

# Graph Generators: State of the Art and Open Challenges

ANGELA BONIFATI, Lyon 1 University, France

IRENA HOLUBOVÁ, Charles University, Prague, Czech Republic

ARNAU PRAT-PÉREZ, DAMA-UPC, Universitat Politècnica de Catalunya, Spain

SHERIF SAKR, University of Tartu, Estonia

The abundance of interconnected data has fueled the design and implementation of graph generators reproducing real-world linking properties, or gauging the effectiveness of graph algorithms, techniques and applications manipulating these data. We consider graph generation across multiple subfields, such as graph databases, Semantic Web, social networks, graph streaming, alongside community detection and graph data analytics. Despite the disparate requirements of modern graph generators throughout these communities, we analyze them under a common umbrella, reaching out the functionalities, the practical usage, and their supported operations. We argue that this classification is serving the need of providing scientists, researchers and practitioners with the right data generator at hand for their work.

This survey provides a comprehensive overview of the state-of-the-art graph generators by focusing on those that are pertinent and suitable for several data-intensive tasks. Finally, we discuss open challenges and missing requirements of current graph generators along with their future extensions to new emerging fields.

CCS Concepts: •**Mathematics of computing** → **Graph theory**; •**Theory of computation** → **Graph algorithms analysis**; •**Information systems** → *Graph-based database models*;

Additional Key Words and Phrases: Big Data management, graph data, generators, benchmarks, synthetic data

## ACM Reference Format:

A. Bonifati, I. Holubová, A. Prat-Pérez, S. Sakr, 2018. Modern Graph Generators: A Survey. *ACM CSUR* 0, 0, Article 0 (2018), 35 pages.

DOI: 0000001.0000001

## 1. INTRODUCTION

Graphs are ubiquitous data structures that span a wide array of scientific disciplines and are a subject of study in several subfields of computer science. Nowadays, due to the dawn of numerous tools and applications manipulating graphs, they are adopted as a rich data model in many data-intensive use cases involving disparate domain knowledge, mathematical foundations and computer science. Interconnected data is often-times used to encode domain-specific use cases, such as recommendation networks, social networks, protein-to-protein interactions, geolocation networks, and fraud detection analysis networks, to name a few.

In general, we can distinguish between two broad types of graph data sets: (1) a single large graph (possibly with several components), such as social networks or Linked

---

Author's addresses: A. Bonifati, Université Claude Bernard Lyon 1, Liris CNRS Campus La Doua 23-25 Avenue Pierre de Coubertin, 69622 Villerubanne, France; I. Holubová, Department of Software Engineering, Faculty of Mathematics and Physics, Charles University, Malostranské nám. 25, 118 00 Praha 1, Czech Republic; A. Prat-Pérez, DAMA-UPC, Universitat Politècnica de Catalunya, Jordi Girona 1-3, 08034 Barcelona, Spain; S. Sakr, University of Tartu, J Liivi 2, Tartu 50409, Estonia.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 ACM. 1539-9087/2018/-ART0 \$15.00

DOI: 0000001.0000001

Data graphs, and (2) a large set of small graphs, such as chemical compounds<sup>1</sup> or linguistics syntax trees<sup>2</sup>. Naturally, the algorithms used in these two classes differ a lot [Sakr and Pardede 2011]. In the former case we can search, e.g., for communities and their features or shortest paths, while in the latter case we usually query for supergraphs, subgraphs or graphs similar to a given graph pattern. In both the cases, as in the other fields, quite often the respective real-world graph data is not publicly available (or simply does not exist when a particular method for data manipulation is proposed). Even in the cases in which real data is abundant, many algorithms and techniques need to be tested on various order of magnitudes of the graph sizes thus leading to the inception of configurable graph generators that reproduce real-world graph properties and provide unique tuning opportunities for algorithms and tools handling such data.

