

Benchmarking graph databases on the problem of community detection

Sotirios Beis, Symeon Papadopoulos, and Yiannis Kompatsiaris

Information Technologies Institute, CERTH, 57001, Thessaloniki, Greece
{sotbeis,papadop,ikom}@iti.gr

Abstract. Thanks to the proliferation of Online Social Networks (OSNs) and Linked Data, graph data have been constantly increasing, reaching massive scales and complexity. Thus, tools to store and manage such data efficiently are absolutely essential. To address this problem, various technologies have been employed, such as relational, object and graph databases. In this paper we present a benchmark that evaluates graph databases with a set of workloads, inspired from OSN mining use case scenarios. In addition to standard network operations, the paper focuses on the problem of community detection and we propose the adaptation of the Louvain method on top of graph databases. The paper reports a comprehensive comparative evaluation between three popular graph databases, Titan, OrientDB and Neo4j. Our experimental results show that, in the current development status, OrientDB is the fastest solution with respect to the Louvain method, while Neo4j performs the query workloads fastest. Moreover, Neo4j and Titan handle better massive and single insertion operations respectively.

1 Introduction

Over the past few years there has been vivid research interest in the study of networks (graphs) arising from various social, technological and scientific activities. Typical examples of social networks are graphs constructed with data from Online Social Networks (OSNs), one of the most famous and widespread Web 2.0 application categories. The rapid growth of OSNs contributes to the creation of high-volume and velocity data, which are modeled with the use of graph structures. The increasing demand for massive graph data management and processing systems has been addressed by the researchers proposing new methods and technologies, such as RDBMS, OODBMS, graph databases, etc. Every solution has its pros and cons so benchmarks to evaluate candidate solutions with respect to specific applications are considered necessary.

Relational databases have been widely used for the storage of a variety of data, including social data, and have proven their reliability. On the other hand RDBMS lack operations to efficiently analyze the relationships among the data points. This led to the development of new systems, such as object and graph databases. More specifically, graph databases are designed to store and manage

effectively big graph data and constitute a powerful tool for graph-like queries, such as “find the friends of a person”.

In this paper we address the problem of comparing graph databases in terms of performance, focusing on the problem of community detection. We implement a clustering workload, which consists of a well-known community detection algorithm for modularity optimization, the Louvain method [1]. We employ cache techniques to take advantage of both graph database capabilities and in-memory execution speed. The use of the clustering workload is the main contribution of this paper, because to our knowledge other existing benchmark frameworks evaluate graph databases in terms of loading time, node creation/deletion or traversal operations, such as “find the friends of a person” or “find the shortest path between two persons”. Furthermore, the benchmark comprises three supplementary workloads that simulate frequently occurring operations in real-world applications, such as the the creation and traversal of the graph. The benchmark implementation is available online as an open-source project¹.

We use the proposed benchmark to evaluate three popular graph databases, *Titan*², *OrientDB*³ and *Neo4j*⁴. For our experiments we used both synthetic and real networks and the comparative evaluation is held with respect to the execution time. Our experimental results show that OrientDB is the most efficient graph database to apply community detection algorithms, in our case the Louvain method. Concerning the supplementary workloads, Neo4j is the fastest alternative, although Titan performs the incremental creation of the graph faster.

The paper is organized as follows. We begin in Section 2 by providing a survey in the area of benchmarks between database systems oriented to store and manage big graph data. In Section 3 we describe the workloads that compose the benchmark. In Section 4 we list some important aspects of the benchmark. Section 5 presents our experimental study, where we describe the datasets used for the evaluation and report the obtained experimental results. Finally, Section 6 concludes the paper and delineates our future work ideas.

