# INSTITUT FÜR INFORMATIK

der Ludwig-Maximilians-Universität München

# Parallel Datalog on Pregel

Bachelor Thesis

Matthias Maiterth

Advisor:                Harald Zauner
Supervising professor:  Prof. Dr. François Bry
Filing date:            October 22, 2012

# Erklärung

Hiermit versichere ich, dass ich die vorliegende Bachelorarbeit selbstständig verfasst habe und keine anderen als die angegebenen Hilfsmittel verwendet habe.

München, den 22. Oktober 2012                    Matthias Maiterth

# Abstract

With the rise of Web 2.0, the efficient processing and analysis of Big Data (e.g. graphs of social networks) has gained strong interest, and still is an active research issue.

For such analysis tasks, Datalog with its concise syntax and clean semantics is a suitable language to express facts and rules concerning Big Data.

Google's MapReduce is the de facto standard for processing Big Data, but for performing recursive computations (as is the case when evaluating Datalog programs), Pregel is better suited.

Pregel was developed for graph processing at Google, and provides a data parallel, vertex-centric approach for computations on a compute cluster.

The focus of this thesis lies on elaborating and implementing a data parallel algorithm for the semi-naive evaluation of Datalog programs using Pregel.

# Acknowledgements

First and foremost, I would like to thank Harald Zauner for offering me a highly interesting topic for this thesis. His support throughout the whole work has been very important and encouraging.

I would like to thank Prof. Dr. François Bry for his helpful advises and insights. The appointments with him throughout the work were an important guideline and built a common thread throughout my work.

I would also like to express my thanks to the other members of the Teaching and Research Unit Programming and Modelling Languages at the Department for Informatics at the LMU.

My greatest thanks to my parents, my friends and the people I love for their support.

# Contents

# Chapter 1

# Motivation: Evaluation of Huge Datalog Programs

With ever growing amounts of data, e.g. graphs of social networks or the web itself, scalable processing of this data has become of utmost importance. Distributed systems have proven capable of processing data of such size.

Datalog [1] provides a concise syntax and clean semantics to define relations in networks of arbitrary size and also allows the definition of logical rules to infer relations that are not yet explicitly represented in the raw data. The versatility and expressive power of Datalog, while at the same time staying simple and understandable, has allowed Datalog not only to stay a hot topic in academia, but also to gain a foothold in industry. Logicblox, for example, uses Datalog to provide cloud based solutions for Big Data enterprise applications in the domains of decision automation, analytics, and planning [2]. SAP is another notable company that uses Datalog to discover, represent, and reconstruct business networks from raw data (mined from a network) [3]. In contrast to primarily interactive systems like the ones mentioned above, we focus on computing the minimal model (or the deductive closure) of a Datalog program, which can be used for further processing later on.

Recent work [4] on the semi-naive evaluation of Datalog based on MapReduce [5] has shown that a query cannot be evaluated by a single MapReduce job or a small number of MapReduce jobs combined. In general, evaluation of a single query requires a large number of MapReduce jobs combined in sequence. This chaining of large numbers of MapReduce jobs results in an overhead that would only be negligible when chaining just a few MapReduce jobs. Consider two MapReduce jobs $A$ and $B$ chained in sequence on the same compute cluster. When $A$ is finished, it writes its output data to a distributed file system. For

the cluster to execute $B$ on $A$'s output, it first has to distribute the map- and reduce-functions of job $B$ and then read and partition $B$'s input from the file system. This interaction with the distributed file system, communication of the user program, and partitioning of data would not pose a problem for a single MapReduce job. But it becomes expensive, when it has to be done repeatedly for a large number of MapReduce jobs.

The overhead of the repeated tasks of interacting with the distributed file system, partitioning the input data, and communicating the user program does not only apply when performing a semi-naive evaluation of Datalog, but more generally to all kinds of complex computations wherever the dependencies within the input data are more complex than what can be expressed by a single instance of the user-defined functions map and reduce.

Since no communication amongst the mappers or reducers is supported by MapReduce, its dataflow model is in general ill-suited for problems with complex dependencies in their input data.

Coming back to evaluating Datalog with MapReduce, especially Datalog programs with recursive rules would in general require a considerably long sequence of MapReduce jobs until no more new knowledge can be derived. Even specialized MapReduce implementations like HaLoop [6], that reduce part of the overhead caused by chaining MapReduce jobs together, cannot eliminate the overhead caused by the mismatch between MapReduce's programming model and input data with complex dependencies (as is the case with most Datalog programs).

In this work, we elaborate an algorithm for evaluating Datalog on a distributed system using Pregel [7], a framework for large-scale graph processing developed at Google. The objective of the algorithm proposed is an efficient computation of the facts implied by a Datalog program relying on data parallelism using a dedicated compute cluster.

In contrast to the dataflow model of MapReduce, Pregel is closely oriented on the bulk-synchronous parallel (BSP) model [8]. A Pregel computation consists of several distributed processes that process data in parallel according to a user program. These processes communicate with each other via messages, and each of them is responsible for a specific partition of the input graph. The expressiveness and flexibility of Pregel's computational model allows us to implement a distributed, semi-naive evaluation algorithm for Datalog as a single Pregel job.

For the above mentioned reasons, the Pregel-based approach does not introduce the overhead, which cannot be avoided by any MapReduce-

based approach.

In Chapter 2, we give a short introduction to Datalog. In Chapter 3, a detailed introduction to Pregel is given including the Apache implementation: Giraph [9]. The concept and implementation of the Datalog evaluation using Pregel is presented in Chapter 4. In Chapter 5, we compare our implementation with DEDALUS [10], another distributed implementation of Datalog. In Chapter 6, we conclude with an outlook on further developments and possible applications.

An experimental evaluation of the algorithm on large Datalog programs utilizing a compute cluster has not been conducted, as it was not part of this thesis. This is to be done in future work.

# Chapter 2

# Datalog

In Section 2.1, we give a short introduction to Datalog. In Section 2.2, the semi-naive evaluation is presented. This evaluation algorithm is adapted for distributed evaluation in Chapter 4.

## 2.1 Introduction to Datalog

Datalog is a language to represent knowledge in a fragment of first-order logic. To be able to combine Datalog and the distributed graph processing framework Pregel, we first introduce Datalog itself as it is described in [1, pp. 96–128] and [11, pp. 11–12]. To understand the structure of Datalog programs, in the following, we examine the program of Figure 2.1.

We begin by introducing atoms: An *atom* consists of a predicate symbol and a list of arguments. Each argument of an atom is either a variable or a constant; function symbols of arities greater than zero are not allowed.

A Datalog program is a finite set of rules. The left part of a rule, before the ":-", is a single atom called the *head* of the rule. The right part, called the *body* of the rule, can consist of an arbitrary number ($\geq 0$) of atoms (separated by &), also referred to as *subgoals*. These atoms can also be simple arithmetic comparisons formed with predi-

```
1   sibling(X,Y):- parent(X,Z) & parent(Y,Z) & X≠Y.
2   parent('Isabella','Ella').
3   parent('Ella','Ben').
4   parent('Daniel','Ben').
```

Figure 2.1: Example Datalog program for computing a `sibling` relation, given `parent` EDB facts.

cates like $=$, $\neq$, and $\geq$, which are referred to as *built-in* predicates. For a rule to be *safe*, each of its variables has to be limited. The limited variables of a rule are defined as follows:

A variable appearing in a rule $r$ is *limited* iff (short for: if and only if) (i) it appears as argument of an ordinary (i.e. not built-in) subgoal of $r$, or (ii) it is equated to a constant (in a subgoal of $r$), or (iii) it is equated to a limited variable of $r$ (in a subgoal of $r$) [1, p.105]. In the following, only safe Datalog programs (i.e. Datalog programs consisting of safe rules only) are considered.

By *facts*, we denote rules with empty body (written as only the head of the rule, leaving out the ":-", see lines 2-4 of Figure 2.1).

A relation can either be an extensional database (EDB) relation, or an intensional database (IDB) relation. An EDB relation is given by a set of facts with the same predicate symbol (which can be seen as the name of the relation), while IDB relations are defined via rules.

Predicate symbols are used for labeling relations. The attributes that name columns in relational algebra are missing, though. The identification of attributes is possible by following the order of the attributes.

For the relation `parent` an EDB relation could look like lines 2-4 of Figure 2.1. This relation can be interpreted as follows: " 'Isabella' has parent 'Ella', 'Ella' has parent 'Ben', and 'Daniel' has parent 'Ben' ".

We use the proof-theoretic interpretation of rules and disallow negation, as in pure Datalog. This means that we derive facts from rules given by a Datalog program or a database, according to [1, p. 97].

If we want to infer facts using a rule, we can use all facts derived from EDB and IDB relations. When replacing the variables with constants and each subgoal of the body is satisfied, the head is true. In other words, a new fact is derived. Let us evaluate rule 1 of the Datalog program in Figure 2.1:

```
sibling(X,Y):- parent(X,Z) & parent(Y,Z) & X≠Y.
```

We can insert the facts from above and infer the following:

```
sibling('Ella','Daniel'):- parent('Ella','Ben') &
        parent('Daniel','Ben') & 'Ella'≠'Daniel'
```

The two predicates for `parent` and the built-in predicate are satisfied. This makes the body true and thus the head is true accordingly.

## 2.2 Semi-naive Evaluation

In this section, we introduce the semi-naive evaluation as defined in Figure 2.2. To begin with, we define the notion of a substitution:

```
1   //Input: F (set of Datalog EDB facts),
2   //       R (set of Datalog rules with non-empty body)
3   Δ_old := F
4   while Δ_old ≠ ∅
5       Δ_new := EVAL-INCR(R, F, Δ_old)
6       F := F ∪ Δ_new
7       Δ_old := Δ_new
8   output F
```

Figure 2.2: Semi-naive evaluation of a Datalog program

A *substitution* $\sigma$ is a function, written in postfix notation, that maps terms to terms and is

- homomorphous, i.e. $f(t_1, \ldots, t_n)\sigma = f(t_1\sigma, \ldots, t_n\sigma)$ for compound terms and $c\sigma = c$ for constants
- identical almost everywhere, i.e.,
  $\{x \mid x \text{ is a variable and } x\sigma \neq x\}$ is finite.

[11, Def. 48]

A substitution $\sigma$ can thus be represented as a finite set $\{x_1 \mapsto t_1, \ldots, x_n \mapsto t_n\}$ where $\{x_1, \ldots, x_n\}$ is the finite set of variables on which $\sigma$ is not identical and $x_i\sigma = t_i$ for $1 \leq i \leq n$.

For the semi-naive evaluation, we define $\texttt{EVAL-INCR}(R, F, \Delta_{old})$, where $R$ is a set of Datalog rules (each with non-empty body), and $F$ and $\Delta_{old}$ are sets of facts, with $\Delta_{old} \subset F$.

$$
\begin{aligned}
\texttt{EVAL-INCR}&(R, F, \Delta_{old}) := \\
&\{h\sigma \mid \exists(h \text{ :- } b_1 \ \& \ \ldots \ \& \ b_n) \in R( \\
&\quad \exists\sigma, \text{ where } \sigma \text{ is a substitution}, (\exists k(1 \leq k \leq n \ \wedge \\
&\quad b_k\sigma \in \Delta_{old} \wedge \forall l(1 \leq l \leq n \wedge l \neq k \Rightarrow b_l\sigma \in F))))\} \\
&\setminus F
\end{aligned}
$$

$\texttt{EVAL-INCR}$ generates a new set of facts, obtained from the heads of rules $r$ of the form '$h$ :- $b_1$ & $\ldots$ & $b_n$' from $R$, where for each rule $r$ all subgoals are satisfied the following way: At least one subgoal is satisfied using a fact from $\Delta_{old}$ and the remaining subgoals are satisfied using facts from $F$. Those facts, that already appear in $F$ are removed from this new set, and the remaining facts are returned as the result of $\texttt{EVAL-INCR}$.

$\texttt{EVAL-INCR}$ can be used in the semi-naive evaluation of Datalog, as shown in Figure 2.2.

7

The idea behind the semi-naive algorithm is evaluating all rules of a Datalog program, using only the most recent addition of facts $\Delta_{old}$ for one subgoal and all facts $F$ for the remaining subgoals to generate a new set of facts $\Delta_{new}$. These new facts are again used for the next evaluation round, until now more new facts can be inferred.

The input for the semi-naive algorithm is a set of Datalog rules $R$ with non-empty body and a set of EDB facts $F$.

- Initially, we initialize $\Delta_{old}$ with the set of facts $F$.

- After that, the tasks from line 5-7 are repeated until $\Delta_{old}$ is empty:

  - $\Delta_{old}$ represents all facts that have been derived in the previous round. (In the first round, these are the input facts $F$.)
  - `EVAL-INCR` is evaluated using $R$, $F$, and $\Delta_{old}$ as arguments, and saving the result in $\Delta_{new}$.
  - $F$ is assigned the union $F \cup \Delta_{new}$. Hence, $F$ contains all facts that have been derived so far.
  - $\Delta_{new}$ is saved to $\Delta_{old}$ since it represents the most recent addition of facts, which is used for evaluating the next round.

- Finally, the set of facts $F$ is returned.

In Chapter 4, we map this evaluation algorithm to Pregel.

# Chapter 3

# Pregel

This chapter presents Pregel [7], a framework for the efficient processing of large (directed) graphs on a compute cluster.

A Pregel computation is a sequence of supersteps, in each of which a vertex can process messages from the previous superstep, send messages to other vertices, and modify its own state, thereby changing the state of the graph as a whole.

