

# Fast, scalable WCOJ graph-pattern matching on in-memory graphs in Spark

Master thesis draft iteration 1.2

Per Fuchs

April 2019

## Abstract

Graph pattern matching with their vast number of cyclic foreign-key joins is a new challenge for data processing systems, like Spark, because their intermediary results grow over linear with regards to the inputs and are materialized by the traditionally used binary joins. Worst-case optimal join algorithms, WCOJ's, are a natural match to tackle this challenge because they do not materialize the aforementioned large intermediary results. We investigate two major open questions regarding WCOJ's. First, we develop a WCOJ specialized to graph-pattern matching, namely self-joins on a relationship with two attributes, and compare its performance with a general WCOJ. Second, we propose a novel method to distribute WCOJ's. We show in our proposal that current methods to distribute WCOJ's, suitable for Spark, do not scale well to bigger graph pattern (five vertices and more). Based on this result we propose to keep the edge relationship cached on all workers but distribute the computation using a logical partitioning. Along the line of argumentation in COST [33], we aim to provide a fast WCOJ implementation in Spark which scales well in use of computational resources by trading off scalability in memory usage. This is a reasonable trade-off as most graphs today fit in main memory [31]. Our thesis provides the first distributed, open-source implementation of a WCOJ in a data processing system widely used in industry, namely Spark.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Spark . . . . .	5
2.1.1	Resilient distributed datasets . . . . .	5
2.1.2	Spark architecture . . . . .	6
2.1.3	Catalyst . . . . .	7
2.1.4	Broadcast variables . . . . .	9
2.2	Graph pattern matching . . . . .	9
2.3	Worst-case optimal join algorithm . . . . .	9
2.3.1	Leapfrog Triejoin . . . . .	10
2.4	Distributed worst-case optimal join in Myria . . . . .	17
2.4.1	Shares . . . . .	17
2.5	Compressed sparse row representation . . . . .	21
2.6	Analysis of public real-world graph datasets . . . . .	22
<b>3</b>	<b>Worst-case optimal join parallelization</b>	<b>23</b>
3.1	Logical Shares . . . . .	23
3.1.1	RangeShares . . . . .	24
3.2	Work-stealing . . . . .	25
<b>4</b>	<b>GraphWCOJ</b>	<b>27</b>
4.0.1	Combining LFTJ with CSR . . . . .	27
4.0.2	Exploiting low average outdegrees . . . . .	29
<b>5</b>	<b>Implementation</b>	<b>31</b>
5.1	General sequential version ( <i>seq</i> ) . . . . .	31
5.1.1	Optimizations . . . . .	31
<b>6</b>	<b>Experiments</b>	<b>34</b>
6.1	Setup . . . . .	34
6.1.1	Hardware and Software . . . . .	34
6.1.2	Algorithms . . . . .	34
6.1.3	Datasets . . . . .	34
6.1.4	Graph patterns . . . . .	35

6.2	Baseline: <code>BroadcastHashJoin</code> vs <code>seq</code> . . . . .	36
6.2.1	Experiment Rationale . . . . .	36
6.2.2	Analysis . . . . .	37
6.3	Scaling of <i>GraphWCOJ</i> . . . . .	38
6.3.1	Results . . . . .	38
6.3.2	Analysis . . . . .	44
<b>7</b>	<b>Related Work</b>	<b>44</b>
7.1	Graphs on Spark . . . . .	44
7.1.1	Fractal a graph pattern mining system on Spark . . . . .	44
7.2	WCOJ on Timely Data Flow . . . . .	44
7.3	Semih’s work on worst-case optimal join for different queries . . . . .	45
7.4	Adaptive Query Exectution . . . . .	45
<b>8</b>	<b>Conclusions</b>	<b>45</b>
8.1	Future work . . . . .	45
8.1.1	Cluster mode . . . . .	45
8.1.2	Deeper integration of Workstealing . . . . .	45
	<b>References</b>	<b>46</b>
<b>A</b>	<b>Experimental Results</b>	<b>49</b>
A.1	<i>GraphWCOJ</i> scaling . . . . .	49

# 1 Introduction

Newly developed worst-case optimal join (WCOJ) algorithms, e.g. Leapfrog Triejoin, turned conventional thinking about join processing on its head because these multi-join algorithms have provably lower complexity than classical binary joins, i.e. join algorithms that join just two tables at-a-time. In the areas of data warehousing and OLAP, this finding does not have much impact, though, since the join patterns most commonly encountered are primary-foreign-key joins, which normally take the form of a tree or snowflake and contain no cycles. The computational complexity of FK-PK joins is by definition linear in size of the inputs. In these "conventional" cases, binary joins, e.g. hash joins, work fine. However, analytical graph queries often use foreign-foreign-key joins, which can grow over linearly in the size of their inputs, and often contain cycles. For these use-cases, worst-case optimal join algorithms excel because matching a pattern consisting of multiple joins causes binary joins to generate a rapidly increasing set of intermediate results, e.g. navigating a social graph with an out-degree in the hundreds, of which many matches are useless and get eliminated by later joins, e.g. the join closing the cycle. These kinds of join patterns are frequently found during graph analysis, e.g. for graph clustering on social network graphs for customer relationship management or recommendation systems and fraud detection in the financial sector [8, 22]. Worst-case optimal join algorithms avoid large result materialization and hence promise to be orders of magnitude faster than binary joins. Therefore, we believe that worst-case optimal join algorithms could be a useful addition to (analytical) graph database systems.

We continue with a short example for a cyclic query and compare how this query is evaluated traditionally and with the new WCOJ's in place. The simplest example of a cyclical join query enumerates all triangles in a graph. This can be formulated as the following datalog query

$$Q(a, b, c) \leftarrow R(a, b), S(b, c), T(c, a) \quad (1)$$

where  $R = S = T$  are aliases for the edge relationship. Traditionally, this would be processed by using multiple binary joins:

$$R \bowtie S \bowtie T \bowtie R \quad (2)$$

Independent of the chosen order, it can be proven that there exists cases where the intermediary result size is in  $\mathcal{O}(n^3)$  with  $n = |R| = |S| = |T|$ . However, it is provable that maximal output of this query is in  $\mathcal{O}(n^{3/2})$  [12, 36]. Hence, binary joins materialize huge intermediary results after processing parts of the query, which are much bigger than the final result after applying all joins. The described problem has been shown to be a fundamental issue with traditional join-at-a-time approaches [12, 36]. Fortunately, worst-case optimal join algorithms can materialize cyclic joins with memory usage linear to their output size by avoiding to produce large intermediary results [leapfrog, 37]. In practice, these algorithms have been shown to be highly beneficial for cyclic queries in analytical graph workloads in an optimized, single machine system [leapfrog, 38] and later in distributed shared-nothing settings [16, 7] - we describe these systems in more detail in section 7.

We identify two challenging, novel directions for our research. First, although, all of the systems cited above focus on queries widely used in graph pattern matching, e.g. clique finding or path queries, they use WCOJ's which are developed for general multi-way joins, however, graph pattern matching uses only self-joins on a relationship with two attributes, namely the edge relationship of the graph. This raises the question if WCOJ's can be optimized by specializing them for graph pattern matching - which is so far the only use-case that has been shown to benefit from WCOJ's in the literature. Second, while the communication costs for worst-case optimal joins in map-reduce like systems (an excellent definition of the term is given in [hypercube]) is well-understood [hypercube, 16], their scalability has not been studied in depth. Given the high complexity of worst-case optimal joins used and the fact that their only integration in a map-reduce like system [16] exhibits a speedup of 8 on 64 nodes (an efficiency of 0.125), leads us to the conclusion that designing a scalable, distributed WCOJ for a map-reduce like system

is an unsolved challenge. We believe it is time to investigate how these algorithms scale in the probably most widely used, general-purpose big data processing engine: Spark. To the best of our knowledge, this would be also the first time a WCOJ is integrated with an industrial-strength cluster computing model. We detail our research questions and goals in ?? and explain how to address challenges in ??.

## 2 Background

TODO introduction

### 2.1 Spark

Spark is the probably most widely used and industry accepted cluster computing model. It improves over former computing models, e.g. MapReduce [20], Hadoop [9] or Haloop [14], by allowing to cache results in memory between multiple queries, using so-called resilient distributed datasets [46]; often abbreviated to RDD.

This section introduces Spark and is organized in four subsections. Section 2.1.1 describes the core data structure of Spark: the RDD's. In section 2.1.2, we explain the different components and processes in a Spark cluster. The query optimizer of Spark, Catalyst, is explained in section 2.1.3. It is the component we integrate our WCOJ with; therefore, it is the module of Spark that is most relevant to this thesis. Finally, in section 2.1.4 we highlight important details about *Broadcast variables* which are used to implement our parallel worst-case optimal join.

#### 2.1.1 Resilient distributed datasets

RDD's form the core of Spark. However, for this thesis, it is not necessary to understand them in great detail. In the next paragraph, we give a short introduction to the relevant aspects of RDD's. For the interested reader, a more in-depth description is given in the original paper [46].

Resilient distributed datasets describe a distributed collection of data items of a single type. In contrast, to other distributed share memory solutions, RDD's do not use fine-grained operations to manipulate single data items but coarse-grained operations which are applied to all data items, e.g. *map* to apply a function to each data item. These operations are called *transformations*. An RDD is built starting from a persistent data source and multiple transformations to apply to this data source. One can store the transformations applied to the input data source as a directed acyclic graph, the so-called *lineage graph*. This graph fully describes the dataset without materializing it because the transformations are deterministic. Hence, the dataset can be computed and recomputed on demand, e.g. when the user asks for the count of all items in the set. Operations which require that the data in the RDD is computed are called *actions*.

RDD's are distributed by organizing their data items into partitions. The partitioning can be chosen by the user or the Spark query optimizer such that it allows to run transformations on all partitions in parallel. For example, one might choose a round-robin partitioning to generate splits of equal size when reading data items from disk or one groups items by hashing a specific key to support parallelizable aggregation on that key per partition. The process of repartitioning an RDD is called a *shuffle*. It is an expensive operation because it involves writing and reading the whole RDD to disk.

Describing datasets as RDD's comes with two main benefits. First, it is resilient because if the dataset or some partitions of it get lost, it is possible to recompute them from persistent storage using lineage graph information. Second, it allows Spark to compute RDD's in parallel.

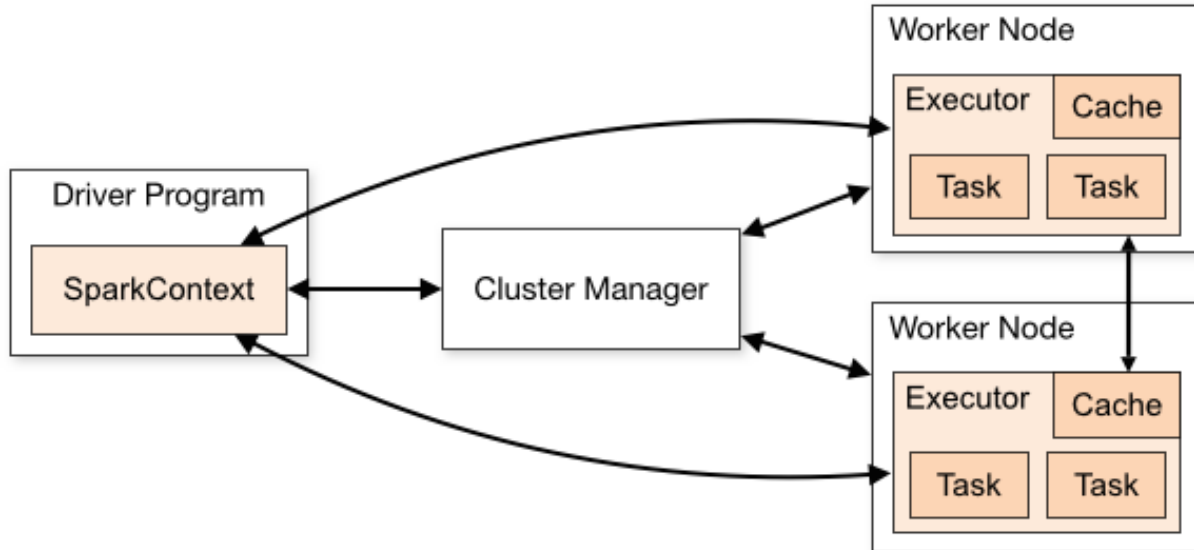


Figure 1: Schematics of a Spark cluster with two workers, each of them with one executor and two threads per executor. Source: Apache Spark Documentation, <https://spark.apache.org/docs/latest/cluster-overview.html>

Spark can parallelize the computation of RDD in two ways. First, by data-parallelism, since different partitions of an RDD can be computed independently from each other. Second, by task parallelism, because some parts of the DAG can be computed without dependence of the others. Indeed, it is possible to compute all parts of an RDD in parallel which are not related in a topological sort of the graph.

### 2.1.2 Spark architecture

Spark allows the user to run his program on a single machine or hundreds of machines organized in a cluster. In this section, we explain the architecture that allows this flexibility. Figure 1 shows the schematics of a Spark cluster setup.

In Spark, each physical machine is called a *worker*. On each worker, Spark starts one or multiple Spark processes in their own JVM instance; each of them is called *executor*. Nowadays, many Spark deployments use a single executor per worker<sup>1</sup>. Each executor runs multiple threads (often one per core on its worker) to execute multiple tasks in parallel. In total, a Spark cluster can run  $\# \text{ workers} \times \# \text{ executors per worker} \times \# \text{ threads per executor}$  tasks in parallel.

Spark uses two kinds of processes to execute an application: a *driver program* and multiple *executors*. When started, the driver program acquires resources from the *cluster manager* for its executor processes. These executors stay alive during the whole Spark application. Then, the driver program continues executing the Spark application. When it encounters parallelizable tasks, it schedules them on the available executors.

All tasks scheduled on the same executor share a cache for in-memory data structures like *Broadcast variables* or persisted RDD partitions. This is important in the context of this thesis because it means that we cache the input graph once per executor; which in many Spark

<sup>1</sup>This is the setup Databricks uses; Databricks is the leading maintainer of the Spark platform and offers professional deployment to many customers.

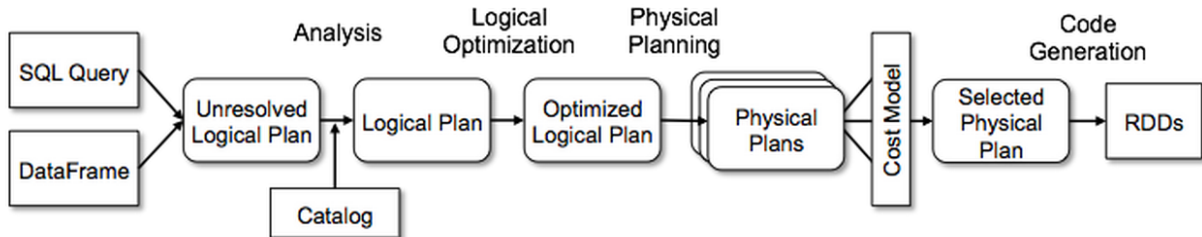


Figure 2: Input and stages of the Catalyst optimizer. Source: Databricks Blog, <https://databricks.com/blog/2015/04/13/deep-dive-into-spark-sqls-catalyst-optimizer.html>

deployments is once per worker or physical machine. This would not be possible if different tasks in the same JVM would not share the same cache.

Spark allows the user to choose a cluster manager to manage resources in the cluster. It comes with good integration for Hadoop YARN [43], Apache Mesos [23] and Kubernetes [27], as well as, a standalone mode where Spark provides its own cluster manager functionality. Finally, one can run Spark on a single machine in *local mode*. In local mode, the driver program and a single executor share a single JVM. The executor uses the cores assigned to Spark to run multiple worker threads. For our experiments, we run Spark purely in local mode.

### 2.1.3 Catalyst

Catalyst [11] is Spark’s query optimizer. It can process queries given as a SQL string or described using the DataFrame API. From a given query it constructs an executable *physical plan*. The query compilation process is organized in multiple stages. Its inputs and stages are shown in fig. 2. Below we explain these in order. We use the triangle given by the datalog rule  $COUNT(triangle(A, B, C)) \leftarrow R(A, B), S(B, C), T(A, C), A < B < C$  as a running example.

The input of Catalyst is a query in the form of a DataFrame or SQL string. From this the optimizer builds a *unresolved logical plan*. This plan can include unresolved attributes, e.g. attribute names which are not matched to a specific data source yet or which have no known type. To resolve this attributes Catalyst uses a *Catalog* of possible bindings which describe the available data sources. This phase is referred to as *Analysis* and results in a *logical plan*. The logical plan represents *what* should be done for the query but not exactly *how*, e.g. it might contain a Join operator but not a Sort-merge join.

We show the logical plan for the triangle query in fig. 3a. As we see, the query is represented as tree where the vertices are operators and the edge indicate dataflow from one operator to another. The leaves of the tree are three aliases of the edge relationship. Two of these source relationships are the input the join between  $R$  and  $S$  via  $B$ . The result of this join and the leaf relationship  $T$  are input to the second join. The tuples produced by this join are filtered to fulfil  $A < B < C$ . Finally, at the root of the tree, there is an aggregation to count all results and report the sum.

The *logical optimization phase* applies batches of rewriting rules until a fixpoint is reached. A simple example of a logical optimization would be rewriting  $2 + 2$  into  $4$ . In the running example of the triangle query, this phase pushes the filters into the two joins. This optimization is called Filter Pushdown. It is efficient because it applies filters earlier within the pipeline reducing the number tuples to process by later operators.

From the *optimized logical plan* the optimizer generates one or multiple *physical plans* by applying so called *Strategies*. They translate a logical operator in one or multiple *physical operators*.

