

DISI - Via Sommarive, 14 - 38123 POVO, Trento - Italy
<http://disi.unitn.it>

AN EVALUATION METHODOLOGY AND EXPERIMENTAL COMPARISON OF GRAPH DATABASES

Matteo Lissandrini, Martin Brugnara, Yannis Velegrakis

April 2017

Technical Report # DISI-17-006

An Evaluation Methodology and Experimental Comparison of Graph Databases

Matteo Lissandrini
University of Trento
ml@disi.unitn.eu

Martin Brugnara
University of Trento
mb@disi.unitn.eu

Yannis Velegrakis
University of Trento
velgias@disi.unitn.eu

ABSTRACT

We are witnessing an increasing interest in graph data. The need for efficient and effective storage and querying of such data has led the development of graph databases. Graph databases represent a relatively new technology, and their requirements and specifications are not yet fully understood by everyone. As such, high heterogeneity can be observed in the functionalities and performances of these systems. In this work we provide a comprehensive study of the existing systems in order to understand their capabilities and limitations. Previous similar efforts have fallen short in providing a complete evaluation of graph databases, and drawing a clear picture on how they compare to each other. We introduce a micro-benchmarking framework for the assessment of the functionalities of the existing systems and provide detailed insights on their performance. We support the broader spectrum of test queries and conduct the evaluations on both synthetic and real data at scales much higher than what has been done so far. We offer a systematic evaluation framework that we have materialized into an evaluation suite. The framework is extensible, allowing the easy inclusion in the evaluation of other datasets, systems or queries.

1. INTRODUCTION

Graph data [56] has become increasingly important nowadays since it can model a wide range of applications, including transportation networks, knowledge graphs [41, 53], and social networks [34]. As the graph datasets are becoming larger and larger, so does the need for their efficient and effective management, analysis, and exploitation. This has led to the development of graph data management systems.

There are two kinds of graph data management systems. One is the graph processing systems [26, 32, 35, 42, 43, 44]. They are systems that analyze graphs with the goal of discovering characteristic properties in their structures, e.g., average degree of connectivity, density, or modularity. They also perform batch analytics at large-scale that implement a number of computationally expensive graph algorithms

like PageRank [50], SVD [29], strongly connected component identification [54], core identification [27], and others. Examples in this category include systems like GraphLab, Giraph, Graph Engine, and GraphX [55]. The second kind of graph management systems comprises the so-called graph databases, or GDB for short [20]. Graph Databases focus on storage and querying tasks where the priority is the high-throughput interrogations of the data, and the execution of transactional operations. Originally, they were implemented by exploiting specialized schemas on relational systems. As the sizes of the graphs was becoming larger and more complex, it became apparent that more dedicated systems were needed. This gave rise to a whole new wave of graph databases, that include Neo4j [11], OrientDB [13], Sparksee [14] (formerly known as DEX), Titan [16], and the more recent, ArangoDB [6] and BlazeGraph [15]. The focus of this work is on this second kind of graph management systems, i.e., the graph databases.

Given the increased popularity that graph databases are enjoying, there is a need for comparative evaluations of their available options. Such evaluations are critically important for practitioners in order to better understand both the capabilities and limitations of each system, as well as the conditions under which perform well, so that they can choose the system that better fits the task at hand. A comparative study is also important for researchers, since they can find where they should invest their future studies. Last but not least, it is of great value for the developers, since it gives them an idea of how graph data management systems compare to the competitors and what parts of their systems need improvement. There is already a number of experimental comparisons on graph databases [28, 38, 39], but they do not provide the kind of complete and exhaustive study needed. They test a limited number of features (i.e., queries), which provide only a partial understanding of each system. Furthermore, existing studies do not perform experiments at large scale, but make prediction on how the systems will perform based on tests on smaller sizes. Apart from the fact that they tend to provide contradictory conclusions, when we performed the experiments at larger scale, results were highly different from those they had predicted. Finally, many of the tests performed are either too simplistic or too generic, to a point that it is not easy to interpret the results and identify the exact limitations of each system.

Given the above motivations, in this work we provide a complete and systematic evaluation of the state-of-the-art graph database systems. Our approach is based not only on the previous works of this kind, but also on the princi-

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing ml@disi.unitn.it
Technical Report, Department of Information Engineering and Computer Science - DISI
University of Trento.

ples that are followed when designing benchmarks [17, 28, 37, 38, 39]. Based on an extensive study of the literature, we have made an effort to cover all the scenarios that have so far been identified. As result, we scrupulously test all the types of insert-select-update-delete queries considered so far, with special attention to the various use-cases, and extend such tests to cover the whole spectrum of tasks, data-types and scale. As an indication of the extent of our work, we test 35 classes of operations with single queries and batch workloads as well (for a total of about 70 different tests) as opposed to 4-13 that existing studies have done, and we scale our experiments up to 28M nodes/ 31M edges, as opposed to the 250K nodes/2.2M edges of existing works. Finally, in the design of the tests, we follow a microbenchmark model [24]. Instead of considering elaborate situations, we have identified primitive operators and we designed tests that provide a clear understanding of each such elementary operator. Complex tasks can be typically decomposed into combinations of basic steps, thus, the performance of more involved tests can be inferred by that of the primitive components. In addition, basic operators are often implemented by opaque components in the system, therefore, by identifying the underperformed operators it is easy to determine system components with limited performance.

The specific contributions of this work are the following: **(i)** We explain the limitations of existing graph database evaluations, and clarify the motives of the current evaluation study (Section 2); **(ii)** We describe the model of a graph database and present the most well-known such systems, both old and new, the features that each one provides, and highlight the implementation choices that characterize them (Section 3); **(iii)** Based on a systematic study of the existing literature, we provide an extensive list of fundamental primitive operations (queries) that graph databases should support (Section 4); **(iv)** We introduce an exhaustive experimental evaluation methodology for graph databases, driven by the micro-benchmarking model. The methodology consists of queries identified previously and a number of synthetic and real-world datasets at different scales, complexity, distribution and other characteristics (Section 5). For fairness across systems, the methodology adopts a standard application interface, which allows the testing of each system using the same set of operations; **(v)** We materialize the methodology into a testing suite based on software containers, which is able to automate the installation and investigation of different graph databases. The suite allows for future extensions with additional tests, and is available online as open source, alongside a number of datasets; **(vi)** We apply this methodology on the state-of-the-art graph databases that are available today, and study the effect that different real and synthetic datasets, from co-citation, biological, knowledge base, and social network domains has on different queries (Section 6), along with a report on our experience with the set-up, loading, and testing of each system. It is important to note that the focus of this work is on single machine installations. We did not evaluate parallel or cluster-based executions. Our goal was, as a first step, to understand how the graph databases perform in a single-machine installation. The question about which system is able to scale-out better, may only come after the understanding of its inherent performance [47, 52]. Multi-machine exploitation is our next step that would naturally complement the current work.

2. EXISTING EVALUATIONS

Since we focus on the existing evaluation of graph databases and not of graph processing systems [26, 35, 43], we do not elaborate further on the latter. For graph databases there are studies, however most of them are incomplete or have become out-dated. In particular, one of the earliest works [20] surveys graph databases in terms of their internal representation and modeling choices. It compares their different data-structures, formats and query languages, but provides no empirical evidence of their effectiveness. Another work [18] compares 9 different systems and distinguishes them into graph databases and graph stores based on their general features, data modeling capabilities and support for specific graph query constructs. Unfortunately, not even this work provides any experimental comparison, and like the previous one, it includes systems that have either evolved considerably since then, or have been discontinued.

A different group of studies [28, 38, 39] has focused on the empirical comparison of the performance of the systems, but even these studies are limited in terms of completeness, consistency, and currency of the results. The first of such works [28] analyzes only 4 systems, two of which are no longer supported. Its experiments are limited both in dataset size as well as in number and type of operations performed. The systems were tested on graphs with at most 1 million nodes, and the operations supported were limited to edge and node insertion, edge-set search based on weights, subgraph search based on 3-hops BFS, and the computation of betweenness centrality. Update operations, graph pattern and path search queries are missing, alongside many scalability tests. A few years later, two empirical works [38, 39] compared almost the same set of graph databases over datasets of comparable small sizes, but agree only partially on the concluded results. In particular, the systems analyzed in the first study [38] were DEX¹, Neo4j, Titan, and OrientDB, while the second study [39] considered also Infinite Graph. The results have shown that for batch insertion DEX¹ is generally the most efficient system, unless properties are attached to elements, in which case Neo4j is the fastest one [39]. For traversal with breadth-first search, both works agree that Neo4j is the most efficient. Nonetheless, the second work claims, but without proving it, that DEX¹ would outperform Neo4j on bigger and denser graphs [39]. In the case of computing unweighted shortest paths between two nodes, Neo4j performs best in both studies, but while Titan ends up being the slowest in the first study [38], it is one of the fastest in the second [39]. For node-search queries, the first work [38] shows that both DEX¹ and OrientDB are the best systems when the selection is based on node identifiers, while the other [39], which implements the search based on a generic attribute, shows Neo4j as the winner. Finally, on update operations, the two experimental comparisons present contradicting results, showing in one study favorable results for DEX¹ and OrientDB, while in the other for Neo4j. Due to these differences, these studies have failed to deliver a consistent picture of the available systems, and also provide no easy way of extending them with additional tests and systems.

The benchmarks proposed in the literature to test the performance of graph databases are also of high significance [19, 21, 31]. Benchmarks typically come with tools to automat-

¹DEX is the old name for the Sparksee system

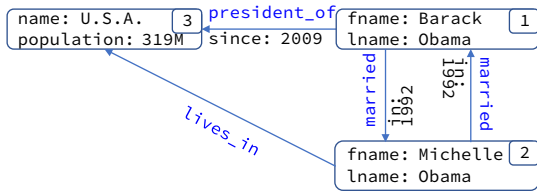


Figure 1: A portion of graph data

ically generate synthetic data, sampling techniques to be used on real data, and query workloads that are designed to pinpoint bottlenecks and shortcomings in the various implementations. These existing benchmarks are domain specific, i.e., RDF focused [19,21] or social network focused [31], but despite the fact that we do not use any of them directly, the design principles and the datasets upon which they have been built have highly influenced our work.

3. GRAPH DATABASES

3.1 Data Model

Graph data are data consisting of nodes (also called vertices) and connections between them called edges. There are various types of graphs depending on the kind of annotations one assumes. In this work we consider generic graphs where every edge has a label and every node or edge has a set of attributes that describes its characteristic properties.

Formally, we axiomatically assume the existence of an infinite set of names \mathcal{N} and an infinite set of values \mathcal{A} . A *property* is an element in the set $\mathcal{N} \times \mathcal{A}$ of name-value pairs.

A *graph* is then a tuple $G = \langle V, E, l, p \rangle$ where V is a set of nodes, E is a set of edges between them, i.e., $E \subseteq V \times V$, $l: E \rightarrow \mathcal{N}$ is an edge labeling function, and $p: \{V \cup E\} \rightarrow 2^{\mathcal{N} \times \mathcal{A}}$ is a property assignment function on edges and nodes.

Note that the above model allows different nodes to have exactly the same properties, and different edges to have exactly the same label and set of properties. To be able to distinguish the different nodes or edges, systems extend the implementation of the above model by means of unique identifiers. In particular, they consider a countable set \mathcal{O} of unique values and a function $id: N \cup E \rightarrow \mathcal{O}$ that assigns to each node and edge a unique value as its identifier. Furthermore, the nodes and edges, as fundamental building blocks of graph data, are typically implemented as atomic *objects* in the systems and are referred to as such.

Figure 1 illustrates a portion of graph data. The annotations containing the colon symbol “ : ” are the properties, while the others are the labels. The number on each node indicates its identifier. For presentation reasons we have omitted the ids on the edges.

3.2 Systems

For a fair comparison we need all systems to support a common access method. Tinkerpop [5], an open source, vendor-agnostic, graph computing framework, is becoming prevalent and the de-facto interface in most graph databases. TinkerPop-enabled system are able to process a common query language: the Gremlin language. Thus, we chose systems that support some version of it through officially recognized implementations. Furthermore, we consider systems

with a licence that permits the publication of experimental comparisons, and also those that were made available to us to run on our server without any fee. Table 1 summarizes the characteristics of the systems we consider in our study. Among others, we show the query languages that these systems support (other than Gremlin). We would’ve also included GraphDB [12] and InfiniteGraph [10], but licensing issues of the first did not allow us to publish any performance verification results, while the second has been discontinued.

3.2.1 ArangoDB.

ArangoDB [6] is a multi-model database. This means that it can work as a document store, a key/value store and a graph database, all at the same time. With this model, objects like nodes, edges or documents, are treated the same and stored into special structures called collections. Apart from Gremlin, it supports its own query language, called AQL, *ArangoDB Query Language*, which is an SQL like dialect that supports multiple data models with single document operations, graph traversals, joins, and transactions. The core, which is open-source (Apache License 2.0), is written in C++, and is integrated with the V8 JavaScript Engine (github.com/v8/v8). That means that it can run user-defined JavaScript code, which will be compiled to native code by V8 on the fly, while AQL primitives are written in C++ and will be executed as such. Nonetheless, the supported way of interacting with the database system is via REST API and HTTP calls, meaning that there is no direct way to embed the server within an application, and that every query will go through a TCP connection.

It supports ACID transactions by storing data modification operations in a write-ahead log, which is a sequence of append-only files containing every write operations executed on the server. While ArangoDB automatically indexes some system attributes (i.e., internal node identifiers), users can also create additional custom indexes. As a consequence, every collection (documents, nodes or edges) has a default primary index, which is an unsorted hash index on object identifiers, and, as such, it can be used neither for non-equality range queries nor for sorting. Furthermore, there exists a default edge index providing for every edge quick access to its source and destination. ArangoDB can serve multiple requests in parallel and supports horizontal scale-out with a cluster deployment using Apache Mesos [4].