2 Related Work

Until now many benchmarks have been proposed, comparing the performance of different databases for graph data. Giatsoglou et al. [2], present a survey of existing solutions to the problem of storing and managing massive graph data. Focusing on the Social Tagging System (STS) use case scenario, they report a comparative study between the Neo4j graph database and two custom storages (H1 and Lucene). Angles et al. [3], considering the category of an OSN as an example of Web 2.0 applications, propose and implement a generator that produces synthetic graphs with OSN characteristics. Using this data and a set of queries that simulate common activities in a social network application, the

¹ <https://github.com/socialsensor/graphdb-benchmarks>

² <http://thinkaurelius.github.io/titan/>

³ <http://www.orienttechnologies.com/>

⁴ <http://www.neo4j.org/>

authors compare two graph databases, one RDF and two relational data management systems. Similarly, in LinkBench [4] a Facebook-like data generator is employed and the performance of a MySQL database system is evaluated. The authors claim that under certain circumstances any database system could be evaluated with LinkBench.

In a recent effort, Grossniklaus et al. [5] define and classify a workload of nine queries, that together cover a wide variety of graph data use cases. Besides graph databases they include RDBMS and OODBMS in their evaluation. Vicknair et al. [6] also present a benchmark that combines different technologies. They implemented a query workload that simulates typical operations performed in provenance systems and they evaluate a graph (Neo4j) and a relational (MySQL) database. Furthermore, the authors describe some objective measures to compare the database systems, such as security, flexibility, etc.

In contrast with the above works, we argue that the most suitable solution to the problem of massive graph storage and management are graph databases, so our research focuses on them. In this direction Bader et al. [7] describe a benchmark that consists of four kernels (operations): (a) bulk load of the data; (b) retrieval of a set of edges that verify a condition (e.g. $\text{weight} > 3$); (c) execution of a k -hops operation; and (d) retrieval of the set of nodes with maximum betweenness centrality. Dominguez et al. [8] report the implementation of this benchmark and a comparative evaluation of four graph database systems (Neo4j, HypergraphDB, Jena and DEX).

Ciglan et al. [9] are based on the ideas proposed in [8] and [10], and extend the discussion focusing primarily on graph traversal operations. They compare five graph databases (Neo4j, DEX, OrientDB, NativeSail and SGDB) by executing some demanding queries, such as “find the most connected component”. Jouili et al. [11] propose a set of workloads similar to [7] and evaluate Neo4j, Titan, OrientDB and DEX. Unlike, previous works they conduct experiments with multiple concurrent users and emphasize the effects of increasing users. Dayarathna et al. [12] implement traversal operation-based workloads to compare four graph databases (Allegrograph, Neo4j, OrientDB and Fuseki). The key difference with other frameworks is that their interest is focused mostly on graph database server and cloud environments.

3 Workload Description

The proposed benchmark is composed of four workloads, Clustering, Massive Insertion, Single Insertion and Query Workload. Every workload has been designed to simulate common operations in graph database systems. Our main contribution is the Clustering workload (CW), however supplementary workloads are employed to achieve a comprehensive comparative evaluation. In this section we describe in more detail the workloads and emphasize their importance by giving some real-world examples.

3.1 Clustering Workload

Until now most community detection algorithms used the main memory to store the graph and perform the required computations. Although, keeping data in memory leads to fast executions times, these implementations have a major drawback: they cannot manage big graph data reliably, which nowadays is a key requirement for big graph processing applications. This motivated this work and more specifically the implementation of the Louvain method on top of three graph databases. We used the Gephi Toolkit⁵ Java implementation of the algorithm as a starting point and applied all necessary modifications to adapt the algorithm to graph databases.

In a first implementation, all the required values for the computations were read directly from the database. The fact that the access of any database (including graph databases) compared to memory is very slow, soon made us realize that the use of cache techniques is necessary. For this purpose we employed the cache implementation of the Guava project⁶. The Guava Cache is configured to evict entries automatically, in order to constrain its memory footprint. Guava provides three basic types of eviction: size-based eviction, time-based eviction, and reference-based eviction. To precisely control the maximum cache size, we utilize the first type of eviction, size-based, and the evaluation was held both between different systems and among different cache sizes. The measurements concern the required time for the algorithm to be completed.

As the authors of the Louvain method mention⁷, the algorithm is a greedy optimization method that attempts to optimize the modularity of a partition of the network. The optimization is performed in two steps. First, the method looks for “small” communities by optimizing modularity locally. Second, it aggregates nodes belonging to the same community and builds a new network whose nodes are the communities. We call those *communities* and *nodeCommunities* respectively. The above steps are repeated in an iterative manner until a maximum of modularity is attained.