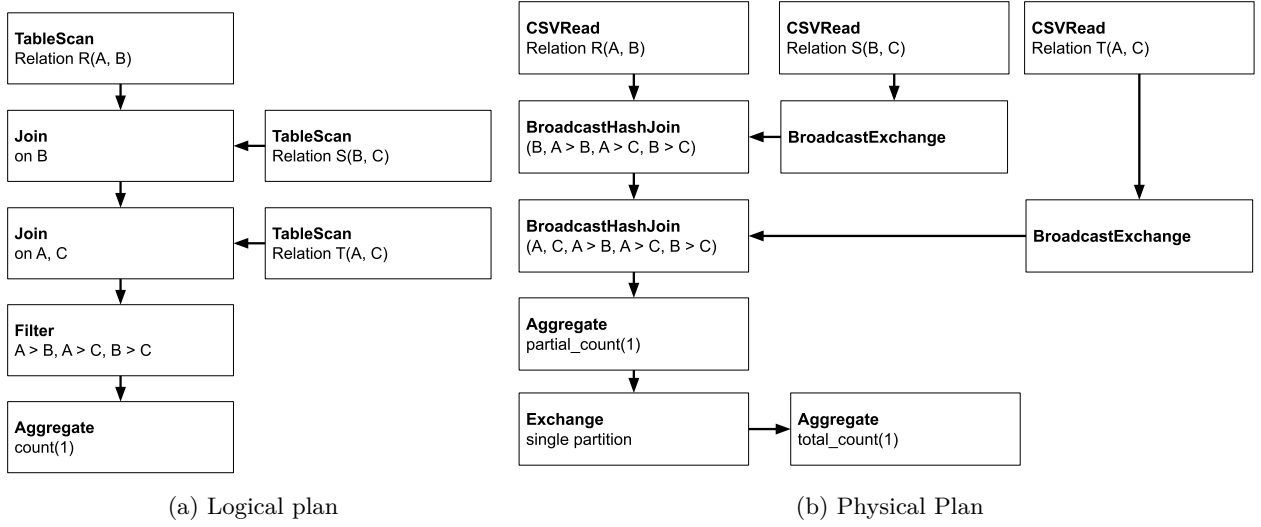


Figure 3: Logical and physical plan for the triangle count query as generated by Catalyst.

*Strategies* are also allowed to return multiple physical plans for a single *logical plan*. In this case, the optimizer selects the best one according to a *cost model*.

The physical plan for the triangle query is shown in fig. 3b. We see multiple examples of translation of a logical operator, which describes what to do, to its physical pendant that also describes how to do it: the *TableScan* becomes a *CSVRead* and the *Joins* are implemented as *BroadcastHashJoins*.

Furthermore, we see the introduction of exchanges. *BroadcastExchanges* precede the *BroadcastHashJoins*. They build a hashtable from their input operators and make them available as a broadcast variable to all executors of the cluster; we explain broadcast variables in depth in section 2.1.4. When an executor is tasked to execute the hash join operator, it acquires the broadcasted hashtable and executes a local hash join of its assigned partitions.

Another exchange operator is introduced for the aggregation. It is broken up into a partial aggregation directly after the last join, an exchange reorganizing all partial counts into a single partition and a second aggregation over that partition to calculate the total count. The last is a good example of Catalyst introducing a shuffle.

To conclude, the translation to a physical plan translates logical operators into concrete implementations of these and adds exchanges to organize the data such that it can be processed independently in partitions.

After generating and choosing a physical plan, Catalyst enters the *code generation* phase in which it compiles Java byte code for some of the physical operators. This code executes often magnitudes faster than interpreted versions of the same operator [11] because it can be specialized towards this particular query, e.g. if a join operates only on integers, code generation can prune all code paths dealing with strings. Indeed, the code generation phase is part of another Spark project called *Tungsten* [45, 5]. In this thesis, we do not build any code generated physical operators. Hence, we do not treat this topic in depth. It is enough to know that all freshly generated Java code is wrapped into a single physical operator. Therefore, it integrates seamlessly with interpreted operators.

Finally, Catalyst arrives at an optimized physical plan which implements the query. The execution of this plan is called *structured query execution* [28]. It translates the plan into RDD operations



implemented by Spark’s core. Hence, the result of Catalysts query compilation is an RDD representing the query. One should note that structured query execution does not materialize the query: the result is an RDD which is a none materialized representation of the operations necessary to generate the result. In this thesis, we are not concerned with the internals of RDD’s. We do not need to introduce any new RDD operations or even touch Spark’s core functionality. Thanks to the extensibility of Catalyst, we can integrate worst-case optimal joins by adding one logical operator, multiple physical operators and a Strategy to translate between them.

#### 2.1.4 Broadcast variables

*Broadcast variables* readonly variables which are accessible by all tasks. They are initialized once by the driver program and should not be changed after initialization. The process of broadcasting them is handled by Spark. It is guaranteed that each broadcast variable is sent only once to each executor and allows it to be spilled to disk if it is not possible to keep the whole value in memory. Furthermore, ‘Spark attempts to distribute broadcast variables using efficient broadcast algorithms to reduce communication costs’ [18]; currently Spark uses a BitTorrent-like communication protocol<sup>2</sup>. Once sent, they are cached once per executor (see also section 2.1.2) and shared by all tasks on this executor. They are cached in deserialized form in memory but can be spilled to disk if they are too big. In this thesis, we use broadcast variables to cache the edge relationship of the graph on all workers.

## 2.2 Graph pattern matching

## 2.3 Worst-case optimal join algorithm

The development of worst-case optimal joins started in 2008 with the discovery that the output size of a relational query is bound by the fractional edge number of its underlying hypergraph [12]. In short, this bound proves that traditional, binary join plans perform asymptotically worse than theoretical possible for the worst-case database instances, e.g. heavily skewed instances. For example, the worst-case runtime of binary joins on the triangle query is in  $\mathcal{O}(N^2)$ , while the AGM bound shows the possibility to solve it in  $\mathcal{O}(N^{3/2})$ . The AGM bound has been treated widely in literature [36, 6, 12]. A particular good explanation is given by Hung Ngo et al in [36]. We refer the reader to these papers for further information. In the next paragraph, we discuss different algorithms matching the AGM bound which are called worst-case optimal joins.

In 2012, Ngo, Porat, Re and Rudra published the first join algorithm matching the AGM bound, called *NPRR* join [37]. In the same year, Veldhuizen proved that the algorithm Leapfrog Triejoin used in LogicBlox, a database system developed by his company, is also worst-case optimal with regards to the fractional edge number bound. We often abbreviate Leapfrog Triejoin to LFJT. Both algorithms have been shown to be instances of a single algorithm, the *Generic Join*, in 2013 by Ngo et al. [36].

Three worst-case optimal join algorithms are known in literature. We choose Leapfrog Triejoin as the basis for our work. The argumentation for this decision is given below. First, we identify the main criteria for this choice. Then, we use them to compare the different algorithms.

The most important argument for our decision is the degree to which the algorithm has been shown to be of practical use. In particular, the number of systems it is used in and openly available data on its performance. If an algorithm is used in academia as well as in industry, we deem this as a advantage. This criteria carries a lot of weight because the first literature on worst-case optimal joins has been rather theoretical but in our work we take a more praxis and

---

<sup>2</sup>See Spark sources: `org.apache.spark.broadcast.TorrentBroadcast`

system oriented perspective.

The practical character of our work also motivates the second dimension which we compare the algorithms in, namely ease of implementation. If two of the three algorithms both have well proven performance, we would like to choose the algorithm that takes less time to implement and is easier to adapt and experiment with. That is, to be able to spend more time on evaluation and optimizations for the graph use-case, instead of, time spent on replicating existing work.

The Leapfrog Triejoin is used in two commercial database solutions: LogicBlox [10] and RelationalAI<sup>3</sup>. Its performance has been reported on in two publications [16, 38]. In particular, it beats various general and graph specific databases for graph pattern matching, i.e. PostgreSQL, MonetDB, NEO4J, graphLab and Virtuoso [38]. The broadest study of its performance uses 15 different datasets and 7 queries [38]. We conclude that the performance of LFTJ is well established by peer reviewed publications as well as industrial usage.

The *NPRR* algorithm has been well analyzed from the theoretical point of view. However, we are not able to find any openly available sources with performance measurements. This disqualifies NPRR as basis for our thesis.

The *Generic Join* is used in at least three academic graph processing engines, namely GraphFlow [24], EmptyHeaded [1] and a unnamed implementation in Timely Dataflow [7]. All three show good performance. However, we are not aware of any commercial systems using GJ.

The comparison of Leapfrog Triejoin, NPRR and *Generic Join* by proven performance rules out NPRR and puts LFTJ and GJ on a similar level. Next, we compare these two algorithm in ease of implementation.

The description of the Leapfrog Triejoin implementation in its original paper [44] is excellent. Furthermore, multiple open source implementation exists [41, 16]. In particular, the implementation of Christian Schroeder for a course at Oxford is helpful because it is standalone and does not require us to understand a whole system<sup>4</sup>.

*Generic Join* is described as a generalization of NPRR and Leapfrog Triejoin in its original paper [36]. Although, well written and algorithmically clear, this explanation is much less practical than the one given for LFTJ which is backed by an executable implementation.

To conclude, we choose Leapfrog Triejoin as basis for our work based on its openly available records of performance, use in academia as well as industrial systems and good description for direct implementation. Furthermore, Peter Boncz (supervisor of this thesis) has direct contact to the inventors of LFTJ giving us access to valuable expertise if necessary.

### 2.3.1 Leapfrog Triejoin

In this section, we described the Leapfrog Triejoin algorithm. In the first paragraph, we give the high-level idea behind the algorithm and some of its requirements. Then we discuss the kind of queries that can be answered with it. The main part of the section, discusses the conceptual algorithm itself. We finish with a short discussion of two implementation problems, namely the data structure to represent the input relationships and the problem of choosing a good variable ordering.

The Leapfrog Triejoin is a variable-oriented join. Given an input queries, it requires a variable ordering. For example in the triangle query  $triangles(a, b, c) \leftarrow R(a, b), S(b, c), T(a, c)$ , the variable ordering could be  $a, b, c$ . Furthermore, the Leapfrog Triejoin requires its input relationships to be sorted by lexicographic, ascending order over the given variable ordering, e.g.  $R$  needs to

---

<sup>3</sup><https://www.relational.ai/>

<sup>4</sup><https://github.com/schroederdewitt/leapfrog-triejoin>

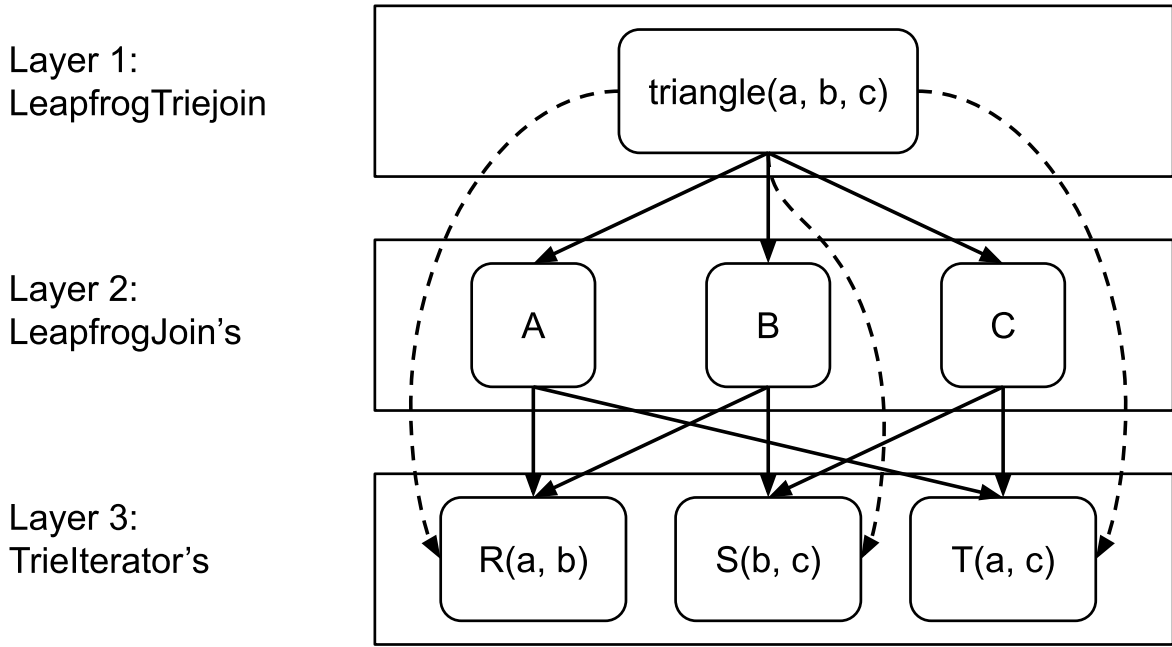


Figure 4: The three layers of the Leapfrog Triejoin algorithm. The configuration for a triangle query is shown: three *TrieIterators* one per input relationship, three *LeapfrogJoins* one per variable and one *LeapfrogTriejoin* component are necessary. The arrows indicate that a component uses another. The *LeapfrogTriejoin* uses all other components but only the vertical part of the *TrieIterators* (dashed arrows). The *LeapfrogJoins* uses the linear part of two *TrieIterators* each.

be sorted by primarily by  $a$  and secondary by  $b$  given the variable ordering  $a, b, c$ . The algorithm is variable-oriented because it fixes one possible binding for  $a$ , one for  $b$  given  $a$  and finally one for  $c$  given  $a$  and  $b$ . This allows it to enumerate the result of the join query without intermediary results. The process can be thought of as a backtracking, depth-first search for possible bindings.

The algorithms implemented in this thesis can process joins of the full conjunctive fragment of first order logic or conjunctive equi-joins in relational algebra terms. Possible extensions to disjunctions, ranges (none-equi joins), negation, projection, functions and scalar operations on join variables are explained in the original Leapfrog Triejoin paper [44]. However, they are not relevant to the core of this work because many interesting graph patterns can be answered using the full conjunctive fragment, e.g. cliques or cycles.

The Leapfrog Triejoin algorithm uses three components which are composed in a layered fashion. The concrete composition used for the triangle query is shown in fig. 4. In this figure, we see three layers each of them made of one or more instances of a component. The components are the *TrieIterator*, *LeapfrogJoin* and *LeapfrogTriejoin*. In the next paragraphs, we explain each layer in order, starting with the lowest layer.

The lowest layer is made of one *TrieIterator* per input relationship. In our example, we have three instances one for  $R$ ,  $S$  and  $T$  each. The *TrieIterator* interface represents the input relationship as a trie with all values for the first attribute on the first level, the values for the second attribute on the second level and so on; an example for this is shown in fig. 5.

The trie contains one level per attribute of the relationship; in the case of the triangle query,

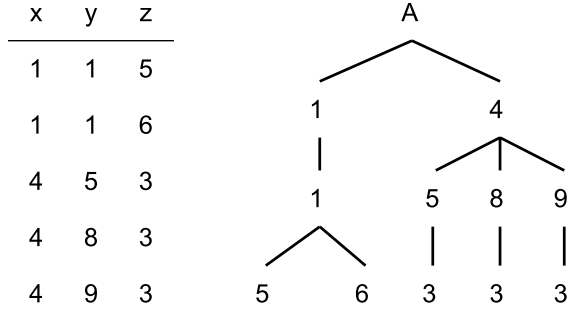


Figure 5: A 3-ary relationship as table (left) and trie (right), to position the iterator at the tuple (1, 1, 5) one calls *open* twice, *key* returns now 5, after a call to *next*, *key* returns 6 and *up* would lead to *key* returning 1.

Method	required complexity
int key()	$\mathcal{O}(1)$
bool atEnd()	$\mathcal{O}(1)$
void up()	$\mathcal{O}(\log_N)$
void open()	$\mathcal{O}(\log_N)$
void next()	$\mathcal{O}(\log_N)$
void seek(key)	$\mathcal{O}(\log_N)$

Table 1: The *TrieIterator* interface with required complexity.  $N$  is the size of relationship represented by the iterator.

there are two levels: one for  $a$  and one for  $b$ . Each level is made of all possible values for its attribute. All tuples of the relationship can be enumerated by a depth-first traversal of the trie.

The *TrieIterator* component offers six methods shown in table 1. The *open* and *up* methods control the level the iterator is positioned at; *open* moves it one level down and *up* moves it one level up. Additionally, *open* places the iterator at the first value for the next level and the *up* method returns to the value of the upper level that was current when the deeper level was opened. We call these two methods the vertical component of the *TrieIterator* interface.

The other four methods are called linear component. All of them operate on the current level of the *TrieIterator*. The *key* function returns the current key (a single integer). The *next* method moves the iterator to the next key on the same level. The *seek(key)* operation finds the least upper bound for its parameter *key*. Finally, the *atEnd* method returns *true* when the iterator is placed behind the last value of the current level.

The middle layer of the Leapfrog Triejoin is made of one *LeapfrogJoin* per variable in the join. This join generates possible bindings for its variable by intersecting the possible values for all input relationship containing the variable. Therefore, it operates on the linear component of all *TrieIterators* of relationships with this variable. Figure 4 for the triangle query shows three *LeapfrogJoin* instances (for  $a, b$  and  $c$ ); each of them uses two *TrieIterators*.

The *LeapfrogJoin* interface has five methods shown in table 2, with their required asymptotical performance. In the following paragraphs, we explain each of them. In short, the join offers an iterator interface over the intersection of its input iterators. This intersection is found by repeatedly seeking the value of the largest input iterator in the smallest input iterator. This

Method	required complexity
int key()	$\mathcal{O}(1)$
bool atEnd()	$\mathcal{O}(1)$
void init()	$\mathcal{O}(\log_N)$
void next()	$\mathcal{O}(\log_N)$
void seek(key)	$\mathcal{O}(\log_N)$

Table 2: The *LeapfrogJoin* interface with required complexity.  $N$  is the size of relationship represented by the iterator.

process resembles a frog taking a leap which gives the join its name. When all iterators point to the same value leapfrogging stops and the value is emitted as part of the intersection.

The *init* operation sorts the input iterator by their current key and finds the first value of the intersection. To find the first value it uses the private method *leapfrogSearch* which is the work-horse of the whole join. The algorithm of this method is shown in algorithm 1. This method loops the process of calling the *seek* method of the its smallest input iterator with the key of the largest input iterator until the smallest and the largest (and therefore all iterators) point to the same value.

```

Data:  iters # sorted array of TrieIterators ;
p # index of the smallest iterator;
Result: Either atEnd is true or key is set to next key of intersection
maxKey ← iters[p % iters.length].key();
while iters[p].key() ≠ maxKey do
    iters[p].seek(maxKey);
    if iters[p].atEnd() then
        atEnd ← true;
        return;
    else
        maxKey ← iters[p].key();
        p ← (p + 1) % k;
    end
end
key ← iters[p].key()

```

**Algorithm 1:** leapfrogSearch()

The *leapfrogNext* method moves the join to its next value. Internally, it uses the *next* function of its smallest iterator and then *leapfrogSearch*.

The operation *leapFrogSeek(key)* first uses the *seek* method of the smallest input iterator to forward it to *key*; then it uses *leapfrogSearch* to either verify that this key is available in all iterators (hence in the intersection) or to find the upper bound of this key.

Finally, the functions *key* and *atEnd* return the current key or if the intersection is complete respectively.