In this survey, we provide a comprehensive overview of the state-of-the-art modern graph generators by focusing on those that are pertinent and suitable for data-intensive tasks. Our aim is to cover a wide range of currently popular areas of graph data processing. In particular, we consider graph generation in the areas of graph databases, Semantic Web, social networks, graph streaming, alongside community detection and graph data analytics. Despite the disparate requirements of modern graph generators throughout these communities, we analyze them under a common umbrella, reaching out the functionalities, the practical usage, and the supported operations of graph generators. The reasons for this scope and classification are as follows:

- (1) Despite the differences of the covered areas, the requirements for modern graph data generators can be similar in particular cases. Reusing or learning from tools in other fields can thus bring new opportunities for both researchers and practitioners.
- (2) The selected classification is serving the need of providing scientists, researchers and practitioners with the right data generator at hand for their work.

To conclude the comparative study and provide a comprehensive view of the field, we also overview the most popular real-world data sets used in the respective covered areas and discuss open challenges and missing requirements of modern graph generators in view of identifying their future extensions to new emerging fields.

*Contributions.* Our survey revisits the representatives of modern graph data generators and summarizes their major features. The comprehensive review and analysis make this paper useful for motivating new graph data generation techniques as well as serving as a technical reference for selecting suitable solutions. In particular, the introductory categorization and comparative study enables the reader to quickly get his/her bearings in the field and identify a subset of generators of his/her interest. Since we do not limit the survey to a particular area of graph data management, the reader can get a broader scope in the field. Hence, while practitioners can find a solution in another previously unexpected area, researchers can identify new target areas or exploit successful results in new fields. Last but not least, to provide a global view of the state-of-the-art, we also provide an overview of most commonly used real-world testing graph datasets. And, finally, we identify general open problems of graph data synthesis and indicate possible solutions to show that it still forms a challenging and promising research area.

*Differences with prior surveys.* To the best of our knowledge, our work is the first to survey the broad landscape of graph data generators spanning different data-intensive

<sup>1</sup><https://pubchem.ncbi.nlm.nih.gov/>

<sup>2</sup><https://catalog.ldc.upenn.edu/ldc99t42>

applications and targeting many computer science subfields. In particular, we cover graph databases, Semantic Web, social networks, graph streaming, community detection, and graph data analytics. However, the literature is still lacking a comprehensive study of graph generators for many of the specific subfields mentioned above.

A limited subset of graph database generators, parallel and distributed graph processing generators, along with a few of the Semantic Web data generators presented in our survey have been discussed in a related book chapter [Bonifati et al. 2018a] while cross-comparing them with respect to input, output, supported workload, data model and query language, along with the distinguished chokepoints. However, the provided classification is inherently database-oriented. In our work, we provide a more comprehensive and broader classification that serves the purpose of letting any researcher or practitioner interested in data generation to be able to make a better choice of the desired graph generator based on its functional and goal-driven features (such as the application domain, the supported operations and the key configuration options). Moreover, in contrast with [Bonifati et al. 2018a], our work encompasses graph generators of several diverse communities, not limiting its scope to a few generators of the database and graph processing communities.

Graph generators matching graph patterns used in data mining have been studied in [Chakrabarti and Faloutsos 2006], focusing on mostly occurring patterns, such as power laws, size of graph diameters and community structure. The considered graph generators are compared in terms of graph type, degree distributions, exponentiality, diameter and community effects. We refer the reader to this survey for taxonomies involving these properties, whereas we provide here a functionality-driven taxonomy across all the categories of graph generators that we consider. We also point out that this survey is outdated as it does not consider the generators that appeared in the last decade. We fill the gap of more recent social network and graph analytics generators respectively in Section 3.4 and 3.5.

Aggarwal and Subbian [Aggarwal and Subbian 2014] have surveyed the evolution of analysis in graphs, by primarily focusing on data mining maintenance methods and on analytical quantification and explanation of the changes of the underlying networks. A brief discussion on evolutionary network data generators is carried out in the paper. The data generation of evolutionary networks is based on Densification Power Law (DPL) and shrinking diameters, community-guided attachment properties [Leskovec et al. 2005b]. Generation of graphs tackling Kronecker recursion with recursive tensor multiplication [Akoglu et al. 2008] is then considered in the above survey. We refer the readers to the aforementioned survey for evolutionary network generators and further discuss the open challenges of evolving graph data in Section 5.