We keep the community and nodeCommunity values stored in the graph database as a property of each node. The implementation is based on three functions that retrieve the required information either by accessing the cache or the database directly. We store this information employing the LoadingCache structure from the Guava Project, which is similar to a ConcurrentMap⁸. More specifically we use the following functions and structures:

- *getNodeNeighbours*: gets the neighbours of a node and stores them to a LoadingCache structure, where the key is the node id and the value is the set of neighbours.

⁵ <https://gephi.org/toolkit/>

⁶ <https://code.google.com/p/guava-libraries/>

⁷ <http://perso.uclouvain.be/vincent.blondel/research/louvain.html>

⁸ <http://docs.oracle.com/javase/7/docs/api/java/util/concurrent/ConcurrentMap.html>

- *getNodesFromCommunity*: gets the nodes from a specific community and stores them to a LoadingCache structure, where the key is the community id and the value is the the set of nodes that the community contains.
- *getNodesFromNodeCommunity*: gets the nodes from a specific nodeCommunity and stores them to a LoadingCache structure, where the key is the nodeCommunity id and the value is the the set of nodes that the nodeCommunity contains.

We use the above information to compute values such as, the degree of a node, the amount of connections a node or a nodeCommunity has with a particular community, the size of a community or a nodeCommunity.

The Clustering Workload is very important due to its numerous applications in OSNs [13]. Some of the most representative examples include topic detection in collaborative tagging systems, such as Flickr or Delicious, tag disambiguation, user profiling, photo clustering, and event detection.

3.2 Supplementary Workloads

In addition to CW, we recognize that a reliable and comprehensive benchmark should contain some supplementary workloads. Here, we list and describe the three additional workloads that constitute the proposed benchmark.

- *Massive Insertion Workload (MIW)*: we create the graph database and configure it for massive loading, then we populate it with a particular dataset. We measure the time for the creation of the whole graph.
- *Single Insertion Workload (SIW)*: we create the graph database and load it with a particular dataset. Every object insertion (node or edge) is committed directly and the graph is constructed incrementally. We measure the insertion time per block, which consists of one thousand edges and the nodes that appear during the insertion of these edges.
- *Query Workload (QW)*: we execute three common queries:
 - FindNeighbours (FN): finds the neighbours of all nodes.
 - FindAdjacentNodes (FA): finds the adjacent nodes of all edges.
 - FindShortestPath (FS): finds the shortest path between the first node and 100 randomly picked nodes.

Here we measure the execution time of each query.

It is obvious that MIW simulates the case study in which graph data are available and we want to load them in batch mode. On the other hand, SIW models a more real-time scenario in which the graph is created progressively. We could claim that the growth of an OSN follows the steps of SIW, by adding more users (nodes) and relationships (edges) between them.

The QW is very important as it applies in most of the existing OSNs. For example with the FN query we can find the friends or followers of a person in Facebook or Twitter respectively, with the FA query we can find whether two users joined a particular Facebook group and with the FS query we can find at which level two users connect with each other in LinkedIn. It is critical for every OSN that these queries can be executed efficiently and in minimal time.

4 Benchmark Description

In this section we discuss some important aspects of the benchmark implementation. The graph database systems selected for the evaluation are Titan (v0.5.0), OrientDB (v2.0-M1) and Neo4j (v2.1.3). The benchmark was implemented in Java 1.7 using the Java API of each database. In order to configure each database, we used the default configuration and the recommendations found in the documentation of the web sites.

For Titan we implement MIW with the BatchGraph interface that enables batch loading of a large number of edges and vertices, while for OrientDB and Neo4j we employ the OrientGraphNoTx and BatchInserter interface respectively, which drop the support for transactions in favor of insertion speed. For all graph databases we implement SIW without using any special configuration. The operations for the QW and CW were implemented using the respective API of each database, except the Titan’s FindShortestPath implementation where we used the Gremlin API.