Pregel has been developed at Google, and is inspired by the bulk-synchronous parallel (BSP) model [8]. According to [8], the BSP model is defined as the combination of the following three attributes:

- A number of *components* for performing computations and memory functions.
- A *router* for delivering messages between components.
- *Supersteps* for synchronizing the computations performed by components. During a superstep, each component performs a task consisting of local computation and message arrivals and transmissions. Once the superstep has been completed by all components, the components proceed to the next superstep (barrier synchronization [8]).

A BSP computation consists of a sequence of such supersteps.

Similar to MapReduce [5], Google's implementation of Pregel is not publicly available, but there exist several open source implementations, e.g. GoldenOrb [12] and Apache Giraph [9].

This chapter is organized as follows: Section 3.1 gives a conceptual overview of Pregel, while Section 3.2 presents more details of the framework. Section 3.3 presents the open source implementation Apache Gi-

raph. In Section 3.4, a general example application is given to visualize Giraph's method of operation.

## 3.1 General Concept of Pregel

As this section presents a conceptual overview of Pregel, it does not address how Pregel computations are actually executed on a distributed system. Hence, the terms "graph" and "vertex" used below are not to be confused with compute nodes in a compute cluster.

The input to a Pregel computation consists of a directed graph and a user program (see Section 3.2) containing a user-defined function called `compute()`.

A Pregel computation consists of a sequence of so-called supersteps. In each superstep, the framework calls the same (user-provided) `compute()`-function in parallel on each (active) vertex of the directed graph. Note that, a Pregel computation has only one `compute()`-function. `compute()` can essentially be used to conduct local computations which generally result in a change of the local state of the vertex or communication with other vertices by sending messages. Once every vertex is finished with the `compute()`-function, the next superstep starts. All messages are synchronized by these supersteps, i.e. a message issued at superstep $S$ will be available at its destination vertex at superstep $S + 1$.

A Pregel computation terminates once all vertices of the directed graph are inactive and no further messages are to be passed. Initially, each vertex is active. If a vertex $V$ calls `voteToHalt()`, $V$ will be inactive in the next superstep (i.e. the framework won't call its `compute()`-function) unless some vertex issued a message to $V$ in the current superstep, since messages always activate the receiving vertex in the next superstep.

The output of a Pregel computation is the resulting directed graph after the computation has terminated.

## 3.2 Presentation of Pregel

This section gives a more detailed presentation of Pregel than was given in the conceptual overview of the last section.

The input to a Pregel computation consists of a directed graph $G$ (consisting of vertices and directed edges) and a user program. The user program contains a vertex class, which includes the `compute()`-

function, and might contain additional classes, e.g. for realizing aggregators and combiners (see details below).

Each vertex of $G$ consists of

- a unique vertex ID,
- a flag indicating whether the vertex is active or not,
- a vertex value,
- a set of incoming messages, and
- a set of outgoing edges.

Each directed edge of $G$ is associated with its source vertex and consists of an edge value and a target vertex ID.

A Pregel computation consists of a sequence of supersteps.

In superstep 0, after the input graph has been read, each vertex is active (i.e. its active flag is set to true) and its set of incoming messages is empty.

In each superstep $S$, the Pregel framework invokes the user-provided `compute()`-function on each active vertex $V$, conceptually in parallel. From within `compute()`, it is possible

- to read each component the vertex $V$ consists of,
- to write $V$'s vertex value,
- to read the number of the current superstep,
- to call `voteToHalt()`, which will set $V$ to inactive in the next superstep,
- to issue messages to vertices that are addressed by their vertex ID,
- to modify the set of outgoing edges and their edge values, and
- to pass values to aggregators (see details below) and to read the aggregated values.

Each superstep in Pregel is concluded by barrier synchronization, which affects both computation and communication as follows:

- Concerning computation, barrier synchronization means that each vertex waits until all other vertices have finished the `compute()`-function.

- Concerning communication, barrier synchronization affects messages, the `voteToHalt()`-function, and aggregators:

  - A message issued at superstep $S$ will not be available at its destination vertex until the next superstep $S + 1$.
  - Calling `voteToHalt()` will not render a vertex $V$ inactive in the current, but in the next superstep (unless a message is available for $V$ in the next superstep).
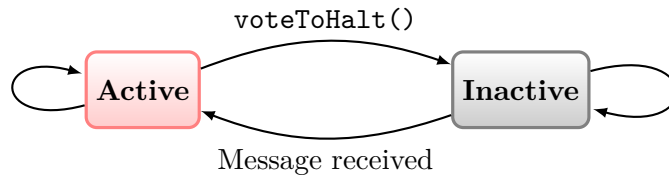
Figure 3.1: Vertex State Machine [7]

> – When using aggregators, the result of aggregating values, which have been passed to an aggregator in the current superstep, is not visible until the next superstep.

A Pregel computation terminates once all vertices of the directed graph are inactive and no further messages are to be passed.

Figure 3.1 illustrates how the state of a vertex can change from active to inactive, and vice versa: A change from active to inactive only occurs, when the vertex itself calls `voteToHalt()`. However, if a vertex receives a message it is guaranteed that the vertex is active in the next superstep. Thus, only when no more messages are to be passed, the existence of no further active vertices can be guaranteed and the computation will terminate. Otherwise, a new superstep follows.

The output of a Pregel computation is the resulting directed graph after the computation has terminated. In order to generate this output, each vertex contributes its vertex ID, its vertex value, and its set of outgoing edges, as specified in the user program.

### Aggregators

Aggregators have already been mentioned above and can be used to introduce a global value that each vertex can read and modify by passing values to the aggregator. For example, an aggregator using the boolean `and` operation to aggregate values can be used to detect if some condition is satisfied for all vertices that pass values to that aggregator.

### Master/Worker Model

In the following, we describe how the vertices of the directed graph $G$ are mapped to a distributed system.

Similar to MapReduce [5], the Pregel framework is realized as master/worker model. When a Pregel computation is executed on a compute cluster, one of its compute nodes is designated as the master,

and all other compute nodes act as workers, which are coordinated by the master. At the beginning of a Pregel computation, the directed graph $G$ is read from a distributed file system.

All vertices of the directed graph $G$ are partitioned via a hash function, which maps each vertex ID to an integer $\in [0, \dots, N-1]$, where $N$ is the number of partitions (which is determined by the master). The master then assigns one or more partitions to each worker.

After the master instructed each worker to perform a superstep, each worker uses one thread per partition and invokes the `compute()`-function on each active vertex of that partition. In addition, each worker delivers the messages, that have been issued in the previous superstep, to their destination vertices before `compute()` is executed.

Each worker notifies the master once it is finished with the superstep, also telling it the number of active vertices in the next superstep. The master then either instructs the workers to perform the next superstep, or to save their partitions to an output file located on a distributed file system when the computation terminates (this is the case iff each worker reported that none of its vertices will be active in the next superstep).

### Combiners, Fault Tolerance, Topology Mutation

The presentation of Pregel is not complete in that it does only mention, but not elaborate on the following features (which were not necessary for the semi-naive, distributed implementation of Datalog presented in Chapter 4):

- Combiners can be used to reduce the message overhead in some cases, if only some aggregated value is of interest, but not the individual value of each single message.

- Pregel offers mechanisms for fault tolerance, e.g. when one or more workers fail.

- The topology of the input graph can be modified (i.e. adding or removing vertices) during the course of a Pregel computation.

Please refer to [7] for further information concerning these features.

### Short comparison to MapReduce

To round out the presentation of Pregel, we concisely compare Pregel and MapReduce [5], with regard to their suitability for processing re-

cursive graph algorithms on large graphs.

In the case of MapReduce, in general a sequence of MapReduce jobs will be needed to process a recursive graph algorithm, see [13]. As a consequence, the entire graph needs to be passed from one MapReduce job to the next. Considerable overhead is caused, since the entire graph is passed from one MapReduce job to the next by means of a distributed file system.

Pregel, in contrast, avoids this overhead by reading the input graph only once, partitioning its vertices and then assigning one or more partitions to each worker. As each vertex has a state and can communicate with other vertices by messages, a recursive graph algorithm can in general be processed by only one Pregel job.

## 3.3   Open Source Implementation: Apache Giraph

In this section, we introduce Giraph [9], an open source implementation of Pregel, which is developed by the Apache Software Foundation.

Technically, Giraph is based on Apache Hadoop [14]. However, only the mapping phase of a single Hadoop job is used to realize a Giraph job. Zookeeper [15] is used for fault tolerance, and the input and output file of a Giraph job are located on the Hadoop Distributed File System (HDFS), which is part of Apache Hadoop.

Conceptually, Giraph is strongly oriented towards the Pregel paper [7]. In the following, we address those concepts of Giraph which can not be found in [7], but are important for the implementation of the distributed evaluation of Datalog with Giraph.

- Worker context: The Pregel paper does not mention any functions that are executed at worker level (recall that `compute()` is executed at vertex level). In Giraph, each worker has a worker context, which includes the following user-defined functions. These can only access values that are not vertex-specific, and can e.g. be used to manage aggregators:

    - Each worker calls `preApplication()` before the beginning of superstep 0. This function can be used to create and initialize (i.e. specify an initial value) aggregators, which can then be accessed from within the worker context (in particular from within the following three functions).

    - After the last superstep has finished but before the result of the Giraph job is written to an output file, each worker calls `postApplication()`. It is mostly used for logging purposes.

14

- Each worker calls `preSuperstep()` before the beginning of each superstep, in particular before it calls `compute()` on any vertex in any of its partitions. At the time `preSuperstep()` is invoked, it is guaranteed that all aggregated values are consistent across all workers – after that, the local value of an aggregator on one worker may differ from that on another worker, as no synchronization takes place until the end of the superstep.

- Each worker calls `postSuperstep()` at the end of each superstep. This function is mostly used for logging purposes, as the values having been passed to an aggregator during the superstep have not been processed yet (i.e. the local value of an aggregator may differ from worker to worker, as explained above).

- Configuration parameters: The Pregel paper does not mention how configuration parameters like the number of workers, the worker context itself, and the format of the input and output files is specified. Details about how Giraph was configured for the implementation of a distributed, semi-naive evaluation of Datalog are given in Section 4.4.

Please visit the Giraph wiki [16] and the Giraph API [17] for references.

## 3.4 Application: Single-Source Shortest Paths

In order to round out the picture of Pregel and Giraph, in the following, we present an example application of Giraph to the single-source shortest path (SSSP) problem, inspired by [7]. The SSSP problem consists in finding a shortest path between a single, distinguished source vertex and every other vertex in a graph, where the (non-negative) cost of an edge is given by its edge value.

Figure 3.2 shows a simplified (i.e. abstracted from technical details) distributed algorithm in the spirit of Dijkstra's sequential algorithm, implemented in Giraph.

The input graph for this algorithm is supposed to be directed, with each edge having a non-negative integer as its value. The vertex IDs are arbitrary strings, and the ID of the distinguished source vertex is **a**. Each vertex value is initially set to $\infty$, which means that the vertex has not yet been reached from the distinguished source vertex.

```
1  public void compute(Iterator<Integer> msgIterator) {
2    int minDist = getVertexId().equals("a")? 0 : ∞;
3
4    while (msgIterator.hasNext())
5      minDist = Math.min(minDist, msgIterator.next());
6
7    if (minDist < getVertexValue()) {
8      setVertexValue(minDist);
9      for (Edge<String,Integer> edge: getOutEdgeMap().values())
10       sendMsg(edge.getDestVertexId(), minDist + edge.getEdgeValue());
11   }
12
13   voteToHalt();
14 }
```

Figure 3.2: (Simplified) Implementation of the SSSP problem, inspired by [7]

In the following, we use the first graph from Figure 3.3 as input graph
and apply the Giraph job given by Figure 3.2 to it. In Figure 3.3, each
active vertex is depicted in red color, and each inactive vertex in gray
color. Furthermore, each vertex is divided into two parts: the upper
part contains the vertex ID, whereas the lower part contains the vertex
value (which is initially set to $\infty$, as explained above).

- At the beginning of superstep 0, the state of the graph is identical
  to the input graph, as no computation has been conducted yet.
  Hence, no messages are received at superstep 0.
  Within superstep 0, the vertex value of the distinguished source
  vertex **a** is updated from $\infty$ to 0, which issues messages along **a**'s
  outgoing edges: one message to **b** with value 1 (as the edge **a**→**b**
  has cost 1), and another message to **c** with value 4 (as the edge
  **a**→**c** has cost 4).

- At the beginning of superstep 1, **a** is inactive (depicted in gray
  color), as each vertex calls `voteToHalt()` at the end of `compute()`
  (see Figure 3.2), and as **a** does not receive any message.
  However, **b** and **c** are active, since they receive the messages that
  **a** issued during superstep 0. For both **b** and **c**, the received value
  (1 and 4, respectively) is smaller than their current vertex value
  ($\infty$ in both cases) – i.e. a path with lower cost than the previously
  known was discovered. Hence, **b** and **c** update their vertex value
  to 1 and 4, respectively, and propagate this update along their
  outgoing edges: **b** issues a message to **c** with value 3 (as **b** can be
  reached from **a** by a path of cost 1, and the edge **b**→**c** has cost
  2), and **c** issues a message to **b** with value 7.

16

- At the beginning of superstep 2, **b** and **c** are active since they receive the messages issued during superstep 1.
  In the case of **c**, the vertex value is updated from 4 to 3, which again is propagated along **c**'s outgoing edges: **c** issues a message to **b** with value 6 (as **c** can now be reached from **a** by a path of cost 3, and the edge **c**→**b** has cost 3).
  In the case of **b**, the vertex value stays the same (hence, no messages are issued by **b**), as **b** could already be reached from **a** by a path of cost 1, and the newly discovered path has cost 7.