*Outline.* The rest of this paper is organized as follows: Section 2 provides the opening categorization and comparison of the generators. Section 3 provides an overview of the existing graph data generators and their main features in the frame of the proposed categories. To provide a complete overview, we list the most commonly used testing data sets for the particular categories in Section 4. In Section 5, we highlight some of the challenges and open problems of graph data synthesis before we conclude in Section 6.

## 2. CLASSIFICATION AND COMPARATIVE STUDY

In order to provide the reader with a general preview of the generators and to enable finding the target solutions easily, we start the survey with a classification and comparative study of the existing tools. In general, there are various ways to classify them. We first provide an overview of the approaches used in related work after which we then introduce our approach. As mentioned in the Introduction, since this survey is

unique especially in terms of scope, our classification and comparative strategy differs as well.

The graph data generators can be classified according to various criteria. For example, [Chakrabarti et al. 2004] introduces two categories – degree-based and procedural generators. In general, *degree-based generators* (e.g., Barabasi-Albert model [Barabasi and Albert 1999]) are commonly attempting to find a graph that matches it, but without providing any information about the graph or attempting to match other graph features (e.g., eigenvalues, small diameter, ..., etc). On the other hand, *procedural generators* (e.g., R-MAT [Chakrabarti et al. 2004]) are commonly targeting simple techniques to produce graphs that are matching the characteristics of the real graphs (e.g., the power law degree distribution).

Paper [Chakrabarti and Faloutsos 2006] introduces five categories of graph models that can be synthesized: (1) *random graph models* (e.g., Erdős-Rényi [Erdos and Renyi 1960]) generated by a random process, (2) *preferential attachment models* (e.g., Barabasi-Albert model) which try to model the power law from the preferential attachment viewpoint, (3) *optimization-based models* (e.g., HOT model [Carlson and Doyle 2000]) resulting from the idea that power laws can result from resource optimizations and tolerance to risks, (4) *tensor-based models* (e.g., R-MAT) targeting a trade-off between low number of model parameters and efficiency, and (5) *internet-specific models* corresponding to hybrids using ideas from the other categories in order to suit the specific features of the graphs.

The type of the generator can also be influenced by the benchmark involving it, whereas we can distinguish, e.g., *domain-specific* benchmarks, *application-specific* benchmarks, *workload-driven* benchmarks, *microbenchmarks* etc.

## 2.1. Classification

At first we classify the generators on the basis of the respective application domains or user communities. In particular we distinguish (1) general graphs, (2) Semantic Web, (3) graph databases, (4) social networks, (5) graph analytics, (6) community detection, and (7) graph data steaming. The selected classes are not rigorously defined (e.g., they are not disjoint as we will show later), but they correspond to the currently most active research areas. Thus we believe that they form a natural first acquaintance for the reader.

In Table I we overview the key characteristics of the data generators clustered according to the respective application domains.<sup>3</sup> In particular, we show:

- Characteristics of the **domain** – its *type* (column “Type”), i.e., fixed, specified using a schema, or extracted from input data and the particular *target* domain, or, in case of a generic tool, the chosen sample domain (column “Target/sample”),
- Characteristics of *read/update operations* (columns “Read” and “Update”), i.e., whether the set of operations is fixed/generated, if it involves operation mixes (i.e., sets/sequences of operations), or if templates of operations are supported.
- Key **configuration** options:
  - whether the generator deals only with structure, or also with *properties* (column “Pro.”) of the graph (Y/N feature),
  - supported types of *distributions* (column “Distributions”) used for generating of the data,
  - *output format* (column “Output”) of the produced graph, and

<sup>3</sup>“GDBs” stands for graph databases, “SNs” stands for social networks, “A” stands for graph analytics, “Co” stands for community detection, and “St” stands for graph streaming. Value “–” means that the information is not available or relevant.

- whether the generator is *distributed* (column “Dis.”) and thus enables more efficient data generation (Y/N feature).

*Number of Generators, Size and Output Format of the Data.* As we can see, the biggest set of available generators can be found in the Semantic Web application domain, probably due to its recent popularity and a number of research groups dealing with this topic. None of these generators is natively implemented in a distributed manner and thus primarily generating output at the Big Data scale. On the other hand, the *Linked Open Data* (LOD) is expected to be large in general, however this is the case of the whole *LOD Cloud*<sup>4</sup>, but not necessarily of the particular data sets forming it.