The last layer of the whole algorithm is a single *LeapfrogTriejoin* instance. It interacts with both lower layers to enumerate all possible bindings for the join. For this it acquires one binding for the first variable from the corresponding *LeapfrogJoin*. Then it moves the *TrieIterators* containing this variable to the next level and finds a binding for the second variable using the next *LeapfrogJoin*. This process continues until all variables are bound and a tuple representing this binding is emitted by the join operator. Then it finds the next possible binding by backtracking.

Algorithm 2 shows the backtracking depth-first traversal. This traversal needs to stop each time when a complete tuple has been found to support the iterator interface of the join. Therefore, it is implemented as a state-machine which stops each time the deepest level is reached and all variables are bound (loop condition in line 35). The next action of the state machine is determined by the outcome of the current action. Hence, we can characterize the state machine by describing each possible action and its possible outcomes. There are three possible actions: *next*, *down* and *up*. We summarize the possible actions, conditions for the next action and if the main loop of the state machine yields the next tuple in table 3 and describe each action below.

The *next* action moves the *LeapfrogJoin* at the current depth to the next possible binding for its variable (line 2). If the *LeapfrogJoin* reached its end, we continue with the *up* action (line 4), otherwise we set the binding and continue by another *next* action, if we are at the deepest level or by moving to the next deeper level by the *down* action (line 7).

The *down* action moves to the next variable in the global variable ordering by opening all related *TrieIterators* and initializing the corresponding *LeapfrogJoin* (line 20 call to *trieJoinOpen*). A *down* can be followed by an *up* if the *LeapfrogJoin* is *atEnd* (line 22), by a *next* action if the trie join is at its lowest level (line 25), or by another *down* action to reach the deepest level.

The *up* action can signal the completion of the join if all bindings for the first variable in the global ordering have been explored, or in other words, the first *LeapfrogJoin* is *atEnd* (condition  $depth == 0 \wedge action == UP$  line 28). Otherwise, all *TrieIterators* corresponding to the current variable are moved upwards by calling *triejoinUp* (line 31) which also updates *depth* and *bindings*. Then, this action is followed by another *up* or a *next* depending on *atEnd* of the current *LeapfrogJoin* (lines 32).

**TrieIterator implementation, backing data structure** While we can implement the *LeapfrogJoin* and *LeapfrogTriejoin* component of the Leapfrog Triejoin from the algorithmic description given above, we are missing some details for a concrete implementation of the *TrieIterator* interface. Mainly, we need to decide for a datastructure to back the *TrieIterator*.

We choose to use sorted arrays as described in [16]. One array is used per column of the input relationship and binary search on these arrays allows us to implement the *TrieIterator* interface with the required asymptotical complexities (see table 1).

**Variable ordering** Finding a good variable ordering for the LFTJ is an interesting research problem in itself. We are aware of two existing approaches.

The first is to create and maintain representative samples for each input relationship and determine the best order based on runs over these samples. This has been implemented in LogicBlox, the first system to use Leapfrog Triejoins [10]. To the best of our knowledge, the exact method of creating the representative samples has not been published.

The second approach is described in great detail by Mhedhbi and Salihoglu in [34]. Its has been implemented in their research graph database Graphflow [24].

They define a novel cost-metric for WCOJs which estimates the costs incurred by constructing the intersections of adjacency lists. The metric takes three factors into account. First, the size adjacency lists.

Second, the number of intermediate matches. The concept of intermediate matches is best understood by a simple example; we see the tailed-triangle query in ?? . Two very different vertice ordering categories exist for this query. The ones that start on  $v_4$  and find all 2-paths of the graph; and vertice orderings that start with  $v_1, v_2, v_3$  in any order which close the triangle first. Clearly, there are more 2-paths in any graph than triangles. Hence, the second category produces far less intermediate matches.

**Data:** *depth* the index of the variable to find a binding for, range from -1 to #variables - 1 \*

```

;
MAX_DEPTH
* the number of variables
* bindings
* array holding the current variable bindings or -1 for no binding
* ;
action
* state of the state machine
* ;
repeat
  switch action do
    case NEXT do
      leapfrogJoins[depth].leapfrogNext() ;
      if leapfrogJoins[depth].atEnd() then
        | action  $\leftarrow$  UP ;
      else
        | bindings(depth)  $\leftarrow$  leapfrogJoins[depth].key() ;
        | if depth == MAX_DEPTH then
          | | action  $\leftarrow$  NEXT
        | end
        | action  $\leftarrow$  DOWN
      end
    end
  end
  case DOWN do
    | depth  $\leftarrow$  depth + 1 ;
    | trieJoinOpen() ;
    | if leapfrogJoins[depth].atEnd() then
      | | action  $\leftarrow$  UP ;
    | else
      | bindings(depth)  $\leftarrow$  leapfrogJoins[depth].key() ;
      | if depth == MAX_DEPTH then
        | | action  $\leftarrow$  NEXT
      | else
        | | action  $\leftarrow$  DOWN
      | end
    | end
  end
  case UP do
    | if (depth == 0) then
      | | atEnd  $\leftarrow$  true ;
    | else
      | depth  $\leftarrow$  depth - 1 ;
      | trieJoinUp() ;
      | if leapfrogJoins[depth].atEnd() then
        | | action  $\leftarrow$  UP ;
      | else
        | | action  $\leftarrow$  NEXT ;
      | end
    | end
  end
end
until !((depth == MAX_DEPTH && bindings[MAX_DEPTH]  $\neq$  -1)) // atEnd;

```

Action	Condition	Next action	Yields
NEXT	$lf.atEnd$	UP	no
	$\neg lf.atEnd \wedge reachedMaxDepth$	NEXT	yes
	$\neg lf.atEnd \wedge \neg reachedMaxDepth$	DOWN	no
DOWN	$lf.atEnd$	UP	no
	$\neg lf.atEnd \wedge reachedMaxDepth$	NEXT	yes
	$\neg lf.atEnd \wedge \neg reachedMaxDepth$	DOWN	no
UP	$depth = 0$ , means highest $lf.atEnd$ is true	– (done)	yes
	$lf.atEnd$	UP	no
	$\neg lf.atEnd$	NEXT	no

Table 3: Summary of actions, conditions for the following action and if a complete tuple has been found. *reachedMaxDepth* is true if we currently find bindings for the last variable in the global order. *lf* abbreviates the *LeapfrogJoin* of the current variable. The columns *Yields* details if the main loop of the state machine yields before computing the next action, this is the case, when all variables have been bound.

Finally, they implement an intersection cache in their system which takes advantage of the fact that some queries can reuse already constructed intersections. So, the last factor taken into account by their cost metrics is the usage of this intersection cache.

They use the described cost metric, a dynamic programming approach to enumerate possible plans and a catalogue of sampled subgraph instances containing the sizes of adjacency lists to intersect and produced intermediate results to estimate the costs for all variable orderings.

Moreover, they implement the ability to change the query ordering adaptively during query execution based on the real adjacency list sizes and intermediate results. They show that adaptive planning can improve the performance of many plans. Furthermore, it makes the query optimizer more robust against choosing bad orderings.

To conclude, the work of Mhedhbi et al. is the most comprehensive study on query vertex orderings for WCOJs currently available; they introduce a cost metric, a query optimizer to use this metric and prove that it is possible and beneficial to compute parts of the results using a different variable order.

In our work, we do not implement an automatic process to choose the best variable order. The order we choose is based on experiments with different orders and intuition of the author.

Integrating the approach of LogixBlox would be possible but require the implementer to find a good sampling strategy because no details are openly available.

The approach of Mhedhbi and Salihoglu is much better documented but also more complex. It consists out of four contribution which build up on each other but could be useful on its own. The cost metric described in their paper applies to our system as well and could be used.

They use this metric for cost estimation in connection with a none-trivial subgraph catalogue. The main challenge in integrating this way of cost estimation with our system is to elegantly integrate catalogue creation in Spark.

Their solution for adaptive variable orderings is helpful because it proves that this technique is beneficial; they also publish performance measurements, so the impact can be evaluated. However, there system employs a *Generic Join* while we use a *Leapfrog Triejoin*. The integration



of adaptive variable orderings into Leapfrog Triejoin is not trivial and it is likely that their implementation is not directly applicable.

Finally, they introduce an intersection cache to make use of repeatedly used intersections. This can be directly applied to our system, e.g. using the decorator pattern around *LeapfrogJoins*. We note that they only cache the last, full n-way-interesection of multiple adjacency lists. It would be interesting to research if the system would benefit from caching partial n-way intersections as well because we noticed that for some queries, e.g. 5-clique, the intersection between the first two lists can be reused more often than the full intersection. This opens the interesting question in which order we should intersect the lists.

We conclude that two concepts to choose a good variable ordering exists which are both (partially) applicable to our system. The LogixBlox approach is simpler and directly integratable but not well documented. The solution used in GraphFlow is far more complex and developed for another WCOJ. Anyhow, the paper describes it in great detail and parts of it could be integrated directly, while others need some engineering effort or need to be redesigned completely.

## 2.4 Distributed worst-case optimal join in Myria

In 2014, a Leapfrog Triejoin variant, dubbed Tributary Join, was used as a distributed join algorithm on a shared-nothing architecture called Myria [16]. They use Tributary Join as a local, serial worst-case optimal join algorithm, combined with the Hypercube shuffle algorithm, also called *Shares*, to partition the data between their machines [2]. The combination of a shuffle algorithm with a WCOJ allows them to distribute an unchanged serial worst-case optimal join version by running it only on a subset of the data on each worker.

This approach is directly applicable to Spark. We could implement a hypercube shuffle for Spark and then choose any WCOJ to run on each partition. However, it is not obvious how well this approach scales because Shares replicates many of its input tuples [16]. The experiments on Myria indicate that the combination of Hypercube shuffles and Tributary Join does not scale well. They report a speedup of 8 on 64 workers compared to the time it takes on 2 nodes, which, although unlikely to be optimal, is not investigated in great detail.

Therefore, we decide to analyse the expected scaling behaviour of Shares for graph pattern matching. Our main concern is that the number of duplicated tuples in the system increases with the query size (number of vertices) and with the number of workers added to the system. We provide a theoretical analysis of the number of duplicated tuples for different query sizes and available workers in a later section of this thesis (??). As result of this investigation, we decide to not physically partition our data but to duplicate it to all workers and seek for alternative strategy to parallelize a worst-case optimal join.

To conclude, the implementation in Myria and in particular the Shares algorithm is the starting point of this thesis. We see it as our baseline for a distributed WCOJ implementation. In the coming section, we explain Shares in detail.

We note that a implementation of a distributed worst-case optimal join on Timely Dataflow exists. However, it is not applicable to Spark. Therefore, we treat it in the related work section (7.2) of this thesis.

### 2.4.1 Shares

Shares partitions the input relationships for a multi-way join over worker nodes, such that, all tuples, which could be joined, end up on the same worker in a single shuffle round. Hence, it allows running any multi-way join algorithm locally after one shuffle round. The output of the join is the union of all local results.



Figure 6: Left: Three aliases of an edge relationship with one triangle. The participating tuples are marked in red, green and blue. Their hypercube coordinates are shown below. Right: Example of a Shares hypercube configuration for the triangle query for 12 workers with three attributes/dimensions of the sizes 3, 2, 2. The tuples marked in red, green and blue end up on the workers with red, green and blue rhombs respectively.

The idea is to organize all workers in a logical hypercube, such that each worker can be addressed by its hypercube coordinate. Then it is straightforward to find a mapping from the attribute values of a tuple to these coordinates, so that joinable tuples arrive on the same worker after one shuffle.

We first explain how to organize the workers in a hypercube and then how to map tuples to these workers. Next, we treat the problem of choosing a good hypercube configuration. Followed, by a summary about optimality of Shares. Finally, we provide analysis of the scaling of Shares for graph pattern matching.

A hypercube is characterized by the number of dimensions and the size of each dimension. Figure 6 shows a hypercube with three dimensions labelled  $a, b$  and  $c$ . They have the size of 3, 2 and 2 for  $a, b$  and  $c$  respectively. It is a possible configuration for the triangle query  $triangle(a, b, c) \leftarrow R(a, b), S(b, c), T(a, c)$  with 12 workers.

Given an input query, Shares builds a hypercube with one dimension per variable in the input. It then chooses the size of each dimension, such that the product is smaller than number of workers. We call  $v$  the numbers of variables in the query and  $p_0, \dots, p_v$  the sizes of each dimension. This allows us to address each worker with a coordinate of the form  $(0..p_0, \dots, 0..p_v)$ . If the product of all dimension sizes is smaller than the number of workers, additional workers are not used. The process of finding the best sizes for the dimensions depends on the input query and the input relationships. We discuss it in a later paragraph of this section.

With this topology in mind, it is straightforward to find a partitioning for all tuples from all relationships such that tuples that could join are sent to the same worker. We choose a hash function for each join variable  $a$  which maps its values in the range of  $[0..p_a]$ . Then each worker determines where to send the tuples it holds by hashing its values. This results in a coordinate in the hypercube which is fixed for all join variables, which occur in the tuple, and unbounded

for join variables which do not occur in the tuple. Then the tuple is sent to all workers with a matching coordinate.

Figure 6 shows how tuples forming a triangle from three relationships are mapped to the workers. The blue, green and red tuple in the relationships form a triangle. The green and the red tuple are sent to 2 workers each and the blue tuple to three workers (marked with small rhombs). They are sent to all workers along the axis where the coordinate is not determined by a value in the tuple. We see that they all end up together on the worker with the coordinate (2, 0, 0). This is where the triangle occurs in the output of the join.

**Finding the best hypercube configuration** The problem of finding the best hypercube configuration is to choose the sizes of its dimensions such that (1) the product of all sizes is smaller than the number of available workers and (2) such that the number of tuples to single worker is minimized. (2) is backed by the assumption that the number of tuples is a good indicator for the amount of work; this assumption is made in all papers discussing the problem [16, 13, 2]. Therefore, we want to minimize the number of tuples on each single worker because the slowest worker determines the run-time. Next, we discuss existing solutions and decide for one of them.

The original Shares paper proposes a none-convex solution which is hard to compute in praxis [2]. Later, Beame et al. define a linear optimization problem which is solvable but leads to fractional hypercube sizes [13]. Hence, it is not possible to use their solution directly. Rounding down would be an obvious solution but as discussed in [16], it can lead to highly suboptimal solutions, in particular with low numbers of workers. Hence, the paper further considers to use a higher number of *virtual* workers and assign these to *physical* workers by a one-to-many mapping. Anyhow, a higher number of workers lead to more replicated tuples. Therefore, this solution does not scale well.

In the end, the paper that integrates Shares and the Tributary join in Myria suggests a practical solution. They enumerate all integral hypercube configurations smaller or equal to the number of available workers. For each configuration they estimate number of assigned tuples, then they choose the configuration with the lowest estimated workload.

They use the following equation to estimate the workload, where  $R$  are all relationships in the query,  $var(r)$  gives the variables occurring in the relationship  $r$ ,  $size(v)$  gives the hypercube size of the dimension for variable  $v$ .

$$workload = \sum_{r \in R} |r| \times \frac{1}{\prod_{v \in var(r)} size(v)} \quad (3)$$

The term  $\prod_{v \in var(r)} size(v)$  gives the numbers of workers that span the hyper-plain over which a relationship  $r$  is partitioned. For example, in fig. 6 the relationship  $(a, b)$  is partitioned over the plain spanned by the dimensions  $a$  and  $b$  with 6 workers. Each tuple in this relationship has a chance of  $\frac{1}{6}$  to be assigned to any of these workers. Hence, the workload caused by  $(a, b)$  is  $|a, b| \times \frac{1}{6}$ .

The paper evaluates this strategy to assign hypercube configurations and finds that it is efficient and practical. We choose to use the same solution for our work.

**Shares is worst-case communication-cost optimal** Shares, as described above, is shown to be worst-case optimal in its communication costs in MapReduce like systems for  $n$ -ary joins using one shuffle round. First, Beame et al. prove that Shares scheme is optimal on databases without skew [13]. Later the same authors are able to give a proof that Shares is also an optimal algorithm for skewed databases if one knows the *heavy-hitter* tuples and splits the join into a skew free part and residual joins for the heavy hitters using different hypercube configurations for each residual join [26].

The implication of these proofs is that it is not possible to find a partitioning scheme for one shuffle round that replicates less data than Shares. This observation is central to our thesis because it is one argument to replicate the graph on all workers instead of using a shuffle algorithm to partition it.

In the rest of this thesis, Shares refers to the original algorithm [2] and not the skew resilient variant *SharesSkew* [4, 13]. This is mainly because even in the presence of skew the original Shares scheme offers good upper bounds, although it can not always match the lowest bound possible [4]. But also because the skew resilient variant requires to know which tuples are *heavy-hitters* (a definition of skew introducing tuples). Finally, while first experiments with SharesSkew exist [4], we are not aware of an extensive study verifying it is possible to integrate SharesSkew into a complete system. Hence, we deem it out of scope for this thesis to attempt a full integration.

Some readers might ask if there are better multi-round algorithms which replicate less data. Indeed, the authors of a Shares related paper raise the same question as future work [26]. They are able to answer this question for specific join queries in [26, 4], e.g. chain-joins and cycle-joins. Later, they present an algorithm which is multi-round optimal for all acyclic queries [3] and one for all queries over binary relationships [25].

The papers about multi-round optimal partitioning schemes are rather theoretical. To the best of our knowledge, only one paper provides practical experiments [4] but has no dedicated implementation section. Also, they have not been shown optimal for general conjunctive join queries but only for special cases. Two of the three papers [c]annot handle clique joins which are important class of joins in our thesis. Additionally, they add additional complexity to the query optimizer, e.g. they require the input query to be represented as generalized hypertree decomposition to calculate their intersection width [3] or to find many different hypercube configurations [25, 4, 26] which is not trivial in praxis and computation intensive as discussed in the last paragraph.

We leave it to future research to investigate the practical application of these algorithms to graph pattern matching. The most interesting paper in this direction is [25]. It develops a multi-round algorithm for n-ary joins on binary relationships like the edge relationship of a graph.

**Analysis of Shares scalability** Next, we analyse the scalability of Shares on growing graph patterns. That is, self-joins over a single relationship which has two variables. In this context, relationships of the join can be seen as the edges of the pattern and variables as vertices.

First, we fix the method to determine the best hypercube configuration  $(p_1 \dots p_k)$ , given a query. For this, we use the method described above and used in [16].

Given the hypercube configuration and a query, we can estimate the workload of each worker by the formula 3. Let  $R$  be the set of all atoms in the join<sup>5</sup>,  $size1(r)$  and  $size2(r)$  be the size of the first respectively second hypercube dimension for the two variables in atom  $r$ . Then, each worker receives  $\sum_{r \in R} \frac{|r|}{size1(r) * size2(r)}$  tuples under the assumption of uniform data distribution and good hash functions. Our argument is that the tuples of each atom  $r$  are divided onto  $size1(r) * size2(r)$  workers; the workers that form the hypercube plain of its two variables.