- At the beginning of superstep 3, only **b** is active as it receives a message with value 6. However, the vertex value of **b** stays unchanged for the same reason as in superstep 2. Hence, as **b** (i) is the only active vertex, (ii) did not issue any message in this superstep, and (iii) calls `voteToHalt()` at the end of `compute()`, the computation terminates (in the next superstep, each vertex would be inactive).

The graph of the last superstep from Figure 3.3 (without the information about active and inactive vertices) represents the result of the computation: the cost of the shortest path from **a** to **a** is 0, from **a** to **b** is 1, and from **a** to **c** is 3.
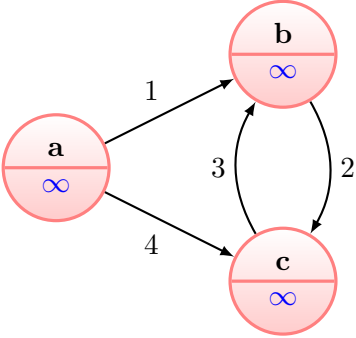
| Superstep | State of the graph at the beginning of the superstep | Messages received |
|---|---|---|
| 0 |  | none |
| 1 |  | – **b** receives 1 from **a**<br>– **c** receives 4 from **a** |
| 2 |  | – **b** receives 7 from **c**<br>– **c** receives 3 from **b** |
| 3 |  | – **b** receives 6 from **c** |

Figure 3.3: Giraph job from Figure 3.2 applied to the graph of superstep 0

# Chapter 4

# Implementing Datalog using Pregel

This chapter represents the main part of this work: Implementing a distributed algorithm for the semi-naive evaluation of Datalog using Pregel.

## 4.1 Motivation on how to implement Datalog using Pregel

This section motivates how the semi-naive evaluation of Datalog (see Figure 2.2) is mapped to Pregel.

In general, there are two ways how problems can be parallelized across multiple processors:

- each processor applies different functions to the same data, which is known as *task parallelism*, and

- each processor applies the same function to different data, which is known as *data parallelism*.

By design (as each active vertex invokes the same `compute()`-function in each superstep, and as the input data is usually distributed among the vertices), Pregel is better suited for data parallelism than for task parallelism. Nevertheless, task parallelism can also be realized in Pregel if each vertex has access to the same data, and if the `compute()`-function consists of several cases, where the case to be executed is determined by e.g. the vertex ID (e.g. realized by an `if-then-else`, or a `switch-case` statement).

According to [18, p. 125], the degree of achievable task parallelism does not grow much with the size of the problem to be solved, while the degree of achievable data parallelism usually does.

For these reasons, a data parallel mapping of the (sequential) semi-naive algorithm for the evaluation of Datalog (see Figure 2.2) to Pregel seems more worthwhile than a task parallel mapping in the case of large Datalog programs with a considerable amount of facts.

Given a Datalog program $P$ that consists of a set of rules $R$ (with non-empty body) and a set of EDB facts $F$, a data parallel mapping to Pregel implies that each Pregel vertex should perform the same computation within each superstep. In the case of the semi-naive evaluation (see Figure 2.2), this computation solely consists of applying the whole set of rules $R$ to the already derived facts. In order for each vertex to conduct this computation, we choose to store the whole set of rules $R$ as part of the vertex value of each vertex.

In addition, applying a data parallel mapping also implies that the remaining part of the Datalog program $P$, namely the EDB facts $F$ are distributed among the vertices. This poses the question of the granularity of such a distribution. If each vertex held all facts of a specific EDB or IDB relation, the degree of parallelism achievable would probably be low when evaluating Datalog programs with few, but large relations. For this reason, we choose a very fine-grained approach for distributing the facts among the vertices: for each constant $c$ appearing in $P$, a vertex with ID $c$ is created that (in its vertex value) stores all those EDB and IDB facts, where $c$ appears as argument. To illustrate, a fact `p(a,b)` would be stored both in the vertex value of vertex `a` and in the vertex value of vertex `b`. This replication may seem disadvantageous at first sight, but it allows to directly address the corresponding vertex when subgoals containing at least one constant are to be processed, e.g. `p(a,X)` and `p(Y,b)` can directly answered by vertex `a` and `b` respectively, without any indirection step being necessary.

To briefly summarize: The data parallel mapping of a Datalog program $P$ to Pregel as presented above consists in

- creating a vertex for every constant $c$ occurring in $P$ and storing (in its vertex value) all those facts from all (EDB and IDB) relations, where $c$ appears as argument, and

- storing the complete set of rules (with non-empty body) $R$ from $P$ in the value of each vertex.

Next, the question is how the function `EVAL-INCR` of Figure 2.2 can be evaluated in a data parallel way, i.e. such that each vertex performs the same computation within each superstep.

In each round, the computation to be performed consists in evaluating `EVAL-INCR`$(R, F, \Delta_{old})$, which for all rules $r \in R$, tries to derive new facts by satisfying at least one subgoal of $r$ with a fact from $\Delta_{old}$, and all remaining subgoals with facts from $F$.

As consequence of the data parallel mapping discussed above, the facts from $F$ are distributed among the vertices. However, in order to satisfy all subgoals of a Datalog rule $r$, in general facts stored at different vertices will have to be taken into account. Hence, some kind of communication between the vertices is necessary. In Pregel, communication between vertices is realized via messages, which are affected by Pregel's barrier synchronization: A message issued at superstep $S$ will not be available at its destination vertex until the next superstep $S+1$ (see Section 3.2). This implies, that in general, more than one superstep is necessary for the computation of one round of the semi-naive evaluation, i.e. for evaluating one call of `EVAL-INCR`.

To solve this issue, rules are processed on a per-subgoal base, i.e. the vertices try to collaboratively satisfy all subgoals, one after the other. To illustrate, consider a rule `r(X,Z):- p(X,Y) & q(Y,Z).` that is processed at vertex `a` whose vertex value contains the fact `p(a,b)`. This fact is used to satisfy and afterwards remove the first subgoal of the rule, leading to a new rule `r(a,Z):- q(b,Z)..` As the remaining subgoal contains the constant `b`, this rule is sent to vertex `b` for further processing (note that vertex `b` stores all facts where `b` appears as argument). Vertex `b` now proceeds in the same way, and tries to satisfy `q(b,Z)` with one of its facts, say `q(b,c)`. The resulting rule `r(a,c):- .` has an empty body, i.e. it is a fact and is stored at the vertex values of vertices `a` and `c`.

Now, consider another rule `r(X,Y):- p(X) & q(Y).` that expresses the cross product of relations `p` and `q`. No matter which vertex satisfies the subgoal `p(X)`, the resulting rule will have `q(Y)` as its only subgoal. `q(Y)` however can be satisfied by any fact having `q` as its predicate symbol, and these facts are stored in a distributed manner across all vertices (e.g. `q(a)` is stored at vertex `a`, `q(b)` is stored at vertex `b`, and so forth). Hence, each vertex needs to able to send rules (i.e. messages) to all vertices, including itself. This is realized by introducing outgoing edges from each vertex to all vertices in Pregel, used for sending messages as broadcast.

```
1   //Input: F (set of Datalog EDB facts),
2   //        R (set of Datalog rules with non-empty body)
3   Δ_old := F
4   while Δ_old ≠ ∅
5       Δ_new := EVAL-INCR-DIST(R, F, Δ_old)
6       F := F ∪ Δ_new
7       Δ_old := Δ_new
8   output F
```

Figure 4.1: Distributed semi-naive evaluation of a Datalog program

Finally, the data parallel mapping of a Datalog program $P$ to Pregel consists in

- creating a vertex for every constant $c$ mentioned in $P$ and storing (in its vertex value) all those facts from all (EDB and IDB) relations, where $c$ appears as argument,

- storing the complete set of rules (with non-empty body) $R$ from $P$ in the vertex value of each vertex,

- introducing outgoing edges from each vertex to all vertices, used for broadcast messages.

In the following section, we present the concept of parallelizing the semi-naive evaluation from Figure 2.2 on a distributed system based on Pregel.

## 4.2 Distributed Semi-naive Evaluation

In this section, we present the concept of the distributed semi-naive evaluation, with reference to the algorithm shown in Figure 4.1.

The input for the distributed semi-naive evaluation is the same as for the non-distributed algorithm: a Datalog program $P$, consisting of a set of rules (with non-empty body) $R$ and a set of EDB facts $F$. Let $C$ be the set of constants appearing in $P$.

As a first step, $P$ has to be distributed, as already described in Section 4.1. Since we use Pregel, we generate a graph consisting of one vertex for each constant $c \in C$ with the following vertex structure:

- Vertex ID: $c$
- State-flag: Active (default, see Section 3.2).
- Vertex value: Consisting of the whole set of rules $R$ and the set of facts $F_c$ from $F$ where $c$ appears as argument. Note that $\bigcup_{c \in C} F_c = F$.

- Incoming messages: Empty set (default, see Section 3.2).
- Outgoing edges: Outgoing edges to all vertices (recall that these are used for broadcast messages)

The state of the vertices changes with each performed superstep and finally represents the result of the evaluation after the last superstep, in a distributed manner.

In the following, we realise the algorithm of Figure 4.1 within Pregel:

- In line 3, $\Delta_{old}$ is initialized with the set of EDB facts $F$. Note that, as explained above, $F$ is stored in a distributed manner across all vertices, by design of the initial vertex values: $F = \bigcup_{c \in C} F_c$. The same applies to $\Delta_{old}$, i.e. a vertex $c$ stores $\Delta_{old,c}$, the set of all facts from $\Delta_{old}$ where $c$ occurs as argument.

- Next, we discuss the `while` loop (lines 4-7):
  - In line 5, `EVAL-INCR-DIST(`$R$`,`$F$`,`$\Delta_{old}$`)` is evaluated, where $R$ is the set of Datalog rules with non-empty body, $F$ is the set of all facts that have been derived in previous rounds, and $\Delta_{old}$ is the set of facts that has been derived in the previous round. `EVAL-INCR-DIST` is computed in a sequence of supersteps using messages for communication between the vertices, which is described in detail below. The result of `EVAL-INCR-DIST` is the set $\Delta_{new}$, consisting of all facts that have been derived in the current round, but not in one of the previous rounds. $\Delta_{new}$ is stored in a distributed manner across all vertices: a vertex $c$ stores $\Delta_{new,c}$, the set of all facts from $\Delta_{new}$ where $c$ occurs as argument.
  - In line 6, the union $F \cup \Delta_{new}$ is formed. This set is stored and computed in a distributed manner, as each vertex forms the union locally: $F_c := F_c \cup \Delta_{new,c}$ (recall that $\bigcup_{c \in C} F_c = F$ and $\bigcup_{c \in C} \Delta_{new,c} = \Delta_{new}$).
  - In line 7, $\Delta_{old}$ is assigned the value of $\Delta_{new}$ for the next round. This assignment is also done in a distributed manner: each vertex $c$ sets $\Delta_{old,c}$ to the value of $\Delta_{new,c}$.

  The `while` loop (see line 4) terminates, when no facts for $\Delta_{new}$ are inferred in `EVAL-INCR-DIST`. This, in consequence, violates the condition of the loop (as $\Delta_{old} = \Delta_{new} = \emptyset$).

- The output $F$ (see line 8) is distributed across all vertices. Collecting the distributed output is handled by design of the Pregel

system: After the last superstep, all vertex values are written into an output file on the distributed file system. The union over all vertex values contains the set of all rules $R$ (with non-empty body), and the desired result $F = \bigcup_{c \in C} F_c$.

Next, we discuss the the local execution of `EVAL-INCR-DIST` on an arbitrary vertex with ID $c$.

We first define a helper function `EVAL-DIST(`$R'$`,`$F'$`)`, which receives a set of rules $R'$ with non-empty body and a set of facts $F'$ as arguments, and returns the following set:

$$\{h\sigma :\text{-}\; b_1\sigma\; \&\; \ldots\; \&\; b_n\sigma \mid \exists(h :\text{-}\; b_1\; \&\; \ldots\; \&\; b_n) \in R'($$
$$\exists k(1 \leq k \leq n\; \wedge$$
$$\exists\sigma, \text{where } \sigma \text{ is a substitution}, (b_k\sigma \in F')))\}$$

In the algorithm of Figure 4.1, one call of `EVAL-INCR-DIST` corresponds to one round of the semi-naive evaluation. As noted above, the evaluation of `EVAL-INCR-DIST` requires a sequence of supersteps $S$. We now examine one arbitrary superstep of $S$:

Initially, we filter the incoming messages for rules with empty body (i.e. facts), which are added to $\Delta_{new,c}$. The remaining messages are denoted by $R_{in,c}$. We have to consider two cases:

- If the current superstep is the first of a round, we define $\Delta'_{old,c}$ to be the set of facts from $\Delta_{old,c}$, where $c$ appears as first argument, and evaluate `EVAL-DIST(`$R, \Delta'_{old,c}$`)`.

- Otherwise (if the current superstep is not the first superstep of a round), we evaluate `EVAL-DIST(`$R_{in,c}, F_c$`)`.

In any case, the result of `EVAL-DIST` is denoted by $R_{ver,c}$.

Let $R_{out,c}$ be an empty set. All rules $r \in R_{ver,c}$ need to be verified. For this, all variable-free subgoals of $r$ are checked:

- If the variable-free subgoal is ordinary and contains $c$ as argument, we check if the subgoal is contained in $F_c$ as fact. If this is the case, the subgoal is removed from $r$, else the rule $r$ is discarded (i.e. not added to $R_{out,c}$).