3.2.2 BlazeGraph.

Blazegraph [15] is open-source and available under GPLv2 or under a commercial licence. It is an RDF-oriented graph database entirely written in Java. Other than Gremlin, it supports SPARQL 1.1, storage and querying of reified statements, and graph analytics.

Storage is provided through a journal file with support for index management against a single backing store, which scales up to 50B triples or quads on a single machine. Full text indexing and search facility are built using a key-range partitioned distributed B+Tree architecture. The database can also be deployed in different modes of replication or distribution. One of them is the federated option that implements a scale-out architecture, using dynamically partitioned indexes to distribute the data across a cluster. While updates on the journal and the replication cluster are ACID, updates on the federation are *shard-wise ACID*. Blazegraph

Table 1: The tested systems

System (versions)	Storage	Protocol	Gremlin	Languages
ArangoDB (2.8)	multi-files	REST (V8 Server)	2.6	AQL, Javascript
BlazeGraph (2.1.4)	journal file	REST, embedded	3.2	Java, SPARQL
Neo4J (1.9, 3.0)	multi-files	REST, WebSocket, embedded	2.6/3.2	Java, Cypher, SPARQL
OrientDB (2.2)	multi-files, in-memory	REST, WebSocket, embedded	2.6	Java, SQL-like
Sparksee (5.1)	multi-files	embedded	2.6	Java, C++, Python, .NET
Titan (0.5, 1.0)	external	REST, embedded	2.6/3.0	Java

uses Multi-Version Concurrency Control (MVCC) for transactions. Transactions are validated upon commit using a unique timestamp for each commit point and transaction. If there is a write-write conflict the transaction aborts. It can operate as an embedded database, or in a client-server architecture using a REST API and a SPARQL end-point.

3.2.3 Neo4J.

Neo4j [11] is another database system implemented in Java. It supports different query methods, i.e., Gremlin, SPARQL, native Java API, and it also provides its own unique language called Cypher. It employs a custom disk-based native storage engine where nodes, relationships, and properties are stored separately on disk. Dynamic pointer compression expands the available address space as needed, allowing the storage of graphs of any size in its latest version, while in older versions it had a limit to 34 billion nodes. Full ACID transactions are supported through an in-memory write-ahead log. A lock manager applies locks on the database objects that are altered during the transaction.

Neo4j has in place a mechanism for fast access to nodes and edges that is based on IDs. The IDs are basically off-sets in one of the store files. Hence, upon the deletion of nodes, the IDs can be reclaimed for new objects. It also supports *schema indexes* on nodes, labels and property values. Finally, it supports full text indexes that are enabled by an external indexing engine (Apache Lucene [3]), which also allows nodes and edges to be viewed and indexed as “key:value” pairs. Other Neo4J features include a REST API access, replication modes and federation for high-availability scenarios.

3.2.4 OrientDB.

OrientDB [13] is a multi-model database, supporting graph, document, key/value, and object data models. It is written in Java and is available under the Apache licence or a Commercial licence. Its multi-model features Object-Oriented concepts with a distinction for classes, documents, document fields, and links. For graph data, a node is a document, and an edge is mapped to a link. Various approaches are provided for interacting with OrientDB, from the native Java API (both document-oriented and graph-oriented), to Gremlin, and extended SQL, which is a SQL-like query language.

OrientDB features 3 storage types: (i) *plocal*, which is a persistent disk-based storage accessed by the same JVM process that runs the queries; (ii) *remote*, which is a network access to a remote storage; and (iii) *memory-based*, where all data is stored into main memory. The disk based storage

(also called Paginated Local Storage) uses a page model and a disk cache. The main components on disk are files called *clusters*. A cluster is a logical portion of disk space where OrientDB stores record data, and each cluster is split into pages, so that each operation is atomic at page level. As we will discuss later (Section 6), the peculiar implementation of this system provides a good performance boost but poses an important limitation to the storing of edge labels.

OrientDB supports ACID transactions through a write ahead log and a *Multiversion Concurrency Control* system where the system keeps the transactions on the client RAM. This means that the size of a transaction is bounded by the JVM available memory. OrientDB also implements SB-Tree indexes (based on B-Trees), hash indexes, and Lucene full text indexes. The system can be deployed with a client-server architecture in a multi-master distributed cluster.

3.2.5 Sparksee.

Sparksee [14, 45], formerly known as DEX [46], is a commercial system written in C++ optimized for out-of-core operations. It provides a native API for Java, C++, Python and .NET platforms, but it does not implement any other query language apart from Gremlin.

It is specifically designed for labeled and attributed multi-graphs. Each vertex and each edge are distinguished by permanent object identifiers. The graph is then split into multiple lists of pairs and the storage of both the structure and the data is partitioned into different clusters of bitmaps for a compact representation. This data organization allows for more than 100 billion vertices and edges to be handled by a single machine. Bitmap clusters are stored in sorted tree structures that are paired with binary logic operations to speedup insertion, removal, and search operations.

Sparksee supports ACID transaction with a *N-readers* and *1-writer* model, enabling multiple read transactions with each write transaction being executed exclusively. Both search and unique indexes are supported for node and edge attributes. In addition a specific neighbor index can also be defined to improve certain traversal operations. Finally, Sparksee provides horizontal scaling, enabling several slave databases to work as replicas of a single master instance.

3.2.6 Titan.

Titan [16] is available under the Apache 2 license. The main part of the system handles data modeling, and query execution, while the data-persistence is handled by a third-party storage and indexing engine to be plugged into it. For storage, it can support an in-memory storage engine (not intended for production use), Cassandra [1], HBase [2], and BerkeleyDB [7]. To store the graph data, Titan adopts the

adjacency list format, where each vertex is stored alongside the list of incident edges. In addition, each vertex property is an entry in the vertex record. Titan adopts Gremlin as its only query language, and Java as the only compatible programming interface. No ACID transactions are supported in general, but are left to the storage layer that is used. Among the three available storage backends only Berkeley DB supports them. Cassandra and HBase provide no serializable isolation, and no multi-row atomic writes.

Titan supports two types of indexes: *graph centric* and *vertex centric*. A graph index is a global structure over the entire graph that facilitates efficient retrieval of vertices or edges by their properties. It supports equality, range, and full-text search. A Vertex-centric index, on the other hand, is local to each specific vertex, and is created based on the label of the adjacent edges and on the properties of the vertex. It is used to speed up edge traversal and filtering, and supports only equality and range search. For more complex indexing external engines like Apache Lucene or Elasticsearch [9] can be used. Due to the ability of Cassandra and HBase to work on a cluster, Titan can also support the same level of parallelization in storage and processing.

4. QUERIES

To generate the set of queries to run on the systems we follow a micro-benchmark approach [24]. The list is the results of an extensive study of the literature and of many practical scenarios. Of the many complex situations we found, we identified the very basic operations of which they were composed. We eliminated repetitions and ended up with a set of common operations that are independent from the schema and the semantics of the underlying data, hence, they enjoy a generic applicability.

In the query list we consider different types of operations. We consider all the “CRUD” kinds, i.e., **C**reations, **R**eads, **U**pdates, **D**eleitions, for nodes, edges, their labels, and for their properties. Specifically for the creation, we treat separately the case of the initial loading of the dataset from the individual object creations. The reason is because the first happens in bulk mode on an empty instance, while the second at run time with data already in the database. We also consider *traversal* operations across nodes and edges, which is characteristic in graph databases. Recall that operations like finding the centrality, or computing strongly connected components are for graph analytic systems and not typical in a graph database. The categorization we follow is aligned to the one found in other similar works [18, 38, 39] and benchmarks [31]. The complete list of queries can be found in Table 2 and is analytically presented next. The syntax is for Gremlin 2.6, but the syntax for gremlin version 3 is quite similar.

4.1 Load Operations

Data loading is a fundamental operation. Given the size of modern datasets, understanding the speed of this operation is crucial for the evaluation of a system. The specific operator (Query 1) reads the graph data from GraphSON² file. In general it’s bound to the speed with which objects are inserted, which will be affected by any index in place and any other consistency check. In some cases GDBs have in place

²A JSON-based format tinkerpop.apache.org/docs/current/reference/#graphson-io-format

special methods or configurations to allow bulk loading, e.g., to deactivate indexing, but in general they are vendor specific, i.e., not found in the Gremlin specifications. Some of them will be described later (Section 7).

4.2 Create Operations

The first category of operations (**C**) includes operators that create new structures in the database. In this group we consider anything that generates new data-entries. Creation operators may be for nodes, edges, or even properties on existing nodes or edges. Often, to create a complex object, e.g., a node with a number of connections to existing nodes, many different such operators may have to be called. Among the others, we also consider a special composite workload where we first insert a node and then also a set of edges connecting it to other existing nodes in the graph.

Insert Node (Query 2) The operator creates a new node in the database with the set of properties that are provided as argument, but without any connection (edges) to other nodes.

Insert Edge (Queries 3, and 4) This operator creates a new edge in the graph between the two nodes specified as arguments, and with the provided label. In a second version, the operator can also take a set of properties as additional argument. In the experiments performed we randomly select nodes among those in the graph, we choose a fresh value as label, and a custom property name and value pair.

Insert Property (Queries 5, and 6) These two operators test the addition of a new property to a specific node and to a specific edge, respectively. The node (or the edge) is explicitly stated, i.e., referred, through its unique id, and, there is no search task involved since the lookup for the object with the specific identifier is performed before the time is measured. In this case the operation are applied directly to the node and edge (**v** and **e**).

Insert Node with Edges (Query 7) This operation requires the insertion of a new node, alongside a number of edges that connect it to other nodes already existing in the database.

4.3 Read Operations

The category of read operations comprises queries that locate and access some part of the graph data stored in the system that satisfy certain conditions. Sometimes, such part may end up being the entire graph.

Graph Statistics (Queries 8, 9, and 10) Many operations often require a scan over the entire graph datasets. Among the queries of this type, three operators were included in the query evaluation set. One that scans and counts all the nodes, one that does the same for all edges, and one that counts the unique labels of the edges. The goal of the last operation is also to stress-test the ability of the system to maintain intermediate information in memory, since it requires to eliminate duplicated before reporting the results.

Search by Property (Queries 11, and 12) These two queries are typical selections. They identify nodes (or edges) that have a specific property. The name and the value of the property are provided as arguments. There may be a unique

Table 2: Test Queries by Category (in Gremlin 2.6)

#	Query	Description	Cat
1.	<code>g.loadGraphSON("/path")</code>	Load dataset into the graph 'g'	L
2.	<code>g.addVertex(p[])</code>	Create new node with properties p	C
3.	<code>g.addEdge(v1 , v2 , l)</code>	Add edge <i>l</i> from <i>v1</i> to <i>v2</i>	
4.	<code>g.addEdge(v1 , v2 , l , p[])</code>	Same as q.3, but with properties <i>p</i>	
5.	<code>v.setProperty(Name, Value)</code>	Add property <i>Name= Value</i> to node <i>v</i>	
6.	<code>e.setProperty(Name, Value)</code>	Add property <i>Name= Value</i> to edge <i>e</i>	
7.	<code>g.addVertex(...); g.addEdge(...)</code>	Add a new node, and then edges to it	
8.	<code>g.V.count()</code>	Total number of nodes	
9.	<code>g.E.count()</code>	Total number of edges	
10.	<code>g.E.label.dedup()</code>	Existing edge labels (no duplicates)	
11.	<code>g.V.has(Name, Value)</code>	Nodes with property <i>Name= Value</i>	
12.	<code>g.E.has(Name, Value)</code>	Edges with property <i>Name= Value</i>	
13.	<code>g.E.has('label',l)</code>	Edges with label <i>l</i>	
14.	<code>g.V(id)</code>	The node with identifier <i>id</i>	
15.	<code>g.E(id)</code>	The edge with identifier <i>id</i>	
16.	<code>v.setProperty(Name, Value)</code>	Update property <i>Name</i> for vertex <i>v</i>	U
17.	<code>e.setProperty(Name, Value)</code>	Update property <i>Name</i> for edge <i>e</i>	
18.	<code>g.removeVertex(id)</code>	Delete node identified by <i>id</i>	D
19.	<code>g.removeEdge(id)</code>	Delete edge identified by <i>id</i>	
20.	<code>v.removeProperty(Name)</code>	Remove node property <i>Name</i> from <i>v</i>	
21.	<code>e.removeProperty(Name)</code>	Remove edge property <i>Name</i> from <i>e</i>	
22.	<code>v.in()</code>	Nodes adjacent to <i>v</i> via incoming edges	T
23.	<code>v.out()</code>	Nodes adjacent to <i>v</i> via outgoing edges	
24.	<code>v.both('l')</code>	Nodes adjacent to <i>v</i> via edges labeled <i>l</i>	
25.	<code>v.inE.label.dedup()</code>	Labels of in coming edges of <i>v</i> (no dupl.)	
26.	<code>v.outE.label.dedup()</code>	Labels of outgoing edges of <i>v</i> (no dupl.)	
27.	<code>v.bothE.label.dedup()</code>	Labels of edges of <i>v</i> (no dupl.)	
28.	<code>g.V.filter{it.inE.count()>=k}</code>	Nodes of at least k-incoming-degree	
29.	<code>g.V.filter{it.outE.count()>=k}</code>	Nodes of at least k-outgoing-degree	
30.	<code>g.V.filter{it.bothE.count()>=k}</code>	Nodes of at least k-degree	
31.	<code>g.V.out.dedup()</code>	Nodes having an incoming edge	
32.	<code>v.as('i').both().except(vs).store(j).loop('i')</code>	Nodes reached via breadth-First traversal from <i>v</i>	
33.	<code>v.as('i').both(*ls).except(j).store(vs).loop('i')</code>	Nodes reached via breadth-First traversal from <i>v</i> on labels <i>ls</i>	
34.	<code>v1.as('i').both().except(j).store(j).loop('i') {!it.object.equals(v2).retain([v2]).path()}</code>	Unweighted Shortest Path from <i>v1</i> to <i>v2</i>	
35.	<code>v1.as('i').both('l').except(j).store(j).loop('i') {!it.object.equals(v2).retain([v2]).path()}</code>	Same as q.34, but only following label <i>l</i>	

* The symbol [] denotes a Hash Map structure

object satisfying the condition of having the specific property, or there may be more than one.

