault tolerance. We also plan experiment with programming languages that are better suited to asynchronous algorithms (e.g. Erlang and Scala) and try different database systems (e.g. MongoDB) as our storage layer.

List of Figures

2.1	Different graph data models (based on [49])	9
2.2	The TinkerPop software stack [10]	11
2.3	Cassandra's ring for data partitioning [5]	13
2.4	Hadoop's architecture [44]	14
2.5	HDFS' architecture	14
2.6	TrainBenchmark subgraph visualized in Neoclipse	15
2.7	Titan graph vertex stored in Cassandra as a row	16
2.8	Using Titan with Cassandra in remote server mode	17
2.9	4store's distributed architecture [37]	18
2.10	Deploying a remote actor in Akka	19
2.11	The metamodel of Ecore	20
2.12	EMF-INCQUERY's architecture	21
3.1	The EMF metamodel of the TrainBenchmark	23
3.2	Graphical representation of the <i>RouteSensor</i> query's pattern. The dashed red arrow defines a negative condition.	24
3.3	INCQUERY-D's architecture on a four-node cluster	26
3.4	A modification on a TrainBenchmark instance model	29
3.5	The Rete net and the partial matches stored in its nodes	30
3.6	The workflow of EMF-INCQUERY (blue) and INCQUERY-D (green)	32
3.7	The EMF model generated from the pattern	34
3.8	The tree editor in INCQUERY-D's tooling	35
3.9	The yFiles viewer in INCQUERY-D's tooling	36
4.1	The non-incremental baseline benchmark's setup	38
4.2	The benchmark scenario	39
4.3	The layout of the distributed Rete net	39
4.4	Batch transformation	41
4.5	Batch validation	42
4.6	<i>check</i> ₀ phase	42
4.7	On-the-fly revalidation (<i>edit</i> and <i>check</i> ₁ phase)	42

4.8	On-the-fly revalidation, <i>edit</i> phase	43
4.9	On-the-fly revalidation, <i>check</i> ₁ phase	43
4.10	<i>Read</i> phase	43
4.11	Batch transformation	44
4.12	On-the-fly revalidation (<i>edit</i> and <i>check</i> ₁ phase)	44
A.1	An example graph based on the TrainBenchmark’s metamodel	55

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] About Data Partitioning in Cassandra. http://www.datastax.com/docs/1.1/cluster_architecture/partitioning, October 2013.
- [6] Akka. <http://akka.io/>, May 2013.
- [7] Apache Cassandra. <http://cassandra.apache.org/>, May 2013.
- [8] Apache Hadoop. <http://hadoop.apache.org/>, May 2013.
- [9] Apache Thrift. <http://thrift.apache.org/>, October 2013.
- [10] Blueprints. <http://blueprints.tinkerpop.com/>, May 2013.
- [11] GraphSON Format. <https://github.com/thinkaurelius/faunus/wiki/GraphSON-Format>, October 2013.
- [12] GraphSON Reader and Writer Library. <https://github.com/tinkerpop/blueprints/wiki/GraphSON-Reader-and-Writer-Library>, October 2013.
- [13] Lustre. <http://lustre.org/>, October 2013.
- [14] Neoclipse. <https://github.com/neo4j/neoclipse>, May 2013.
- [15] NoSQL Databases. <http://nosql-database.org/>, May 2013.
- [16] Objectivity – InfiniteGraph. <http://www.objectivity.com/infinitegraph>, October 2013.

- [17] OrientDB Graph-Document NoSQL DBMS. <http://www.orientdb.org/>, October 2013.
- [18] Planet Cassandra – Companies. <http://planetcassandra.org/Company/ViewCompany>, October 2013.
- [19] Protocol Buffers – Google’s data interchange format. <https://code.google.com/p/protobuf/>, October 2013.
- [20] Sparsity-technologies: DEX high-performance graph database. <http://www.sparsity-technologies.com/dex>, October 2013.
- [21] The GraphML File Format. <http://graphml.graphdrawing.org/>, October 2013.
- [22] The R Project for Statistical Computing. <http://www.r-project.org/>, October 2013.
- [23] TinkerPop. <http://www.tinkerpop.com/>, May 2013.
- [24] Xtext – Language Development Made Easy! <http://www.eclipse.org/Xtext/>, October 2013.
- [25] Acharya, A. et al. Implementation of production systems on message-passing computers. *IEEE Trans. Parallel Distr. Syst.*, 3(4):477–487, July 1992.
- [26] Don Batory. The LEAPS Algorithm. Technical report, Austin, TX, USA, 1994.
- [27] 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.
- [28] Bergmann, Gábor et al. Incremental evaluation of model queries over EMF models. In *MODELS*, volume 6394 of *LNCS*. Springer, 2010.
- [29] Bin Cao, Jianwei Yin, Qi Zhang, and Yanming Ye. A MapReduce-Based Architecture for Rule Matching in Production System. In *CLOUDCOM*. IEEE, 2010.
- [30] E. F. Codd. A relational model of data for large shared data banks. *Commun. ACM*, 13(6):377–387, June 1970.
- [31] Jeffrey Dean and Sanjay Ghemawat. MapReduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, January 2008.
- [32] Eclipse MDT Project. Eclipse OCL, 2011. <http://eclipse.org/modeling/mdt/?project=ocl>.