- If the variable-free subgoal is ordinary and does not contain $c$ as argument, $r$ is not altered.

- If the variable-free subgoal is built-in and evaluates to `true`, the subgoal is removed from $r$. If it evaluates to `false`, $r$ is discarded (i.e. not added to $R_{out,c}$).

24

```
1   path(X,Y) :- edge(X,Y).
2   path(X,Z) :- path(X,Y) & path(Y,Z).
3   edge(a,b).
4   edge(b,c).
5   edge(c,a).
```

Figure 4.2: Example Datalog program for computing the transitive closure over a cyclic base relation

After all variable-free subgoals of $r$ have been checked, $r$ is added to $R_{out,c}$, unless it was discarded.

As a last step, the rules in $R_{out,c}$ are issued as messages. The recipients for the messages are obtained in the following way: If the body of a rule is empty, the recipients are all vertices whose ID appears as argument in the head of the rule. Otherwise, the ordinary subgoals are searched for constants. If constants were found, the rule is sent to all vertices whose ID is equal to one of the constants. Otherwise (if the ordinary subgoals contain no constants), the rule is broadcasted, as any vertex might be able to satisfy one of its subgoals. After issuing the messages, the superstep ends.

The next superstep of $S$ starts. `EVAL-INCR-DIST` terminates when no more messages are issued in a superstep. After the last superstep of $S$, $\Delta_{new}$ is stored in a distributed manner across all vertices: $\Delta_{new} = \bigcup_{c \in C} \Delta_{new,c}$.

This concludes the concept of the distributed semi-naive evaluation using Pregel. In the next section, we illustrate our approach with a simple example, followed by a section, which presents the Java implementation of the algorithm.

## 4.3   Evaluation by Example

In this section, we carry out our algorithm for the distributed evaluation of Datalog programs on a simple, yet illustrative example program shown in Figure 4.2.

First, we begin by creating the vertex structure. The only constants appearing in the Datalog program in Figure 4.2 are a, b, and c, hence we create vertices with IDs a, b, and c. Each vertex stores all those facts from Figure 4.2, where its ID appears as argument. Therefore, vertex a stores the facts edge(a,b) and edge(c,a), vertex b stores the facts

`edge(a,b)` and `edge(b,c)`, and vertex `c` stores the facts `edge(b,c)` and `edge(c,a)`. Furthermore, each vertex stores the two rules shown in lines 1 and 2 of Figure 4.2, and each vertex has outgoing edges to all vertices.

In the following tables, the left column depicts the state of the graph, and the right column contains all messages that are available at the corresponding superstep. For reasons of clarity, the rules in lines 1 and 2 of Figure 4.2 are omitted from the vertices in the left column – nevertheless, they are part of the vertex value of each vertex. Each fact that is stored in a vertex value has an index indicating the number of the round, in which this fact has been derived (e.g. $\texttt{path(a,a)}_3$ means that the fact `path(a,a)` has been derived in round 3). Note that in the tables below, the facts of $\Delta_{old}$ are always those facts indexed with the number of the previous round.

In the following descriptions, we use the notation introduced in Section 4.2. Additionally, by "rule 1" and "rule 2", we denote the rule in line 1 and 2 of Figure 4.2, respectively.

**Superstep 0 (round 1)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | none |

The first superstep, superstep 0, is the first superstep of round 1. In this case, each vertex $v$ just uses those facts from $\Delta_{old,v}$ where its vertex ID appears as first argument, e.g. vertex `a` just uses `edge(a,b)`, but does not use `edge(c,a)`. Let's continue to describe the computation that vertex `a` performs: Regarding rule 1, the subgoal `edge(X,Y)` can be satisfied by `edge(a,b)`, leading to the rule `path(a,b):- edge(a,b).` which is added to $R_{ver,\texttt{a}}$. Regarding rule 2, no subgoal can be satisfied, as no `path`-facts have been derived yet. Verifying the only rule in $R_{ver,\texttt{a}}$ leads to the fact `path(a,b):- .`, which is added to $R_{out,\texttt{a}}$, and sent to vertices `a` and `b`.

The computation at the remaining vertices `b` and `c` is performed in an analogous manner and results in the fact `path(b,c):- .` and `path(c,a):- .`, respectively. Each of these two facts is also sent to its corresponding vertices.

## Superstep 1 (round 1)

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(a,b). path(c,a).`<br><br>**b:** `path(a,b). path(b,c).`<br><br>**c:** `path(b,c). path(c,a).` |

Superstep 1 also belongs to round 1, as messages have been issued in superstep 0. These messages are now available at their destination vertices, as depicted in the right column of the above table. Let's again look at vertex `a`: it receives the facts `path(a,b)` and `path(c,a)` and adds these to $\Delta_{new,\mathtt{a}}$, which is stored in `a`'s vertex value. No rules with non-empty subgoal have been received, so $R_{in,\mathtt{a}}$ is empty, therefore `EVAL-DIST` produces an empty output and no messages are issued.

Analogously, vertex `b` and `c` receive only facts and store them in $F_\mathtt{b}$ and $F_\mathtt{c}$, respectively.

## Superstep 2 (round 2)

| State of the graph at the beginning | Messages received |
|---|---|
|  | none |

As none of the vertices issued a message in superstep 1, `EVAL-INCR-DIST` terminates and each vertex $v$ adds all facts from $\Delta_{new,c}$ to $F_c$ and sets $\Delta_{old,c}$ to the value of $\Delta_{new,c}$. Since $\Delta_{old} \neq \emptyset$, a new round starts (round 2). As superstep 2 is the first superstep of this new round, for each vertex $v$ only those facts in $\Delta_{old,v}$ are considered, where $v$ appears as first argument.

For vertex `a`, the computation proceeds as follows: `path(a,b)` is the only fact in $\Delta_{old,\mathtt{a}}$ with `a` as first argument. This fact can't satisfy the subgoal of rule 1, but it can satisfy both subgoals `path(X,Y)` and `path(Y,Z)` of rule 2. For the first subgoal `path(X,Y)` the resulting rule is `path(a,Z):- path(a,b) & path(b,Z).`, and for the second subgoal it is `path(X,b):- path(X,a) & path(a,b).`. In the verification step, the satisfied subgoal `path(a,b)` is removed from both rules,

which results in the rules `path(a,Z):- path(b,Z).` and `path(X,b):-`
`path(X,a)..` These are sent to vertex `b` and `a`, respectively, for further
processing.

Again, the computation at the remaining vertices `b` and `c` is per-
formed in an analogous manner.

**Superstep 3 (round 2)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(c,Z):- path(a,Z).`<br>    `path(X,b):- path(X,a).`<br><br>**b:** `path(a,Z):- path(b,Z).`<br>    `path(X,c):- path(X,b).`<br><br>**c:** `path(b,Z):- path(c,Z).`<br>    `path(X,a):- path(X,c).` |

As messages have been issued in superstep 2, superstep 3 still belongs
to the same round as superstep 2 (round 2). In contrast to superstep
1, the messages received at the current superstep 3 are no facts, but
rules with non-empty body.

In the case of vertex `a`, the computation proceeds as follows: As
both incoming messages `path(c,Z):- path(a,Z).` and `path(X,b):-`
`path(X,a).` have non-empty body, they are added to $R_{in,\mathsf{a}}$. Vertex `a`
can use all facts from $F_\mathsf{a}$ for satisfying the subgoals of these two rules.

Concerning the first rule, the subgoal `path(a,Z)` can be only sat-
isfied by the fact `path(a,b)`. Concerning the second rule, the subgoal
`path(X,a)` can be only satisfied by `path(c,a)`. After verification, both
rules result in the fact `path(c,b):- .`, which is sent to vertices `b` and
`c`.

For the vertices `b` and `c`, the derived facts are `path(a,c)` and
`path(b,a)`, respectively.

## Superstep 4 (round 2)

| State of the graph at the beginning | Messages received |
|---|---|
| $a$<br>$edge(a,b)_0$<br>$edge(c,a)_0$<br>$path(a,b)_1$<br>$path(c,a)_1$<br><br>$b$<br>$edge(a,b)_0$<br>$edge(b,c)_0$<br>$path(a,b)_1$<br>$path(b,c)_1$<br><br>$c$<br>$edge(b,c)_0$<br>$edge(c,a)_0$<br>$path(b,c)_1$<br>$path(c,a)_1$ | **a:** `path(a,c). path(b,a).`<br><br>**b:** `path(b,a). path(c,b).`<br><br>**c:** `path(a,c). path(c,b).` |

The computation in superstep 4 is similar to that of superstep 1: As messages have been issued in superstep 3, superstep 4 still belongs to round 2 and just facts are received as incoming messages, each of which is stored in the set $\Delta_{new,v}$ in the vertex value of the corresponding vertex $v$.

## Superstep 5 (round 3)

| State of the graph at the beginning | Messages received |
|---|---|
| $a$<br>$edge(a,b)_0$<br>$edge(c,a)_0$<br>$path(a,b)_1$<br>$path(c,a)_1$<br>$path(a,c)_2$<br>$path(b,a)_2$<br><br>$b$<br>$edge(a,b)_0$<br>$edge(b,c)_0$<br>$path(a,b)_1$<br>$path(b,c)_1$<br>$path(b,a)_2$<br>$path(c,b)_2$<br><br>$c$<br>$edge(b,c)_0$<br>$edge(c,a)_0$<br>$path(b,c)_1$<br>$path(c,a)_1$<br>$path(a,c)_2$<br>$path(c,b)_2$ | none |

As none of the vertices issued a message in superstep 4, `EVAL-INCR-DIST` terminates and each vertex $v$ adds all facts from $\Delta_{new,c}$ to $F_c$ and sets $\Delta_{old,c}$ to the value of $\Delta_{new,c}$. Since $\Delta_{old} \neq \emptyset$, a new round is started (round 3), and hence superstep 5 is the first superstep in this new round. This implies, that for each vertex $v$ only those facts in $\Delta_{old,v}$ are considered, where $v$ appears as first argument.

In the case of vertex `a`, the computation proceeds as follows. The only fact in $\Delta_{old,a}$ with `a` as first argument is `path(a,c)`. This fact can't satisfy the subgoal of rule 1, but it can satisfy both subgoals `path(X,Y)` and `path(Y,Z)` of rule 2. For the first subgoal `path(X,Y)`, the resulting rule is `path(a,Z):- path(a,c) & path(c,Z).`, and for the second subgoal, it is `path(X,c):- path(X,a) & path(a,c)..` In the verification step, the satisfied subgoal `path(a,c)` is removed from

29

both rules, which results in the rules `path(a,Z):- path(c,Z).` and
`path(X,c):- path(X,a)..` These are sent to vertex `c` and `a` for further processing, respectively.

**Superstep 6 (round 3)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(b,Z):- path(a,Z).`<br>    `path(X,c):- path(X,a).`<br><br>**b:** `path(c,Z):- path(b,Z).`<br>    `path(X,a):- path(X,b).`<br><br>**c:** `path(a,Z):- path(c,Z).`<br>    `path(X,b):- path(X,c).` |

The computation of superstep 6 is similar to that of superstep 3: each vertex $v$ uses facts from $F_v$ to satisfy a subgoal in the incoming rules.

Let's again look at vertex `a`. Its first incoming rule `path(b,Z):- path(a,Z).` leads to the facts `path(b,b)` and `path(b,c)`, and its second incoming rule leads to the facts `path(c,c)` and `path(b,c)`. All of these facts are sent to the corresponding vertices.

**Superstep 7 (round 3)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(a,a).   path(a,b).`<br>    `path(c,a).`<br><br>**b:** `path(a,b).   path(b,b).`<br>    `path(b,c).`<br><br>**c:** `path(b,c).   path(c,a).`<br>    `path(c,c).` |

The computation in superstep 7 is again similar to that of superstep 1: superstep 7 still belongs to round 3 and just facts are received as incoming messages, each of which is stored in the set $\Delta_{new,v}$ in the vertex value of the corresponding vertex $v$.

**Superstep 8 (round 4)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | none |

As no message was issued in superstep 7, `EVAL-INCR-DIST` terminates and each vertex $v$ adds all facts from $\Delta_{new,c}$ to $F_c$ and sets $\Delta_{old,c}$ to the value of $\Delta_{new,c}$. Since $\Delta_{old} \neq \emptyset$, a new round is started (round 4), and hence superstep 8 is the first superstep in this new round. This implies, that for each vertex $v$ only those facts in $\Delta_{old,v}$ are considered, where $v$ appears as first argument.

In the case of vertex `a`, `p(a,a)` is the only fact that can be used to satisfy the subgoals of rule 2, which finally generates the rules `path(a,Z):- path(a,Z).` and `path(X,a):- path(X,a).` and which both are sent to vertex `a` for further processing.

As all rules issued in this superstep are not able to derive new facts (since the only subgoal of the body is equal to the head), these rules could be discarded in an optimized version of the algorithm.

**Superstep 9 (round 4)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(a,Z):- path(a,Z).`<br>     `path(X,a):- path(X,a).`<br><br>**b:** `path(b,Z):- path(b,Z).`<br>     `path(X,b):- path(X,b).`<br><br>**c:** `path(c,Z):- path(c,Z).`<br>     `path(X,c):- path(X,c).` |

The computation of superstep 9 is again similar to that of superstep 3:

each vertex $v$ uses facts from $F_v$ to satisfy a subgoal in each incoming rule.