Search by Label (Query 13) The search by label task is similar to the search by property, but has only one operator since labels are only on edges. Labels are fundamental components of almost every graph dataset, and this is probably the reason why the syntax in Gremlin 3.x distinguishes between labels and properties with a special provision, while in 2.6, they were treated equally. In a graph database edge labels have a primary role, also usually, labels are not optional and are immutable, hence searching edges based on a specific label should receive special attention.

Search by Id (Queries 14, and 15) As it happens in almost

any other kind of database, a fundamental search operation is the one done by reference to a key, i.e., ID. Those are *system defined*, and in some cases based on the internal data organization of the system. These two queries have been included, to retrieve a node and an edge via their unique identifier.

4.4 Update Operations

Data update operators are typical of dynamic data, and graph data is no exception. Since edges are first class citizens of the system, an update of the structure of the graph, i.e., on the connectivity of two or more nodes, requires either the creation of new edges or deletion of existing. In contrast, updates on the properties of the objects are possible without

deletion/insertion, as it happens in other similar databases. Thus, we have included Queries 16, and 17 to test the ability of a system to change the value of a property of a specific node or an edge. In this case, as above, we do not consider the time required to first retrieve the object to be updated.

4.5 Delete Operations

To test how easily and efficiently data can be removed from a graph database, we included three types of deletions: one for a node, one for an edge and one for a property.

Delete Node (Query 18) Deleting a specific node requires the elimination of all its properties, all its edges, as well as the node itself. It may result to a very costly operation when many different data-structures are involved.

Delete Edge (Query 19) Similarly to the node case, deleting an edge requires the prior removal of its properties. This operation is probably one of the most common delete operations in continuously evolving graphs.

Delete Property (Queries 20, and 21) The last two queries eliminate a property from an node or an edge, respectively. As the structure of a node or edge is not fixed, it may happen that either element lose a property.

4.6 Traversals

The ability to conveniently perform traversal operations is one of the main reason why graph models are preferred to others. A traversal means moving across different nodes that are connected in a consecutive fashion through edges.

Direct Neighbors (Queries 22, 23) A popular operation is the one that, given a node, retrieves those directly reachable from it (1-hop), i.e., those that can be found by following either an incoming or an outgoing edge.

Filter Direct Neighbors (Query 24) The specific query performs a traversal of only one hop, and for edges having a specific label. The reason why it is considered separately from other traversals is because it is very frequent and involves no recursion, and as such, it is often subject to separate efficient implementation by the various systems.

Node Edge-Labels (Queries 25, 26, and 27) Given a node, there is often the need to know the labels of the incoming, outgoing, or both edges. These three kinds of retrieval is exactly what this set of three queries perform, respectively.

K-Degree Search (Queries 28, 29, 30, and 31) For many real application scenarios there is a need to identify nodes with many connections, i.e., edges, since this is an indicator of the importance of a node. The number of edges a node has is called the degree of the node, and nodes with high degree are usually hubs in the network. The first three queries identify and retrieve nodes with at least k edges. They differ from each other in considering only incoming edges, only outgoing, or both. The fourth query identifies nodes with at least one incoming edge and is often used when a hierarchy needs to be retrieved.

Breadth-First Search (Queries 32, and 33) A number of search operations give preference to nodes found in close proximity, and they are better implemented with a breadth-first search from a given node. This ability is tested with these two queries, with the second being a special case of

the first that considers only edges with a specific label.

Shortest Path (Queries 34, and 35) Another traditional operation on graph data is the identification of the path between two nodes that contain the smallest number of edges. For this we included these two queries, with the second query being a special case of the first that considers only edges with a specific label. In particular, given an unweighted graph, they determine the shortest path between two nodes via a BFS-like traversal.

5. EVALUATION METHODOLOGY

Fairness, reproducibility, and extensibility have been three fundamental principles in our evaluation of the different systems. In particular, a common query language and input format for the data has been adopted for all the systems. For the query executions, it has been ensured that they have been performed in isolation so that they have not been affected by external factors. Any random selection made in one system (e.g., a random selection of a node in order to query it) has been maintained the same across the other systems. Furthermore, all experiments have been performed on the same machine to avoid any effect caused by hardware variations. Both real and synthetic datasets have been used, especially on large volumes in order for the experiments to be able to highlight the differences across the systems. Finally, our evaluation methodology has been materialized in a test suite and made available online. It contains scripts, data and queries, and is extensible to new systems and queries.

Common Query Language. We have opted for a common query language across all the systems to ensure that the semantics of the queries we run are interpreted in the same way by the different systems. In particular, we selected as application layer the Apache TinkerPop [5] framework and the expressive query language Gremlin [51], which is becoming the de-facto standard for graph databases. In the context of graph databases, TinkerPop acts as a database-independent connectivity layer, while Gremlin is the analogous to SQL in relational databases [36]. All the graph databases we tested have official adapters for Gremlin already implemented.

Software Containers. To ensure full control over the environment in which each system runs, and to facilitate reproducibility, we opted for installing and running each graph database within a dedicated software container [23]. A popular solution is Docker [8], an open source software that creates a level of “soft” virtualization of the operating system, which allows an application within the environment to access machine resources directly without the overhead of interacting with an actual virtual machine. Furthermore, thanks to the so called *overlay file-system* (AUFS [49]), it is possible to create a “snapshot” of a system and its files, and then share the entire computational environment. This allows the sharing of our one-click installation scripts for all the databases and our testing environment, so that the experiments can be repeated elsewhere.

Hardware. For the experiments we used a machine with a 24-core CPU, an Intel Xeon E5-2420, 1.90GHz processor, 128 GB of RAM, 2TB hard disk with 20000 rpm, Ubuntu 14.04.4 operating system, and with Docker 1.13, configured to use AUFS on *ext4*. Each graph database was configured to be free to use all the available machine resources, e.g., for the JVM we used the option `-Xmx120GB`. For other parameters

we used the settings recommended by the vendor. The latter applies also to Apache Cassandra that was serving as the back-end for Titan.

Evaluation Approach. The Gremlin queries are called for execution via Groovy³ scripts. For the systems supporting different major versions of Gremlin, we tested both. Note that Gremlin has no specification for indexes. Some systems create indexes automatically in a way agnostic to the user while others require explicit commands in their own language. We opted for the default behavior of not taking any action and letting the system go with its default indexing policy. We will consider explicit indexing in a future work.

In the case of queries with parameters, for fair comparisons, the parameter values are decided in advance and kept the same for all the systems. For instance, query 14 needs the ID of the node to retrieve. If a different node is retrieved in every system, then it won't be possible to compare them. For this reason, when we need to evaluate a query, we first decide the parameters to use and then start the executions on the different systems. The same applies on queries that need to start from a node or an edge, e.g. query 22 needs to know the node v . For these queries, the node is decided first and then the query is run for that same node in all the systems. Naturally, the time needed to identify the node (or edge) that will be used in the query and retrieve its id, is not part of the time reported as execution time for the respective queries. A similar approach is followed also for the multi-fold evaluation. When we perform k runs of the same query on the same system and dataset (to select the average response time as the indicative performance), we first select k parameter values, nodes, or edges to query (usually through random sampling), and then perform each of the k runs with one of these k parameters (or nodes, or edges).

In the scalability studies of queries 11 and 12 that are performing selection based on a property value, it is important that the performance variation observed when running the same query on datasets of different sizes is due to the size of the data and not due to the cardinality variation of the answer set. To achieve this goal, we select to use properties that not only exist in all the datasets of different sizes, but also have the same cardinality in all of them. In case such properties do not exist, we create and assign at loading time two different properties with predefined names to 10 random edges and 10 random nodes in each dataset and then use these property names for queries 11 and 12. (The different case of the same type of query run on the same dataset producing results of different cardinalities is covered by the different parameter values that are decided in the k -fold experiments.)

Unfortunately, almost all the databases, when loading the data, create and assign their own internal identifiers. This creates a problem when we later need to run across all the systems the same query that is using the identifier as a parameter. For this reason, when loading the data, we add to every node a property $\langle\langle\textit{objectID}, id\rangle\rangle$ where the id is the node identifier in the original dataset. As a result, even if the system decides to replace the identifier of the node with an internal one, we still have a way to find the node using the $\textit{objectID}$ property. So before starting the evaluation of query $g.V(id)$, for instance, on the graph database system S , where id is the node identifier in the original dataset, we first

search in the system S and retrieve the internal identifier iid of the node with the attribute $\langle\langle\textit{objectID}, id\rangle\rangle$. Then, we execute the query $g.V(iid)$ instead of the $g.V(id)$, and report its execution time as the time for the evaluation of the query $g.V(id)$.

Each query is executed 10 times. If the query has parameters, then a different parameter is provided to it on each iteration. These parameters have all been decided in advance as explained previously. The 10 times that a query execution is repeated are performed first in isolation and then in batch mode. For the isolation, we turn the system on, run the single query with one of the parameters, then shut the system off, and reset the file-system. Then repeat again with the next parameter. In this way, each run is unaffected by what has run before. In batch mode, we turn the system on, run the query with the first parameter, then with the second, then the third, and so forth. At the end we shut down the system. The isolation mode makes no sense to be repeated for the queries 8, 9, 10, 28, 29, 30 and 31 since they have no graph-dependent parameters, thus, every isolation mode repetition will be identical to the others. Thus, these queries are evaluated only once in isolation and not in batch. In queries 28, 29 and 30, the k has been considered a threshold and not a parameter, and the fixed value $k=50$ has been considered throughout the experiments. In total, for every evaluation of a specific system with a specific dataset, 337 query executions are taking place.

Test Suite. We have materialized the evaluation procedure into a software package (a test suite) and have made it available on-line⁴, enabling repeatability and extensibility. The suite contains the scripts for installing and configuring each database in the Docker environment and for loading the datasets. The queries themselves are also contained in the suite. There is also a python script that instantiates the Docker container and provides the parameters required by each query. To test a new query it suffices to write it into a dedicated script, while in order to perform the tests on a new dataset one only needs to place the dataset in GraphSON format in a JSON file in the directory from where the data are loaded.

Datasets. We have tested our system on both real and synthetic datasets. One dataset (*MiCo*) describes co-authorship information crawled from the CS Microsoft Academic portal [30]. Nodes represent authors while edges represent co-authorship between two authors and have as a label the number of co-authored papers. Another dataset (*Yeast*) is a protein interaction network [22]. Nodes represent budding yeast proteins (*S.cerevisiae*) [25] and have as labels the short name, a long name, a description, and a label based on its putative function class. Edges represent protein-to-protein interactions and have as label the two classes of the proteins involved. A third real dataset is Freebase [33], which is one of the largest knowledge bases nowadays. Nodes represent entities or events, and edges model relationships between them. We have taken the latest snapshot [40, 48] and have considered four subgraphs of it of different sizes. One subgraph (*Frb-O*) was created by considering only the nodes related to the topics of organization, business, government, finance, geography and military, alongside their respective edges. Furthermore, we randomly selected 0.1%, 1%, and

³A superset of Java: groovy-lang.org

⁴<https://disi.unitn.it/~lissandrini/gdb.html>

10% of the edges from the complete graph, which alongside the nodes at their endpoints created 3 graph datasets, the *Frb-S*, *Frb-M*, and *Frb-L*, respectively.

For a synthetic dataset we used the data generator⁵ provided by the Linked Data Benchmark Council⁶ (LDBC) [31] to produce a graph that mimics the characteristics of a real social network with power-law structure, and real-word characteristics like assortativity based on interests or preferences (*ldbc*). The generator was instructed to produce a dataset simulating the activity of 1000 users over a period of 3 years. The *ldbc* is the only dataset with properties on both edges and nodes. The others have properties only on the nodes.

Table 3 provides the characteristics of the aforementioned datasets. It reports the number of nodes ($|V|$), edges ($|E|$), labels ($|L|$), connected components ($\#$), the size of the maximum connected component (Maxim), the graph density (Density), the network modularity (Modularity), the average degree of connectivity (Avg), the max degree of connectivity (Max), and the diameter (Δ).

As shown in the table, the *MiCo* and the *Frb* are sparse, while the *ldbc* and *Yeast* are one order of magnitude denser, which reflects their nature. The *ldbc* is the only dataset with a single component, while the *Frb* datasets are the most fragmented. The average and maximum degree are reported because large hubs become bottleneck in traversals.

Evaluation Metrics. For the evaluation we consider the disk space, the data loading time, the query execution time, but we also comment on the experience with installing and running each system.

6. EXPERIMENTAL RESULTS

In the tests we run we noticed that *MiCo* and *ldbc* were giving results similar to the *Frb-M* and *Frb-O*. The *Yeast* was so small that didn't highlight any particular issue, especially when compared to the results of *Frb-S*. Thus, in what follows, we will talk mainly about the results of the *Frb-S*, *Frb-O*, *Frb-M*, and *Frb-L* and only when there is some different behavior from the others we will be mentioning it. Additional details about the experimental results on the other datasets can be found below (Section 6.8).

Regarding the documentation, Neo4J and OrientDB provide in-depth information. Sparksee, Titan and ArangoDB are limited in some aspects, yet clear for basic installation, configurations and operational needs. The Titan documentation is the less self-contained, especially on how to be configured with Cassandra. Finally, the BlazeGraph documentation is largely outdated.