To ensure that a benchmark provides meaningful and trustworthy results, it is necessary to guarantee its fairness and accuracy. There are many aspects that can influence the measurements, such as the system overhead. It is really important that the results do not come from time periods with different system status (e.g. different number of processes in the background), so we execute MIW, SIW and QW sequentially for each database. In addition to this, we execute them in every possible combination for each database, in order to minimize the possibility that the results are affected by the order of execution. We report the mean value of all measurements.

Regarding the CW, in order to eliminate the cold cache effects we execute it twice and keep always the second value. Moreover, as we described in the previous section to get an acceptable execution time, cache techniques are necessary. The cache size is defined as a percentage of total nodes. For our experiments we use six different cache sizes (5%, 10%, 15%, 20%, 25%, 30%) and we report the respective improvements.

5 Experimental Study

In this section we present the experimental study. At first we describe the datasets used for the evaluation. We include a table with some important statistics of each dataset. Then we report and discuss the results.

5.1 Datasets

The right choice of datasets that will be used for running database benchmarks is important to obtain representative and meaningful results. It is necessary to test the databases on a sufficient number of datasets of different sizes and complexity to get an approximation of the database scaling properties.

For our evaluation we use both synthetic and real data. More specifically, we execute MIW, SIW and QW with real networks derived from the SNAP dataset collection⁹. On the other hand, with the CW we use synthetic data generated with the LFR-Benchmark generator [1] that produces networks with power-law degree distribution and implanted communities within the network. The Table 1 presents the summary statistics of the datasets.

Table 1: Datasets used in the experiments

Dataset	Nodes	Edges	max. κ	$\langle \kappa \rangle$	$\langle cc \rangle$
Graph1k	1,000	7,578	150	15.156	0.404
Graph5k	5,000	74,097	450	29.639	0.445
Graph10k	10,000	180,314	750	36.063	0.446
Graph20k	20,000	389,448	1,320	38.945	0.420
Graph30k	30,000	666,008	1,750	44.401	0.410
Graph40k	40,000	1,006,945	2,000	50.347	0.395
Graph50k	50,000	1,256,044	2,750	50.242	0.436
Enron (EN)	36,692	367,662	1,383	20.041	0.497
Amazon (AM)	334,863	925,872	168	5.530	0.398
Youtube (YT)	1,134,890	2,987,624	28,576	5.265	0.081
Livejournal (LJ)	3,997,962	34,681,189	14,703	17.349	0.045

5.2 Benchmark Results

In this section we report and discuss the performance of Titan, OrientDB and Neo4j employing the proposed benchmark. Table 2 lists the required time for the execution of MIW and QW, while Figure 1 illustrates the experimental results of SIW. Table 3 and Figure 2 depict the measurements of CW. Note that in every table we mark the best performance with bold. All experiments were run on an Intel Core i7 at 3.5Ghz with 16GB of main memory and a 1.4 TB hard disk, the OS being Ubuntu Linux 12.04 (64bit).

Table 2 summarizes the measurements of the MIW and QW for all the benchmarked graph databases with respect to each real dataset. According to the benchmark results, we observe that Neo4j handles the massive insertion of the data more efficiently from its competitors. Titan is also an effective alternative, while OrientDB could not load the data in a comparable time.

Concerning the QW, Table 2 indicates that Neo4j performs queries more effectively than the other candidates. More specifically, although OrientDB has slightly smaller execution time comparing to Neo4j in the FN query load for the Enron dataset, Neo4j is considerably faster in all other cases. It is worth

⁹ <http://snap.stanford.edu/data/index.html>

mentioning that the shortest path search is limited to paths of depth 6, because with larger depth, the FS query workload in Titan cannot be executed in a reasonable amount of time.