**Superstep 10 (round 4)**

| State of the graph at the beginning | Messages received |
|---|---|
|  | **a:** `path(a,a).` `path(a,b).` `path(a,c).` `path(b,a).` `path(c,a).`<br><br>**b:** `path(a,b).` `path(b,a).` `path(b,b).` `path(b,c).` `path(c,b).`<br><br>**c:** `path(a,c).` `path(b,c).` `path(c,a).` `path(c,b).` `path(c,c).` |

The computation in superstep 10 is again similar to that of superstep 1: superstep 10 still belongs to round 4, and just facts are received as incoming messages. However, none of these facts is "new", i.e. all received facts have already been derived in previous rounds and are already contained in $F$. Hence, `EVAL-INCR-DIST` returns the empty set, and in consequence, the `while` loop terminates as its condition is not satisfied anymore.

The result of the computation is contained in the vertex values of `a`, `b`, and `c`, and represents the minimal model of the Datalog program from Figure 4.2.

## 4.4  Java Implementation

In this section, we give an overview of the Java implementation of the algorithm in Figure 4.1. The complete source-code can be found in the Appendix.

The implementation is based on Apache Giraph version 0.1 [9] and Hadoop version 0.20.203.0 [14]. In addition to that, we make use of the parser component of IRIS [19], an open source Datalog reasoner.

The `main()`-function is located in the class `PDoP` (Appendix A), and calls `run()`. In the `run()`-function, the Giraph job is configured and started. The configuration consists of the following tasks:

- Setting the worker context
  (class `PDoPWorkerContext`, Appendix B)

- Specifying the input format
  (class `PDoPInputFormat`, Appendix F)

- Specifying the output format
  (class `PDoPOutputFormat`, Appendix G)

- Setting the worker configuration (Specifying the number of workers, and fault tolerance)

**Setting the worker context**

The worker context (class `PDoPWorkerContext`, Appendix B) specifies the following four functions `preApplication()`, `postApplication()`, `preSuperstep()`, and `postSuperstep()`:

- `preApplication()`: In this function, aggregators are registered and initialized. For our application, we use four aggregators:

  - `VoteAggregator` is used to check if all vertices are finished with the current round, and whether a new round can start. The `VoteAggregator` is realized by a boolean aggregator that uses the boolean `and` operation as method of aggregation (class `AndAggregator`, Appendix C).
  - `LastTimeAggregator` is used to check if a new round was started in the previous superstep, which indicates the first superstep of a round. The `LastTimeAggregator` is realized by a boolean aggregator that uses the boolean `or` operation as method of aggregation (class `OrAggregator`, Appendix D).
  - `TerminateAggregator` is used to detect when the algorithm is finished. In simple Giraph jobs, `voteToHalt()` is called at the end of every superstep, thus making the job stop when no more messages are issued, but in our case no more messages just indicate that a new round can start. For this reason we need this aggregator to indicate when to terminate (utilizing the `OrAggregator`, Appendix D).
  - `RoundAggregator` (Appendix E) is used as the round counter. The round counter indicates how often the `while` loop of the evaluation algorithm has been entered so far. This is useful if we want to identify in which round a fact was derived, and gives the possibility to differentiate facts of $\Delta_{old}$, $\Delta_{new}$, and $F$ according to a specific round. (The `RoundAggregator` is

33

useful for debugging purposes, but the algorithm could be implemented without it.)

- `preSuperstep()`: This function contains the logic for evaluating the aggregated values. Placing the logic in `preSuperstep()` ensures, that every worker has the same aggregator values before calling `compute()` on its active vertices.

- `postSuperstep()` and `postApplication()` are not used for the implementation of our algorithm.

**Specifying the input and output format**

The input and output format are specified by the following classes:

- The class `PDoPInputFormat` uses the class `PDoPReader` to read the input file from the HDFS, which contains one line for each vertex represented in the following structure:

$$[\texttt{<ID> <RULES><FACTS> <EDGES>}]$$

where `<ID>` is the vertex ID, `<RULES>` are all rules with non-empty body, `<FACTS>` are only those facts concerning the specific vertex (as specified in Section 4.2), and `<EDGES>` are all vertex IDs (separated by commas) indicating the destination for the edges. Rules and facts are given in Datalog syntax, as specified by IRIS.

- The class `PDoPOutputFormat` uses the class `PDoPWriter` to write the results of each vertex to an output file on the HDFS.

Having finished the configuration of the Giraph job, we can start the evaluation, which brings us to the central function: `compute()`.

**The `compute()`-function**

The structure of the `compute()`-function is as follows:

1. If the value of the `TerminateAggregator` is `true`, the vertex calls `voteToHalt()` and the computation terminates.

2. The vertex value is parsed using the class `RuleFactQuerySplitter` (Appendix J).

3. The incoming messages are parsed and filtered using the `RuleFactQuerySplitter` and the `Filter` classes (Appendix K). `Filter` uses the interface `Condition` (Appendix L) together with

34

its implementations for the type of facts to be generated (Appendices: M,N,O). With `Filter`, we can generate $\Delta_{old,v}$, $\Delta_{new,v}$, and $F_v$ from all facts that are available to a vertex $v$, as needed (see Section 4.2).

4. Next, the evaluation using the appropriate rules according to Section 4.2 is executed by utilizing the class `PartialEvaluation` (Appendix P). If the evaluation generated new rules, these are sent as messages.

5. As a last step, the `VoteAggregator` is called (with a positive vote, if no messages were sent), and the vertex value is updated.

Finally, in order to represent relations, the class `SimplePDoPRelation` (Appendix Q) has been created, based on the class `SimpleRelation` from the IRIS reasoner, as the latter does not offer a public constructor.

For further details refer to the Appendix, as this section gave just an overview of our implementation.

# Chapter 5

# Related Work

Besides this work, another approach to the distributed evaluation of Datalog is DEDALUS [10]. In the following, we give a short introduction to its key ideas.

DEDALUS is a language that builds upon Datalog and includes negation, aggregation, and a choice construct. In order to capture time, DEDALUS uses a `successor` predicate with `successor(X,Y)` being true iff $X = Y + 1$. In a well-formed DEDALUS rule, every subgoal must use the same variable $\mathcal{T}$ at its rightmost attribute, and the variable $\mathcal{S}$ at its rightmost head attribute. $\mathcal{S}$ and $\mathcal{T}$ describe time, and $\mathcal{S}$ must be restricted in one of the following two ways:

1. A rule is deductive if $\mathcal{S}$ is bound to the value $\mathcal{T}$; that is, the body contains the subgoal $\mathcal{S} = \mathcal{T}$.

2. A rule is inductive if $\mathcal{S}$ is the successor of $\mathcal{T}$; that is, the body contains the subgoal `successor(`$\mathcal{T}$`, `$\mathcal{S}$`)`.

[10]

With these supplements, the two key features of DEDALUS are achieved: mutable state and asynchronous processing and communication.

**Mutable State**

Mutable state is achieved by introducing predicates `p_pos` and `p_neg` for each EDB predicate `p` in a DEDALUS program, together with the following rules

$$p\_pos(\texttt{A}_1,\texttt{A}_2,\ldots,\texttt{A}_n,\mathcal{S}) \;\texttt{:-}\; p(\texttt{A}_1,\texttt{A}_2,\ldots,\texttt{A}_n,\mathcal{T}), \;\; \mathcal{S}=\mathcal{T}. \quad (5.1)$$

$$\begin{aligned}
p\_pos(\texttt{A}_1,\texttt{A}_2,\ldots,\texttt{A}_n,\mathcal{S}) \;\texttt{:-}\;\; & p\_pos(\texttt{A}_1,\texttt{A}_2,\ldots,\texttt{A}_n,\mathcal{T}), \\
& \neg p\_neg(\texttt{A}_1,\texttt{A}_2,\ldots,\texttt{A}_n,\mathcal{T}), \\
& \texttt{successor}(\mathcal{T},\ \mathcal{S}). \qquad\qquad (5.2)
\end{aligned}$$

Rule 5.1 is a deductive rule that derives a `p_pos`-fact for each `p`-fact, but also allows rule 5.2 to derive additional `p_pos`-facts. Rule 5.2 is an inductive rule that captures the notion of mutable state: as long as no corresponding `p_neg`-fact exists, a `p_pos`-fact at time $\mathcal{T}$ will also be true at the next point in time $\mathcal{S}$.

## Asynchronous Processing and Communication

By introducing a `choice` construct, asynchronous processing and communication is possible. `choice` allows to model non-determinism, e.g. as it occurs in unreliable networks with delay, loss of or out-of-order messages. A `choose` subgoal looks as follows:

$$\texttt{choose(X,Y)},$$

where $\texttt{X} := (\texttt{X}_1,\texttt{X}_2,\ldots,\texttt{X}_n)$ and $\texttt{Y} := (\texttt{Y}_1,\texttt{Y}_2,\ldots,\texttt{Y}_m)$. This subgoal enforces the functional dependency $\texttt{X} \to \texttt{Y}$: for each value of $\texttt{X}$, a value for $\texttt{Y}$ has to be "chosen" non-deterministically.

With these additions, implementing Lamport clocks and reliable broadcasts is possible in DEDALUS (which is not possible in pure Datalog), see [10].

## Comparison of Dedalus with our implementation

DEDALUS and our implementation are similar in the following aspects:

- Concerning the distribution, both in our implementation and in DEDALUS, each compute node stores the whole set of rules, while the facts are distributed across the compute nodes.

- Concerning the evaluation, DEDALUS also makes use of the semi-naive algorithm for efficiently evaluating a single timestep.

However, they differ in the following aspects:

- Our implementation focuses on pure Datalog, while DEDALUS is an extension thereof, including (among others) an explicit notion of time, i.e. a fact $p(\texttt{A}_1,\ldots,\texttt{A}_n,\mathcal{T})$ is interpreted as the fact $p(\texttt{A}_1,\ldots,\texttt{A}_n)$ being true at timestep $\mathcal{T}$.

- DEDALUS is not synchronized by design, i.e. it may happen that facts from the future are used to derive facts in the past. Our implementation, however, is synchronized by design of the underlying framework Pregel.

# Chapter 6

# Conclusion and Future Work

Not only in the case of social networks, companies are interested in data they can use to generate profit. Thus, as a logical consequence, collecting this data (e.g. user activity) in data warehouses has become common practise. However, the efficient evaluation of this data is not yet that well established [20]. Knowledge discovery in big data sets is gaining ever growing importance, with the side effect of a renewed interest in Datalog. Recent work [21] has successfully connected these three cornerstones: Datalog, Machine Learning, and Big Data. This gives us sufficient claims that the efficient evaluation of arbitrary Datalog programs on a compute cluster is of great value, not only to academia.

In this work, we developed a data parallel, semi-naive evaluation algorithm for Datalog programs that utilizes the distributed graph processing framework Pregel. We have shown the method of operation of this distributed algorithm, and are fairly confident that this approach scales well for Big Data applications – even if an experimental evaluation remains to be done, as it is outside the scope of this bachelor's thesis.

When querying a data set multiple times, computing the minimal model just once and using it for query answering is more efficient than evaluating each query separately, e.g. in a top-down manner. With our approach, we are able to generate the minimal model using Pregel, which was explicitly designed for practical computing problems concerning large graphs [7].

In future work, the distributed, semi-naive evaluation algorithm for Datalog developed in this thesis can be extended to Datalog¬ by adding stratified negation, which can be easily realized. As this implementation can be seen as a proof of concept, no special emphasis has been put on optimization, yet. For example, Pregel combiners (see Section 3.2)

can be used in further work to reduce the messaging overhead caused by some redundant messages. In addition, messages that are not able to derive new knowledge (e.g. `p(X):- p(X).`) need not to be considered at all.