In terms of system configuration Neo4J doesn't require any specific set-up. OrientDB instead supports a default maximum number of edge labels equal to 32676 divided by the number of CPU cores, and requires disabling a special feature for supporting more. ArangoDB requires two configurations, one for the engine, and one for the V8 javascript server for logging. With only default values this system generated 40 GB of log files in about 24 hours of activity, with a single active client and it is not able to allocate more than 4GB of memory. For Titan instead the most important configurations are for the JVM Garbage Collection and for the

Cassandra backend. Moreover, for large datasets, it is necessary to disable automatic schema creation, and create it manually before data loading.

Finally, the systems based on Java, namely, BlazeGraph, Neo4J, OrientDB and Titan, are sensitive to the JVM garbage collection, especially for very large datasets that require large amount of main-memory. As a general rule, the option `-XX:+UseG1GC` for the *Garbage First* (G1) garbage collector is strongly recommended.

6.1 Data Loading

The Task. For many systems, loading the data simply by executing the Gremlin query 1 was causing system failures or was taking days. For OrientDB and ArangoDB we are forced to load the data using their native routines. With Gremlin, ArangoDB sends each node and edge insertion instruction separately to the server via a HTTP call making it prohibitively slow even for small datasets. For OrientDB, limited edge-label storing features and long loading times required us to pass through some server-side implementation-specific commands in order to load the datasets. BlazeGraph required the explicit activation of the *bulk loading* feature otherwise we were facing loading times in the order of days. Titan, for any medium to large size dataset requires disabling the automatic schema creation during loading, otherwise its storage back-end (Cassandra) would get swamped with extra consistency check operations. This means that the complete schema of the graph should be known to the system prior to the insertion of the data and is immutable unless implementation specific directives are issued to update it. In the Gremlin implementation in all other systems those operations are transparent to the user. As a result, only Neo4J and Sparksee managed the loading through the gremlin API with no issues, and they did so in times comparable to those achieved by the built-in scripts provided by ArangoDB. Consequently, since (for the loading alone) the different systems did not go through exactly the same procedures, discussions regarding the loading times need to be taken with this information in mind.

The time. For the *Yeast*, which is the smallest dataset, loading times vary from a couple of seconds (with ArangoDB) to a minute (with Titan (v.1.0)). With the *Frb-S* dataset, loading times range from 16 seconds (with ArangoDB) to 16 minutes (with BlazeGraph). Titan and OrientDB are the second slowest, requiring around 5 minutes. Neo4J is usually the second fastest in all loadings tasks being only ten seconds slower than ArangoDB. This ranking stays similar when using the *MiCo* and *ldbc* datasets.

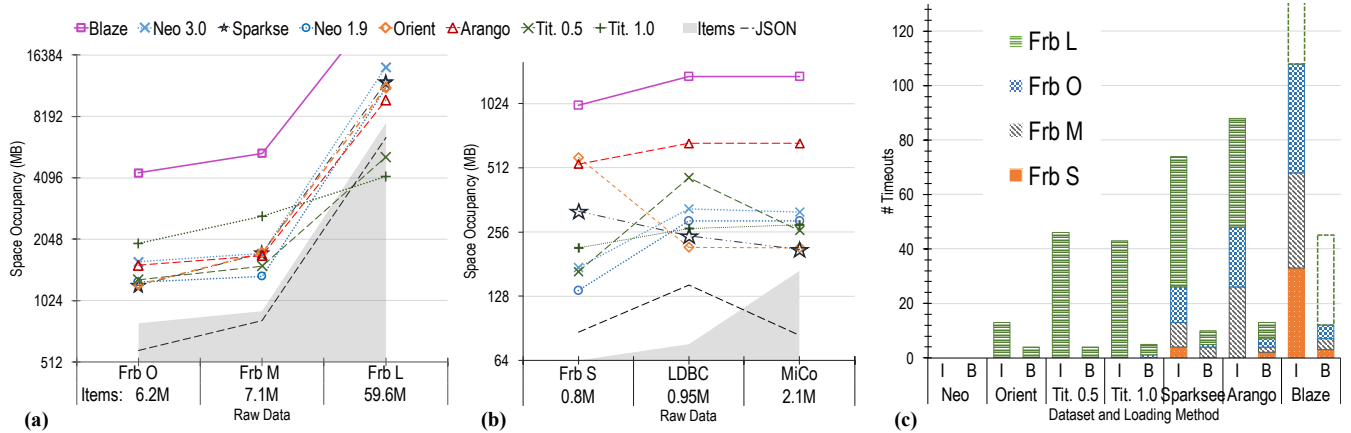
Using the *Frb-O*, *Frb-M*, *Frb-L*, we observed that loading time increases proportionally to the number of elements (nodes and edges) within each dataset. With the largest dataset (*Frb-L*) ArangoDB has the fastest loading time (~ 19 min) and only Neo4J (v.3.0) is just few minutes slower, followed by Neo4J (v.1.9) (~ 38 min), and Sparksee (~ 48 min). OrientDB, instead, took almost 3 hours, while both versions of Titan approximately 4.5 hours. BlazeGraph instead took almost 4.45 hours to load *Frb-M* and was not able to load *Frb-L* at all (we stopped the process after 96 hours). Hence we are not able to show results for BlazeGraph on this particular dataset. Nonetheless we will see in the following that BlazeGraph had almost consistently the worst performance

⁵github.com/ldbc/ldbc_snb_datagen

⁶ldbccouncil.org

Table 3: Dataset Characteristics

	V	E	L	Connected Component		Density	Modularity	Degree		Δ
				#	Maxim			Avg	Max	
<i>Yeast</i>	2.3K	7.1K	167	101	2.2K	$1.34 \cdot 10^{-3}$	$3.66 \cdot 10^{-2}$	6.1	66	11
<i>MiCo</i>	100K	1.1M	106	1.3K	93K	$1.10 \cdot 10^{-6}$	$5.45 \cdot 10^{-3}$	21.6	1.3K	23
<i>Frb-O</i>	1.9M	4.3M	424	133K	1.6M	$1.19 \cdot 10^{-6}$	$9.82 \cdot 10^{-1}$	4.3	92K	48
<i>Frb-S</i>	0.5M	0.3M	1814	0.16M	20K	$1.20 \cdot 10^{-6}$	$9.91 \cdot 10^{-1}$	1.3	13K	4
<i>Frb-M</i>	4M	3.1M	2912	1.1M	1.4M	$1.94 \cdot 10^{-7}$	$7.97 \cdot 10^{-1}$	1.5	139K	37
<i>Frb-L</i>	28.4M	31.2M	3821	2M	23M	$3.87 \cdot 10^{-8}$	$2.12 \cdot 10^{-1}$	2.2	1.4M	33
<i>ldbc</i>	184K	1.5M	15	1	184K	$4.43 \cdot 10^{-5}$	0	16.6	48K	10


Figure 2: Space occupancy on disk required by the systems on the various datasets compared to the size of the original Json file and number of elements in the dataset ((left) and (center)) and number of Time-Outs for Isolation (I) and Batch (B) modes (right)

on the various datasets we were able to test.

The Space. We exploited the docker utilities to measure the disk size occupied by the data in each system. The results are illustrated in Figures 2(a) and 2(b). For each system, we obtained the size of the docker image with the system installed and its required libraries, then we measured again the size of said image after the data loading step. The difference gives a precise account of all files that the loading phase has generated.

Loading *Yeast* that is small, and not reported in figure, leaves the image size almost unchanged for both Neo4J (v.1.9) and Titan (v.0.5), and only 10, 20, 30 and 70MB are added for Neo4J (v.3.0), Sparksee, Titan (v.1.0), and OrientDB, respectively. Instead, ArangoDB generates almost 150MB of additional disk space, and BlazeGraph more than 830MB, the latter due to the size of journal and automatic-indexing that, when generated, are multiples of a fixed size.

With the *Frb-O* dataset, as Figure 2 illustrates, Sparksee, OrientDB, Neo4J (v.1.9), and Titan (v.0.5) are all equally *compact*, with a delta on the disk image size of about 1.2GB. For *Frb-M*, though, Neo4J (v.1.9) and Titan (v.0.5) are equally effective in disk size and a little better than the others, requiring respectively 1.3GB and 1.5GB to store a dataset of 816MB and 7.1 million elements. Titan (v.1.0) has, on both medium size datasets (*Frb-O* and *Frb-M*), the

second worst performance (the worst being BlazeGraph), with three to four times the space consumption of the original dataset in plain text. Instead, for the *Frb-L*, Titan (v.1.0) scales much better, compressing 6.4GB of raw data into 4.1GB, followed by Titan (v.0.5) taking 5.1GB. The remaining databases are almost equivalent, taking from 10 to 14GB. Exception is the BlazeGraph, on all the datasets, requiring on average three times the size of any other system. Note that for BlazeGraph on *Frb-L* the reported size is at the time-out. This shows that the compression strategy of Titan is the most compact at larger scales.

The comparison between the disk space required by the systems to store *Frb-S*, *MiCo* and *ldbc* (Figure 2(b)) reveals a peculiar behavior for Sparksee and OrientDB, where the space occupied on disk is smaller for the two larger datasets. As a matter of fact, the *ldbc* dataset stored as plain text file occupies twice more space on disk than the *Frb-S* file, and contains 2 hundred thousands more elements. Nonetheless Sparksee requires for *ldbc* about 25% less space, and OrientDB less than half the space occupied on disk by the corresponding image with *Frb-S*. *MiCo* has comparable size, in plain text, to *Frb-S*, and contains twice the objects of *Frb-S*, but still the respective docker images of OrientDB and Sparksee for *MiCo* are almost half the size of their images containing the *Frb-S*. These disproportions can be explained

by the fact that *Frb-S* contains almost $\sim 2K$ different edge labels, while *MiCo* 106, and *ldbc* only 15, and apparently these systems are sensitive to the number of edge labels.

It is important to note here that we have tried also much larger datasets, but we were not able to load them on a large number of systems so we could not have comparison across all the systems and we have not reported them.

6.2 Completion Rate

Since graph databases are often used for on-line applications, ensuring that queries terminate in a reasonable time is important. We count the queries that could not complete within 2 hours, in isolation or in batch mode, and illustrate the results in Figure 2(c). Note that, if one instantiation of one query fails to run within the allotted amount of time in isolation, when executed in a batch it will cause the failure of the entire batch as well.

Neo4J, in both version, is the only system which completed successfully all tests with all parameters on all datasets. OrientDB is the second best, with just few timeouts on the large *Frb-L*. BlazeGraph is at the other end of the spectrum, collecting the highest number of timeouts. It reaches the time limit even in some batch executions on *Yeast*, and also, even though we cannot report timeouts with *Frb-L*, we can safely assume they should happen for the same queries which fail on *Frb-M*. In general the most problematic queries are those that have to scan or filter the entire graph, i.e., queries Q.9 and Q.10. Some shortest-path searches, and some bread first traversal with depth 3 or more in most databases reach the timeout on *Frb-O*, *Frb-M* and *Frb-L*. Filtering of nodes based on their degree (Q.28, Q.29, and Q.30) and the search for nodes with at least one incoming edge (Q.31) are proved to be extremely problematic almost for all databases apart from Neo4J and Titan (v.1.0). In particular for Sparksee these queries cause the system to exhaust the entire available RAM and swap space on all Freebase subsamples. BlazeGraph fails also these last queries on all the Freebase datasets, while ArangoDB fails it only on *Frb-M* and *Frb-L*, and OrientDB instead only on *Frb-L*.

6.3 Insertions, Updates and Deletions

For operations that add new objects (nodes, edges, or properties), tests show extremely fast performances for Sparksee, Neo4J (v.1.9), and ArangoDB, with times below 100ms, with Sparksee being generally the fastest (Figure 3(a)). Moreover, with the only exception of BlazeGraph, all databases are almost unaffected by the size of the dataset. We attribute this to the use of write-ahead logs, and the internal configuration of the adopted data-structures. BlazeGraph is instead the slowest with times between 10 seconds and more than a minute. Both versions of Titan are the second slowest with times around 7 seconds for insertion of nodes, and 3 seconds for insertion of edges or properties, while for the insertion of a node with all the edges (Q.7) it takes more than 30 seconds. Sparksee, ArangoDB, OrientDB, and Neo4J (v.1.9) complete the task in less than a second. OrientDB is among the fastest for insertions of nodes (Q.2) and properties on both nodes and edges (Q.5 and Q.6), but is much slower, with inconsistent behavior, for insertion of edges. Neo4J (v.3.0), is more than an order of magnitude slower than its previous version, with a fluctuating behavior that does not depend on the size of the dataset.

Similar results are obtained for the update of properties on both nodes and edges (Q.16, and Q.17), and for the deletion of properties on edges (Q.21).

The performance of node removal (Q.18) for both OrientDB and Sparksee seems highly affected by the structure and size of the graphs (Figure 3(b)). On the other hand, ArangoDB and Neo4J (v.1.9) remain almost constantly below the 100ms threshold, while Neo4J (v.3.0) completes all the deletions between 0.5 and 2 seconds. Finally, for the removal of nodes, edges, and node properties, Titan shows almost one order of magnitude improvement.

For creations, updates and deletions, as a whole, the fastest are Neo4J (v.1.9), with constant times below 100ms, and then Sparksee, but with quite a scale-sensitive behavior for edge-deletion, that is shared with OrientDB. ArangoDB is also consistently among the fastest, but its interaction through REST calls, and the fact that it does not support transactions, constitutes a bias on those results in its favor since the time is measured on the client side.