In the special case of graph pattern matching where all atoms of the query are pointing to the same relationship, we can optimize the hypercube shuffle such that a tuple is only sent once to a worker, although it might be assigned to it via multiple atoms.

If we apply this optimization, we can predict the probability with which each tuple is assigned to a worker using the Poisson binomial distribution. The Poisson binomial distribution

---

<sup>5</sup>An atom in a datalog join is the reference to a relationship, e.g.  $triangle(a, b, c) \leftarrow R(a, b), S(b, c), T(a, c)$  has three atoms named  $R, S$  and  $T$ . In this section, we differentiate atoms and relationships because multiple atoms can point the same underlying relationship which becomes of particular importance.

Pattern	Edges	workload [64]/[128]
Triangle	3	0.18 / 0.12
4-clique	6	0.59 / 0.44
5-clique	10	0.90 / 0.82
House	5	0.42 / 0.32
Diamond	8	0.76 / 0.67

Table 4: Workload on 64 and 128 workers in percentage of tuples of the edge table assigned to each worker estimated by using Poisson binomial distribution to estimate the workload and the method from [16] to determine the optimal shares configuration.

$\Pr(n, k, u_0, \dots, u_n)$  allows us to calculate the likelihood that  $k$  out of  $n$  independent, binary and differently distributed trials succeed, under the condition that the  $i$  the trial succeeds with a probability of  $u_i$ . We use  $n = |R|$ ,  $k = 0$  and  $u_i = 1/(size1(r_i) * size2(r_i))$  to calculate the probability that a tuple is not assigned to an arbitrary, fixed worker. This allows us to predict the number of tuples assigned to each worker by  $|E| * (1 - \Pr(|R|, 0, u_0, \dots, u_{|R|}))$  with  $E$  being the edge relationship.

Table 4 shows the expected percentage of tuples from the edge relationship assigned to each worker for graph patterns of different sizes calculated using Poisson binomial distribution and optimal shares assignments according to the method used in [16]. As we can see in this table, the number of tuples assigned to each worker grows over linear in the size of the graph pattern. Furthermore, doubling the number of workers is inefficient to counter this growth.

In particular, already small clique queries of four vertices replicate over half of the tuples on all 64 workers. 5-clique queries require nearly a full broadcast with each worker holding 82% or 90% of all tuples with 128 respectively 64 workers. The diamond query used in practice by the Twitter recommendation engine has to replicate far more than half of the tuples to all workers.

The second observation has two reasons. First, doubling the number of workers does not allow us to double the dimensions of the hypercube because a hypercube always needs product of all dimension sizes to be built. Second, the number of replicated tuples increases with a growing hypercube because each tuple is replicated to more workers; namely  $\prod_{r \in R/r} size1(r) * size2(r)$  workers. This is because each tuple binds only two out of all variables. Hence, it is replicated over many dimensions.

In light of the numbers presented in table 4 and in line [7], we conclude that the communication costs for Shares converge towards a full broadcast for bigger graph patterns and scaling becomes increasingly inefficient. By this observation and the fact that hypercube shuffling is an optimal scheme (see the last paragraph), we decide against using any partitioning scheme in our work but replicate the edge relationship on all workers.

## 2.5 Compressed sparse row representation

Compressed sparse row representation (short CSR) is a well known, low-memory representation for static graphs [15, 42]. To ease its explanation, we assume that the graph’s vertices are identified by the numbers from 0 to  $|V| - 1$ . However, our implementation allows the use of arbitrary vertex identifiers in  $\mathcal{N}$  by storing the translation in an additional array of size  $|V|$ .

CSR uses two arrays to represent the edge relationship of the graph: one of size  $|E|$  which is a projection of the edge relationship onto the *dst* attribute and a second of size  $|V| + 1$  which stores indices into the first array. To find all destinations directly reachable from a source *src*

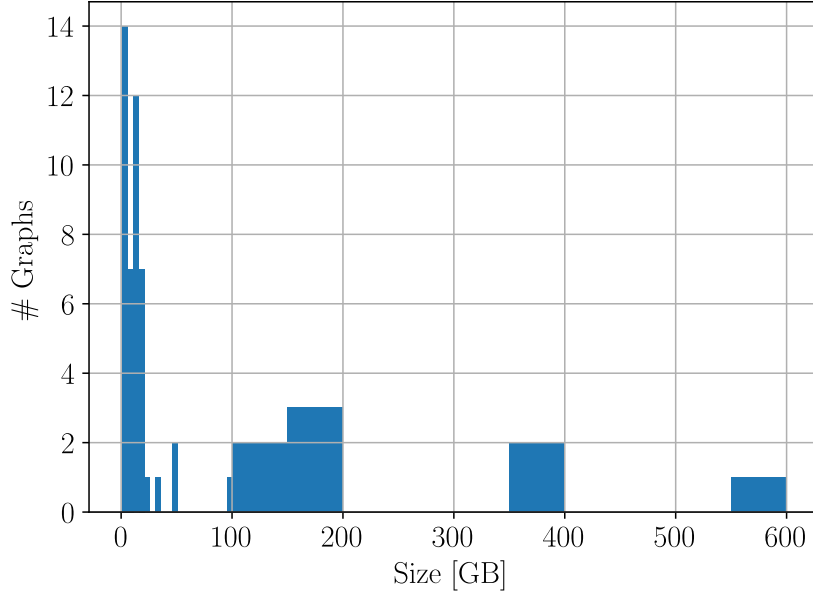


Figure 7: Sizes of all graphs from the SNAP and Laboratory of Web Algorithms dataset collection in giga bytes. The histogram shows graphs up to 100 GB in buckets of 5 GB and in buckets of 50 GB after. In total, we see data collected about 157 graphs.

$\in V$ , one accesses the second array at *src* for the correct index into the first array for a list of destinations.

The CSR format has two beneficial properties in the context of this thesis. First, it allows locating all destinations for a source vertex by one array lookup; hence, in constant time. Second, the representation is only, roughly, half as big than a simple columnar representation. A uncompressed columnar representation needs  $2 \times |E|$  while CSR uses only  $|V| + 1 + |E|$ , note that for most real-world graph  $|V| \ll |E|$  holds (see section 2.6).

## 2.6 Analysis of public real-world graph datasets

In this section, we present a short analysis of the sizes of real world graph datasets. For this, we collect data about all graphs from the SNAP and Laboratory of Web Algorithms dataset collection [**snap-datasets**, **laboratory-of-web-algorithms-datasets**]. The graphs in the Snap dataset are a bit older; they have been collected between 2000 and 2010. All Laboratory of Web Algorithms graph have been collected between 2007 and 2018. Both dataset collections are heavily used and cited in academia [**fractal**, **longbin**, 7, 38, 16]. Two of these papers are from 2019.

For our size calculation we assume that the graph is stored in compressed sparse row representation section 2.5 using integers for the vertex ID's. Then, we determine the storage size in bits by the formulae  $32 \times |V| + 32 \times |E|$  with 32 the size of an integer in bits,  $V$  the set of all vertices in the graph and  $E$  the set of all edges in the graph.

Figure 7 shows a histogram of sizes for all 157 graphs from the two datasets. 104 of these graphs are smaller than 1 GB and only 8 graphs are bigger than 100 GB. The biggest graph is the friendship graph of Facebook from the year 2017 with 552.2 GB.

We conclude that even the biggest graph can be fitted in main memory of many cluster machines today. The vast majority could be fitted in the main memory of a simple desktop machine or laptop.

### 3 Worst-case optimal join parallelization

Based on the fact that Shares is an optimal partitioning scheme for n-ary joins in MapReduce like systems [2] and our analysis that Shares converges to a full broadcast of the graph edges (see ??), we decided to forego physical partitioning of the graph. We cache the graph in memory such that each Spark task can access the whole graph. Then, we experiment with multiple *logical* partitioning schemes which ensure that each task processes only some parts of the graph. This design has a big advantage over physical partitionings. Each worker holds the full edge relationship, therefore, it can answer any possible query without needing to shuffle data or materializing new data structures for the *LeapfrogTriejoin*, e.g. sorted arrays or CSR representations. Arranging the data into suitable data structures and shuffling data is a one-off action on system startup.

This design allows us to implement a new flavour of the Shares partitioning in which we filter the vertices of the graph on-the-fly while processing it with our *GraphWCOJ* algorithm. We describe this contribution in section 3.1.

Furthermore, we consider a work-stealing based partitioning which does not replicate any work and produces less skew than Shares. This comes at the price of implementing work-stealing on Spark. The design of work-stealing in Spark is described in section 3.2

#### 3.1 Logical Shares

Shares has been developed as an optimal shuffle for n-ary joins on MapReduce like systems. So, it is used to physical partition the tables participating in the join overall workers of the system. Then, each worker works only on the tuples it holds in its partition. This has been implemented in Myria to be used with a WCOJ [16] and for Hadoop [TODO]. We describe the Shares and Myria in more detail in ?? and assume that the reader is familiar with this section.

The idea of Shares can also be used for a *logical* partitioning scheme. Instead, of partitioning the graph before computing the join, we determine if a tuple should be considered by the join on-the-fly. We do so by assigning a coordinate of a hypercube to each worker. Then each worker is responsible for the tuples which match its coordinate as in the original Shares. However, a huge difference to a physical Shares partitioning exists: while physical Shares keeps one prefiltered, materialized partition of the edge relationship per relationship of the join in memory, we keep only a single copy of the graph edge relationship in memory which can be used by all *TrieIterator* for all edge table aliases. Once, the edge relationship is broadcasted and cached, we do not need to materialize prefiltered partitions of it before every query.

Filtering tuples on-the-fly in the LFTJ comes with a challenge: in the *LeapfrogTriejoin* we do not consider whole tuples but only single attributes of a tuple at the time, e.g. a *LeapfrogJoin* only considers one attribute and cannot determine the whole tuple to which this attribute belongs. Fortunately, a tuple matches only if all attributes match the coordinate of the worker. Hence, we can filter out a tuple if any of its attributes do not match. For example, we can exclude a value in a *LeapfrogJoin* without knowing the whole tuple.

Integrating Shares and LFTJ comes with two important design decisions. First, the *LeapfrogTriejoin* operates on a complete copy of the edge relationship. Hence, we need to filter out the values that do not match the coordinate of the worker. Second, we need to compute the optimal Hypercube configuration. We describe our solutions below.

The first design decision is where to filter the values. The *LeapfrogTriejoin* consists out of multiple components which are composed as layers upon each other. On top we have the *LeapfrogTriejoin* which operates on one *LeapfrogJoin* per attribute. The *LeapfrogJoins* uses multiple *TrieIterators*. Our first instinct is to push the filter as deep as possible into these layers.

We built a *TrieIterator* that never returns a value which hash does not match the coordinate. This is implemented by changing the *next* and *seek* methods such that they linearly consider further values until they find a matching value if the return value of the original function does not match. However, the resulting LFTJ was so slow that we abandoned this idea immediately. We hypothesize that this is the case because the original *next* and *seek* method is now followed by a linear search for a matching value. Furthermore, many of these values are later dropped in the intersection of the *LeapfrogJoin* which can also be seen as a filter over the values of the *TrieIterators*. As we know from ??, the *LeapfrogJoin* is a rather selective filter. It does not make sense to push a less selective filter below a more selective filter.

With this idea in mind, we build a logical Shares implementation that filters the return values of the *leapfrogNext* method. This is implemented as a decorator pattern around the original *LeapfrogJoin*. The use of the decorator pattern allows us to easily integrate Shares with the LFTJ while keeping it decoupled enough to use other partitioning schemes.

The second design decision is how and where to compute the best hypercube configuration. The how has been discussed extensively in former literature [2, 16]. We implement the exhaustive search algorithm used in the Myria system [16].

In the interest of a simple solution, we compute the best configuration on the master before starting the Spark tasks for the join. We note that the exhaustive algorithm could be optimized easily and it would be worthwhile to introduce a cache for common configurations. Due to time constraints, we leave this to future work and keep our focus on the scaling behaviour of Shares.

To conclude, we succeeded to integrate Shares with *LeapfrogTriejoin* and report our results in ??. We cannot improve on the main weakness of Shares that it duplicates a lot of work. Indeed, our design filters tuples only after the *LeapfrogJoin*. Therefore, all tuples are considered in the *TrieIterator* and their binary search of the first variable. This does not influence scaling much because only the correct logical partition of values for the first variable are used as bindings in the *LeapfrogTriejoin*. This means they are still filtered early enough before most of the work happens. We improve over a physical Shares by using the same CSR data structure for all *TrieIterator*. Therefore, we do not need to materialize a prefiltered data structure for each *TrieIterator* and query which saves time and memory if the partitions become large for bigger queries.

### 3.1.1 RangeShares

In the last section, we raised the point that our Shares implementation only filters out value after the *LeapfrogJoins*. This is because a hash-based filter needs to consider single values one-by-one. In this section, we explore the possibility to use range based filters which can be pushed into the *TrieIterators*. However, we warn the reader that this is a negative result. It leads to high skew which hinders good scaling of this idea.

We observe that the general idea behind Shares is to introduce a mapping per attribute from the value space into the space of possible hypercube coordinates, e.g. so far all Shares variants use a hash function per attribute to map the values onto the hypercube. We investigate the possibility to use ranges as mapping function, e.g. in a three-dimensional hypercube with three workers per dimension, we could divide the value space into three ranges; a value matches a coordinate if it is in the correct range. Contrary, to hash-based mappings which are checked value by value until one matches, a range check is a single conditional after each *seek* and *next* function call. Furthermore, this conditional is predictable for the processor because, for all but one call, the value is in range and returned. So, contrary to hash-based filter we can push a range based filter



into the *TrieIterators*.

We implement this idea by dividing the vertice ids per attribute into as many ranges as the size of the corresponding hypercube dimension. For example, assume we have edge ids from 0 to 899, three attributes and the hypercube dimension have the size 3, 2 and 2. Then, we choose the ranges [0,300), [300,600) and [600,900) for the first attribute and the ranges [0,450) and [450,900) for the other two attributes. The worker with the coordinate (0, 0, 0) is then assigned the ranges [0,300), [0,450) and [0,450). It configures its *TrieIterators* accordingly such that they are limited to these ranges.

We run first experiments to evaluate this idea. We expect it to scale better than a hash-based Shares because it saves intersection work in the *LeapfrogJoins*. However, we find that high skew between the workers leads to much worse performance than a hash-based Shares. The explanation is that if a worker is assigned the same range multiple times and this range turns out to take long to compute, it takes much longer than all other workers.

To mitigate this problem, we break down the vertice ids into more ranges than there are workers in the hypercube dimension corresponding to the attributes. Then, we assign multiple ranges to each *TrieIterator* in such a way that the overlap on the first two attributes equals the overlap of a hash-based implementation and assign the ranges of the later attributes randomly such that all combinations are covered. However, experiments still show a high skew: some workers find many more instances of the searched pattern in their ranges than others. For the triangle query on *LiveJournal*, we find that the fastest worker outputs only 0.4 times the triangles than the slowest worker. We conclude that the pattern instances are unevenly distributed over the ranges of vertice ids which leads to high skew in a range based solution. We stopped our investigation in this direction.

### 3.2 Work-stealing

Normally, Spark uses static, physical partitioning of the data. As we learned in the last section that can lead to a trade-off between the ability to handle skew and duplicated work. A standard approach to handle skew and unbalanced workloads is work-stealing. For this, the work is not statically partitioned before-hand but organized in many smaller tasks which can be solved by all workers. Workers are either assigned an equal split of tasks and steal tasks from other workers when they are out of work or all tasks are arranged in a queue accessible for all workers, so that workers can poll tasks from it whenever they are out of work. In either way, this results in a situation where no task is guaranteed to be solved by a single worker and each worker only finishes when no free tasks are left in the system. Hence, the maximum amount of skew roughly the size of the smallest task. There is no duplicated work because different tasks should not overlap and each task exists only once in the whole system.

This leaves us with two design decisions. First, how to organize the workload of a *LeapfrogTriejoin* into tasks. Second, how do workers get their tasks? We address these questions in order by first describing our preferred solution and then their integration with Spark. We conclude the section with an evaluation of the limitations of this integration.

Before, we remember that we built our solution for the Spark *local-mode* (see section 2.1.2). Hence, the Spark master and all executors are threads within the same JVM process. This is important because it allows us to share data structures between multiple Spark tasks as normal JVM objects. We discuss how the design can be extended Spark's *cluster-mode* in future work (??) and related work (??).

The first design decision is the definition of a work-stealing task. It is not necessary to define the tasks such that they have all the same size. However, it is important to choose the task size small enough to avoid skew. Furthermore, the tasks should not overlap so that work is not duplicated. We choose to define a task as the work necessary to find all possible tuples for a

single binding of the first join variable. This is none overlapping. The task size can vary widely and is query dependent. However, given the huge amount of tasks (as many as vertices in the graph), we believe this to be small enough. We plan to evaluate this during our experiments.

The second design choice in work-stealing is how to hand tasks to workers. For simplicity, we choose to use a shared, thread-safe queue that holds all tasks. The main drawback of this solution is that the access to the queue has to be synchronized between all workers. If there are too many workers contending for the critical section of polling a job from the queue, they can slow each other down. However, the critical section is short because it includes only the call to the poll method of the queue. Additionally, we decide to implement a batching scheme such that a single poll can assign multiple tasks to a worker. This allows us to fine-tune how often a worker needs to return to the queue for new tasks.

It turns out that the work-stealing scheme as described above is straightforward to integrate into Spark. We choose a Scala object<sup>6</sup> to hold a dictionary which associates an ID for each query with a thread-safe queue instance. This queue can be accessed by each Spark thread. Due to the association between query and queue, it is possible to run multiple queries in parallel without interference.

The queue for a query is filled by the first Spark task that accesses it. This can be implemented by a short synchronized code section at the beginning of all tasks. It checks if the queue is empty and if so pushes one task per possible binding (all graph vertices) or batch of possible bindings. The synchronized section is fast and only runs once when the tasks start. Hence, it comes at neglectable performance costs.

Once the queue is filled, we run our normal *LeapfrogTrieJoin* with filtered *Leapfrog join* for the first attribute. This filter is implemented as a decorator around the original *Leapfrog join*. The *leapfrogNext* method of this decorator returns only values that are polled from the work-stealing queue before.