# List of Figures

# Appendix

**Appendix A:**
pdop.PDoP.java

```java
package pdop;

import java.io.IOException;
import java.util.HashMap;
import java.util.HashSet;
import java.util.Iterator;
import java.util.LinkedList;
import java.util.List;
import java.util.Map;
import java.util.Set;

import org.apache.giraph.graph.EdgeListVertex;
import org.apache.giraph.graph.GiraphJob;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.BooleanWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;
import org.apache.mahout.math.Arrays;
import org.deri.iris.api.basics.IAtom;
import org.deri.iris.api.basics.IPredicate;
import org.deri.iris.api.basics.IRule;
import org.deri.iris.api.basics.ITuple;
import org.deri.iris.api.terms.ITerm;
import org.deri.iris.factory.Factory;
import org.deri.iris.storage.IRelation;

import pdop.giraph.graph.PDoPWorkerContext;
import pdop.giraph.graph.aggregator.AndAggregator;
import pdop.giraph.graph.aggregator.OrAggregator;
import pdop.giraph.graph.aggregator.RoundAggregator;
import pdop.giraph.io.PDoPInputFormat;
import pdop.giraph.io.PDoPOutputFormat;
import pdop.iris.PartialEvaluation;
```

```java
37 import pdop.iris.RuleFactQuerySplitter;
38 import pdop.iris.SimplePDoPRelation;
39 import pdop.iris.filter.Filter;
40 import pdop.iris.filter.conditions.EqualCondition;
41 import pdop.iris.filter.conditions.LessEqualCondition;
42 import pdop.iris.filter.conditions.TrueCondition;
43
44 import com.google.common.base.Preconditions;
45
46 /*
47  * Parallel Datalog on Pregel
48  * (Giraph implementation after SimpleShortestPathsVertex
       example)
49  */
50
51 public class PDoP extends EdgeListVertex<Text, Text, Text,
     Text> implements
52     Tool {
53
54   @Override
55   public void compute(Iterator<Text> msgIterator) throws
        IOException {
56
57     RoundAggregator roundAggregator = (RoundAggregator)
          getAggregator("RoundAggregator");
58     AndAggregator voteAggregator = (AndAggregator)
          getAggregator("VoteAggregator");
59     OrAggregator lastTimeAggregator = (OrAggregator)
          getAggregator("LastTimeAggregator");
60     OrAggregator terminateAggregator = (OrAggregator)
          getAggregator("TerminateAggregator");
61
62     if (terminateAggregator.getAggregatedValue().get()) {
63       System.out.println("Finished␣after␣" + getSuperstep()
64           + "␣Supersteps.␣In␣Round␣"
65           + roundAggregator.getAggregatedValue().get());
66       voteToHalt();
67     } else {
68       boolean voteForNextRound = true;
69
70       RuleFactQuerySplitter splitter = new
            RuleFactQuerySplitter(
71         getVertexValue().toString());
72       List<IRule> rules = splitter.getRules();
73       Map<IPredicate, IRelation> facts = splitter.getFacts()
          ;
74
75       if (getSuperstep() == 0) {
76         // add timestamps to first facts
77         Map<IPredicate, IRelation> factsWtime = new HashMap<
            IPredicate, IRelation>();
78         for (IPredicate pred : facts.keySet()) {
79           factsWtime.put(
80               Factory.BASIC.createPredicate(
```

```
81                pred.getPredicateSymbol() + "_0",
82                pred.getArity()), facts.get(pred));
83      }
84    facts = factsWtime;
85    /*
86     * all facts look like this now: a_0('a','b'). These
              are set to
87     * VertexValue at the end of compute
88     */
89    }
90
91    StringBuilder ruleBuilder = new StringBuilder();
92    while (msgIterator.hasNext()) {
93      ruleBuilder.append(msgIterator.next().toString());
94    }
95
96    RuleFactQuerySplitter msgSplitter = new
          RuleFactQuerySplitter(
97        ruleBuilder.toString());
98
99    List<IRule> msgRules = msgSplitter.getRules();
100   List<IRule> partialRules = new LinkedList<IRule>();
101
102   Map<IPredicate, IRelation> msgedFacts = new HashMap<
          IPredicate, IRelation>();
103   Map<IPredicate, IRelation> truncatedFacts = Filter.
          filter(facts,
104       new TrueCondition());
105   long timeStamp = roundAggregator.getAggregatedValue().
          get();
106
107   // check msgedRules for real facts and partial Rules
108   for (IRule msgRule : msgRules) {
109     if (msgRule.getBody().size() == 0) { // add as fact
110       IAtom headAtom = msgRule.getHead().get(0).getAtom
              ();
111       ITuple headTuple = headAtom.getTuple();
112       IPredicate factPred = headAtom.getPredicate();
113       IRelation tuplesOfPredicate = truncatedFacts.get(
              factPred);
114       if (tuplesOfPredicate == null
115           || !(tuplesOfPredicate.contains(headTuple))) {
116         // fact not there: add.
117         // otherwise fact is there: do nothing.
118         String factPredSymbol = factPred.
              getPredicateSymbol()
119           + "_" + timeStamp;
120         IPredicate newFactPred = Factory.BASIC.
              createPredicate(
121           factPredSymbol, headTuple.size());
122         IRelation factRel;
123         if (msgedFacts.containsKey(newFactPred)) {
124           factRel = msgedFacts.get(newFactPred);
125         } else {
```

```
126              factRel = new SimplePDoPRelation ();
127            }
128            factRel.add(headTuple);
129            msgedFacts.put(newFactPred, factRel);
130          }
131        } else { // add rule as partialRule
132          partialRules.add(msgRule);
133        }
134      }
135      facts.putAll(msgedFacts);
136
137      /*
138       * Vertex value is set with these rules and facts.
                Following only
139       * messages are generated with these rules, but no
                local Rules or
140       * facts are changed.
141       */
142
143      if (getSuperstep() == 0) {
144        // First evaluation of rules
145        voteForNextRound = evaluateThis(truncatedFacts,
              rules);
146      } else {
147        // if last time was increased use deltafacts and
                evalutate with
148        // all rules!
149        if (lastTimeAggregator.getAggregatedValue().get()) {
150          Map<IPredicate, IRelation> deltaFacts = Filter.
                filter(
151            facts, new EqualCondition(roundAggregator
152              .getAggregatedValue().get() - 1));
153          // remove facts with vertex ID not appearing as
                first
154          // argument from deltaFacts! (not yet implemented)
155          voteForNextRound = evaluateThis(deltaFacts, rules)
                ;
156        } else {
157          // if last time was not increased use all facts
                from the
158          // previous rounds and evaluate only with partial
                Rules!
159          Map<IPredicate, IRelation> previousFacts = Filter.
                filter(
160            facts, new LessEqualCondition(roundAggregator
161              .getAggregatedValue().get() - 1));
162          voteForNextRound = evaluateThis(previousFacts,
              partialRules);
163        }
164      }
165      voteAggregator.aggregate(new BooleanWritable(
          voteForNextRound));
166
167      // Update vertex value
```

```java
168        StringBuilder builder = new StringBuilder();
169        for (IRule rule : rules) {
170          builder.append(rule);
171        }
172        for (IPredicate pred : facts.keySet()) {
173          for (int i = 0; i < facts.get(pred).size(); i++) {
174            builder.append(pred);
175            builder.append(facts.get(pred).get(i));
176            builder.append(".");
177          }
178        }
179        setVertexValue(new Text(builder.toString()));
180      }
181    }
182
183    /*
184     * evaluates the Rules with the given facts, collects all
             the messages and
185     * eliminates duplicates Sends them out in the end.
186     */
187    public boolean evaluateThis(Map<IPredicate, IRelation>
             facts,
188        List<IRule> rules) {
189      boolean voteForNextRound = false; // only set this to
             true, when no
190                         // message was added
191      Map<IRule, Set<ITerm>> cleanedMessages = new HashMap<
             IRule, Set<ITerm>>();
192      Set<IRule> broadcasts = new HashSet<IRule>();
193      for (IRule rule : rules) {
194        PartialEvaluation eval = new PartialEvaluation(facts,
             rule);
195        Map<IRule, Set<ITerm>> partialRules = eval.
             getPartialRules();
196        for (IRule msg : partialRules.keySet()) {
197          if (!cleanedMessages.containsKey(msg)) {
198            cleanedMessages.put(msg, new HashSet<ITerm>());
199          }
200          cleanedMessages.get(msg).addAll(partialRules.get(msg
             ));
201        }
202        broadcasts.addAll(eval.getBroadcastRules());
203      }
204      for (IRule msg : cleanedMessages.keySet()) {
205        for (ITerm recipient : cleanedMessages.get(msg)) {
206          sendMsg(new Text(recipient.getValue().toString()),
207              new Text(msg.toString()));
208        }
209      }
210      for (IRule broadcast : broadcasts) {
211        sendMsgToAllEdges(new Text(broadcast.toString()));
212      }
213      if (cleanedMessages.isEmpty() && broadcasts.isEmpty()) {
214        voteForNextRound = true; // no messages were sent
```

```java
215        }
216      return voteForNextRound ;
217    }
218
219    @Override
220    public int run ( String [] argArray ) throws Exception {
221      System . out . println ( Arrays . toString ( argArray ));
222
223      Preconditions . checkArgument ( argArray . length >= 3 ,
224          "run :␣Must␣have␣3␣arguments␣<input␣path >␣<output␣
              path >␣"
225            + "<#␣of␣workers >" );
226
227      GiraphJob job = new GiraphJob ( getConf () , getClass () .
             getName ());
228      job . setVertexClass ( getClass ());
229      job . setWorkerContextClass ( PDoPWorkerContext . class );
230      job . setVertexInputFormatClass ( PDoPInputFormat . class );
231      job . setVertexOutputFormatClass ( PDoPOutputFormat . class );
232      // hadoop paths !
233      FileInputFormat . addInputPath ( job , new Path ( argArray [0]))
             ;
234      FileOutputFormat . setOutputPath ( job , new Path ( argArray
             [1]));
235      job . setWorkerConfiguration ( Integer . parseInt ( argArray [2])
             ,
236          Integer . parseInt ( argArray [2]) , 100.0 f );
237      return job . run ( true ) ? 0 : -1;
238    }
239
240    /*
241     * Use from commandline with the following : > inputfolder
             outputfolder
242     * number_of_workers
243     */
244    public static void main ( String [] args ) throws Exception {
245      System . exit ( ToolRunner . run ( new PDoP () , args ));
246    }
247 }
```

## Appendix B:
pdop.giraph.graph.PDoPWorkerContext.java

```java
1 package  pdop . giraph . graph ;
2
3 import  org . apache . giraph . graph . WorkerContext ;
4 import  org . apache . hadoop . io . BooleanWritable ;
5 import  org . apache . hadoop . io . LongWritable ;
6
7 import  pdop . giraph . graph . aggregator . AndAggregator ;
```

```java
8  import pdop.giraph.graph.aggregator.RoundAggregator;
9  import pdop.giraph.graph.aggregator.OrAggregator;
10
11 /**
12  * Worker Context needed for Aggregators
13  */
14 public class PDoPWorkerContext extends WorkerContext {
15
16   @Override
17   public void preApplication() throws InstantiationException
        ,
18       IllegalAccessException {
19     registerAggregator("RoundAggregator", RoundAggregator.
          class);
20     registerAggregator("VoteAggregator", AndAggregator.class
          );
21     registerAggregator("LastTimeAggregator", OrAggregator.
          class);
22     registerAggregator("TerminateAggregator", OrAggregator.
          class);
23
24     RoundAggregator roundAggregator = (RoundAggregator)
          getAggregator("RoundAggregator");
25     AndAggregator voteAggregator = (AndAggregator)
          getAggregator("VoteAggregator");
26     OrAggregator lastTimeAggregator = (OrAggregator)
          getAggregator("LastTimeAggregator");
27     OrAggregator terminateAggregator = (OrAggregator)
          getAggregator("TerminateAggregator");
28
29     roundAggregator.setAggregatedValue(new LongWritable(0));
30     voteAggregator.setAggregatedValue(new BooleanWritable(
          true));
31     lastTimeAggregator.setAggregatedValue(new
          BooleanWritable(false));
32     terminateAggregator.setAggregatedValue(new
          BooleanWritable(false));
33   }
34
35   @Override
36   public void postApplication() {
37   }
38
39   /*
40    * Logic for rounds in preSuperstep, since values in
          aggregators are not
41    * synchronized if used in postSuperstep.
42    *
43    * Use postSuperstep for logging only!
44    */
45   @Override
46   public void preSuperstep() {
47
48     RoundAggregator roundAggregator = (RoundAggregator)
```

```
            getAggregator("RoundAggregator");
49      AndAggregator voteAggregator = (AndAggregator)
            getAggregator("VoteAggregator");
50      OrAggregator lastTimeAggregator = (OrAggregator)
            getAggregator("LastTimeAggregator");
51      OrAggregator terminateAggregator = (OrAggregator)
            getAggregator("TerminateAggregator");
52
53      if (voteAggregator.getAggregatedValue().get()) {
54        roundAggregator.setAggregatedValue(new LongWritable(
            roundAggregator
55          .getAggregatedValue().get() + 1));
56        if (lastTimeAggregator.getAggregatedValue().get()) {
57          terminateAggregator
58              .setAggregatedValue(new BooleanWritable(true));
59        } else {
60          lastTimeAggregator
61              .setAggregatedValue(new BooleanWritable(true));
62        }
63      } else {
64        lastTimeAggregator.setAggregatedValue(new
            BooleanWritable(false));
65        voteAggregator.setAggregatedValue(new BooleanWritable(
            true));
66      }
67
68      useAggregator("RoundAggregator");
69      useAggregator("VoteAggregator");
70      useAggregator("LastTimeAggregator");
71      useAggregator("TerminateAggregator");
72    }
73
74    @Override
75    public void postSuperstep() {
76    }
77 }
```

## Appendix C:

pdop.giraph.graph.aggregator.AndAggregator.java

```
1 package pdop.giraph.graph.aggregator;
2
3 import org.apache.giraph.graph.Aggregator;
4 import org.apache.hadoop.io.BooleanWritable;
5
6 public class AndAggregator implements Aggregator<
    BooleanWritable> {
7
8   boolean voteAggregator;
9
```

```
10    @Override
11    public void aggregate(BooleanWritable value) {
12      voteAggregator = voteAggregator && value.get();
13    }
14
15    @Override
16    public void setAggregatedValue(BooleanWritable value) {
17      voteAggregator = value.get();
18
19    }
20
21    @Override
22    public BooleanWritable getAggregatedValue() {
23      return new BooleanWritable(voteAggregator);
24    }
25
26    @Override
27    public BooleanWritable createAggregatedValue() {
28      return new BooleanWritable();
29    }
30 }
```

## Appendix D:
pdop.giraph.graph.aggregator.OrAggregator.java

```
1 package pdop.giraph.graph.aggregator;
2
3 import org.apache.giraph.graph.Aggregator;
4 import org.apache.hadoop.io.BooleanWritable;
5
6 public class OrAggregator implements Aggregator<
     BooleanWritable> {
7
8    boolean voteAggregator;
9
10    @Override
11    public void aggregate(BooleanWritable value) {
12      voteAggregator = voteAggregator || value.get();
13    }
14
15    @Override
16    public void setAggregatedValue(BooleanWritable value) {
17      voteAggregator = value.get();
18
19    }
20
21    @Override
22    public BooleanWritable getAggregatedValue() {
23      return new BooleanWritable(voteAggregator);
24    }
```

```
25
26    @Override
27    public BooleanWritable createAggregatedValue() {
28      return new BooleanWritable();
29    }
30  }
```