6.4 General Selections

With *read* queries, some heterogeneous behaviours start to show up. The search by ID (Figure 4(b)) differs significantly from all other queries in this class. BlazeGraph is again the slowest, with performances around 500ms for the search of nodes, and instead 4 seconds or more for the search of edges. All other systems take less than 400ms to satisfy both queries, with Titan the slowest among them. Here Sparksee, OrientDB and Neo4J (v.1.9) return in about 10ms, hinting to the fact that, given the ID, they are able to jump immediately to the right position on disk where to find it.

In counting nodes and edges (Q.8, and Q.9), Sparksee has the best performance followed by Neo4J (v.3.0). As a matter of fact Sparksee and Neo4J (v.3.0) complete the two tasks in less than 10 seconds on all sizes of Freebase, while Neo4J (v.1.9) take more than an minute on the *Frb-L*. For BlazeGraph and ArangoDB, node counting is one of the few queries in this category that complete before timeout. In particular in Q.8 BlazeGraph is faster than ArangoDB, but then it hits the time limit for Q.9 on all Freebase subsamples, while ArangoDB, at least for *Frb-S* it's able to get the answer in time also on the other queries. Edge iteration, on the other hand, seems hard for ArangoDB that rarely completes within 2 hours for the Freebase datasets.

Computing the set of unique labels (Q.10) changes a little the ranking. Here, the two versions of Neo4J are the fastest databases, while Sparksee gets a little slower. The search for nodes (Q.11) and edges (Q.12) based on property values performs similar to the search for edges based on labels (Q.13), for almost all databases. In these three queries, Neo4J (v.3.0) gives the shortest time, with the Q.13 performing slightly faster than the others, getting an answer in a little more than 10 seconds on the larger dataset, while Neo4J (v.1.9), Sparksee, and OrientDB are at least one order of magnitude slower. Only for Sparksee we notice relevant differences between Q.12 and Q.13. Hence, equality search on edge labels is not really optimized in the various systems.

6.5 Traversals

As mentioned above, the most important class of queries that GDBs are called to satisfy regards the *traversal* of graph structures. In the performance of traversal queries

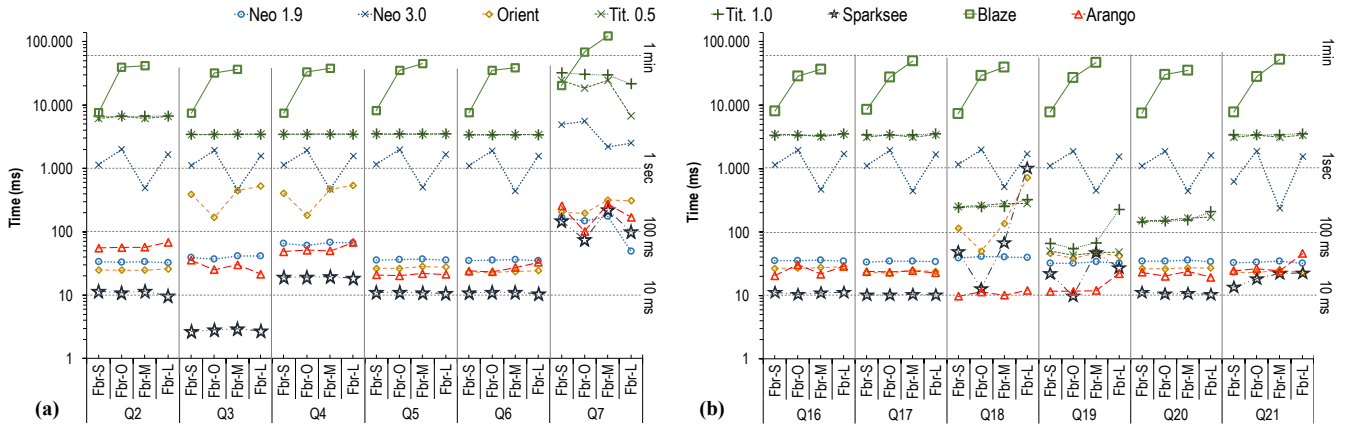


Figure 3: Time required for (a) insertions and (b) updates and deletions.

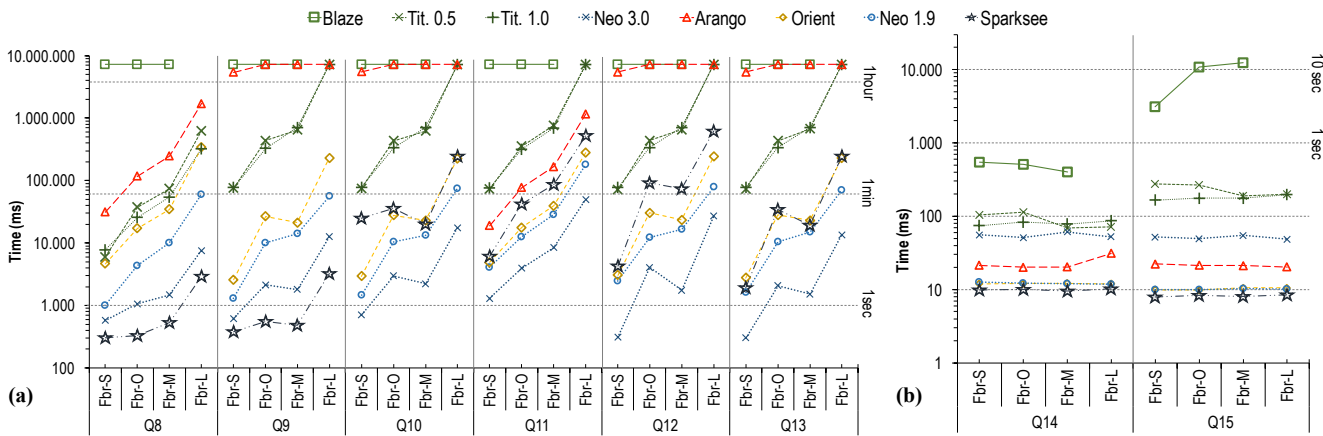


Figure 4: Selection Queries: The Id-based (right) perform orders of magnitude better than the rest (left)

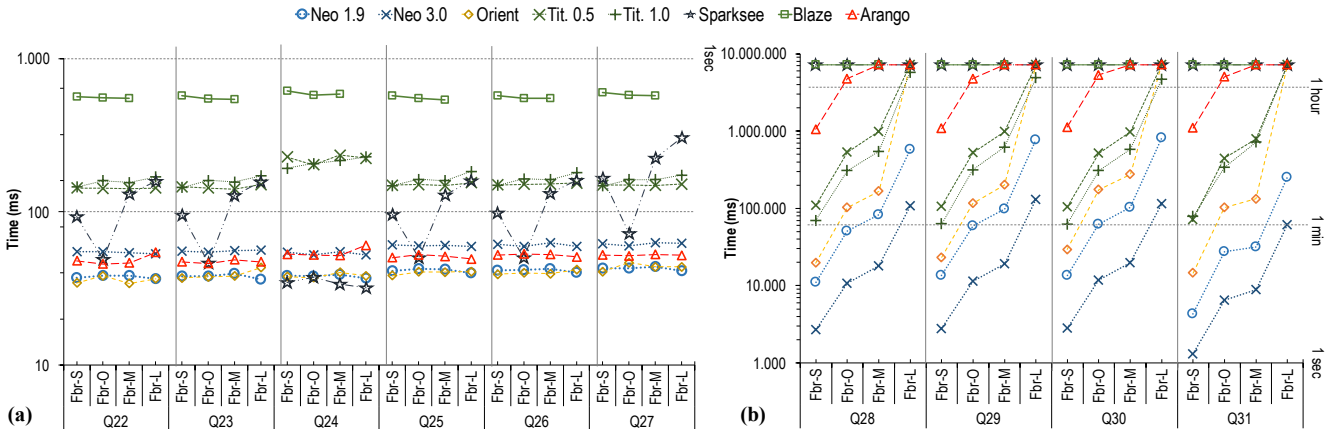


Figure 5: Time required for traversal operations: (a) local access to node edges, and (b) filtering on all nodes

that access the direct neighborhood of a specific node (Q.22 to Q.27), we observe (Figure 5(a)) that OrientDB, Neo4J (v.1.9), ArangoDB, and then Neo4J (v.3.0) are the fastest, with response times below the 60ms, and being robust to the size and structure of the dataset. Sparksee seems to be

more sensitive to the structure and size of the graph, requiring around 150ms on *Frb-L*. The only exception for Sparksee is when performing a visit of the direct neighborhood of a node filtered by the edge labels, in which case it is on par with the former systems. BlazeGraph is again an order of

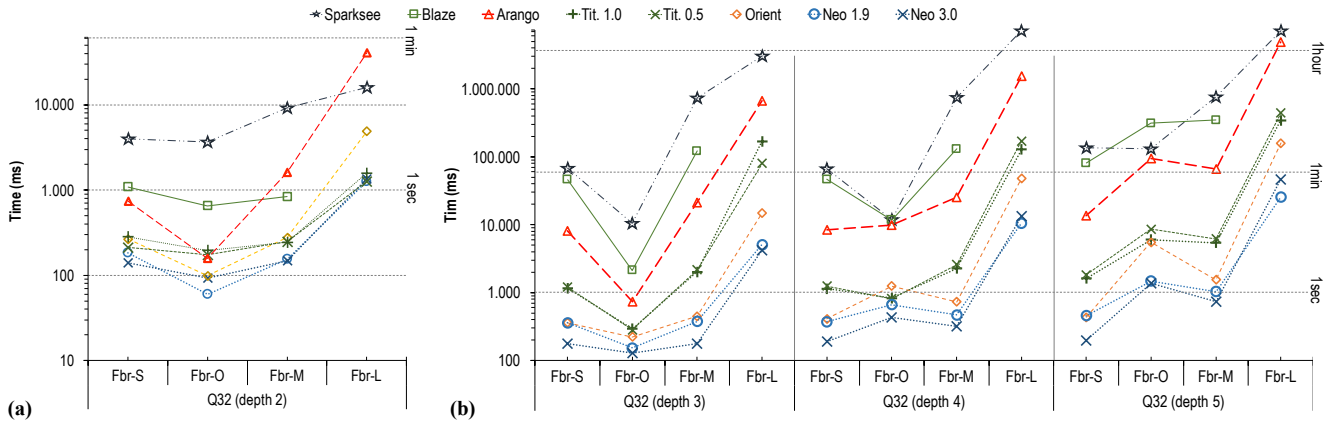


Figure 6: Time required for breadth-first traversal (a) at depth=2, and (b) at depth>=3

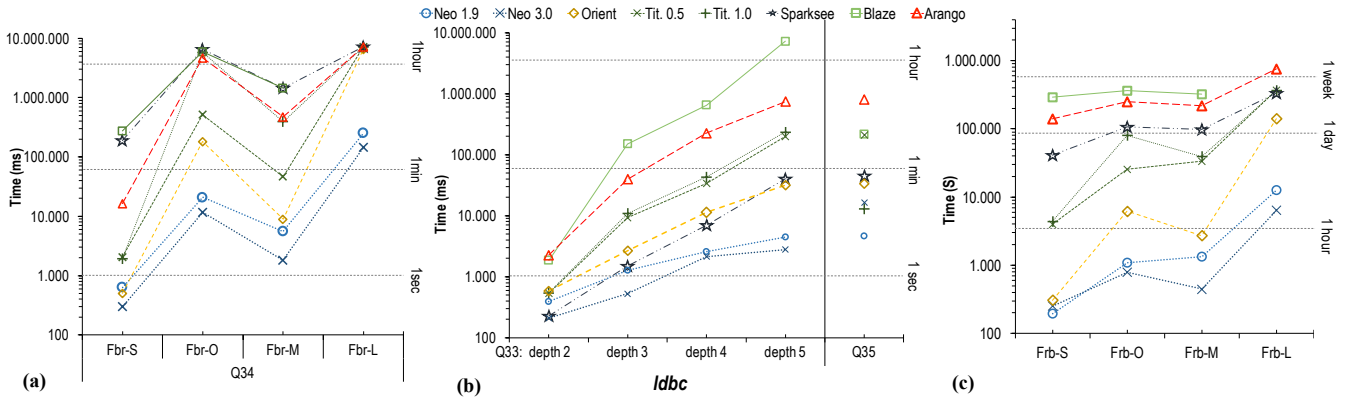


Figure 7: Performance of (a) Shortest Path, (b) label-constrained BFS and Shortest Path, and (c) Overall

magnitude slower (~600ms) preceded by Titan (~160ms). When comparing the performance of queries Q.28 to Q.31 that traverse the entire graph filtering nodes based on the edges around them, as shown in Figure 5(b), the clear winner is Neo4J (v.3.0), with its older version being the second fastest. Those two are also the only two engines that complete the query on all datasets. In particular Neo4J (v.3.0) completed each query on *Frb-L* in less than two minutes on average, while Neo4J (v.1.9) took at least 10 minutes for the same dataset. All tested systems are obviously affected by the number of nodes and edges to inspect. Sparksee is unable to complete any of these queries on Freebase due to the exhaustion of the available memory, indicating probably a problem in the implementation, as this never happens in any other case. BlazeGraph as well hits the timeout limit on all samples, while ArangoDB is able to complete only on *Frb-S* and *Frb-O*. Nevertheless, they all complete the task on *Yeast*, *ldbc* and *MiCo*.

We study breadth-first searches (Q.32 and Q.33) and shortest path searches (Q.34 and Q.35) separately from the other traversal operations. The performance of the unlabeled version of breadth-first-search, shown in Figure 6, highlights once more the good scalability of both versions of Neo4J at all depths. OrientDB and Titan give the second fastest times for depth 2, with times 50% slower than those of Neo4J. For depth 3 and higher, as Figure 6(b) illustrates, OrientDB is a

little faster than Titan. On the other hand, in these queries we observe that Sparksee is actually the slowest engine, even slower than BlazeGraph. For query Q.33 in Figure 7(a), which is the shortest path with no label constraint, the performance of the system is similar to the above, apart from BlazeGraph and Sparksee that are in this case very similar.