The second large group of generators also corresponds to a popular application domain – social networks. In this case the size of the output graph is important if we require realistic features of the result. However, most of the proposals do not provide an implementation at all and, thus, unfortunately this aspect cannot be discussed reliably.

In case of the other application-specific domains the amount of generators is relatively small. As we will show in Section 2.2, the situation is not so critical. Some of the application-specific generators can be re-used also in other application domains or the general graph generators can be used.

Considering the output format of the data, as expected, the generators produce data in a standard format (e.g., RDF) or in a reasonable graph-related form (e.g., edge list).

*Domain.* If we consider the features of the particular domain of the generators, in most cases it is expectably fixed. It is the simplest option, but it does mean that the generator is simple too – it can still focus on other complex aspects of the output. The particular target (in case of fixed) or sample (in case of schema-driven or extracted) domains do not provide very rich areas. They either correspond to the respective application domains (like in the case of social networks), or they are based on well-known and commonly used data sets (such as, e.g., DBLP [dbl 2016] or DBpedia [Bizer et al. 2009]). Except for general graphs without any domain, in all other cases we can find a flexible representative with schema-driven/pattern-driven domain or a domain extracted from sample data.

*Operations.* In the case of operations, most commonly the respective generators are accompanied with a set of fixed read operations or sequences of operations (query mixes) representing typical behavior of a user. In some cases this aspect is more flexible – either query templates are used or the operations are generated, e.g., in order to access most of the data in the generated graph. On the other hand, update operations are provided only in a small amount of cases. As we will discuss in Section 5.6, the more general problem of evolving graphs is still an open issue.

*Configuration.* A natural feature of the generators is to provide as realistic result as possible. Hence, most of them focus not only on the structure of the output graph, but also the properties. For the purpose of generating graphs with near real-world characteristics various distributions are used, such as power-law, Zipfian etc. Especially interesting are distributions extracted from real-world data (such as, e.g., Facebook in case of the social network domain).

<sup>4</sup><https://lod-cloud.net/>

Table 1. Key characteristics of the generators

		Domain		Operations		Configuration		
Generator	Type	Target/sample	Read	Update	Pro.	Distributions	Output	Dis.
<b>General</b>	Preferential attachment	–	N	N	Y	power-law	edge-list	N
	R-MAT	–	N	N	Y	power-law	edge-list	N
	GraphGen	–	N	N	Y	user-defined	node/edge-list	N
	BTER	–	N	N	N	user-defined	edge-list	Y
	Darwin	–	N	N	N	user-defined	edge-list	Y
	LUBM	fixed	fixed	N	Y	random (LCCG)	RDF	N
<b>Semantic web</b>	LBBM	extracted	N	N	Y	Monte Carlo	RDF	N
	UOBM	fixed	fixed	N	Y	random	RDF	N
	IIMB	fixed	N	N	Y	random	RDF	N
	BSBM	fixed	fixed	N	Y	mostly normal	RDF, relational	N
	SP2Bench	fixed	fixed	N	Y	based on DBLP	RDF	N
	[Duan et al. 2011]	extracted	N	N	Y	–	RDF	N
	DBPSB	extracted	templates	N	Y	random	RDF	N
	LODIB	fixed	N	N	Y	44 types	RDF	N
	SIB	fixed	social network	mixes	Y	from real-world data	RDF	N
	Geographica	fixed	OpenStreetMap	fixed + templates	Y	–	RDF	N
	WatDiv	schema-driven	user-defined	templates	N	uniform, normal, Zipfian	RDF	N
	RBench	extracted	DBLP Yago	templates	N	from real-world data	RDF	N
	LDBC SPB	fixed	media	mixes	N	power-law, skewed values, value correlation	RDF	N
	LinkGen	schema-driven	user-defined	templates	N	Gaussian, Zipfian	RDF	N
<b>GDBs</b>	XGDBench	fixed	social network	generated	Y	power-law	MAG	Y
	gMark	schema-driven	user-defined	generated	N	uniform, normal, Zipfian	N-triples	N
	graphGen	pattern-driven	user-defined	–	Y	–	Graphson, CypherQueries	N
<b>SNs</b>	[Barrett et al. 2009]	fixed	N	N	Y	simulation-driven	impl. NA	–
	[Yao et al. 2011]	fixed	N	N	N	power-law	impl. NA	–
	LinkBench	fixed	social network	generated	Y	Facebook	impl. NA	–
	S3G2	fixed	social network	N	Y	Facebook	CSV, RDF	Y
	[Ali et al. 2014]	schema-driven	N	N	Y	power-law	CSV	N
	LDBC SNB	fixed	social network	generated	Y	Facebook	CSV, RDF	Y
<b>A</b>	[Nettleton 2016]	schema-driven	N	generated	N	power-law	impl. NA	–
	HPC Scal. Graph Anal.	fixed	fixed	N	N	uniform	edge-list	N
<b>Co</b>	Danon et al.	–	–	N	N	uniform	edge-list	N
	LFR	–	–	N	N	power-law	edge-list	N
	LFR-Overlapping	–	–	N	N	power-law	edge-list	N
<b>St</b>	S2Gen	schema-driven	social network	Y	N	user-defined	RDF	N
	RSPLab	schema-driven	agnostic	Y	N	user-defined	RDF	N