Table 2: MIW and QW results (sec)

Graph	Workload	Titan	OrientDB	Neo4j
EN	MIW	9.36	62.77	6.77
AM	MIW	34.00	97.00	10.61
YT	MIW	104.27	252.15	24.69
LJ	MIW	663.03	9416.74	349.55
EN	QW-FN	1.87	0.56	0.95
AM	QW-FN	6.47	3.50	1.85
YT	QW-FN	20.71	9.34	4.51
LJ	QW-FN	213.41	303.09	47.07
EN	QW-FA	3.78	0.71	0.16
AM	QW-FA	13.77	2.30	0.36
YT	QW-FA	42.82	6.15	1.46
LJ	QW-FA	460.25	518.12	47.07
EN	QW-FS	1.63	3.09	0.16
AM	QW-FS	0.12	83.29	0.302
YT	QW-FS	24.87	23.47	0.08
LJ	QW-FS	123.50	86.87	18.13

The results of SIW for each real dataset are illustrated in Figure 1. Each sub-figure includes three diagrams, one for every graph database, that plot the required time for the insertion of a block. As we described in Section 3, a block consists of 1,000 edges and the nodes that appear during the insertion of these edges. In order to present more readable diagrams for the three technologies we used a logarithmic scale for the time axis. It appears that Titan is the most efficient solution for single insertion of data. Moreover, we observe that the performance of OrientDB and Neo4j is comparable, however OrientDB seems to perform much better.

The experimental results of CW are reported in Table 3. We observe that OrientDB is considerably faster than its competitors. Moreover Table 3 indicates that while Titan has comparable execution times with Neo4j for small graphs, it does not scale as good as Neo4j. Thus, for graphs with >1,000 nodes, Neo4j is much faster.

Additionally, Table 3 points out the positive impact of increasing the cache size. We observe that for all graph databases regardless of the graph size, as the cache size increases the execution time decreases. We wrap up the comparative evaluation with Figure 2 that depicts the scalability of each database when the

CW is executed. Every sub-figure contains six diagrams, one for each cache value, that plot the required time for the convergence of the algorithm for the respective synthetic graph. For better representation we used a logarithmic scale for the time axis. We can deduce that since the diagrams are linear in logarithmic scale, the actual execution time should grow exponentially with the graph size.

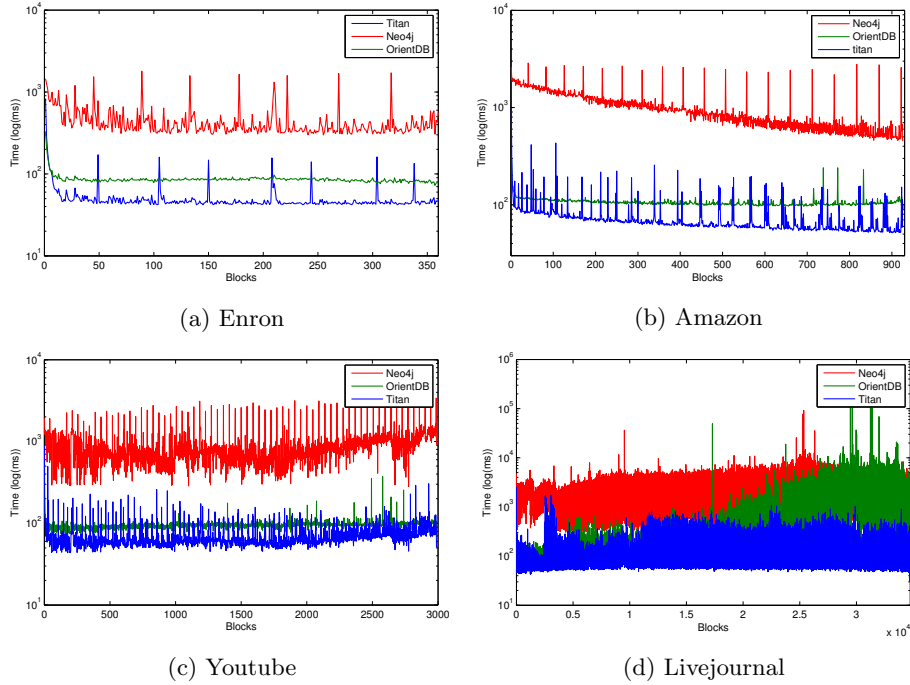


Fig. 1: SIW benchmark results

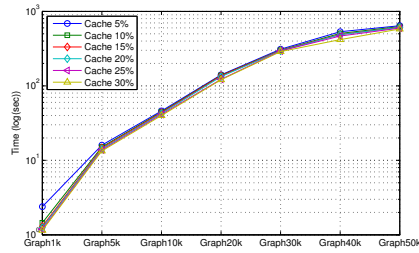
In summary, we found that OrientDB is clearly the best candidate to apply community detections algorithms (the Louvain method in our case) and Neo4j is the most efficient solution for the MIW and QW. On the other hand, Titan is the fastest alternative for the incremental creation of a graph database (SIW). Titan also has competitive performance in the MIW, but does not scale very well compared to its two competitors.