The label-filtered version of both the breadth first search and the shortest path query on the Freebase samples (not shown in a figure) were extremely fast for all datasets because the filter on edge labels cause the exploration to stop almost immediately. Running the same queries on *ldbc* we still observe (Figure 7(b)) that Neo4J is the fastest engine, but also Sparksee is the second fastest in par with OrientDB for the breadth-first search, while on the shortest path search filtered on labels, Titan (v.1.0) gets the second place.

6.6 Overall Performance

To sum up the evaluation we can compare the cumulative time taken by each system to complete the entire set of queries in both single and batch executions (Figure 7(c)). Overall Neo4J is the fastest engine. The newer version has been updated to handle graphs with arbitrary number of nodes and edges, while the older version supported “only” some billions. The more complex data structure put in place in the new version is probably the cause of the slower and fluctuating times recored in the **C**, **U**, **D** class of queries.

Nonetheless, on the most time-consuming queries for class **R** and **T** Neo4J (v.3.0) is usually the fastest, and Neo4J (v.1.9) the runner-up. Pretty good running times have also been recorded for OrientDB, which is often on par with Neo4J, and in cases is better than one of its two versions. It does not, however, do well in cases where large portions of the graph have to be processed and kept in memory, e.g., with *Frb-L*. Titan results quite often one order of magnitude slower than the best engine. It shows difficulties in create and update operations, however, it is much better in deletions, most likely due to the *tombstone* mechanism, where it marks an item as removed instead of actually removing it. Nonetheless, this method seems to result slower than the write ahead log (WAL) adopted by the others. Sparksee gives almost consistently the best times in the operations for creating, updating, and deleting objects. Although it is not very fast with deletions of nodes having a lot of edges, it is still better than others. It also performs best in edge and node counts, as well as in retrieval of nodes and edges by ID, thanks to its internal compressed data-structures. Nevertheless, it performs worse than others for the other retrieval queries and is the worst in most traversals, showing good performances only when a filter on edge labels was applied. Finally, it gives a lot of timeouts on the degree-based node search queries. ArangoDB excels only in few queries. For creation, updates and deletes, it ranks among the fastest. For retrievals its performance is in general poor, except when searching by ID, while for traversals it has a narrow lead over Sparksee and BlazeGraph. Finally BlazeGraph results in a generally poor performance. The indexes it builds automatically do not seem to help much. Most likely, it is optimized for SPARQL queries only and not for a generic graph management.

6.7 Single vs Batch Execution.

We looked at the times differences between single executions (run in isolation) and batch. We report times for each batch execution for *Frb-S*, *Frb-O*, *Frb-M*, and *Frb-L* in Figures 8, 9, 10, and 11. Running the queries in batch mode does not create any major changes in the way the systems compare to each other. For the retrieval queries, the batch requests of the 10 queries were taking exactly 10 times the time of one iteration, i.e., no benefit obtained from the batch execution. Exception is for queries 14 and 15 (Figure 9 b), here times to retrieve 10 nodes by their internal IDs are almost exactly the same as for retrieving one single node (see Figure 4 above). Such behavior suggests that the systems load the data into main memory at the first call, and then retrieves everything from there.

Instead, for the create, update and delete operations, the batch is less than 10 times the time needed for one iteration, meaning that in single mode most of the time we measure is some initiation set-up time for the operation. For traversal queries the batch executions only stressed the differences between faster and slower databases.

6.8 *Yeast*, *MiCo*, and *ldbc*

In the following we report on the results of the tests performed on the *Yeast*, *MiCo*, and *ldbc* datasets, which are generally smaller than the Freebase samples, and also have a much smaller number of edge labels. Results for queries in isolation mode are reported in Figure 12, 14, 16, and 18, while results for the batch mode execution are in Figure 13, 15,

17, and 19. Experiments on these datasets, as noted above, show again similar relative performances compared to the results on the Freebase samples described earlier. In general we see Sparksee performing among the fastest databases more often. ArangoDB’s performance as well is much more similar to the other systems. BlazeGraph instead is usually the slowest also on those datasets. As a matter of fact, even in tests with *Yeast*, BlazeGraph is not always able to terminate queries within the timeout limit, which indicates some serious implementation problems for some of the selection queries (Figure 14).

7. EXPERIENCES

In general our experiences cover a large spectrum of issues, technical challenges that we faced, and areas of improvement that are related to the usability of the various systems.

Installation, Configuration, Documentation and Support.

First we stress that the only 2 systems that we were able to install and run as expected were Neo4J and Sparksee. For those, after downloading the relevant binaries and following the instructions provided on the respective websites, we were almost immediately able to load our datasets, at all sizes, and run some queries. For the others, as mentioned earlier (and below) we had to overcome some difficulties in importing the datasets, configuring the systems properly, and understanding the errors raised when running some of the queries. As a result, for those system that are open-source and hosted on a public repository, we reported those problems and bugs found as issues. In total we issued 8 support request (comprising bug issues) for ArangoDB, 4 for OrientDB, 2 for Titan, and 1 for BlazeGraph.

For ArangoDB and OrientDB some of those bugs have been fixed in official releases of the software or have been included in the development road-map. Instead those regarding Titan and BlazeGraph didn’t receive any reply from the developers (in many months) and, where possible, were either fixed or circumvented in our local installs. This also describes the level of support received by the respective development teams.

Regarding the documentation, we note that Neo4J and OrientDB are provided with pretty in-depth informations for developers. Sparksee, Titan and ArangoDB have some documentation, limited in some aspects, but still clear for basic installation, configurations and operational needs. Among those Titan manual contains a lot of confusion among the various existing software versions, and in some cases, the provided instructions and example-code are not actually self-contained. Also, given the reliance on Cassandra for the storage, it is to note the reduced amount of information on how to properly configure this system and how to tackle the various problems arising with it. BlazeGraph’s documentation, instead, is largely outdated. Also, even though the system relies a lot on the user for proper configuration, the information provided is generally cryptic.

Regarding the configuration of the other systems, we report that Neo4J doesn’t require any specific configuration. OrientDB instead supports by default a number of edge labels at most equal to 32676 divided by the number of cores in the machine (e.g., 4084 edge labels on a 8 cores machine), for supporting more labels, it requires a special feature to be disabled. ArangoDB requires two configurations, one for the engine, and one for the V8 javascript server, the second

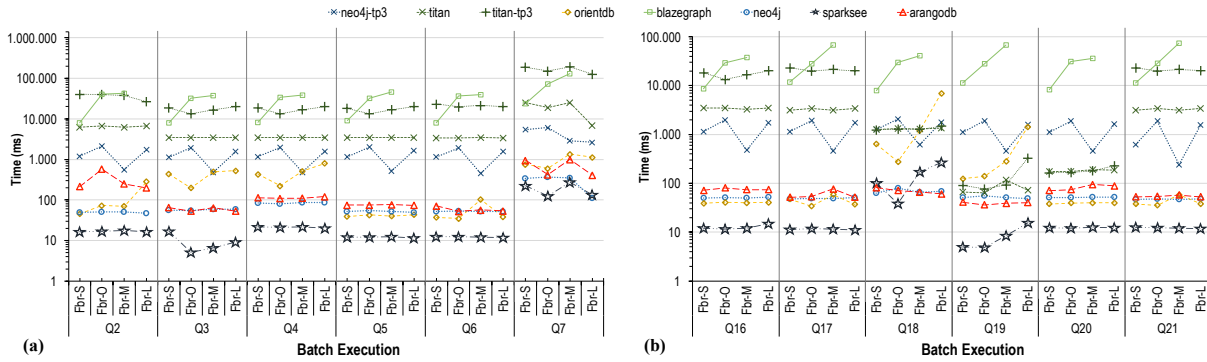


Figure 8: Time required in Batch Mode for (a) insertions and (b) updates and deletions.

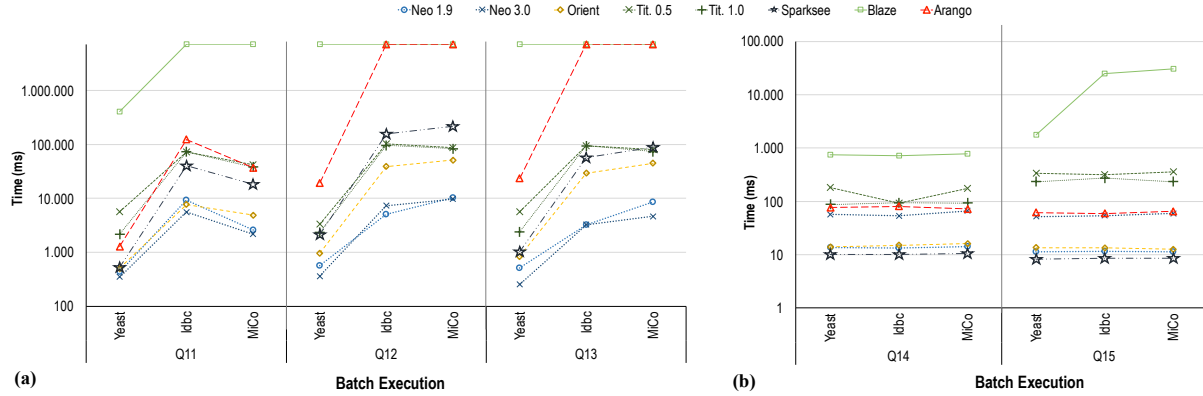


Figure 9: Selection Queries in Batch Mode: The Id-based (right) perform orders of magnitude better than the rest (left), and compared to the isolation mode they take the same amount of time

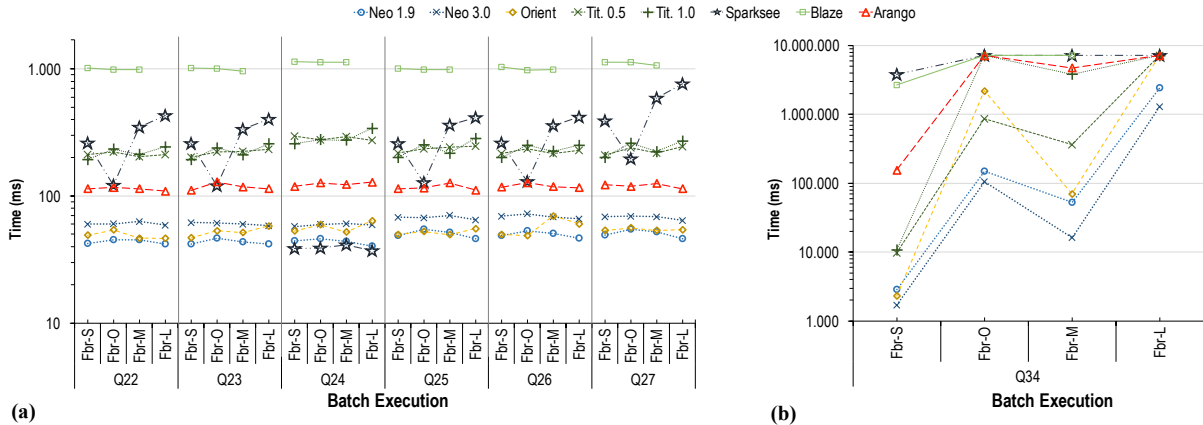


Figure 10: Time required for traversal operations in Batch Mode: (a) local access to node edges, and (b) for shortest path search

regards the level of logging of the system. Without proper configuration (with only default values) this system generated 40 GB of log files in about 24 hours of activity, with a single active client. For Titan instead the most important configurations are for the JVM Garbage Collection and for the Cassandra backend. Additionally, with large datasets, it is necessary to disable automatic schema creation, and to create instead the schema manually before loading the data.

All systems based on java, were also extremely sensitive to the effect of the garbage collection routines. When dealing with data-intensive applications and a large amount of main-memory, it is necessary to provide a customized configuration to the JVM, yet, none of the systems provide clear instructions on how to tune it properly for their needs, but they only propose generic advices.

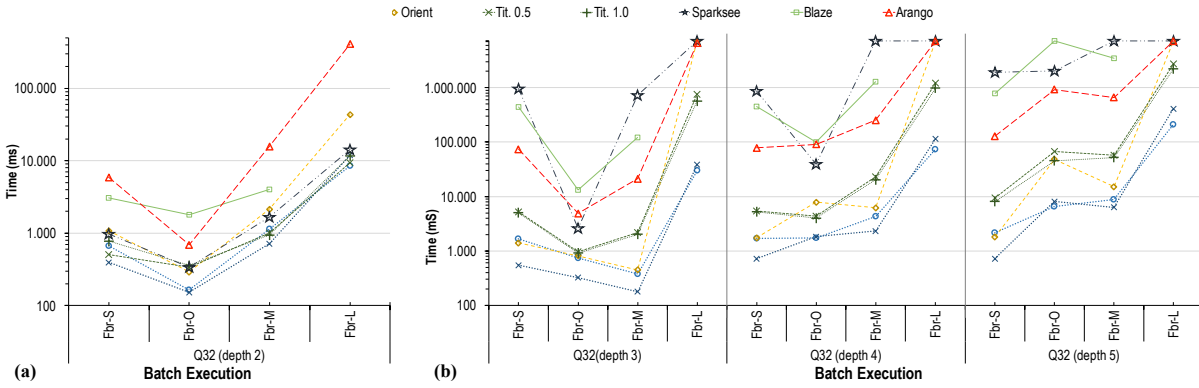


Figure 11: Time required for breadth-first traversal in batch mode (a) at depth=2, and (b) at depth ≥ 3

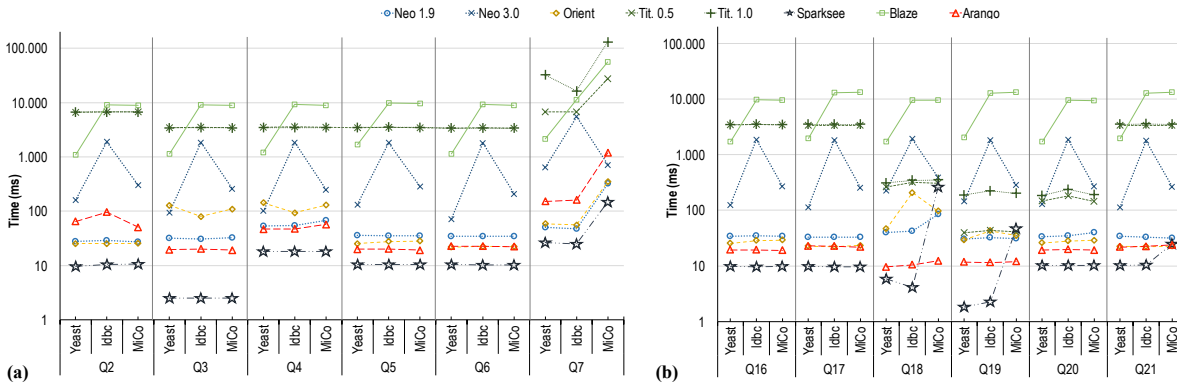


Figure 12: Time required on *Yeast*, *ldbc*, and *MiCo* for (a) insertions and (b) updates and deletions in isolation mode.