Table II. Overlapping of classes of generators

	Generator	General	Semantic Web	Graph databases	Social networks	Analytics	Community detection	Steaming
General	Preferential attachment	x				x		
	R-MAT	x				x		
	GraphGen	x		x				
	BTER	x						
	Darwini	x						
Semantic web	LUBM		x					
	LBBM		x					
	UOBM		x					
	IIMB		x					
	BSBM		x					
	SP <sup>2</sup> Bench		x					
	[Duan et al. 2011]		x					
	DBPSB		x					
	LODIB		x					
	SIB		x					
	Geographica		x					
	WatDiv		x					
	RBench		x					
	LDBC SPB		x					
	LinkGen		x					
GDBs	XGDBench			x	x			
	gMark		x	x	x			
	graphGen			x				
SNs	[Barrett et al. 2009]				x			
	[Yao et al. 2011]				x			
	LinkBench			x	x			
	S3G2		x	x	x			
	[Ali et al. 2014]				x			
	LDBC SNB		x	x	x			
	[Nettleton 2016]				x			
A	HPC Scal. Graph Anal.					x		
Co	Danon et al.	x				x	x	
	LFR	x				x	x	
	LFR-Overlapping	x				x	x	
St	S2Gen							x
	RSPLab							x

## 2.2. Overlapping

As we have mentioned, the basic classification that we have used in this paper is based on a relatively vague division of the generators on the basis of the current application domains or research areas. In addition, some of the generators are either general or have features applicable in other domains. So the classes can overlap, as depicted in Table II.

For example, many general or domain-agnostic graph generators, such as the preferential attachment [Barabasi and Albert 1999] or R-MAT, are typically used to test graph analytics frameworks when large real graphs are not available.

Similarly, some social network graph generators such as LDBC-SNB, S3G2 or LinkBench, can be used to test graph databases. In the case of the first two, even though they are designed not to be specific to any type of technology, the graph databases are their main target. Additionally, they also provide serializers for RDF, thus they can also be used to test RDF systems.

In the case of LinkBench, nothing prevents the user to load the generated graph in a graph database (Facebook uses MySQL in that paper) to test a workload similar to Facebook and extend and complement it with more graph queries like those in LDBC-SNB.

Generators for community detection aim, in general, at creating graphs with a more realistic structure (graphs with communities of nodes where the density of edges is larger internally than externally). Even though these generators are, in general, used to test community detection algorithms (they generate also the expected communities in the graph), some studies also use them for general graph purposes or to test graph analytics algorithms besides community detection.

### 3. GRAPH DATA GENERATORS

In this section, we discuss in more detail the different graph data generators using the classification introduced before. For each category we briefly describe the key features of each of the representative examples. The aim is to provide the readers with a detailed look at each of the tools in the context of its competitors from the same domain.