6 Conclusions and Future Work

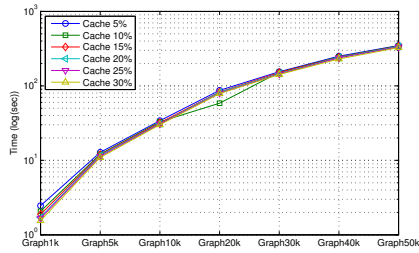
In this paper we proposed a benchmark framework for the comparative evaluation of database systems oriented to store and manage graph data. The benchmark consists of four workloads, Massive Insertion, Single Insertion, Query and Clustering Workload. For the Clustering Workload we implemented a well-known

Table 3: CW results (sec)

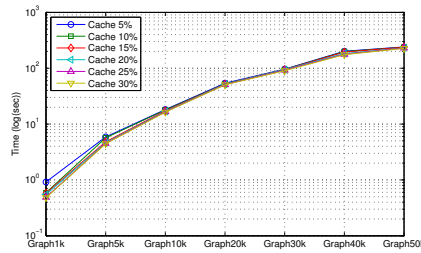
Graph-Cache	Titan	OrientDB	Neo4j
Graph1k-5%	2.39	0.92	2.46
Graph1k-10%	1.45	0.59	2.07
Graph1k-15%	1.30	0.58	1.88
Graph1k-20%	1.25	0.55	1.72
Graph1k-25%	1.19	0.49	1.67
Graph1k-30%	1.15	0.48	1.55
Graph5k-5%	16.01	5.88	12.8
Graph5k-10%	15.10	5.67	12.13
Graph5k-15%	14.63	4.81	11.91
Graph5k-20%	14.16	4.62	11.68
Graph5k-25%	13.76	4.51	11.31
Graph5k-30%	13.38	4.45	10.94
Graph10k-5%	46.06	18.20	34.05
Graph10k-10%	44.59	17.92	32.88
Graph10k-15%	43.68	17.31	31.91
Graph10k-20%	42.48	16.88	31.01
Graph10k-25%	41.32	16.58	30.74
Graph10k-30%	39.98	16.34	30.13
Graph20k-5%	140.46	54.01	87.04
Graph20k-10%	138.10	52.51	85.49
Graph20k-15%	137.25	52.12	82.88
Graph20k-20%	133.11	51.68	82.16
Graph20k-25%	122.48	50.79	79.87
Graph20k-30%	120.94	50.49	78.81
Graph30k-5%	310.25	69.38	154.60
Graph30k-10%	301.80	94.98	151.81
Graph30k-15%	299.27	94.85	151.12
Graph30k-20%	296.43	94.67	146.25
Graph30k-25%	294.33	92.62	144.08
Graph30k-30%	288.50	90.13	142.33
Graph40k-5%	533.29	201.19	250.79
Graph40k-10%	505.91	199.18	244.79
Graph40k-15%	490.39	194.34	242.55
Graph40k-20%	487.31	183.14	241.47
Graph40k-25%	467.18	177.55	237.29
Graph40k-30%	418.07	174.65	229.65
Graph50k-5%	642.42	240.58	348.33
Graph50k-10%	624.36	238.35	344.06
Graph50k-15%	611.70	237.65	340.20
Graph50k-20%	610.40	230.76	337.36
Graph50k-25%	596.29	230.03	332.01
Graph50k-30%	580.44	226.31	325.88



(a) Titan



(b) Neo4j



(c) OrientDB

Fig. 2: CW benchmark results

community detection algorithm, the Louvain method, on top of three graph databases. Employing the proposed benchmark we evaluated the selected graph databases, Titan, OrientDB and Neo4j using both synthetic and real networks.