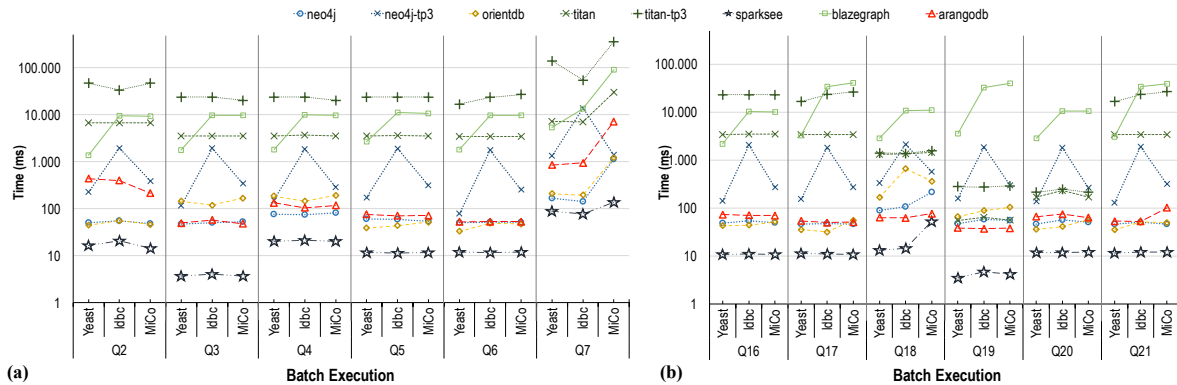


Figure 13: Time required for (a) insertions and (b) updates and deletions in batch mode for *Yeast*, *ldbc*, and *MiCo*.

Finally we report on the Tinkerpop/Gremlin documentation. For version 2.6 the list of supported methods with some examples are provided ⁷, for version 3 the official manuals are much more extended ⁸, although not to the benefit of clarity. In this sense, we also hope that the code of the

queries implemented in this study serve as more concrete tutorial for understanding the basics of the Gremlin language.

Loading problems. As already mentioned, we encountered a great deal of issues when trying to load the datasets in some of the databases tested. ArangoDB in particular, when using Gremlin for loading, sends each node and edge insertion instruction separately to the server (in a HTTP call). This method results too slow, even for small datasets, so that we were forced to use some routines provided by the

⁷gremlindocs.spmallete.documentup.com

⁸<http://tinkerpop.apache.org/docs/current/reference/>

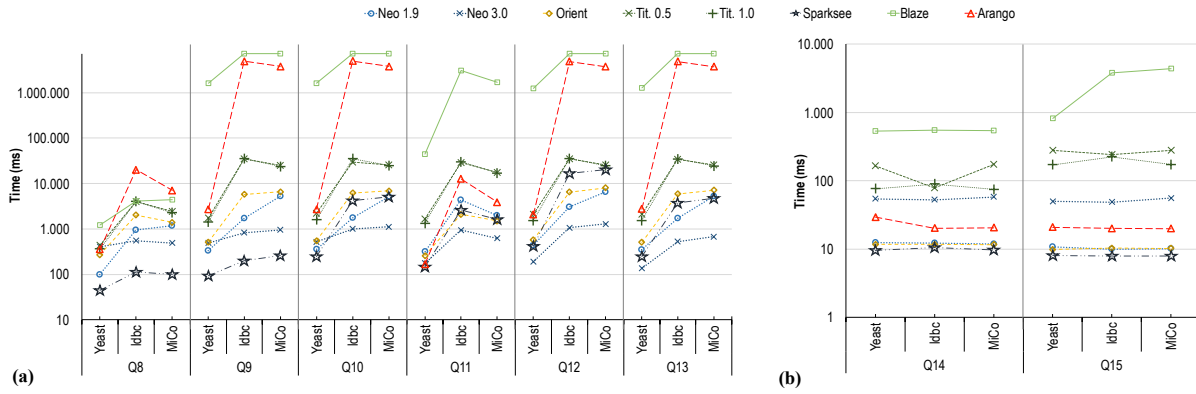


Figure 14: Selection Queries in Batch Mode for *Yeast*, *ldbc*, and *MiCo*: The Id-based (right) perform orders of magnitude better than the rest (left), and compared to the isolation mode they take the same amount of time

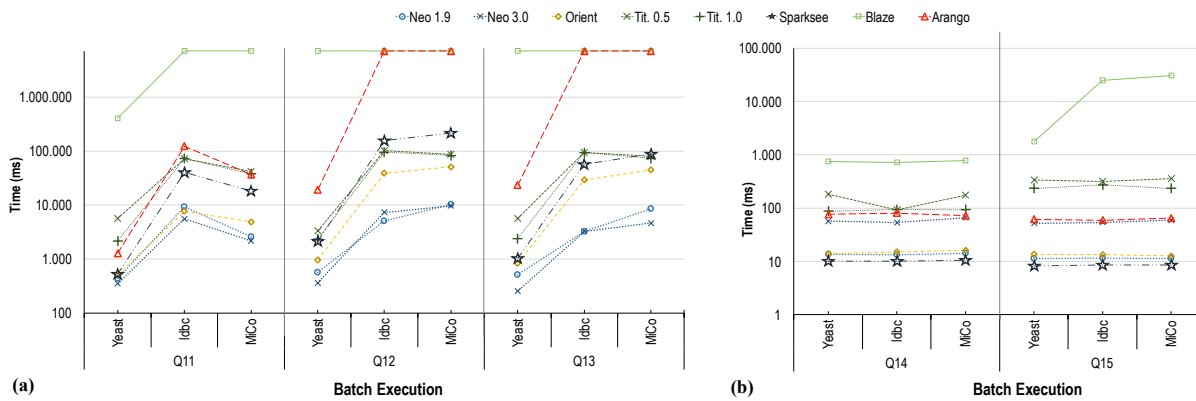


Figure 15: Selection Queries in Batch Mode for *Yeast*, *ldbc*, and *MiCo*: The Id-based (right) perform orders of magnitude better than the rest (left), and compared to the isolation mode they take the same amount of time

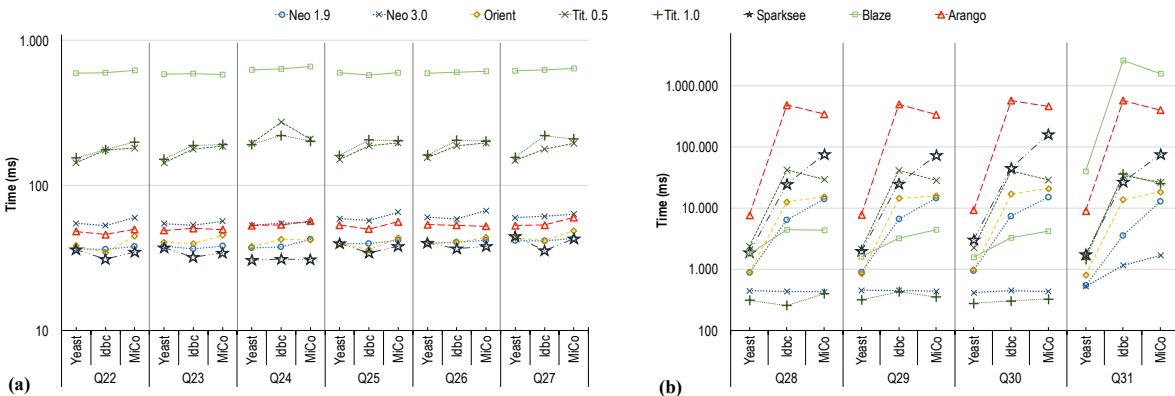


Figure 16: Time required for traversal operations in batch mode for *Yeast*, *ldbc*, and *MiCo*: (a) local access to node edges, and (b) for shortest path search

back-end system itself. For BlazeGraph, with the exception of the smallest datasets, we had to activate a specific *bulk loading* feature otherwise we were facing loading times in the order of days. Still, as mentioned, this was not enough when we tried to load the *Frb-L* sample. OrientDB as well required us to pass through some server-side implementation-specific

commands in order to load the datasets. In particular, it didn't support non-alphanumeric characters in edge-label, and for the Freebase samples we had to disable some features that were limiting the maximum number of edge-labels. Finally, Titan (in both versions) for any medium to large sized datasets requires disabling the automatic schema creation

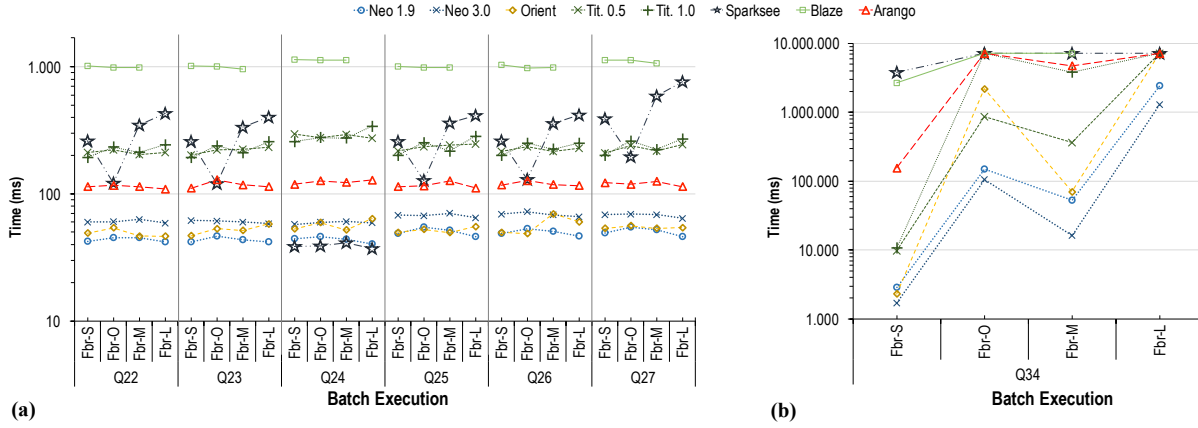


Figure 17: Time required for traversal operations in batch mode for *Yeast*, *ldbc*, and *MiCo*: (a) local access to node edges, and (b) for shortest path search

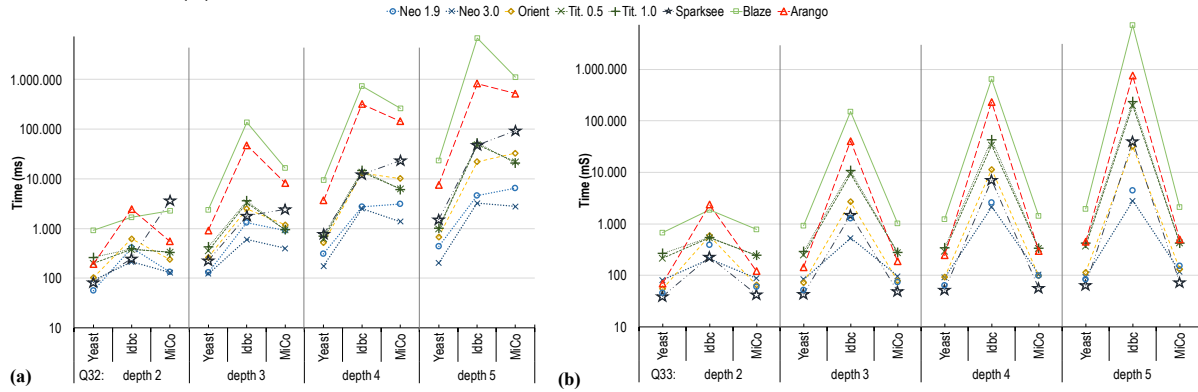


Figure 18: Time required for breadth-first traversal in batch mode for *Yeast*, *ldbc*, and *MiCo* (a) at depth=2, and (b) at depth ≥ 3

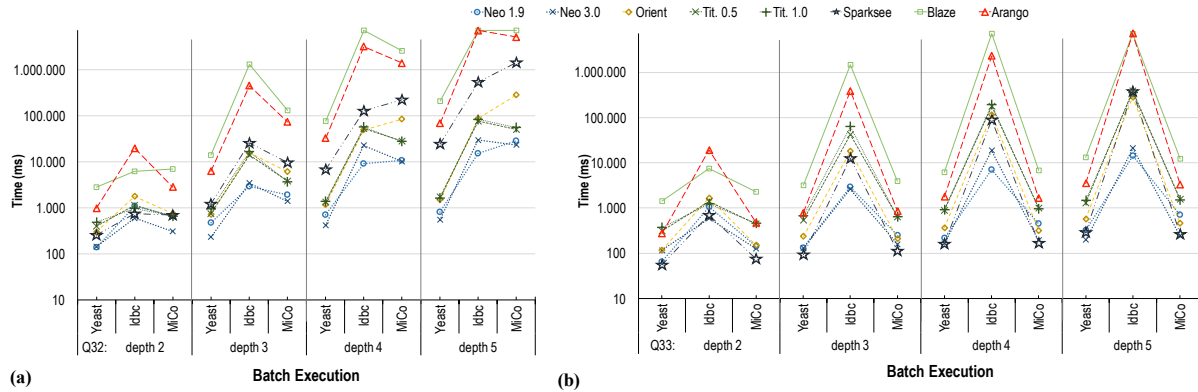


Figure 19: Time required for breadth-first traversal in batch mode for *Yeast*, *ldbc*, and *MiCo* (a) at depth=2, and (b) at depth ≥ 3

during loading, otherwise its storage back-end (Cassandra) would get swamped with extra consistency check operations. This means that the complete schema of the graph, in terms of node and edge labels and properties, should be known to the system prior to the insertion of the data, the same way one should declare the schema in a relational database before loading any data. This required us to issue a set of instructions, before loading the data, to create such schema.