- [33] Charles Forgy. Rete: A fast algorithm for the many pattern/many object pattern match problem. *Artificial Intelligences*, 19(1):17–37, 1982.
- [34] 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.
- [35] Thomas Goldschmidt and Axel Uhl. Efficient OCL impact analysis, 2011.
- [36] RDF Core Working Group. Resource Description Framework (RDF). <http://www.w3.org/RDF/>, 2004.
- [37] S. Harris, N. Lamb, and N. Shadbolt. 4store: The Design and Implementation of a Clustered RDF Store. In *5th International Workshop on Scalable Semantic Web Knowledge Base Systems (SSWS2009)*, 2009.
- [38] 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%.
- [39] 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.
- [40] 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.
- [41] 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.
- [42] Adam Jacobs. The pathologies of big data. *Commun. ACM*, 52(8):36–44, August 2009.
- [43] Avinash Lakshman and Prashant Malik. Cassandra: a decentralized structured storage system. *SIGOPS Oper. Syst. Rev.*, 44(2):35–40, April 2010.

- [44] 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.
- [45] 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.
- [46] Miranker, Daniel P et al. Diamond: A SPARQL query engine, for linked data based on the Rete match. *AIMWD*, 2012.
- [47] Neo Technology. Neo4j. <http://neo4j.org/>, 2013.
- [48] Mikko Rinne. SPARQL update for complex event processing. In *ISWC'12*, volume 7650 of *LNCS*. 2012.
- [49] Marko A. Rodriguez and Peter Neubauer. Constructions from dots and lines. *CoRR*, abs/1006.2361, 2010.
- [50] Markus Scheidgen. How Big are Models – An Estimation. Technical report, Department of Computer Science, Humboldt Universität zu Berlin, 2012.
- [51] Sherif Sakr. Supply cloud-level data scalability with NoSQL databases. <http://www.ibm.com/developerworks/cloud/library/cl-nosqldatabase/index.html>, March 2013.
- [52] Andrés Taylor and Alistair Jones. Cypher Query Language. <http://www.slideshare.net/graphdevroom/cypher-query-language>, 2012.
- [53] The Eclipse Project. Eclipse Modeling Framework. <http://www.eclipse.org/emf>, October 2012.
- [54] W3C. SPARQL Query Language for RDF. <http://www.w3.org/TR/rdf-sparql-query/>.

Appendix A

Graph formats

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

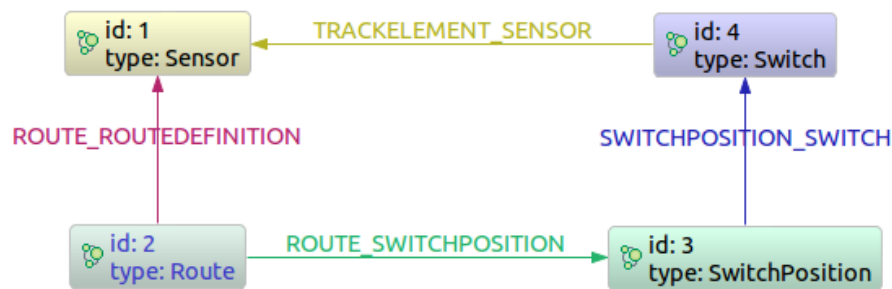


Figure A.1: An example graph based on the TrainBenchmark's metamodel

A.1 Property graph formats

A.1.1 GraphML

```
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 TrainBenchmark's metamodel stored in GraphML format

A.1.2 Blueprints GraphSON

```

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 TrainBenchmark's metamodel stored in Blueprints GraphSON format

A.1.3 Faunus GraphSON

```

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 TrainBenchmark's metamodel stored in Faunus GraphSON format

A.2 Semantic graph formats

A.2.1 RDF

```

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"/>

```

```

23 </rdf:Description>
24
25 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#1">
26   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#Sensor"/>
27 </rdf:Description>
28
29 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#2">
30   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#Route"/>
31 </rdf:Description>
32
33 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#3">
34   <rdf:type rdf:resource="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#Switch"/>
35 </rdf:Description>
36
37 <rdf:Description rdf:about="http://www.semanticweb.org/ontologies/2011/1/
    TrainRequirementOntology.owl#4">
38   <rdf:type rdf: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 TrainBenchmark's metamodel stored in RDF format.