The experimental results demonstrate that in most cases the measurements are comparable when processing small graphs. But when the size of the datasets grows significantly, Neo4j appears to be the most efficient solution for storing and querying graph data. On the other hand, when there is a need for successive local queries, like the ones that the Louvain method employs, OrientDB is the best candidate. Last, Titan seems to be the best alternative for single insertion operations.

In the future we hope to investigate the performance gain after we parallelize the operations of the graph databases. Moreover, it would be interesting to run the benchmark employing the distributed implementations of Titan and OrientDB in order to examine their horizontal and vertical scalability properties. Also, we intend to improve the performance of the implemented community detection algorithm and test it on graphs of much larger size.

Acknowledgments

This work was supported by the SocialSensor FP7 project, partially funded by the EC under grant agreement 287975.

References

1. Blondel, V.D., Guillaume, J.L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* **2008**(10) (2008) P10008
2. Giatsoglou, M., Papadopoulos, S., Vakali, A.: Massive graph management for the web and web 2.0. In Vakali, A., Jain, L., eds.: *New Directions in Web Data Management 1*. Volume 331 of *Studies in Computational Intelligence*. Springer Berlin Heidelberg (2011) 19–58
3. Angles, R., Prat-Pérez, A., Dominguez-Sal, D., Larriba-Pey, J.L.: Benchmarking database systems for social network applications. In: *First International Workshop on Graph Data Management Experiences and Systems*. GRADES '13, New York, NY, USA, ACM (2013) 15:1–15:7
4. Armstrong, T.G., Ponnkantti, V., Borthakur, D., Callaghan, M.: Linkbench: a database benchmark based on the facebook social graph. (2013)
5. Grossniklaus, M., Leone, S., Zäschke, T.: Towards a benchmark for graph data management and processing. (2013)
6. Vicknair, C., Macias, M., Zhao, Z., Nan, X., Chen, Y., Wilkins, D.: A comparison of a graph database and a relational database: A data provenance perspective. In: *Proceedings of the 48th Annual Southeast Regional Conference*. ACM SE '10, New York, NY, USA, ACM (2010) 42:1–42:6
7. Bader, D.A., Feo, J., Gilbert, J., Kepner, J., Koester, D., Loh, E., Madduri, K., Mann, B., Meuse, T., Robinson, E.: Hpc scalable graph analysis benchmark (2009)
8. Dominguez-Sal, D., Urbn-Bayes, P., Gimnez-Va, A., Gmez-Villamor, S., Martinez-Bazan, N., Larriba-Pey, J.: Survey of graph database performance on the hpc scalable graph analysis benchmark. In Shen, H., Pei, J., zsu, M., Zou, L., Lu, J., Ling, T.W., Yu, G., Zhuang, Y., Shao, J., eds.: *Web-Age Information Management*. Volume 6185 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg (2010) 37–48
9. Ciglan, M., Averbuch, A., Hluchy, L.: Benchmarking traversal operations over graph databases. In: *Data Engineering Workshops (ICDEW)*, 2012 IEEE 28th International Conference on. (April 2012) 186–189
10. Dominguez-Sal, D., Martinez-Bazan, N., Munes-Mulero, V., Baleta, P., Larriba-Pey, J.: A discussion on the design of graph database benchmarks. In Nambiar, R., Poess, M., eds.: *Performance Evaluation, Measurement and Characterization of Complex Systems*. Volume 6417 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg (2011) 25–40
11. Jouili, S., Vansteenbergh, V.: An empirical comparison of graph databases. In: *Social Computing (SocialCom)*, 2013 International Conference on. (Sept 2013) 708–715
12. Dayarathna, M., Suzumura, T.: Xgdbench: A benchmarking platform for graph stores in exascale clouds. In: *Cloud Computing Technology and Science (Cloud-Com)*, 2012 IEEE 4th International Conference on. (Dec 2012) 363–370
13. Papadopoulos, S., Kompatsiaris, Y., Vakali, A., Spyridonos, P.: Community detection in social media. *Data Mining and Knowledge Discovery* **24**(3) (2012) 515–554