Queries, Groovy, and Gremlin. Last, we report that using Groovy as support language for Gremlin was quite problematic in some cases. As a matter of fact the Groovy

language has dynamic types, and uses type inference along with peculiar handling of variable scope. As a result, explicit type casting is needed when providing the values to queries in some systems, especially with numbers. For example, in Sparksee if one attribute is of *Long* type and size (i.e., larger than a 32 bit number), then all values for the attributes with the same name need to be passed and queried as *Long* values, otherwise values compatible with the *Integer* type will be treated as such, and the search will result in a mismatch, independently of the value they represent. For Titan, instead, when not provided by the schema declared a

priory, each value should be inserted as the smaller available type, i.e., if a number is within the integer range, it should be converted to the integer type. With the other systems instead, types are handled transparently for the user, and work without explicit type casts.

A second problem with Gremlin is the lack of explicit operations for pattern-matching queries and shortest paths queries. Both types of queries could be implemented with the composition of basic constructors (although for weighted shortest path the implementation would be extremely hard), while would be better to have an abstract operator in the language and leave to the engine the implementation of advanced and optimized algorithms.

Finally, Gremlin doesn't provide a way to handle indexes, this as well is a limitation of the language that requires for the user to access directly the back-end system.

8. CONCLUSIONS

We performed an extensive experimental evaluation of the state-of-the-art graph databases. We scaled to levels that have not been considered before, and included systems that have not been previously considered. Furthermore, we provided a principled and systematic evaluation methodology based on micro-benchmarking that contains 35 different operations. We also described the challenges we faced in loading the large datasets and running the queries, and how we overcame these challenges. We materialized our methodology into a suite that we made available on-line.⁹ It includes, scripts, datasets, and queries, among any other interesting material. To the best of our knowledge, our study is the most complete and up-to-date study of graph databases available nowadays. Apart from the direct benefits, our work can complement studies on the different (but highly related) graph analytic systems.

One of the features not tested yet is parallelism, which is part of our future work. In fact all our experiments were conducted on a single machine, not exploiting any of the parallel features that almost all system provide. Nevertheless, it is important to notice that many systems advertise their ability to scale to multiple machines more than other features, but they seemed unable to exploit at best the resources of a single machine. In particular, in some cases even simple queries for relatively small db sizes were taking 2 or more hours to complete.

References

[1] Apache Cassandra. <http://cassandra.apache.org>.
 [2] Apache Hbase. <http://hbase.apache.org>.
 [3] Apache Lucene. <http://lucene.apache.org>.
 [4] Apache Mesos. <http://mesos.apache.org>.
 [5] Apache tinkerspop. <http://tinkerspop.apache.org/>.
 [6] Arangodb. <https://www.arangodb.com/>.
 [7] BerkeleyDB. <http://www.oracle.com/technetwork/products/berkeleydb>.
 [8] Docker inc., docker. <https://www.docker.com/>.
 [9] Elasticsearch. <http://www.elastic.co/products/elasticsearch>.
 [10] Infinitegraph. <http://www.objectivity.com/products/infinitegraph>.
 [11] Neo technology, inc., neo4j. <http://neo4j.com>.
 [12] Ontotext graphdb. <http://graphdb.ontotext.com/>.
 [13] Orient technologies, orientdb. <http://orientdb.com/orientdb/>.
 [14] Sparsity technologies, sparksee. <http://www.sparsity-technologies.com/>.

[15] Systap, llc., blazegraph. <https://www.blazegraph.com/>.
 [16] Thinkaurelius, titan. <http://titan.thinkaurelius.com/>.
 [17] B. Alexe, W. C. Tan, and Y. Velegrakis. Stbenchmark: towards a benchmark for mapping systems. *PVLDB*, 1(1):230–244, 2008.
 [18] R. Angles. A comparison of current graph database models. In *Proceedings of the 2012 IEEE 28th International Conference on Data Engineering Workshops, ICDEW '12*, pages 171–177, Washington, DC, USA, 2012. IEEE Computer Society.
 [19] R. Angles, P. Boncz, J. Larriba-Pey, I. Fundulaki, T. Neumann, O. Erling, P. Neubauer, N. Martínez-Bazan, V. Kotsev, and I. Toma. The linked data benchmark council: A graph and rdf industry benchmarking effort. *SIGMOD Rec.*, 43(1):27–31, May 2014.
 [20] R. Angles and C. Gutierrez. Survey of graph database models. *ACM Comput. Surv.*, 40(1):1:1–1:39, Feb. 2008.
 [21] R. Angles, A. Prat-Pérez, D. Dominguez-Sal, and J.-L. Larriba-Pey. Benchmarking database systems for social network applications. In *First International Workshop on Graph Data Management Experiences and Systems, GRADES '13*, pages 15:1–15:7, New York, NY, USA, 2013. ACM.
 [22] V. Batagelj and A. Mrvar. Yeast, pajek dataset. <http://vlado.fmf.uni-lj.si/pub/networks/data/>, 2006. <http://vlado.fmf.uni-lj.si/pub/networks/data/>.
 [23] C. Boettiger. An introduction to docker for reproducible research. *SIGOPS Oper. Syst. Rev.*, 49(1):71–79, Jan. 2015.
 [24] H. Boral and D. J. Dewitt. A methodology for database system performance evaluation. In *Proceedings of the International Conference on Management of Data*, pages 176–185, 1984.
 [25] D. Bu, Y. Zhao, L. Cai, H. Xue, X. Zhu, H. Lu, J. Zhang, S. Sun, L. Ling, N. Zhang, et al. Topological structure analysis of the protein–protein interaction network in budding yeast. *Nucleic acids research*, 31(9):2443–2450, 2003.
 [26] M. Capotá, T. Hegeman, A. Iosup, A. Prat-Pérez, O. Erling, and P. Boncz. Graphalytics: A big data benchmark for graph-processing platforms. In *Proceedings of the GRADES'15, GRADES'15*, pages 7:1–7:6, New York, NY, USA, 2015. ACM.
 [27] J. Cheng, Y. Ke, S. Chu, and M. T. Oszu. Efficient core decomposition in massive networks. In *Proceedings of the 2011 IEEE 27th International Conference on Data Engineering, ICDE '11*, pages 51–62, Washington, DC, USA, 2011. IEEE Computer Society.
 [28] D. Dominguez-Sal, P. Urbón-Bayes, A. Giménez-Vañó, S. Gómez-Villamor, N. Martínez-Bazán, and J. L. Larriba-Pey. Survey of graph database performance on the hpc scalable graph analysis benchmark. In *Proceedings of the 2010 International Conference on Web-age Information Management, WAIM'10*, pages 37–48, Berlin, Heidelberg, 2010. Springer-Verlag.
 [29] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Mach. Learn.*, 56(1-3):9–33, June 2004.
 [30] M. Elseidy, E. Abdelhamid, S. Skiadopoulos, and P. Kalnis. Grami: Frequent subgraph and pattern mining in a single large graph. *Proc. VLDB Endow.*, 7(7):517–528, Mar. 2014.
 [31] O. Erling, A. Averbuch, J. Larriba-Pey, H. Chafi, A. Gubichev, A. Prat, M.-D. Pham, and P. Boncz. The ldbc social network benchmark: Interactive workload. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, SIGMOD '15*, pages 619–630, New York, NY, USA, 2015. ACM.
 [32] J. Fan, A. G. S. Raj, and J. M. Patel. The case against specialized graph analytics engines. In *CIDR 2015, Seventh Biennial Conference on Innovative Data Systems Research, Asilomar, CA, USA, January 4-7, 2015, Online Proceedings*, 2015.
 [33] Google. Freebase data dumps. <https://developers.google.com/freebase/data>, 2015.
 [34] O. Goonetilleke, S. Sathe, T. Sellis, and X. Zhang. Microblogging queries on graph databases: An introspection. In *Proceedings of the GRADES'15, GRADES'15*, pages 5:1–5:6, New York, NY, USA, 2015. ACM.
 [35] M. Han, K. Daudjee, K. Ammar, M. T. Oszu, X. Wang, and T. Jin. An experimental comparison of pregel-like graph processing systems. *Proc. VLDB Endow.*, 7(12):1047–1058, Aug. 2014.
 [36] F. Holzschuher and R. Peinl. Performance of graph query languages: Comparison of cypher, gremlin and native access in neo4j. In *Proceedings of the Joint EDBT/ICDT 2013 Work-*

⁹<https://disi.unitn.it/~lissandrini/exemplar.html>

- shops*, EDBT '13, pages 195–204, New York, NY, USA, 2013. ACM.
- [37] E. Ioannou, N. Rassadko, and Y. Velegrakis. On generating benchmark data for entity matching. *J. Data Semantics*, 2(1):37–56, 2013.
- [38] S. Jouili and V. Vansteenbergh. An empirical comparison of graph databases. In *Proceedings of the 2013 International Conference on Social Computing, SOCIALCOM '13*, pages 708–715, Washington, DC, USA, 2013. IEEE Computer Society.
- [39] V. Kolomičenko, M. Svoboda, and I. H. Mlýnková. Experimental comparison of graph databases. In *Proceedings of International Conference on Information Integration and Web-based Applications & Services, IIWAS '13*, pages 115:115–115:124, New York, NY, USA, 2013. ACM.
- [40] M. Lissandrini. Freebase exq data dump. <https://disi.unitn.it/~lissandrini/notes/freebase-data-dump.html>, 2017.
- [41] M. Lissandrini, D. Mottin, T. Palpanas, D. Papadimitriou, and Y. Velegrakis. Unleashing the power of information graphs. *SIGMOD Rec.*, 43(4):21–26, Feb. 2015.
- [42] Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein. Distributed graphlab: A framework for machine learning and data mining in the cloud. *Proc. VLDB Endow.*, 5(8):716–727, Apr. 2012.
- [43] Y. Lu, J. Cheng, D. Yan, and H. Wu. Large-scale distributed graph computing systems: An experimental evaluation. *Proc. VLDB Endow.*, 8(3):281–292, Nov. 2014.
- [44] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: A system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, SIGMOD '10, pages 135–146, New York, NY, USA, 2010. ACM.
- [45] N. Martínez-Bazan, M. A. Águila Lorente, V. Muntés-Mulero, D. Domínguez-Sal, S. Gómez-Villamor, and J.-L. Larriba-Pey. Efficient graph management based on bitmap indices. In *Proceedings of the 16th International Database Engineering & Applications Symposium, IDEAS '12*, pages 110–119, New York, NY, USA, 2012. ACM.
- [46] N. Martínez-Bazan, S. Gómez-Villamor, and F. Escalé-Claveras. Dex: A high-performance graph database management system. In *Proceedings of the 2011 IEEE 27th International Conference on Data Engineering Workshops, ICDEW '11*, pages 124–127, Washington, DC, USA, 2011. IEEE Computer Society.
- [47] F. McSherry, M. Isard, and D. G. Murray. Scalability! but at what cost? In *15th Workshop on Hot Topics in Operating Systems (HotOS XV)*, Kartause Ittingen, Switzerland, 2015. USENIX Association.
- [48] D. Mottin, M. Lissandrini, Y. Velegrakis, and T. Palpanas. Exemplar queries: A new way of searching. *The VLDB Journal*, 25(6):741–765, Dec. 2016.
- [49] J. Okajima. aufs: Advanced multi layered unification filesystem. <http://aufs.sourceforge.net/>.
- [50] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. TR 1999-66, Stanford InfoLab, Nov.
- [51] M. A. Rodríguez. The gremlin graph traversal machine and language (invited talk). In *Proceedings of the 15th Symposium on Database Programming Languages, DBPL 2015*, pages 1–10, New York, NY, USA, 2015. ACM.
- [52] A. Rowstron, D. Narayanan, A. Donnelly, G. O’Shea, and A. Douglas. Nobody ever got fired for using hadoop on a cluster. In *Proceedings of the 1st International Workshop on Hot Topics in Cloud Data Processing, HotCDP '12*, pages 2:1–2:5, New York, NY, USA, 2012. ACM.
- [53] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web*, pages 697–706. ACM, 2007.
- [54] R. Tarjan. Depth first search and linear graph algorithms. *SIAM JOURNAL ON COMPUTING*, 1(2), 1972.
- [55] D. Yan, Y. Bu, Y. Tian, A. Deshpande, and J. Cheng. Big graph analytics systems. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD '16*, pages 2241–2243, New York, NY, USA, 2016. ACM.
- [56] M. Yannakakis. Graph-theoretic methods in database theory. In *Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS '90*, pages 230–242, New York, NY, USA, 1990. ACM.