Our integration of work-stealing in Spark comes with some limitations. We see it more as a proof-of-concept that work-stealing is a good choice for the parallelization of worst-case optimal joins in Spark than as a solid implementation of work-stealing in Spark. The later is not possible within the time-frame of this thesis. In the following, we discuss the constraints of our integration.

Work-stealing leads to an unforeseeable partitioning of the results: it is not possible to foresee which bindings end up in a certain partition nor can we guarantee a specific partition size. If the user relies on any specific partitioning, he needs to repartition the results after. Moreover, we cannot guarantee to construct an equal partitioning over multiple runs of the same query. If the user depends on a stable partitioning per query, he should cache the query after the worst-case optimal join execution.

We do not integrate our work-stealing scheme into the Spark scheduler but we provide a best-effort implementation because we use all resources assigned to us as soon as they are assigned to us. We can handle all scheduler decisions. The first worker assigned fills the queue. The worker who takes the last element from the queue sets a boolean that this query has been completed. Hence, tasks that are started after the query has been computed, do not recompute the query.

We do not provide a fault-tolerant system. We see two possibilities to make our system fault-tolerant. First, one can stop all tasks if a single task fails and restart the computation with the last cached results before the worst-case optimal join. Second, one could extend the critical section of polling a queue value by the *LeapfrogJoin* by two more operations: we peek at the value from the queue without removing it, log the value in a set of values per task and then poll it and remove it from the queue. With these operations, it is guaranteed that the master can reconstruct all values that a failed worker thread considered. So, after a task failure, the master

---

<sup>6</sup>Methods and fields defined on a Scala object are the Scala equivalent to static methods and fields in Java. Most importantly they are shared between all threads of the same JVM.

can add these values to the work queue again such that other tasks will redo the computations.

## 4 GraphWCOJ

### 4.0.1 Combining LFTJ with CSR

For our graph pattern matching specialized *LeapfrogTriejoin* version we choose CSR (see section 2.5) as backing data structure. This data structure is typically used for static graphs and we show that it is a great match for LFTJ. In this section, we shortly describe the implementation of a CSR based *TrieIterator*, point out the differences between this new version and a column based *TrieIterator* (as described in section 5.1<sup>7</sup>) and conclude with an experiment demonstrating the power of this optimization.

The implementation of a CSR based *TrieIterator* is straightforward except for one design change: instead of using the vertice identifier from the graph directly, we use their indices in the CSR representation. This change is rather minor because it can be contained at any level by using a hash map for translation, e.g. in the *TrieIterator* itself, in the *LeapfrogTriejoin* or at the end of the query by an additional mapping operation. In the current system, the translation is performed by the LFTJ implementation to allow easy integration into other projects. However, it is possible to work on the indices throughout the whole system to save the translation costs. We now outline how to implement each of the *TrieIterator* methods, under the assumption that all vertices have outgoing edges. Then, we drop this assumption and explain the necessary changes. The creation of the CSR data structure itself is described in ??<sup>8</sup>.

We use the variables *depth*, *srcPosition* and *dstPosition* to store if we are operating on the source or destination level and the positions of the iterator on the respective level. We call the two arrays of the CSR datastructure *dst* for the array storing all edge destinations and *dstIndices* for the array storing indices into *dst* (refer to ?? for a complete explanation of these arrays). The *open* function does nothing if we open the first level; if we open the second level, it sets *dstPosition* to *dstIndices[srcPosition]*. Additionally, it updates *depth*. The *up* function only updates *depth*. The *next* method increases *srcPosition* or *dstPosition* by one depending on *depth*. The *seek(key)* method sets *srcPosition* to *key* if *depth* equals 0 or uses binary search to find *key* in *dst.slice(dstPosition, dstIndices[srcPosition + 1])* and then sets *dstPosition* accordingly. The *key* function returns *srcPosition* on the first level and *dst[dstPosition]* if *depth* equals 1. The *atEnd* function is: *if (depth == 0) srcPosition == dstIndices.length - 1 else dstPosition == dstIndices[srcPosition + 1]*.

<Alternative representation of the *TrieIterator* implementation as table. Let me know what you prefer. I think the text is better after seeing both.>

---

<sup>7</sup>To be rewritten and integrated

<sup>8</sup>To be written.

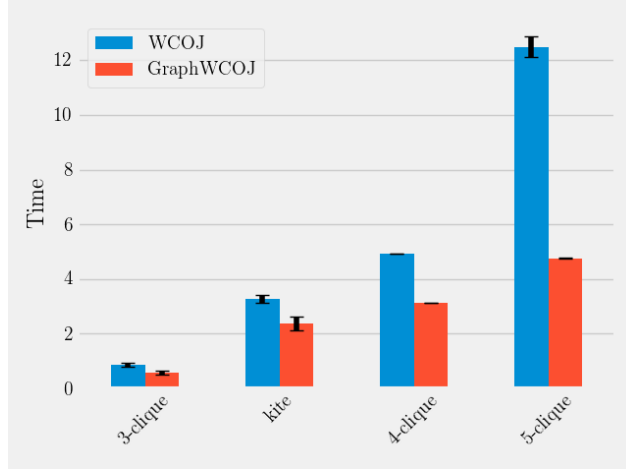


Figure 8: Barchart comparing join time of WCOJ and GraphWCOJ on multiple queries on SNB-sf1

Method	1st level	2nd level
open	update depth	$\text{dstPosition} \leftarrow \text{dstIndices}[\text{srcPosition}]; \text{update depth}$
up	update depth	update depth
next	$\text{srcPosition}++$	$\text{dstPosition}++$
seek(k)	$\text{srcPosition} \leftarrow k$	$\text{dstPosition} \leftarrow \text{lub}(k, \text{dst}[\text{dstPosition}:\text{dstIndices}[\text{srcPosition} + 1]])$
key	$\text{srcPosition}$	$\text{dst}[\text{dstPosition}]$
atEnd	$\text{srcPosition} = \text{dstIndices.length} - 1$	$\text{dstPosition} = \text{dstIndices}(\text{srcPosition} + 1)$

Table 5: Tabular summary of *TrieIterator* implementation based on CSR. *lub* is *leastUpperBound*.  $\text{dst}[\langle \text{start} \rangle : \langle \text{end} \rangle]$  slices the array.

To resolve the assumption of no empty outgoing adjacency lists, we adapt *open*, *next* and *seek* to skip source positions without outgoing edges. This is easy to detect because then  $\text{dstIndices}(x) == \text{dstIndices}(x + 1)$ . We can skip these cases by simple linear search until we find a valid position. This solution is sufficient because there are only a few vertices with no outgoing edges in real-world graphs.

The *TrieIterator* implementation based on CSR is much faster than the column based iterator; mainly due to the fact that the *seek* method on the first level can be implemented in  $\mathcal{O}(1)$ , instead of  $\mathcal{O}(\log n)$ . This optimization has huge potential because these searches are the most costly operations for a column based *TrieIterator* [16]. Note that searches on the second level are fast, due to the fact that most graphs have a low outdegree (see section 2.6). Additionally to this advantage, CSR based *TrieIterator* do less bookkeeping because they support only 2 levels and spent nearly no time on processing *atEnd* for the second level, while a column based *TrieIterator* needs to calculate the number of outgoing edges for each source vertex in its *open* method, to allow a fast *atEnd* method).

We conclude that CSR based *TrieIterator*’s are a promising match for LFTJ and graph pattern matching. The improvements of this optimization can be seen in fig. 8. It demonstrates an up to 2.6 speedup over a column based LFTJ. We also see that the optimization has a stronger impact on queries with more edges and vertices, e.g. 5-clique. For a more thorough evaluation refer to the experiment section 6.

#### 4.0.2 Exploiting low average outdegrees

In our study of real-world graphs (??), we show that most real-world graphs have a small, average outdegree. The outdegree over all graphs is TODO and the maximum average outdegree of a single graph is TODO at TODO. These results lead to the hypothesis that the intersection of multiple adjacency lists is small, e.g. below 10 in many cases. We can exploit this fact by materializing the intersections in the *Leapfrog joins* directly in one go; instead of, generating one value at-the-time in an iterator like fashion as described in the original paper [leapfrog]. This is beneficial because it makes better use of data locality; we elaborate this statement in a later paragraph.

We structure the remainder of this section as follows. First, we shortly reiterate the most important facts about *Leapfrog joins* for this chapter (TODO just organize both sections behind each other?). Second, we analysis the intersection workload in terms of input sizes and result size to confirm our hypothesis and gain valuable insights to choose the best intersection algorithm. Third, we explain the algorithm we chose based on the analysis. Fourth, we point out differences to the original Leapfrog Triejoin. Finally, we present a short experiment showing the performance gains of this optimization.

*Leapfrog joins* build the intersection between multiple adjacency lists. This is done in an iterator-like fashion in their *leapfrog\_next* method by repeatedly finding the upper-bound for the largest value in the lowest iterator. This algorithm is asymptotically optimal for the problem of n-way intersections. However, we claim that it is (1) to complex for small intersections and (2) should generate all values at once instead of one-by-one to improve performance on real-world adjacency lists.

To determine the best algorithms to build the n-way intersection in the *Leapfrog joins*, we run some experiments to characterize the workload. Towards this goal, we log the size of the full intersection, the size of the smallest iterator participating and the size of the largest intersection between the smallest iterator and any other iterator on 5-clique queries on **SNB-sf-1**. Figure 9 depicts these metrics as cumulative histograms. In the next paragraphs, we point out the most important observations in each of these graphs.

Figure 9c shows the size distribution of the smallest iterator, as to be expected for a social network graph, the outdegree is between 1 and 200. We do not see the long-tail distribution typical for power-law graphs because we choose the smallest iterator out of 5 and even though there are vertices with a much higher outdegree, the chance of encountering 5 of these in a single intersection is small. We note that in 80% of all cases the smallest iterator has a size lower than 80 and above that the distribution slowly increases to 100%

Figure 9b illustrates the size distribution of intersecting the smallest iterator with any other iterator, such that the intersection is maximal. We choose this specific metric to motivate one of our design choices later on. As for the smallest iterator, some of these intersections are as big as 200 but most of them are much smaller. However, unlike for the smallest iterator metric, 80% of the intersections contain less than 21 elements and the frequency increases to 100% in a steep curve. This last observation is even stronger for the size of the total intersection (fig. 9a): the size is less than 5 in 80% of all intersections and increases similarly steep to 100%. The maximum is a little lower than 200.

These observations confirm our hypothesis that the size of the intersections is small (below 5) and do not show the same long-tail distribution as the whole social network graph. Hence, we can materialize them without running at risk of building big intermediary results. Furthermore, the experiment shows that optimizing by taking iterator sizes into account is worthwhile but only for the smallest iterator because once we start with the smallest iterator the further intersections are small (below 21) in the vast majority of all instances.

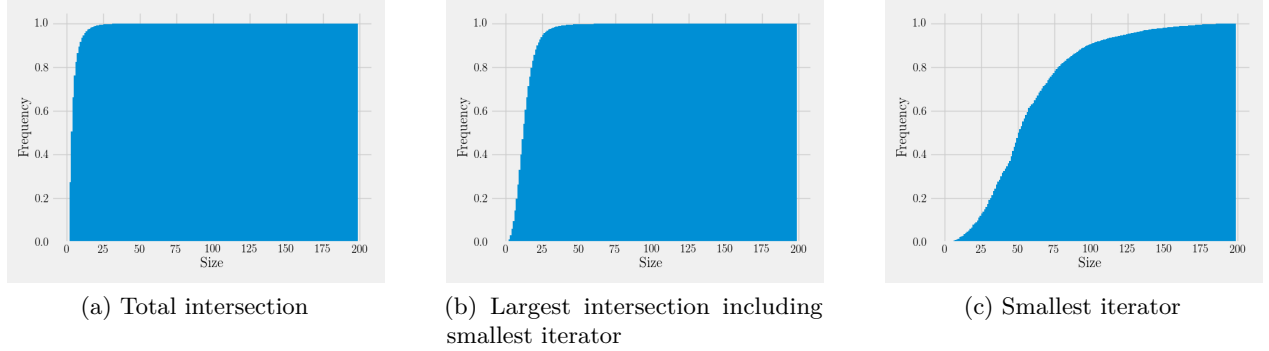


Figure 9: Cumulative histograms of total intersection sizes, largest intersection with the smallest iterator and any other, and size of the smallest iterator participating in 5-clique on SNB-sf-1.

In the coming paragraphs, we detail how to build the  $n$ -way intersection of multiple adjacency list such that we gain performance by better use of data-locality than original the *Leapfrog join*. We choose to use pairwise intersections over multi-way intersection algorithms for their simpler, linear memory access patterns.

From our analysis, we conclude that the final intersection size is strongly dependent on the smallest iterator and that the intersection of the smallest iterator with any other iterator is close to the final size. These insights translate into two design decisions. First, we start with the smallest iterator<sup>9</sup>. However, we do not take the sizes of any other iterators into account because the effort for sorting the iterators by size would not pay off.

Second, we use two different tactics to build the pairwise intersections. The first intersection between two iterators is built *in-tandem*, where we seek the upper bound of the higher value in the smaller iterator. This algorithm guarantees to find the intersection between two iterators in the asymptotical fastest way. After this first intersection, the intermediary result is quite small. Therefore, we use the simpler scheme of linearly iterating the intermediary and probing the iterator by binary search with fall-back to linear search.

Finally, we point out a few exceptional cases and pitfalls for implementors:

- If all iterators of the *Leapfrog join* are on their first level, the intersection is near  $|V|$ . In this case, we fall back to the original *Leapfrog join*.
- We use an array to materialize the intersections because Scala collections are slow. Instead of deleting elements, we replace them with a special value.
- Allocating a new array for every *Leapfrog join* initialization is costly. We estimate the size of the intersection by the size of the smallest iterator and reuse the array whenever possible. We use a sentry element to mark the end of the array.

The carefully chosen operations explained above are faster than the original *Leapfrog join* algorithm for two reasons. First, although, the original algorithm uses an asymptotical better  $n$ -way intersection, multiple binary intersections are preferable with the rather small adjacency lists of graph workloads. In particular, in our situation where the size of the first binary intersection is already very close to the size of the final intersection.

Second, the *Leapfrog join* generates a single value, then yields control to other parts of the *LeapfrogTriejoin* algorithm and later touches the same adjacency lists again to generate the next value. Our approach touches the adjacency lists exactly once per *Leapfrog* initialization and

---

<sup>9</sup>We take advantage of the fact that CSR allows us to determine the size of iterators cheaply (see ??)

condenses the intersection into a much smaller array. This array is more likely to stay cached while the other parts of the *LeapfrogTriejoin* do their work.

Figure 10 shows that a materialized *Leapfrog join* out-performance better than the original algorithm. As expected, the optimization is more powerful for bigger queries because they work with more adjacency lists per *Leapfrog join*, e.g. for a triangle each *Leapfrog join* intersects only two iterators while for a 5-clique each join handles 5 adjacency lists. Anyhow, even for the triangle query, we see a small but clear improvement and a lower error due to better cache use. We refer the reader to our experiment section (6) for further experiments on our graph specialized WCOJ and detailed descriptions of the datasets and queries used.

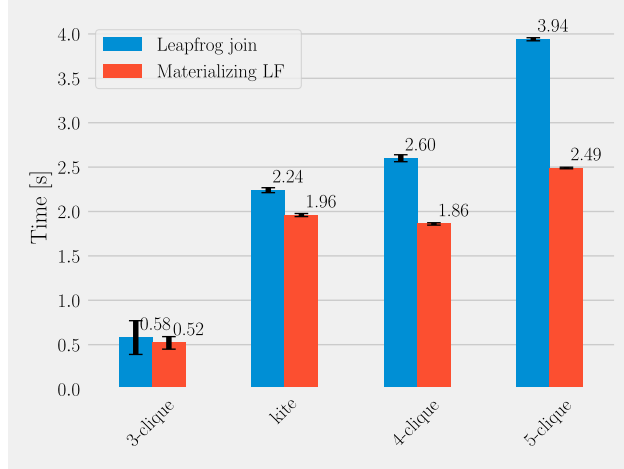


Figure 10: Barchart showing GraphWCOJ with and without *Leapfrog join* materialization enabled for different queries on SNB-sf1.

## 5 Implementation

### 5.1 General sequential version (*seq*)

#### 5.1.1 Optimizations

A simple, idiomatic Scala implementation of the Tributary join is not able to beat Spark’s **BroadcastHashjoin** on any other query than the triangle query. Hence, we spent roughly 2 weeks to optimize our first implementation. After, we are able to beat Spark’s **BroadcastHashjoin** on nearly all queries and datasets. In this section, we discuss the implemented optimization and give a rough estimate of how important each of these is. In total, we improved the WCOJ running time from 248.2 seconds to 44.5 seconds on the unfiltered 5-clique query on the **Amazon-0602** dataset. We list all optimizations in table 6 and label them ‘very important’, ‘important’ and ‘minor’ based on the performance improvement directly after applying it.

It is not helpful to give more detailed information on the effect of single optimization because they are not independent of each other. Hence, they might have a hugely different effect when applied in a different order, e.g. we first applied an optimization to the binary search and then optimized the *LeapfrogJoin.next* method to avoid many searches. Hence, giving detailed runtime measurements for the binary search optimization would overestimate its value. It is out of the scope of this work to study the dependency and order of the optimization to gain correct runtime measurements.

Category	Optimization	Impact
<b>LFTJ</b>	<i>LeapfrogJoin.init</i> avoid sorting iterators	very important
	<i>ArrayTrieIterable.next</i> in $\mathcal{O}(1)$ for deepest level	very important
<b>Binary search</b>	linear search for short search spaces	important
	avoiding unnecessary conditions	important
<b>Spark</b>	direct use of arrays instead of <code>ColumnVector</code>	important
<b>Scala</b>	use <i>while</i> instead of <i>map</i> , <i>foreach</i> , <i>exists</i> , etc	very important
	use <i>Array</i> instead of Scala’s collections	very important
	use of <i>private[this]</i>	minor
	enable compiler optimization	minor
<b>General</b>	remove array lookups from the critical path $column(depth)(position) \rightarrow currentColumn(position)$	very important
	use <i>Array</i> instead of <i>Map</i> if keys are integers and dense	important
	strength reduction $(i + 1) \% 5 \rightarrow \text{if } (i == 4) 0 \text{ else } i + 1$	important

Table 6: Summary of all optimizations used for *seq* and an estimate of their impact.

We discuss the optimization in categories: Leapfrog Triejoin specific, binary search specific, Spark related, Scala related and general. We conclude the section with some changes we tried that do not improve performance.