## Appendix E:
pdop.giraph.graph.aggregator.RoundAggregator.java

```
1  package pdop.giraph.graph.aggregator;
2
3  import org.apache.giraph.graph.Aggregator;
4  import org.apache.hadoop.io.LongWritable;
5
6  public class RoundAggregator implements Aggregator<
       LongWritable> {
7
8    long roundcounter;
9
10   /**
11    * used as aggregate(newvalue);
12    */
13   @Override
14   public void aggregate(LongWritable value) {
15     roundcounter = value.get();
16   }
17
18   @Override
19   public void setAggregatedValue(LongWritable value) {
20     roundcounter = value.get();
21
22   }
23
24   @Override
25   public LongWritable getAggregatedValue() {
26     return new LongWritable(roundcounter);
27   }
28
29   @Override
30   public LongWritable createAggregatedValue() {
31     return new LongWritable();
32   }
33 }
```

## Appendix F:
pdop.giraph.io.PDoPInputFormat.java

```java
package pdop.giraph.io;

import java.io.IOException;

import org.apache.giraph.graph.VertexReader;
import org.apache.giraph.lib.TextVertexInputFormat;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.InputSplit;
import org.apache.hadoop.mapreduce.TaskAttemptContext;

public class PDoPInputFormat extends
    TextVertexInputFormat<Text, Text, Text, Text> {

  @Override
  public VertexReader<Text, Text, Text, Text>
      createVertexReader(
      InputSplit split, TaskAttemptContext context) throws
          IOException {
    return new PDoPReader(
        textInputFormat.createRecordReader(split, context));
  }
}
```

## Appendix G:
pdop.giraph.io.PDoPOutputFormat.java

```java
package pdop.giraph.io;

import java.io.IOException;

import org.apache.giraph.graph.VertexWriter;
import org.apache.giraph.lib.TextVertexOutputFormat;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.RecordWriter;
import org.apache.hadoop.mapreduce.TaskAttemptContext;

public class PDoPOutputFormat extends TextVertexOutputFormat
    <Text, Text, Text> {

  @Override
  public VertexWriter<Text, Text, Text> createVertexWriter(
      TaskAttemptContext context) throws IOException,
      InterruptedException {
    RecordWriter<Text, Text> recordWriter = textOutputFormat
        .getRecordWriter(context);
    return new PDoPWriter(recordWriter);
  }
```

```
21 }
```

## Appendix H:
pdop.giraph.io.PDoPReader.java

```java
 1 package pdop.giraph.io;
 2
 3 import java.io.IOException;
 4 import java.util.HashMap;
 5 import java.util.LinkedList;
 6 import java.util.Map;
 7
 8 import org.apache.giraph.graph.BasicVertex;
 9 import org.apache.giraph.graph.BspUtils;
10 import org.apache.giraph.lib.TextVertexInputFormat.
      TextVertexReader;
11 import org.apache.hadoop.io.LongWritable;
12 import org.apache.hadoop.io.Text;
13 import org.apache.hadoop.mapreduce.RecordReader;
14
15 import pdop.PDoP;
16
17 /**
18  * for {@link PDoP} should support input files given in the
      following form:
19  * [VertexID \tab rules.facts.]
20  */
21 public class PDoPReader extends TextVertexReader<Text, Text,
      Text, Text> {
22
23   /**
24    * Constructor with the line record reader.
25    *
26    * @param lineRecordReader
27    *            Will read from this line.
28    */
29   public PDoPReader(RecordReader<LongWritable, Text>
        lineRecordReader) {
30     super(lineRecordReader);
31   }
32
33   @Override
34   public BasicVertex<Text, Text, Text, Text>
        getCurrentVertex()
35       throws IOException, InterruptedException {
36
37     BasicVertex<Text, Text, Text, Text> vertex = BspUtils
38         .<Text, Text, Text, Text> createVertex(getContext()
39             .getConfiguration());
40
```

```
41    String line = getRecordReader ().getCurrentValue ().
          toString ();
42    StringBuilder builder = new StringBuilder (line);
43    builder.deleteCharAt (0);
44    builder.deleteCharAt (builder.length () - 1);
45
46    String [] lineSegments = builder.toString ().split ("\t");
47    try {
48      Text vertexId = new Text (lineSegments [0]);
49      Text vertexVal = new Text ();
50      vertexVal.set (lineSegments [1]);
51
52      Map <Text , Text > vertexEdge = new HashMap <Text , Text >()
            ;
53      String [] edgeVertices = lineSegments [2].split (",");
54      for (String edgeVertex : edgeVertices) {
55        vertexEdge.put (new Text (edgeVertex), new Text ());
56      }
57      LinkedList <Text > vertexMsg = null;
58      vertex.initialize (vertexId , vertexVal , vertexEdge ,
            vertexMsg );
59
60    } catch (IllegalArgumentException e) {
61      throw new IllegalArgumentException (
62          "next:␣Couldn't␣get␣vertex␣from␣line" + line, e);
63    }
64    return vertex ;
65  }
66
67  @Override
68  public boolean nextVertex () throws IOException ,
        InterruptedException {
69    return getRecordReader ().nextKeyValue ();
70  }
71 }
```

## Appendix I:
pdop.giraph.io.PDoPWriter.java

```
1 package pdop.giraph.io;
2
3 import java.io.IOException ;
4
5 import org.apache.giraph.graph.BasicVertex ;
6 import org.apache.giraph.lib.TextVertexOutputFormat .
      TextVertexWriter ;
7 import org.apache.hadoop.io.Text ;
8 import org.apache.hadoop.mapreduce.RecordWriter ;
9
10 public class PDoPWriter extends TextVertexWriter <Text , Text ,
```

```
      Text > {
11   public PDoPWriter ( RecordWriter < Text , Text >
        lineRecordWriter ) {
12     super ( lineRecordWriter );
13   }
14
15   /**
16    * Output should be in the following form again ! [ VertexID
          \tab
17    * rules.facts .]
18    */
19   @Override
20   public void writeVertex ( BasicVertex < Text , Text , Text , ?>
        vertex )
21      throws IOException , InterruptedException {
22     StringBuilder builder = new StringBuilder ();
23     builder . append ("[");
24     builder . append ( vertex . getVertexId (). toString ());
25     builder . append ("\t");
26     builder . append ( vertex . getVertexValue (). toString ());
27     builder . append ("]");
28
29     getRecordWriter (). write ( new Text ( builder . toString ()) ,
          null );
30   }
31 }
```

## Appendix J:
pdop.iris.RuleFactQuerySplitter.java

```
1 package pdop . iris ;
2
3 import java . util . HashMap ;
4 import java . util . LinkedList ;
5 import java . util . List ;
6 import java . util . Map ;
7
8 import org . deri . iris . api . basics . IPredicate ;
9 import org . deri . iris . api . basics . IQuery ;
10 import org . deri . iris . api . basics . IRule ;
11 import org . deri . iris . compiler . Parser ;
12 import org . deri . iris . storage . IRelation ;
13
14 public class RuleFactQuerySplitter {
15
16   private Map < IPredicate , IRelation > facts = new HashMap <
        IPredicate , IRelation >();
17   private List < IRule > rules = new LinkedList < IRule >();
18   private List < IQuery > queries = new LinkedList < IQuery >();
19
```

```
20    /**
21     * The program is split into rules and facts and querries
           in the
22     * constructor. Use getters to obtain these afterwards.
23     */
24    public RuleFactQuerySplitter(String program) {
25      try {
26        Parser parser = new Parser();
27        parser.parse(program);
28        facts = parser.getFacts();
29        rules = parser.getRules();
30        queries = parser.getQueries();
31      } catch (Exception e) {
32        System.out.println("Parser␣Error␣in␣RuleFactSplitter:␣
              " + e);
33      }
34    }
35
36    public Map<IPredicate, IRelation> getFacts() {
37      return facts;
38    }
39
40    public List<IRule> getRules() {
41      return rules;
42    }
43
44    public List<IQuery> getQueries() {
45      return queries;
46    }
47 }
```

## Appendix K:
pdop.iris.filter.Filter.java

```
1 package pdop.iris.filter;
2
3 import java.util.HashMap;
4 import java.util.Map;
5
6 import org.deri.iris.api.basics.IPredicate;
7 import org.deri.iris.factory.Factory;
8 import org.deri.iris.storage.IRelation;
9
10 import pdop.iris.SimplePDoPRelation;
11 import pdop.iris.filter.conditions.Condition;
12
13 public class Filter {
14   /**
15    * Filters facts according to given roundcounter,
          timestamp of the facts and
```

```
16    * condition. Used to create DeltaFacts, "oldFacts" and
          trunkated Facts.
17    * (all of these are trunkated by default)
18    *
19    * @param roundcounter
20    * @param facts
21    * @param condition
22    * @return filtered facts according to condition (
          trunkated)
23    */
24
25   public static Map<IPredicate, IRelation> filter(
26       Map<IPredicate, IRelation> facts, Condition c) {
27
28     Map<IPredicate, IRelation> filteredFacts = new HashMap<
          IPredicate, IRelation>();
29
30     for (IPredicate factPred : facts.keySet()) {
31       String pred = factPred.getPredicateSymbol();
32       String[] predWOTime = pred.split("_");
33       filteredFacts.put(
34           Factory.BASIC.createPredicate(predWOTime[0],
35               factPred.getArity()), new SimplePDoPRelation()
                  );
36     }
37
38     for (IPredicate filteredPred : filteredFacts.keySet()) {
39       SimplePDoPRelation relAcc = new SimplePDoPRelation();
40       for (IPredicate factPred : facts.keySet()) {
41         String[] split = factPred.getPredicateSymbol().split
              ("_");
42         if (split[0].equals(filteredPred.getPredicateSymbol
              ())) {
43           if (c.isSatisfied(Long.valueOf(split[1]))) {
44             IRelation factRel = facts.get(factPred);
45             for (int i = 0; i < factRel.size(); i++) {
46               relAcc.add(factRel.get(i));
47             }
48           }
49         }
50       }
51       filteredFacts.put(filteredPred, relAcc);
52     }
53     return filteredFacts;
54   }
55 }
```

## Appendix L:
pdop.iris.filter.conditions.Condition.java

```java
package pdop.iris.filter.conditions;

public interface Condition {

  public boolean isSatisfied(long round);

}
```

## Appendix M:
pdop.iris.filter.conditions.EqualCondition.java

```java
package pdop.iris.filter.conditions;

public class EqualCondition implements Condition {

  long value;

  public EqualCondition(long v) {
    value = v;
  }

  public boolean isSatisfied(long round) {
    return round == value;
  }
}
```

## Appendix N:
pdop.iris.filter.conditions.LessEqualCondition.java

```java
package pdop.iris.filter.conditions;

public class LessEqualCondition implements Condition {

  long value;

  public LessEqualCondition(long v) {
    value = v;
  }

  public boolean isSatisfied(long round) {
    return round <= value;
  }
}
```

## Appendix O:
pdop.iris.filter.conditions.TrueCondition.java

```java
1 package pdop.iris.filter.conditions;
2
3 public class TrueCondition implements Condition {
4   public boolean isSatisfied(long round) {
5     return true;
6   }
7 }
```