#### 3.1. General Graphs

We start by focusing on approaches that have been designed for dealing with the generation of general graph data that is not aimed at a particular application domain. Currently, there exists a number of tools which involve a kind of general graph data generator, such as gtools from projects nauty and Traces [McKay and Piperno 2014] or the Stanford GraphBase [Knuth 2008]. We, however, focus on primarily generating/benchmarking projects targeting the Big Data world.

*Preferential Attachment.* Barabasi and Albert [Barabasi and Albert 1999] introduced a graph generation model that relies on two main mechanisms. The first mechanism is continuously expanding the graphs by adding new vertices. The second mechanism is to preferentially attach the new vertices to the nodes/regions that are already well connected. So, in this approach, the generation of large graphs is governed by standard, robust self-organizing mechanisms that go beyond the characteristics of individual applications.

*R-MAT.* R-MAT (*Recursive Matrix*) is a procedural synthetic graph generator which is designed to generate power-law degree distributions [Chakrabarti et al. 2004]. The generator is recursive and employs a fairly small number of parameters. In principle, the strategy of this generator is to achieve simple means to produce graphs that are matching the properties of the real graphs. In particular, the design goal of R-MAT is to produce graphs which mimics the degree distributions, imitate a community structure and have a small diameter. R-MAT can generate weighted, directed and bipartite graphs.

*GraphGen.* For the purpose of testing the scalability of an indexing technique called FG-index [Cheng et al. 2007] on the size of the database of graphs, their average size and average density, the authors have also implemented a synthetic generator called



GraphGen<sup>5</sup>. It relies on data generation code for associations and sequential patterns provided by IBM<sup>6</sup>. GraphGen yields a collection of undirected, labeled and connected graphs. It addresses the performance evaluation of frequent subgraph mining algorithms and graph query processing algorithms. The result is represented as a list of graphs, each consisting of a list of nodes along with a list of edges.

**BTER.** BTER (Block Two-Level Erdős-Rényi) [Kolda et al. 2014] is a graph generator based on the creation of multiple Erdős-Rényi graphs with different connection probabilities of which they are connected randomly between them. As the main feature, BTER is able to reproduce input degree distributions and average clustering coefficient per degree values. The generator starts by grouping the vertices by degree  $d$ , and forming groups of size  $d + 1$  of nodes with degree  $d$ . Then, these groups are assigned an internal edge probability in order to match the observed average clustering coefficient of the nodes of such degree. Based on this probability, for each node, the excess degree (i.e., the degree that in expectation will not be realized internally in the group) is computed and used to connect nodes from different groups at random. The authors describe a highly scalable MapReduce based implementation that is capable of generating large graphs (with billions of nodes) in a reasonable amounts of time.

**Darwini.** Darwini [Edunov et al. 2016] is an extension of BTER designed to run on Vertex Centric computing frameworks like Pregel [Malewicz et al. 2010] or Apache Giraph [Ching et al. 2015], with the additional feature that it is more accurate when reproducing the clustering coefficient of the input graph. Instead of just focusing on the average clustering coefficient for each degree, Darwini is able to model the clustering coefficient distribution per degree. It achieves this by gathering the nodes of the graph into buckets based on the expected number of closed triangles that they need to close in order to attain the expected clustering coefficient. The latter is sampled from the input distributions. Then, the vertices in each bucket are connected randomly with a probability that would produce the expected desired number of triangles for each bucket. Then, as in BTER, the excess degree is used to connect the different buckets. The authors report that Darwini is able to generate graphs with billions and even trillions of edges.

### 3.2. Semantic Web

With the dawn of the concept of Linked Data it is a natural development that there would emerge respective benchmarks involving both synthetic data and real-world data sets with real-world characteristics. The used data sets correspond to RDF representation of relational-like data [Guo et al. 2005; Bizer and Schultz 2009], social network-like data [Schmidt et al. 2010], or specific and significantly more complex data structures such as biological data [Wu et al. 2014b]. In this section, we provide an overview of benchmarking systems involving a kind of graph-based RDF data generator or data modifier.