Binary search specific optimizations become a category on its own because the sorted search is the most expensive operation in the Tributary join. According to profiler sessions, the join spends more than 70% of its time in this method. This result is in line with the observation that ‘in the Tributary join algorithm, the most expensive step is the binary search’ from [16].

We applied two Tributary join specific optimizations. The first in the class `leapfrogTriejoin.LeapfrogJoin` (see also ??) and the second in the `leapfrogTriejoin.ArrayTrieIterable.TrieIteratorImpl`.

The *LeapfrogJoin.init* method is originally described in [leapfrog] to sort its *TrieIterators*. However, the method can be improved by avoiding to sort the *TrieIterators* (line 11). We can start moving the *TrieIterator* without sorting them and arrive at an ordered array in  $\mathcal{O}(n)$  steps -  $n$  defined as the size of *iterators*. This approach improves over the original algorithm in two ways: (1) it starts moving the *TrieIterators* to their next intersection immediately without sorting them first and (2) orders the array in fewer steps than traditional sorting algorithms.

To implement this we find the maximum *key* value in all iterators and store the index of this *TrieIterator* in  $p$ . Then we move the *TrieIterator* at  $p + 1$  to the least upper bound of this *max* (by calling *seek*) and store the result as the new maximum. We proceed with this process - wrapping  $p$  around when it reaches *iterators.length* - until  $p$  equals the original maximum index. Now, we are either in a state in which all *TrieIterators* point to the same value and we are done - the *LeapfrogJoin* is initialized - or we arrived at a state in which the *iterators* array is sorted according to *key* and can proceed as in the original *LeapfrogJoin.init* method. To apply this optimization one replaces the call to *sort* in line 11 with the procedure explained above<sup>10</sup>.

The second Leapfrog Triejoin specific optimization is to change the `ArrayTrieIterable.TrieIteratorImpl.next` method. This method moves the iterator to the next value on the same level of the trie. Hence, it generally runs in  $\mathcal{O}(\log n)$ ,  $n$  being the number of tuples in the relationship, because it needs

<sup>10</sup>The implementation of Scala’s array sort for objects is slow because it copies the array twice and casts the values to *Java.Object* such that it can use Java’s sorting methods. Before we applied the sorting optimization above, we replaced Scala’s sort method with an optimized insertion sort, which was faster than Scala’s sorting method - the *iterator* array contains normally at most 20 items.



to find the least upper bound of  $key + 1$ . However, under the assumption that all tuples are unique - which is fulfilled for the use-case of an edge relationship - the last level of the trie is unique. Hence, we can move to the next value by simply increasing the position by one, which is an operation in  $\mathcal{O}(1)$ .

The binary search is the most expensive operation of the Leapfrog Triejoin. Hence, special attention needs to be paid while implementing it. Our most important optimization is to change to a linear search once we narrowed the search space to a certain threshold - currently at 60 values. We experimented with values from 0 to 400 and found that 60 was optimal but even going as high as 120 values would not change the performance much.

Another important optimization is to avoid unnecessary if-statements in the loop of the binary search, e.g. the implementation on Wikipedia and many other example implementations use an if-statement with three branches for smaller, bigger and equal but two branches for greater than and less-or-equal suffice for a least upper bound search.

A similar optimization can be applied to a linear search on a sorted array: intuitively one would use the while-loop condition  $array(i) > key \wedge i < end$  with  $key$  being the key to find the least upper bound for,  $i$  the loop invariant and  $end$  the exclusive end of the search space. Anyhow, it is faster to check for  $key > array(end - 1)$  once before the loop and return if this is the case because the value cannot be found in the search space. This obviously circumvents the main loop of the linear search; additionally, it simplifies the loop condition to  $array(i) > key$ .

The Spark infrastructure uses the interface `ColumnVector` to represent columns of relationships. The implementation `OnHeapColumnVector` is a simple wrapper around an array of the correct type with support for *null* values and *append* operations. First, we used this data structure to represent our columns but we could see a clear increase in performance by replacing it by an implementation that exposes the array to allow the binary search to run on the array directly. This is likely due to saving virtual function calls in the hottest part of our code. The implementation is straightforward and can be found in our repository in `leapfrogTriejoin.ExposedArrayColumnVector`; we implemented it only for the `Long` datatype.

We found many standard optimizations and Scala specific optimizations to be really useful. Most likely these are the optimizations that brought the biggest performance improvements. However, they are well-known, so we mention them only in the table 6. For Scala specific optimizations one can find good explanations at [17].

Apart from the aforementioned very useful optimizations, we investigated multiple other avenues in hope for performance improvements which did not succeed, we list these approaches here to save others the work of investigating:

- reimplement in Java
- use of a Galloping search before the binary search
- unrolling the while-loop in *LeapfrogTriejoin.moveToNextTuple*
- predicating the *action* variable in *LeapfrogTriejoin.moveToNextTuple*

Finally, we believe that code generation for specific queries that combines the functionality of *LeapfrogTriejoin*, *LeapfrogJoin* and *ArrayTrieIterator* into one query specific function would lead to noticeable performance improvements. The reason for this belief is that our implementation takes about 3.46 seconds for a triangle query on the Twitter social circle dataset while a triangle query specific Julia implementation, of a colleague of ours, needs only half a second. The main difference between our implementation and his are: the language used (Julia is a high-performance, compiled language) and the fact that his implementation has no query interpretation overhead but cannot handle any other query than the triangle query.

However, a code generated Leapfrog Triejoin is out of scope for this thesis, also, we are aware of efforts by RelationalAI to write a paper about this specific topic. We are looking forward to

seeing their results.

## 6 Experiments

TODO introduction

### 6.1 Setup

#### 6.1.1 Hardware and Software

We run our experiments on machines of the type `diamond` of the Scilens cluster owned by the CWI Database Architecture research group. These machines feature 4 Intel Xeon E5-4657Lv2 processors with 12 cores each and hyperthreading of 2 (48 cores / 96 threads) Each core has 32 KB of 1st level cache, 32KB 2nd level cache. The 3rd level cache are 30 MB shared between 12 cores. The main memory consists of 1 TB of RAM DDR-3 memory.

The machines run a Fedora version 30 Linux system with the 5.0.17-300.fc30.x86\_64 kernel. We use Spark 2.4.0 with Scala 2.11.12 on Java openJDK 1.8. In the majority of our experiments, we use Spark in its standard configuration with enabled code generation. We also tune the parameters for driver and executor memory usage (`spark.driver.memory` and `spark.executor.memory`) to fit all necessary data into main memory.

#### 6.1.2 Algorithms

In our experiments we use 4 different join algorithms. Two of them are worst-case optimal joins. That is our Leapfrog Triejoin implementation, *LFTJ*, and a graph-pattern matching specialized Leapfrog Triejoin developed in this thesis: *GraphWCOJ*. *LFTJ* is only run as sequential algorithm as a baseline against *GraphWCOJ*. We compare these to algorithms in ??.

The other two algorithms are Spark’s versions of *BroadcastHashJoin* and *SortmergeJoin*. We compare them against the sequential version of *LFTJ* and *GraphWCOJ* in ?? and their scaling with *GraphWCOJ* in ?. We adjust the "`spark.sql.autoBroadcastJoinThreshold`" parameter to control if Spark is using a *BroadcastHashJoin* or a *SortMergeJoin*.

#### 6.1.3 Datasets

We run the majority of our experiments on two datasets from different use-cases, social networks and product co-purchase. We motivate our choice in the next paragraph. Table 7 includes a list of all graph datasets mentioned throughout the thesis.

TODO add Vertices and edges numbers

The SNB benchmark [29] generates data emulating the posts, messages and friendships in a social network. For our experiments, we only use the friendships relationship (`person_knows_person.csv`) which is an undirected relationship. After generation only edges of the kind  $src < dst$  exist, we generate the opposing edges before loading the dataset, such that the edge table becomes truly undirected. The benchmark comes with an extensively parameterizable graph generation engine which allows us to experiment with sizes as small as 1GB and up to 1TB for big experiments and different levels of selectivity. The different sizes are called scale-factor or `sf`, e.g. `SNB-sf1` refers to a Social network benchmark dataset generated with default parameters

Name	Variant	Vertices	Edges	Source
<b>SNB</b>	sf1		453,032	[29]
<b>Amazon</b>	0302	262,111	1,234,877	[30]
	0601	403,394	3,387,388	[30]
<b>Twitter</b>	sc-d	81,306	1,768,135	[30]
	sc-u	TODO	TODO	[30]
<b>LiveJournal</b>		4,847,571	68,993,773	[30]
<b>Orkut</b>		3,072,441	117,185,083	[30]

Table 7: A summary of all datasets mentioned in the thesis. Explanation of them and for the variants is given in running text.

and scale-factor 1. We include the exact parameter used for generation in our repository under `experiments/snb/params.txt`.

The Amazon co-purchasing network contains edges between products that have been purchased together and hence are closely related to each other [30]. This is a directed relationship from the product purchased first to the product purchased second, both directions of an edge can exist if the order in which products have been purchased varies. The Snap dataset collection contains multiple Amazon co-purchase datasets, each of them containing a single day of purchases. We choose the smallest and biggest dataset from the 2nd of March and the 1st of June 2003, we call them **Amazon-0302** and **Amazon-0601**. We pick co-purchase datasets for evaluation because former work often concentrated on social networks and web crawl based graphs [16, 7] but [40] points out that the biggest graphs are actually graphs like the aforementioned Amazon graph containing purchase information.

To allow comparisons with former work, we run a subset of our experiments on the Twitter social circle network from [30]. This dataset includes the follower relationship of one thousand Twitter users; each of these follows 10 to 4.964 other users and relationships between these are included. The graph is originally directed but for some experiments, we add reversed edges to make the graph undirected - again for comparison with former work. We call this graph **Twitter-sc-d** and **Twitter-sc-u** for the directed respectively undirected variant.

The *LiveJournal* graph represents the friendship relationship of a medium sized social network.

#### 6.1.4 Graph patterns

In this section, we detail the graph patterns used throughout our experiments. Most of the queries are cyclic because that has been shown to be the primary use-case for WCOJ in former research [38, 16]. WCOJ’s also have been successfully applied to selective path queries in [38]; however, this result have not been reproduced by any other paper.

To most of our queries, we apply a filter to make them more realistic, e.g. a clique query does make more sense if it is combined with a smaller-than filter, which requires that the attributes are bound such that  $a$  smaller than  $b$ , smaller than  $c$ . Because otherwise, one gets the same clique in all possible orders in the output, which not only takes much more time but is also most likely not the result a user would want. We ensure that filters can be pushed down through or in the join by Spark as well as by the WCOJ to compare both algorithms on an equal basis. A complete list of all queries and filters used is shown in table 8. The less known queries are also detailed in text. Patterns and filters might be used in all possible combinations, we name the resulting query  $\langle pattern \rangle - \langle filter \rangle$ , e.g. *triangle-lt*.

Name	Parameters	Vertices	Edges	Example pattern
<b>triangle</b>	NA	3	3	$a \rightarrow b; a \rightarrow c; b \rightarrow c$
<b>n-clique</b>	# vertices	$n$	$1/2 \times n \times (n - 1)$	see above
<b>n-cycle</b>	# vertices	$n$	$n$	$a \rightarrow b; b \rightarrow c; c \rightarrow z; z \rightarrow a$
<b>n-s-path</b>	# edges / selectivity	$n$	$n - 1$	$a \rightarrow b; b \rightarrow c; c \rightarrow z$
<b>house</b>	NA	5	9	$a \rightarrow b; a \rightarrow c; a \rightarrow d; b \rightarrow c; b \rightarrow d; c \rightarrow d; c \rightarrow e; d \rightarrow e$
<b>Filters</b>				
<b>distinct</b>				$a \neq b; a \neq c; a \neq d; b \neq c; \dots$
<b>lt</b>				$a < b; b < c; c < d; \dots$

Table 8: Summary of patterns and filters used.

For a selective path query, we first select two sets of nodes with respect to the *selectivity* parameter. Then we search for all path of a certain length according to the *edges* parameter, e.g. **4-0.1-path** finds all paths between two randomly selected, fixed sets of vertices of length 4 - the sets of nodes contain roughly 10% of all input nodes and are not guaranteed to be intersection free.

## 6.2 Baseline: BroadcastHashJoin vs seq

In this experiment, we compare the runtime of our sequential Leapfrog Triejoin implementation, **seq**, with the runtime of Spark’s **BroadcastHashjoin**. Towards, this goal we ran all queries from table 8 on our three main datasets: **ama-0302**, **ama-0601** and **snb-sf1**. The clique patterns are combined with the less-than filter, **n-clique-lt**, and the cycle pattern with the distinct filter, **n-cycle-distinct**. These seem to be the most realistic setups because cliques are fully symmetric and one wants to avoid redundant results. For cycles, the less-than filter is too restrictive because it excludes cycles for which  $a < b > c$ . We show our results in table 9 and in barcharts (fig. 11, 12 and 13). Section 6.2.2 analyzes the results.

Our experiment measures the time it takes to perform a **count** on the cached dataset using **BroadcastHashjoin** and **seq**. For **BroadcastHashjoin**, the time to run the whole query is reported. For **seq**, we report setup time and the time, it takes to run the join, separately. Setup time includes the sorting, materialization and copying the results of our join from a Scala **Array** into the **UnsafeInternalRow** format expected by Spark. TODO sorting not yet (data is presorted in files) This section is focused on comparing the runtimes excluding the setup time - rational given in section 6.2.1.

### 6.2.1 Experiment Rationale

**Question:** Why do we compare against Spark’s **BroadcastHashjoin** instead of **SortMergeJoin**?

**Answer:** Because even when all data is arranged in a single partition, for simple sequential processing, Spark schedules its **SortMergeJoin** to use a shuffle. A shuffle writes and reads data to and from disk. Hence, **SortMergeJoin** is much slower than a **BroadcastHashJoin**. We compared the algorithms on the **Amazon-0601** dataset for the **triangle** (8.1 seconds vs 58.9 seconds) and **5-clique** pattern (32.9 seconds vs 850.9 seconds). We assume that Spark is able to optimize its broadcasts when **local[1]** is used to start the Spark session because then Spark uses the driver as executor.

**Question:** Why do we exclude setup times from the WCOJ times?

**Answer:** Because our final implementation `dist` is meant to cache the readily sorted and formatted edge tables and reuse it for multiple queries. We anticipate that this is necessary to benefit from WCOJ’s in general. Furthermore, we optimize the setup times in later implementation. Hence, the current setup code is much slower than the one we expect to use for later implementations.

**Question:** Why is the time to copy results into `UnsafeInternalRow` format for WCOJ counted as setup time?

**Answer:** It is time solely spent for integration with Spark and not Leapfrog Triejoin specific. Furthermore, it could be avoided by working directly on the `UnsafeInternalRow` format within our `seq` implementation. However, this would require us to work with unmanaged memory (the `UnsafeInternalRow` interface is slower than working on `Arrays`) and we deem this as an unnecessary engineering overhead.

**Question:** Is Spark’s code generation a huge advantage for the `BroadcastHashjoin`? **Answer:** Yes, we ran Spark without code generation for comparison on the `Amazon-0302` dataset for `triangle` and `5-clique`: with code generation Spark takes 3.1 and 4.2 seconds without 14 and 16.

### 6.2.2 Analysis

For now, we settle to simply point out the most important observations and postpone deeper analysis, e.g. influence of dataset size and characteristics, to experiments run against our later implementations, i.e. `seq-graph-pattern` and `dist`.

We are able to beat Spark’s `BroadcastHashjoin` on all datasets and queries except `5-clique-lt` on `Amazon-0602`. Generally, we see that for `n-clique` patterns the speedup over Spark decreases for bigger  $n$ . This is due to the fact that many binary joins in a `n-clique` are actually semi-joins which do not increase but decrease the size of intermediary results, e.g. for `5-clique` on `Amazon-0302` only 3 out of 9 joins lead to a bigger intermediary result.

The cycle query results are highly interesting because we see an increasing speedup for higher  $n$  on `Amazon-0602` but a decreasing speedup on `Amazon-0302`. Unfortunately, we are (TODO currently) not able to provide `n-cycle` results for the `SNB-sf1` dataset, due to the fact that `BroadcastHashjoin`’s take more than 22 hours for the `6-cycle` which blocked our experiments. `5-cycle` runs at the moment.

The `House` and `5-clique` pattern seem to be quite similar - the `House` is a `5-clique` with two missing edges. However, as the count of their results indicates these two edges lead to dramatically different outcomes. Hence, their different timing and speedup behaviour.

The `Kite` pattern produces consistently the second highest speedup after the `3-clique`. Most likely due to the fact that a `Kite` is two triangles back-to-back.

The path query shows very different behaviour on the `Amazon` and the `SNB` datasets. This might be due to the different selectivity; it is extremely high on the co-purchase datasets and rather low on the social network benchmark. This different in selectivity is not surprising given that the `SNB` network fulfills the small world property, while the `Amazon` dataset relates products purchased together which naturally leads to multiple loosely connected, denser components. We will run the path queries with a different selectivity on the two input vertex sets to confirm this hypothesis.