## Appendix P:
pdop.iris.PartialEvaluation.java

```java
1 package pdop.iris;
2
3 import java.util.HashMap;
4 import java.util.HashSet;
5 import java.util.LinkedList;
6 import java.util.List;
7 import java.util.Map;
8 import java.util.Set;
9
10 import org.deri.iris.api.basics.IAtom;
11 import org.deri.iris.api.basics.ILiteral;
12 import org.deri.iris.api.basics.IPredicate;
13 import org.deri.iris.api.basics.IRule;
14 import org.deri.iris.api.basics.ITuple;
15 import org.deri.iris.api.terms.ITerm;
16 import org.deri.iris.api.terms.IVariable;
17 import org.deri.iris.factory.Factory;
18 import org.deri.iris.storage.IRelation;
19
20 public class PartialEvaluation {
21
22   private Map<IRule, Set<ITerm>> partialRules;
23
24   // partial Rules with rule and recipients
25   public Map<IRule, Set<ITerm>> getPartialRules() {
26     return partialRules;
27   }
28
29   private void addPartialRule(IRule rule, Set<ITerm>
        recipients) {
30     partialRules.put(rule, recipients);
31   }
32
33   // partial Rule for broadcasting
34   private Set<IRule> broadcastRules;
35
```

```
36    public Set<IRule> getBroadcastRules() {
37      return broadcastRules;
38    }
39
40    private void addBroadcastRule(IRule rule) {
41      broadcastRules.add(rule);
42    }
43
44    /**
45     * evalutate rule with given facts. (Partial evaluation)
46     *
47     * @param facts
48     * @param rule
49     *           obtain partial rules with getter!
50     */
51    public PartialEvaluation(Map<IPredicate, IRelation> facts,
          IRule rule) {
52      partialRules = new HashMap<IRule, Set<ITerm>>();
53      broadcastRules = new HashSet<IRule>();
54      for (ILiteral literal : rule.getBody()) { // for every
          Literal in the
55                                // body
56        IRelation values = facts.get(literal.getAtom().
            getPredicate()); // get
57                                              // facts
58                                              // for
59                                              // that
60                                              // exact
61                                              // literal
62        if (values != null) { // if empty stop
63          for (int i = 0; i < values.size(); i++) {
64            IRule interRule = changeRule(values.get(i),
                literal, rule);
65            if (interRule != null) {
66              interRule = checkRule(facts, interRule);
67              if (interRule != null) {
68                if (!rule.equals(interRule)) {
69                  Set<ITerm> recipients =
                      getRecipientsFromRule(interRule);
70                  if (recipients.isEmpty()) {
71                    addBroadcastRule(interRule);
72                  } else {
73                    addPartialRule(interRule, recipients);
74                  }
75                }
76              }
77            }
78          }
79        }
80      }
81    }
82
83    public static Set<ITerm> getRecipientsFromRule(IRule rule)
        {
```

```
84    Set < ITerm > recipients = new HashSet < ITerm >();
85    if ( rule.getBody().isEmpty()) {
86      recipients.addAll(checkSegmentForRecipients(rule.
          getHead()));
87    } else {
88      recipients.addAll(checkSegmentForRecipients(rule.
          getBody()));
89    }
90    return recipients;
91  }
92
93  public static Set < ITerm > checkSegmentForRecipients(List<
        ILiteral > segment) {
94    Set < ITerm > recipients = new HashSet < ITerm >();
95    for ( ILiteral literal : segment) {
96      for ( ITerm term : literal.getAtom().getTuple()) {
97        if (!(term instanceof IVariable)) {
98          recipients.add(term);
99        }
100       }
101     }
102   return recipients;
103 }
104
105 private static IRule checkRule ( Map < IPredicate , IRelation >
        facts , IRule rule) {
106   List < ILiteral > changedBody = new LinkedList < ILiteral >();
107   // check rules for facts from and remove parts that are
          true( aka
108   // represented in the facts)
109   for (int i = 0; i < rule.getBody().size(); i++) {
110     IAtom atomToCheck = rule.getBody().get(i).getAtom();
111     boolean builtIn = isBuiltIn(atomToCheck.getPredicate()
          );
112     if ( builtIn ) {
113       // check for Build in Predicates: so far implemented
              ::
114       // NOT_EQUAL
115       if (checkBuiltIn(changedBody , atomToCheck) == true)
            {
116         return null;
117       }
118     } else { // atomToCheck is no built - in predicate ,
            continue checking
119             // facts
120       IRelation edb = facts.get(atomToCheck.getPredicate()
            );
121       if (edb != null && edb.size() != 0) { // entries in
              edb , check!
122         boolean anyEDBTrue = false;
123         for (int j = 0; j < edb.size(); j++) {
124           boolean edbTrue = true;
125           for (int k = 0; k < atomToCheck.getTuple().size
                (); k++) {
```

```
126                    if ( atomToCheck
127                        . getTuple ()
128                        . get (k)
129                        . getValue ()
130                        . toString ()
131                        . equals ( edb . get (j) . get (k) . getValue ()
132                            . toString ())) {
133                      edbTrue &= true ;
134                    } else {
135                      edbTrue &= false ;
136                      break ;
137                    }
138                  }
139                  if ( edbTrue ) {
140                    anyEDBTrue = true ;// true => can be removed
141                    break ;
142                  }
143                }
144                if (! anyEDBTrue ) { // no EDB value verifies , add
                     for others
145                          // Vertices to Check.
146                  changedBody . add ( Factory . BASIC . createLiteral (true
                       ,
147                      atomToCheck ));
148                }
149
150            } else { // no entries !
151              if ( atomToCheck . getTuple () . getAllVariables () . size
                   () == 0) {
152                // p(d):-t(a). still possible ! -> check constant
                     =|= ID -> broadcast ( not yet implemented ),
                     else return null
153                return null ; // filled false.
154              } else { // add partial for others to check.
155                changedBody . add ( Factory . BASIC . createLiteral (true
                     ,
156                    atomToCheck ));
157              }
158            }
159          }
160        }
161    return Factory . BASIC . createRule ( rule . getHead () ,
           changedBody );
162  }
163
164  private static IRule changeRule ( ITuple fact , ILiteral
         literal , IRule rule ) {
165    Map < IVariable , ITerm > varMap = new HashMap < IVariable ,
         ITerm >(); // map
166                                        // for
167                                        // ( Variable , Constant )
                                              for
168                                        // replacement
169    // add every Variable to the Map and save the string
```

```
                  that matches given
170     // the facts.
171     for (int i = 0; i < literal.getAtom().getTuple().size();
            i++) { // for
172                                         // every
173                                         // Element
174                                         // in
175                                         // the
176                                         // Tuple
177       ITerm term = literal.getAtom().getTuple().get(i);
178       if (term instanceof IVariable) { // check if it is a
            IVariable
179         varMap.put((IVariable) term, fact.get(i)); // and
              add it to the
180                                 // map
181       } else { // term is a constant
182         if (!term.equals(fact.get(i))) {
183           return null; // the constant in the rule is not
                equal to the
184                   // constant in fact, abort!
185         }
186       }
187     }// Map finished, change rule now!
188     List<ILiteral> changedHead = changeSegment(varMap, rule.
            getHead());
189     List<ILiteral> changedBody = changeSegment(varMap, rule.
            getBody());
190     return Factory.BASIC.createRule(changedHead, changedBody
            );
191   }
192
193   /**
194    * Replaces the variables in literals with the appropriate
           constants from
195    * the varMap
196    *
197    * @param varMap
198    * @param literals
199    * @return new literals with substitutions
200    */
201   private static List<ILiteral> changeSegment(Map<IVariable,
          ITerm> varMap,
202       List<ILiteral> literals) {
203     List<ILiteral> newLiterals = new LinkedList<ILiteral>();
204
205     for (int i = 0; i < literals.size(); i++) {
206       ILiteral literal = literals.get(i);
207       ITuple tuple = literal.getAtom().getTuple();
208       List<ITerm> newTerms = new LinkedList<ITerm>();
209       for (int j = 0; j < tuple.size(); j++) {
210         ITerm newTerm = varMap.get(tuple.get(j));
211         if (newTerm == null) {
212           newTerms.add(tuple.get(j));
213         } else {
```

```
214          newTerms.add(newTerm);
215        }
216      }
217      ILiteral newLiteral = Factory.BASIC.createLiteral(true
             , literal
218          .getAtom().getPredicate(), Factory.BASIC
219          .createTuple(newTerms));
220      newLiterals.add(newLiteral);
221    }
222    return newLiterals;
223  }
224
225  /**
226   * checks if atomToCheck is a verifiable with a built-in
         function
227   *
228   * @param changedBody
229   * @param atomToCheck
230   * @return true, if verified; false, if others need to
         verify it.
231   */
232  private static boolean checkBuiltIn(List<ILiteral>
        changedBody,
233      IAtom atomToCheck) {
234
235    if (atomToCheck.getPredicate().getPredicateSymbol().
         equals("NOT_EQUAL")) {
236      if (atomToCheck
237          .getTuple()
238          .get(0)
239          .getValue()
240          .toString()
241          .equals(atomToCheck.getTuple().get(1).getValue().
              toString())) {
242        return true; // both equal, drop.
243      } else {
244        if (atomToCheck.getTuple().getAllVariables().size()
             != 0) {
245        // still variables to be filled => add
246        changedBody.add(Factory.BASIC.createLiteral(true,
247            atomToCheck));
248        }
249      }
250    }
251    return false;
252  }
253
254  private static boolean isBuiltIn(IPredicate pred) {
255    return pred.getPredicateSymbol().equals("NOT_EQUAL");
256  }
257 }
```

## Appendix Q:
pdop.iris.SimplePDoPRelation.java

```java
1  package pdop.iris;
2
3  import java.util.List;
4
5  import org.deri.iris.api.basics.ITuple;
6  import org.deri.iris.storage.IRelation;
7  import org.deri.iris.utils.UniqueList;
8
9  public class SimplePDoPRelation implements IRelation {
10
11    /*
12     * SimpleRelation from Iris with explicit public
           constructor
13     */
14
15    /**
16     * Constructor. For performance reasons where the user of
           the class can
17     * enforce uniqueness (or does not require it), uniqueness
           enforcement can
18     * be turned off.
19     *
20     * @param forceUniqueness
21     *             true, if this object should enforce
           uniqueness.
22     */
23    public SimplePDoPRelation() {
24      mTuples = new UniqueList<ITuple>();
25    }
26
27    public boolean add(ITuple tuple) {
28      assert mTuples.isEmpty() || (mTuples.get(0).size() ==
           tuple.size());
29
30      return mTuples.add(tuple);
31    }
32
33    public boolean addAll(IRelation relation) {
34      boolean added = false;
35
36      for (int i = 0; i < relation.size(); ++i)
37        if (add(relation.get(i)))
38          added = true;
39
40      return added;
41    }
42
43    public ITuple get(int index) {
44      return mTuples.get(index);
45    }
```

```java
46
47     public int size() {
48       return mTuples.size();
49     }
50
51     public boolean contains(ITuple tuple) {
52       return mTuples.contains(tuple);
53     }
54
55     @Override
56     public String toString() {
57       return mTuples.toString();
58     }
59
60     /** The array list (or unique list) of tuples. */
61     private final List<ITuple> mTuples;
62
63 }
```

# Bibliography

[1] Jeffrey D. Ullman. *Principles of Database and Knowledge-Base Systems, Volume I.* Computer Science Press, 1988.

[2] Inc LogicBlox. Datalog for Enterprise Applications: from Industrial Applications to Research. `http://www.logicblox.com/presentations/arefdatalog20.pdf`, March 2010. accessed on: September 20, 2012.

[3] Daniel Ritter and Till Westmann. Business Network Reconstruction Using Datalog. In Barceló and Pichler [22], pages 148–152.

[4] Marianne Shaw, Paraschos Koutris, Bill Howe, and Dan Suciu. Optimizing Large-Scale Semi-Naïve Datalog Evaluation in Hadoop. In Barceló and Pichler [22], pages 165–176.

[5] Jeffrey Dean and Sanjay Ghemawat. MapReduce: simplified data processing on large clusters. In *Proceedings of the 6th conference on Symposium on Operating Systems Design & Implementation - Volume 6*, OSDI'04, pages 10–10, Berkeley, CA, USA, 2004. USENIX Association.

[6] Yingyi Bu and Bill Howe. HaLoop. `http://code.google.com/p/haloop/`, 2012. accessed on: September 24, 2012.

[7] Grzegorz Malewicz, Matthew H. Austern, Aart J.C Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, SIGMOD '10, pages 135–146, New York, NY, USA, 2010. ACM.

[8] Leslie G. Valiant. A bridging model for parallel computation. *Commun. ACM*, 33(8):103–111, August 1990.

[9] Apache Software Foundation. Apache Incubator Giraph. `http://incubator.apache.org/giraph/`, February 2012. accessed on: September 24, 2012.

[10] Peter Alvaro, William Marczak, Neil Conway, Joseph M. Hellerstein, David Maier, and Russell C Sears. Dedalus: Datalog in Time and Space. Technical Report UCB/EECS-2009-173, EECS Department, University of California, Berkeley, Dec 2009.

[11] François Bry, Norbert Eisinger, Thomas Eiter, Tim Furche, Georg Gottlob, Clemens Ley, Benedikt Linse, Reinhard Pichler, and Fang Wei. Foundations of Rule-Based Query Answering. In *Reasoning Web, Third International Summer School 2007*, volume 4636 of *LNCS*. Springer-Verlag, 2007.

[12] Ravel. GoldenOrb. `http://goldenorbos.org/`, 2012. accessed on: October 15, 2012.

[13] Jimmy Lin and Michael Schatz. Design patterns for efficient graph algorithms in MapReduce. In *Proceedings of the Eighth Workshop on Mining and Learning with Graphs*, MLG '10, pages 78–85, New York, NY, USA, 2010. ACM.

[14] Apache Software Foundation. Apache Hadoop. `http://hadoop.apache.org/`, September 2012. accessed on: October 15, 2012.

[15] Apache Software Foundation. Apache ZooKeeper. `http://zookeeper.apache.org/`, September 2012. accessed on: October 15, 2012.

[16] Apache Software Foundation. Apache Incubator Giraph Wiki. `https://cwiki.apache.org/confluence/display/GIRAPH/Index`, February 2012. accessed on: September 24, 2012.

[17] Apache Software Foundation. Apache Incubator Giraph API. `http://incubator.apache.org/giraph/apidocs/index.html`, February 2012. accessed on: September 24, 2012.

[18] David E. Culler, Anoop Gupta, and Jaswinder Pal Singh. *Parallel Computer Architecture: A Hardware/Software Approach*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1st edition, 1997.

[19] Barry Bishop, Florian Fischer, Uwe Keller, Nathalie Steinmetz, Christoph 'Gigi' Fuchs, Matthias Pressnig, Darko Anicic, Cristina Feier, Holger Lausen, and Richard Pöttler. IRIS Reasoner. `http://iris-reasoner.org/`, February 2010. accessed on: October 15, 2012.

[20] Avanade® Research Insights. Global Survey: Is Big Data Producing Big Returns? `http://www.avanade.com/Documents/Research%20and%20Insights/avanade-big-data-executive-summary-2012.pdf`, June 2012. accessed on: October 11, 2012.

[21] Yingyi Bu, Vinayak R. Borkar, Michael J. Carey, Joshua Rosen, Neoklis Polyzotis, Tyson Condie, Markus Weimer, and Raghu Ramakrishnan. Scaling Datalog for Machine Learning on Big Data. *CoRR*, abs/1203.0160, 2012.

[22] Pablo Barceló and Reinhard Pichler, editors. *Datalog in Academia and Industry - Second International Workshop, Datalog 2.0, Vienna, Austria, September 11-13, 2012. Proceedings*, volume 7494 of *Lecture Notes in Computer Science*. Springer, 2012.