Finally, we observe that all three datasets lead to quite different results which are most likely not comparable to each other without deeper research in the characteristics of the datasets themselves. In particular, it becomes clear that co-purchase datasets and social network datasets must have very different characteristics. Although, `SNB-sf1` is much smaller than `Amazon-0601`, queries on it take a similar or even much more time, e.g. `5-clique-lt` takes 14.21 seconds on

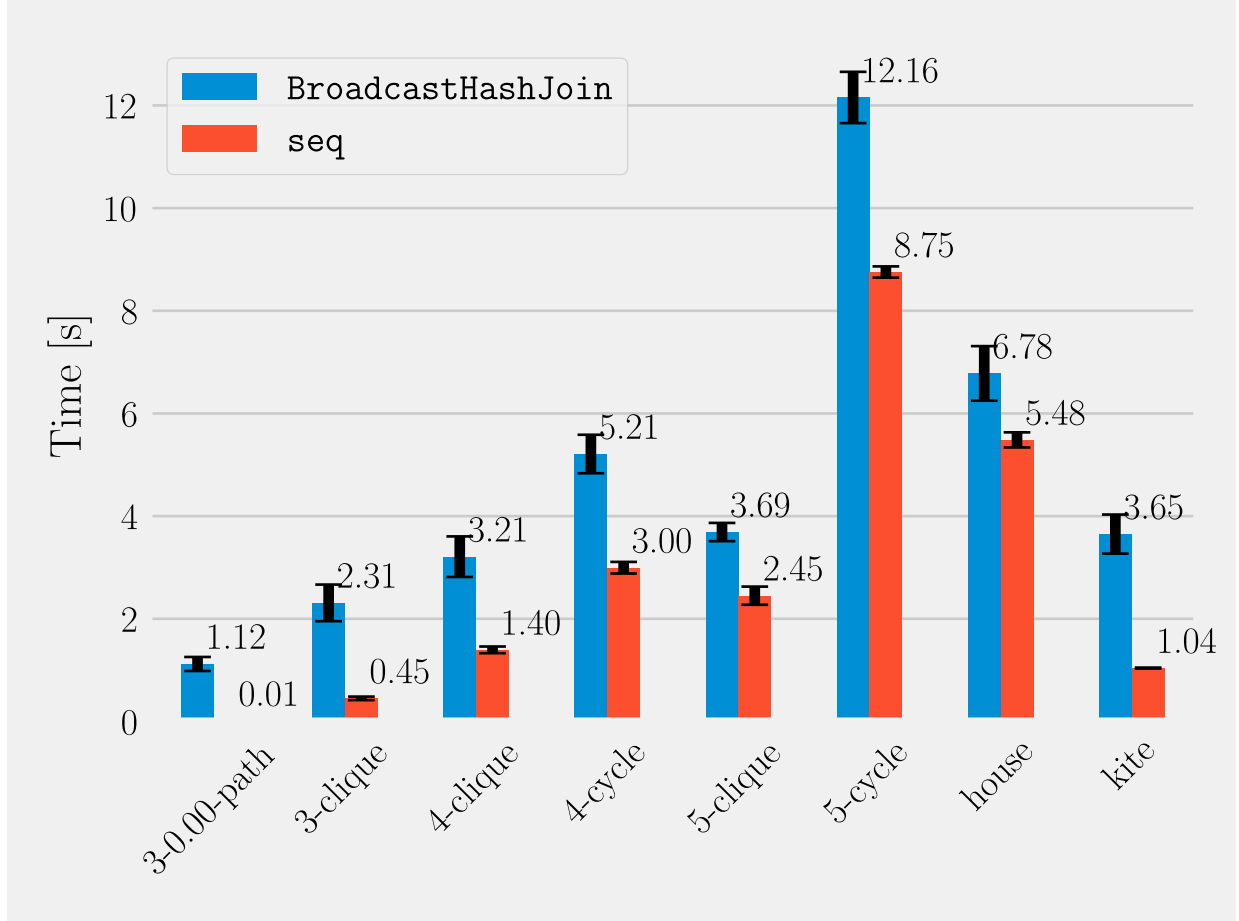


Figure 11: seq vs BroadcastHashJoin on Amazon-0302

the bigger dataset and 12.65 seconds smaller, even though, the result set is much smaller on **SNB-sf1**; **4-cycles-distinct** takes roughly 8 times longer on the small dataset and has a much bigger result set. In general, we see a higher speedup on **SNB-sf1**

### 6.3 Scaling of *GraphWCOJ*

In this section, we aim to analyse and compare the scaling of *GraphWCOJ* using different partitioning schemes. Towards this goal, we run *GraphWCOJ* on datasets of different size namely Twitter, LiveJournal and Orkut. We compare two partitioning schemes: Shares and *work-stealing*. These are the two most promising schemes identified in [10]. The experiment is performed on 3-clique and 5-clique. 3-clique is the smallest of our queries. Therefore, it is most difficult to scale. 5-clique takes much longer than 3-clique. Hence, it shows how query size influences the scaling. Also, it increases the job size for the *work-stealing* partitioning scheme.

#### 6.3.1 Results

We first describe our expectations of the experiment outcome. We assume that scaling improves with the dataset size. Hence, we should see the highest speedups for Orkut, then LiveJournal and the lowest speedups for Twitter. Also, we expect the scaling to improve with the query size. Both hypothesis are grounded in the fact that more work to distribute often leads to stronger

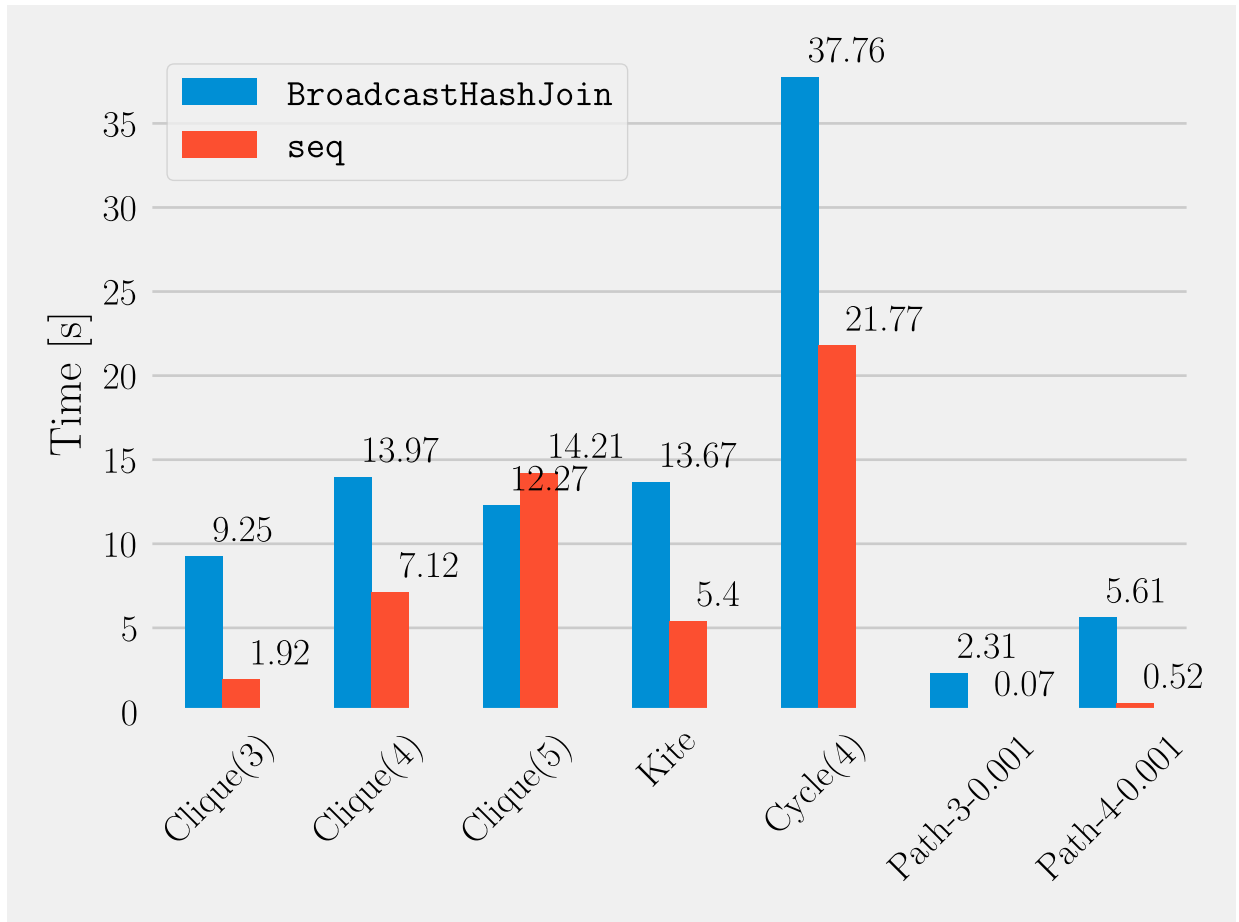


Figure 12: seq vs BroadcastHashJoin on Amazon-0601

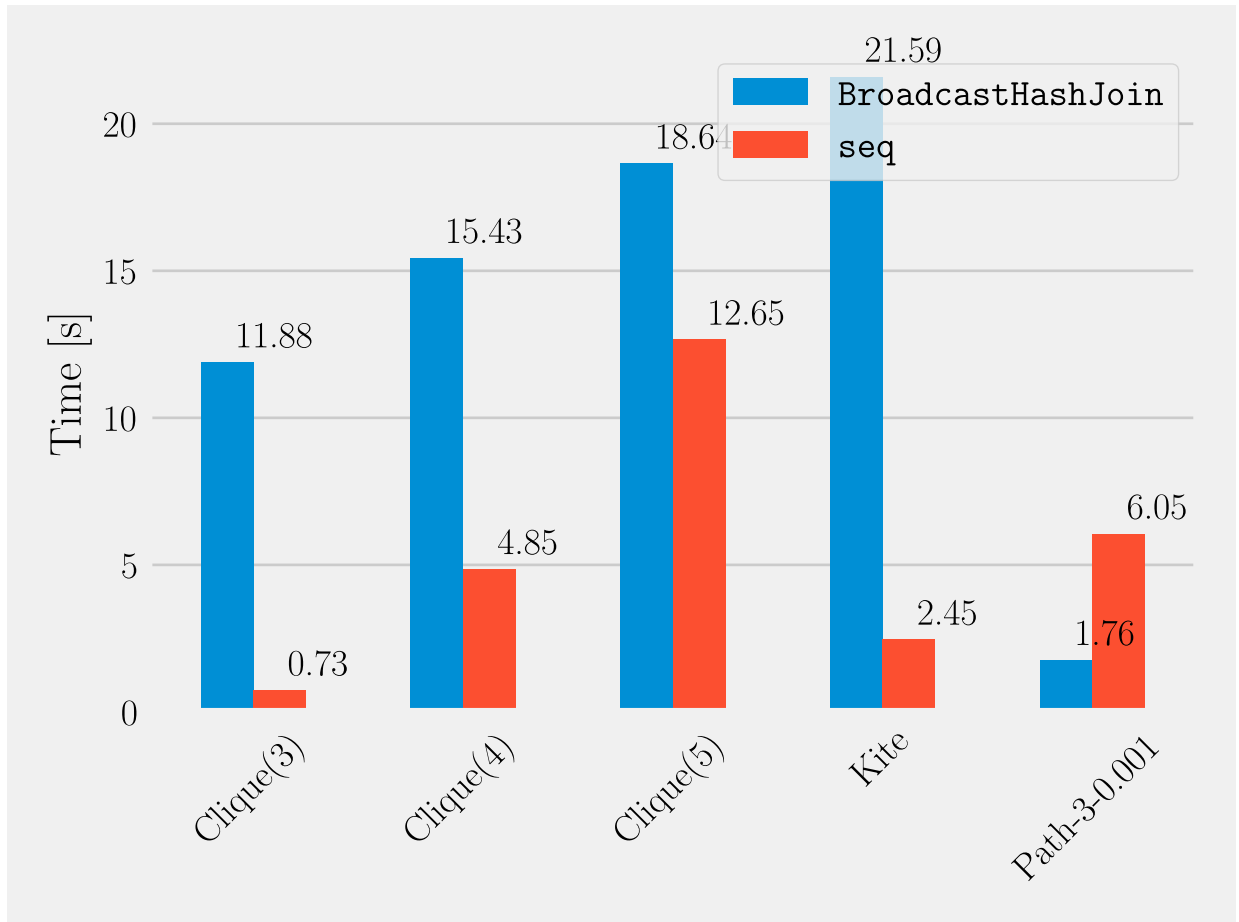


Figure 13: seq vs BroadcastHashJoin on SNB-sf1

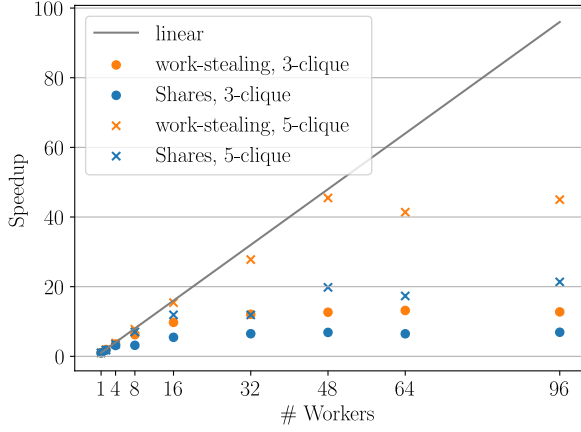


Query	# Result	BroadcastHashJoin	seq	setup	Speedup
Clique(3)	188,705	3.08	0.49	2.51	6.3
Clique(4)	47,824	3.85	1.25	4.55	3.1
Clique(5)	7,064	4.22	2.32	7.60	1.8
House	2,941,664	7.87	5.66	6.18	1.4
Kite	220,637	3.82	1.04	6.73	3.7
Cycle(4)	3,076,324	4.90	2.75	3.40	1.8
Cycle(5)	6,389,425	13.15	8.28	4.87	1.6
Cycle(6)	11,682,732	38.00	27.62	6.79	1.4
Path-3-0.001	2,092	1.36	0.01	0.89	136.0
Path-4-0.001	8,829	3.11	0.03	1.58	103.7

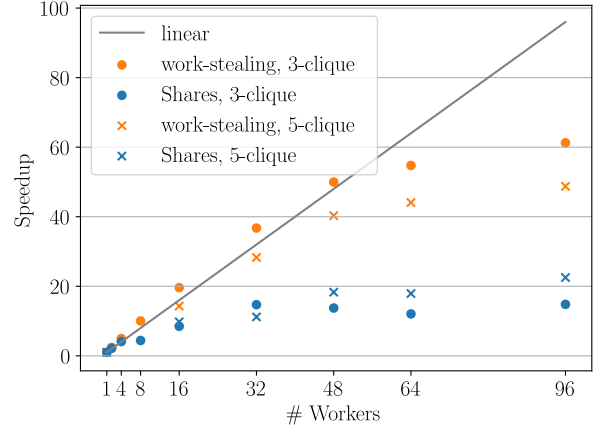
Query	# Result	BroadcastHashJoin	seq	setup	Speedup
Clique(3)	1,137,579	9.25	1.92	6.23	4.8
Clique(4)	911,548	13.97	7.12	11.74	2.0
Clique(5)	585,171	12.27	14.21	19.23	0.9
House	177,745,510	66.53	57.37	23.45	1.2
Kite	3,672,375	13.67	5.40	15.93	2.5
Cycle(4)	45,175,984	37.76	21.77	14.12	1.7
Cycle(5)	263,436,965	227.97	115.60	24.96	2.0
Cycle(6)	1,484,438,088	2020.04	872.62	83.12	2.3
Path-3-0.001	98,929	2.31	0.07	2.30	33.0
Path-4-0.001	922,888	5.61	0.52	4.25	10.8

Query	# Result	BroadcastHashJoin	seq	setup	Speedup
Clique(3)	540,225	11.88	0.73	1.09	16.3
Clique(4)	260,786	15.43	4.85	2.33	3.2
Clique(5)	40,137	18.64	12.65	3.08	1.5
House	100,206,468	300.13	85.18	6.58	3.5
Kite	1,780,235	21.59	2.45	2.20	8.8
Cycle(4)	232,636,376	636.77	162.89	10.68	3.9
Path-3-0.001	65,290,479	1.76	6.05	3.28	0.3
Path-4-0.001	7,235,522,926	111.73	655.52	293.27	0.2

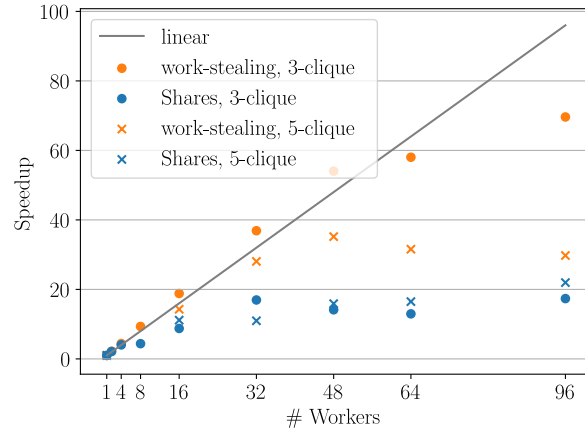
Table 9: Runtimes for **BroadcastHashJoin** and **seq**. The speedup is calculated between join times and excludes setup. From top to bottom for dataset: **ama-0302**, **ama-0601** and **snb-sf1**. All times in seconds.



(a) Twitter dataset



(b) LiveJournal dataset



(c) Orkut dataset

Figure 14: Scaling behaviour of Shares and work-stealing on three different datasets and two different queries. The batch size parameter for *work-stealing* is chosen for balance between lock contention and worker skew: 50 for Twitter and 3-clique on LiveJournal, 1 for 5-clique on LiveJournal and 20 on the Orkut dataset. Measurements for 5-clique for low levels of parallelism are missing for LiveJournal and Orkut due to the time it would take to collect the results.

scaling. Additionally, we believe that *work-stealing* shows better scaling than Shares because it does not duplicate work. Finally, we have no clear cut expectations to the scaling behaviour of *work-stealing*. Theoretically, we could expect linear scaling for it because no work is duplicated, synchronization overhead is minimal and work balance should be given by the scheme. However, we measure on a quite complex hardware platform which complicates scaling behaviour.

First of all, we work on a machine with 4 sockets. This can influence scaling positively and negatively. Positively because adding more sockets means to add significantly more L3 cache (30 MB shared per socket). If we do not use all cores on a socket, each used core can use a bigger share of this cache. Negatively because each socket is in a different NUMA zone and the graph is not guaranteed to be cached in all NUMA zone. Indeed, Spark shares the broadcasts for all tasks on a single executor. So there is only one copy in memory.

Additionally, we run on an Intel processor with hyperthreading. Hence, we can not expect linear speedup above 48 workers because after multiple threads will share resources and cannot be expected to reach the same performance as two cores.

To conclude, we expect sub-linear speedup for Shares and better but still sub-linear speedup for *work-stealing*. Anyhow, it super-linear scaling in MapReduce like systems is not unheard of and could be possible on our machines.

We describe our observations per dataset; starting with Twitter. As expected, both partitioning schemes scale better when we increase the query size. For 5-clique, *work-stealing* exhibits near linear scaling up to 48 workers, while clique-3 reaches the maximum speedup of 6.22 for 8 workers. The highest speedup for 5-clique is 45 on 96 workers; clique-3 reaches its highest speedup with 13.2 on 64 workers. Shares lacks behind in scaling for both queries and all levels of parallelism. The best observed speedup is 21.3 for 5-clique and 96 workers.

The experiment on LiveJournal confirms our hypothesis that bigger datasets lead to better speedups; the highest observed speedup is 61.2 for *work-stealing* on 3-clique and 36.81 for Shares on 5-clique each with 96 workers. Also, we can confirm that Shares scales better on 5-clique than on 3-clique; with the exception of 32 workers. However, this is not the case for *work-stealing*. *work-stealing* shows better speedups on clique-3 than on clique-5. Nevertheless, *work-stealing* beats Shares on both queries and all levels of parallelism.

Additionally, we see two strange scaling behaviours for LiveJournal. First, super-linear scaling for 3-clique and *work-stealing*. We hypothesize that this is the fact because if the 32 processes are distributed over all 4 sockets they share in total 120 MB of L3 cache while a single process can use only 30 MB of L3 cache. To confirm this we rerun the experiment with 1, 8, 16 and 32 workers while using *taskset* to bind the application to the first 8, 16 or 32 cores. This rules out the use of more than 1, 2 or 3 sockets respectively. In this experiment, we measure speedup of 8.6, 16.6 and 32.9 for 8, 16 and 32 workers. This is significantly lower than the speedup measurements without *taskset*. We conclude that this confirms our hypothesis and believe that the slight super-linear scaling that remains arises from the bigger amount of L1 and L2 cache in the system.

Second, Shares exhibits lower speedup of 12.1 for 64 workers which is lower than for 32 workers (14.7) and 14.8 for 96 workers. This can be explained by the chosen Shares configuration. For 32 workers, the best configuration is given by the hypercube of the sizes 4, 4, 2. For 64 workers, we get the hypercube with 4 workers on each axis. Hence, although we are doubling the number of workers, we use the new workers only to partition work along the last axis, in the case of 3-clique along the C attribute axis. Partitioning work along the last axis leads to a high amount of duplicated work on the first two axis. Additionally, with 64 workers at least 12 of these workers are not exclusive cores but cores shared by two hyperthreads. In total, we get a lower speedup. This changes slightly for 96 workers because the optimal hypercube configuration here is 6, 4, 4 which adds more workers along the first axis. However, the scaling only increases marginally by 0.1 from 32 workers which is quite disappointing given that the number of threads increased by a threefold.

One could argue that we should use a different definition of *best* hypercube configuration. As we see, it is not necessarily efficient to distribute the computation along the last axis. We implemented version of the configuration finder that considers only the first  $i$  axes and call this partitioning scheme *i-prefixShares*. TODO report results However for time constraints, we do not investigate this issue further and do not include *i-prefixShares* in our further experiments.

The Orkut and LiveJournal datasets lead to highly similar scaling results: strong linear scaling for *work-stealing* up to 32 workers, *work-stealing* scales significantly better than Shares for 3-clique but not for 5-clique and Shares exhibits less speedup for 64 workers than for 32 and 64 workers.

### 6.3.2 Analysis

## 7 Related Work

### 7.1 Graphs on Spark

Due to its generality, widespread acceptance in the industry, the ability to use cloud hardware and its fault tolerance by design, it is an attractive target for big graph processing. For example, GraphFrames [19], GraphX [21] (a Pregel [32] implementation) or graph query languages as G-CORE [8] and openCypher with their ‘Cypher for Apache Spark’ [39] all aim to ease graph processing on Spark. The last two technologies translate their graph specific operations to the relational interface of Spark (SparkSQL) to profit from Spark’s relational query optimizer Catalyst [11]. Hence, we believe that the WCOJ’s, with their efficiency for analytical graph queries, are a valuable addition to Spark’s built-in join algorithms in general and these graph-on-spark systems in particular.

#### 7.1.1 Fractal a graph pattern mining system on Spark

### 7.2 WCOJ on Timely Data Flow

A second distributed version of worst-case optimal joins was published in 2018 based on a Timely Data Flow system [7, 35]. In Timely Data Flow, shuffling is a streaming, asynchronous operation which requires no disk access<sup>11</sup>. Therefore, the number of shuffle operations is less important than in Hadoop or Spark. The authors take advantage of this fact by using a uniform, non-replicating partitioning scheme for their relationships. Then they implement a worst-case optimal join using distributed data flow operators [35], e.g. min and intersection. Similar to us, the authors focus on scalability and efficiency in their work but due to the use of a streaming, in memory shuffles and a fine-grained batched processing scheme their approach is unlikely to be successful in Spark. Hence, their and our research share the same goals, however, we aim to achieve it in a more restrictive, but widely used and industrial accepted, system.

---

<sup>11</sup>The most commonly used cluster computing engines (Hadoop and Spark) implement shuffling as a synchronizing operation that requires to write and read all tuples from disk.

7.3 Semih's work on worst-case optimal join for different queries

7.4 Adaptive Query Execution

## 8 Conclusions

8.1 Future work

8.1.1 Cluster mode

8.1.2 Deeper integration of Workstealing

## References

- [1] Christopher R Aberger et al. “Emptyheaded: A relational engine for graph processing.” In: *ACM Transactions on Database Systems (TODS)* 42.4 (2017), p. 20.
- [2] Foto N Afrati and Jeffrey D Ullman. “Optimizing multiway joins in a map-reduce environment.” In: *IEEE Transactions on Knowledge and Data Engineering* 23.9 (2011), pp. 1282–1298.
- [3] Foto N Afrati et al. “GYM: A multi-round distributed join algorithm.” In: *20th International Conference on Database Theory (ICDT 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. 2017.
- [4] Foto N Afrati et al. “Sharesskew: An algorithm to handle skew for joins in map-reduce.” In: *Information Systems* 77 (2018), pp. 129–150.
- [5] Sameer Agarwal, Davies Liu, and Reynold Xin. *Apache Spark as a Compiler: Joining a Billion Rows per Second on a Laptop, Deep dive into the new Tungsten execution engine*. 2016. URL: <https://databricks.com/blog/2016/05/23/apache-spark-as-a-compiler-joining-a-billion-rows-per-second-on-a-laptop.html> (visited on 09/26/2019).
- [6] Andreas Amler. “Evaluation of Worst-Case Optimal Join Algorithms.” Master’s thesis. Technische Universität München, 2017.
- [7] Khaled Ammar et al. “Distributed evaluation of subgraph queries using worst-case optimal low-memory dataflows.” In: *Proceedings of the VLDB Endowment* 11.6 (2018), pp. 691–704.
- [8] Renzo Angles et al. “G-CORE: A core for future graph query languages.” In: *Proceedings of the 2018 International Conference on Management of Data*. ACM. 2018, pp. 1421–1432.
- [9] *Apache hadoop*. URL: <http://hadoop.apache.org> (visited on 09/26/2019).
- [10] Molham Aref et al. “Design and implementation of the LogicBlox system.” In: *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM. 2015, pp. 1371–1382.
- [11] Michael Armbrust et al. “Spark sql: Relational data processing in spark.” In: *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM. 2015, pp. 1383–1394.
- [12] Albert Atserias, Martin Grohe, and Dániel Marx. “Size bounds and query plans for relational joins.” In: *Foundations of Computer Science, 2008. FOCS’08. IEEE 49th Annual IEEE Symposium on*. IEEE. 2008, pp. 739–748.
- [13] Paul Beame, Paraschos Koutris, and Dan Suciu. “Skew in parallel query processing.” In: *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. ACM. 2014, pp. 212–223.
- [14] Yingyi Bu et al. “HaLoop: efficient iterative data processing on large clusters.” In: *Proceedings of the VLDB Endowment* 3.1-2 (2010), pp. 285–296.
- [15] Aydin Buluç et al. “Parallel sparse matrix-vector and matrix-transpose-vector multiplication using compressed sparse blocks.” In: *Proceedings of the twenty-first annual symposium on Parallelism in algorithms and architectures*. ACM. 2009, pp. 233–244.

- [16] Shumo Chu, Magdalena Balazinska, and Dan Suciu. “From theory to practice: Efficient join query evaluation in a parallel database system.” In: *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM. 2015, pp. 63–78.
- [17] Databricks. *Databricks Scala Guide*. 2018. URL: <https://github.com/databricks/scala-style-guide/blob/7eb5477781c11f9a75a2d8d6ef773ca6965f4ea0/README.md> (visited on 05/25/2019).
- [18] databricks. *Apache Spark Documentation: RDD Programming Guide*. URL: <https://spark.apache.org/docs/2.2.3/rdd-programming-guide.html> (visited on 09/26/2019).
- [19] Ankur Dave et al. *GraphFrame*. 2016. URL: <https://databricks.com/blog/2016/03/03/introducing-graphframes.html> (visited on 02/18/2019).
- [20] Jeffrey Dean and Sanjay Ghemawat. “MapReduce: Simplified Data Processing on Large Clusters.” In: *OSDI’04: Sixth Symposium on Operating System Design and Implementation*. San Francisco, CA, 2004, pp. 137–150.
- [21] Joseph E Gonzalez et al. “GraphX: Graph Processing in a Distributed Dataflow Framework.” In: *OSDI*. Vol. 14. 2014, pp. 599–613.
- [22] Pankaj Gupta et al. “Real-time twitter recommendation: Online motif detection in large dynamic graphs.” In: *Proceedings of the VLDB Endowment* 7.13 (2014), pp. 1379–1380.
- [23] Benjamin Hindman et al. “Mesos: A platform for fine-grained resource sharing in the data center.” In: *NSDI*. Vol. 11. 2011. 2011, pp. 22–22.
- [24] Chathura Kankanamge et al. “Graphflow: An active graph database.” In: *Proceedings of the 2017 ACM International Conference on Management of Data*. ACM. 2017, pp. 1695–1698.
- [25] Bas Ketsman and Dan Suciu. “A worst-case optimal multi-round algorithm for parallel computation of conjunctive queries.” In: *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. ACM. 2017, pp. 417–428.
- [26] Paraschos Koutris, Paul Beame, and Dan Suciu. “Worst-case optimal algorithms for parallel query processing.” In: *19th International Conference on Database Theory (ICDT 2016)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. 2016.
- [27] *kubernetes*. URL: <https://kubernetes.io> (visited on 09/26/2019).
- [28] Jacek Laskowski. *The Internals of Spark SQL, QueryExecution - Structured Query Execution Pipeline*. URL: <https://jaceklaskowski.gitbooks.io/mastering-spark-sql/spark-sql-QueryExecution.html> (visited on 09/26/2019).
- [29] LDBC. *LDBC SNB Documentation*. 2017. URL: [https://github.com/ldbc/ldbc\\_snb\\_docs](https://github.com/ldbc/ldbc_snb_docs) (visited on 04/10/2019).
- [30] Jure Leskovec and Andrej Krevl. *SNAP Datasets: Stanford Large Network Dataset Collection*. <http://snap.stanford.edu/data>. June 2014.
- [31] Jure Leskovec and Rok Sosič. “Snap: A general-purpose network analysis and graph-mining library.” In: *ACM Transactions on Intelligent Systems and Technology (TIST)* 8.1 (2016), p. 1.

- [32] Grzegorz Malewicz et al. “Pregel: a system for large-scale graph processing.” In: *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM. 2010, pp. 135–146.
- [33] Frank McSherry, Michael Isard, and Derek G Murray. “Scalability! But at what {COST}?” In: *15th Workshop on Hot Topics in Operating Systems (HotOS {XV})*. 2015.
- [34] Amine Mhedhbi and Semih Salihoglu. “Optimizing Subgraph Queries by Combining Binary and Worst-Case Optimal Joins.” In: *CoRR* abs/1903.02076 (2019). arXiv: 1903.02076. URL: <http://arxiv.org/abs/1903.02076>.
- [35] Derek G Murray et al. “Naiad: a timely dataflow system.” In: *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*. ACM. 2013, pp. 439–455.
- [36] Hung Q Ngo, Christopher Ré, and Atri Rudra. “Skew strikes back: New developments in the theory of join algorithms.” In: *arXiv preprint arXiv:1310.3314* (2013).
- [37] Hung Q Ngo et al. “Worst-case optimal join algorithms.” In: *Proceedings of the 31st ACM SIGMOD-SIGACT-SIGAI symposium on Principles of Database Systems*. ACM. 2012, pp. 37–48.
- [38] Dung Nguyen et al. “Join processing for graph patterns: An old dog with new tricks.” In: *Proceedings of the GRADES’15*. ACM. 2015, p. 2.
- [39] openCypher Project. *CAPS: Cypher for Apache Spark*. 2016. URL: <https://github.com/opencypher/cypher-for-apache-spark> (visited on 02/18/2019).
- [40] Semih Salihoglu and M Tamer Özsu. “Response to “Scale Up or Scale Out for Graph Processing”.” In: *IEEE Internet Computing* 22.5 (2018), pp. 18–24.
- [41] Christian Schroeder dewitt. *Leapfrog Triejoin implementation for ‘Database Systems and Implementation’ at Oxford University*. 2012. URL: <https://github.com/schroeder-dewitt/leapfrog-triejoin> (visited on 03/14/2019).
- [42] William F Tinney and John W Walker. “Direct solutions of sparse network equations by optimally ordered triangular factorization.” In: *Proceedings of the IEEE* 55.11 (1967), pp. 1801–1809.
- [43] Vinod Kumar Vavilapalli et al. “Apache hadoop yarn: Yet another resource negotiator.” In: *Proceedings of the 4th annual Symposium on Cloud Computing*. ACM. 2013, p. 5.
- [44] Todd L Veldhuizen. “Leapfrog triejoin: A simple, worst-case optimal join algorithm.” In: *arXiv preprint arXiv:1210.0481* (2012).
- [45] Reynold Xin and Josh Rosen. *Project Tungsten: Bringing Apache Spark Closer to Bare Metal*. 2015. URL: <https://databricks.com/blog/2015/04/28/project-tungsten-bringing-spark-closer-to-bare-metal.html> (visited on 09/26/2019).
- [46] Matei Zaharia et al. “Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing.” In: *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*. USENIX Association. 2012, pp. 2–2.



## A Experimental Results

### A.1 *Graph*WCOJ scaling

Partitioning	Query	Parallelism	Time	Speedup
Shares	3-clique	1	0.0	1.0
Shares	3-clique	2	0.0	1.8
Shares	3-clique	4	0.0	3.1
Shares	3-clique	8	0.0	3.2
Shares	3-clique	16	0.0	5.5
Shares	3-clique	32	0.0	6.5
Shares	3-clique	48	0.0	6.9
Shares	3-clique	64	0.0	6.5
Shares	3-clique	96	0.0	6.9
Shares	5-clique	1	2.6	1.0
Shares	5-clique	2	1.5	1.8
Shares	5-clique	4	0.8	3.5
Shares	5-clique	8	0.4	7.0
Shares	5-clique	16	0.2	11.9
Shares	5-clique	32	0.2	11.9
Shares	5-clique	48	0.1	19.8
Shares	5-clique	64	0.2	17.3
Shares	5-clique	96	0.1	21.4
work-stealing	3-clique	1	0.0	1.0
work-stealing	3-clique	2	0.0	2.0
work-stealing	3-clique	4	0.0	3.6
work-stealing	3-clique	8	0.0	6.2
work-stealing	3-clique	16	0.0	9.7
work-stealing	3-clique	32	0.0	12.1
work-stealing	3-clique	48	0.0	12.7
work-stealing	3-clique	64	0.0	13.2
work-stealing	3-clique	96	0.0	12.8
work-stealing	5-clique	1	2.6	1.0
work-stealing	5-clique	2	1.4	1.9
work-stealing	5-clique	4	0.7	3.8
work-stealing	5-clique	8	0.3	7.8
work-stealing	5-clique	16	0.2	15.4
work-stealing	5-clique	32	0.1	27.8
work-stealing	5-clique	48	0.1	45.5
work-stealing	5-clique	64	0.1	41.4
work-stealing	5-clique	96	0.1	45.0

Table 10: Table showing the speedup of *GraphWCOJ* for 3-clique and 5-clique on the Twitter dataset. Time is shown in minutes.

Partitioning	Query	Parallelism	Time	Speedup
Shares	3-clique	1	1.7	1.0
Shares	3-clique	2	0.8	2.2
Shares	3-clique	4	0.4	4.1
Shares	3-clique	8	0.4	4.4
Shares	3-clique	16	0.2	8.5
Shares	3-clique	32	0.1	14.7
Shares	3-clique	48	0.1	13.7
Shares	3-clique	64	0.1	12.1
Shares	3-clique	96	0.1	14.8
Shares	5-clique	1	531.3	1.0
Shares	5-clique	16	54.3	9.8
Shares	5-clique	32	47.5	11.2
Shares	5-clique	48	29.0	18.3
Shares	5-clique	64	29.7	17.9
Shares	5-clique	96	23.6	22.5
work-stealing	3-clique	1	1.7	1.0
work-stealing	3-clique	2	0.7	2.4
work-stealing	3-clique	4	0.3	5.0
work-stealing	3-clique	8	0.2	10.0
work-stealing	3-clique	16	0.1	19.6
work-stealing	3-clique	32	0.0	36.7
work-stealing	3-clique	48	0.0	50.0
work-stealing	3-clique	64	0.0	54.7
work-stealing	3-clique	96	0.0	61.3
work-stealing	5-clique	1	531.3	1.0
work-stealing	5-clique	16	37.1	14.3
work-stealing	5-clique	32	18.8	28.3
work-stealing	5-clique	48	13.2	40.3
work-stealing	5-clique	64	12.1	44.1
work-stealing	5-clique	96	10.9	48.7

Table 11: Table showing the speedup of *GraphWCOJ* for 3-clique and 5-clique on the LiveJournal dataset. Time is shown in minutes.

Partitioning	Query	Parallelism	Time	Speedup
Shares	3-clique	1	16.3	1.0
Shares	3-clique	2	7.7	2.1
Shares	3-clique	4	4.0	4.0
Shares	3-clique	8	3.7	4.4
Shares	3-clique	16	1.9	8.8
Shares	3-clique	32	1.0	17.0
Shares	3-clique	48	1.2	14.1
Shares	3-clique	64	1.3	13.0
Shares	3-clique	96	0.9	17.4
Shares	5-clique	1	1112.1	1.0
Shares	5-clique	16	99.7	11.2
Shares	5-clique	32	101.4	11.0
Shares	5-clique	48	70.0	15.9
Shares	5-clique	64	67.5	16.5
Shares	5-clique	96	50.6	22.0
work-stealing	3-clique	1	16.3	1.0
work-stealing	3-clique	2	7.4	2.2
work-stealing	3-clique	4	3.7	4.4
work-stealing	3-clique	8	1.7	9.3
work-stealing	3-clique	16	0.9	18.8
work-stealing	3-clique	32	0.4	36.9
work-stealing	3-clique	48	0.3	54.0
work-stealing	3-clique	64	0.3	58.0
work-stealing	3-clique	96	0.2	69.6
work-stealing	5-clique	1	1112.1	1.0
work-stealing	5-clique	16	77.8	14.3
work-stealing	5-clique	32	39.6	28.1
work-stealing	5-clique	48	31.6	35.2
work-stealing	5-clique	64	35.2	31.6
work-stealing	5-clique	96	37.4	29.8

Table 12: Table showing the speedup of *GraphWCOJ* for 3-clique and 5-clique on the Orkut dataset. Time is shown in minutes.