**LUBM.** The use-case driven Lehigh University Benchmark (LUBM)<sup>7</sup> considers the university domain. The ontology defines 43 classes and 32 properties [Guo et al. 2005]. In addition, the LUBM benchmark provides 14 test queries. In particular, the benchmark focuses on *extensional* queries, i.e., queries which target the instance data over ontologies, as an opposite to *intentional* queries, i.e., queries which target the proper-

<sup>5</sup><https://www.cse.ust.hk/graphgen/>

<sup>6</sup>From 1996, no longer available at <http://www.almaden.ibm.com/cs/projects/iis/hdb/Projects/datamining/mining.shtml>

<sup>7</sup><http://swat.cse.lehigh.edu/projects/lubm/>

ties and classes of the ontology. The Univ-Bench Artificial (UBA) data generator features repeatable and random data generation (exploiting classical linear congruential generator, LCG, of numbers). In particular, the data which is produced by the generator are assigned zero-based indexes (i.e., the first university is named *University0* and so on), thus they are reproducible at any time with the same indexes. The generator allows to specify a seed for random number generation, along with the desired number of universities, and the starting index of the universities.

An extension of LUBM, the Lehigh BibTeX Benchmark (LBBM) [Wang et al. 2005], enables generating synthetic data for different ontologies. The data generation process is managed through two main phases: (1) the property-discovery phase, and (2) the data generation phase. In particular, LBBM provides a probabilistic model that can emulate the properties of the data of a particular domain and generate synthetic data exhibiting similar properties. A Monte Carlo algorithm is employed to output synthetic data. The approach is demonstrated on the Lehigh University BibTeX ontology which consists of 28 classes along with 80 properties. The LUBM benchmark includes 12 test queries that were designed for the benchmark data. Another extension of LUBM, the University Ontology Benchmark (UOBM)<sup>8</sup>, focuses on two aspects: (1) usage of all constructs of OWL Lite and OWL DL [owl 2004] and (2) lack of necessary links between the generated data which thus form isolated graphs [Ma et al. 2006]. In the former case the original ontology is replaced by the two types of extended versions from which the user can choose. In the latter case cross-university and cross-department links are added to create a more complex graph.

**IIMB.** Contrary to the previous work, Ferrara et al. [Ferrara et al. 2008] proposed the ISLab Instance Matching Benchmark (IIMB)<sup>9</sup> for the problem of instance matching. For any two objects  $o_1$  and  $o_2$  adhering to the same ontology or to different ontologies, instance matching is specified in the form of a function  $Om(o_1, o_2) \rightarrow \{0, 1\}$ , where  $o_1$  and  $o_2$  are linked to the same real-world object (in which case the function maps to 1) or  $o_1$  and  $o_2$  are representing different objects (in which case the function maps to 0). It targets the domain of movie data which contains 15 named classes, along with 5 object properties and 13 datatype properties. The data are extracted from IMDb<sup>10</sup>. The data generator corresponds to a data modifier which simulates differences between the data. In particular it involves data value differences (such as typographical errors or usage of different standard formats, e.g., for names), structural heterogeneity (represented by different levels of depth for properties, diverse aggregation criteria for properties, or missing values specification) and logical heterogeneity (such as instantiation on different subclasses of the same superclass or instantiation on disjoint classes).

**BSBM.** The Berlin SPARQL Benchmark (BSBM)<sup>11</sup>, is centered around an e-commerce application domain with object types such as *Customer*, *Vendor*, *Product* and *Offer* in addition to the relationship among them [Bizer and Schultz 2009]. The benchmark provides a workload that has 12 queries with 2 types of query workloads (i.e., 2 sequences of the 12 queries) emulating the navigation pattern and search of a consumer seeking a product. The data generator is capable of producing arbitrarily scalable datasets by controlling the number of products ( $n$ ) as a scale factor. The scale factor also impacts other data characteristics, such as, e.g., the depth of type hierarchy of products, branching factor, the number of product features, etc. BSBM can output

<sup>8</sup><https://www.cs.ox.ac.uk/isg/tools/UOBMGenerator/>

<sup>9</sup><http://www.ics.forth.gr/isl/BenchmarksTutorial/>

<sup>10</sup><http://www.imdb.com/>

<sup>11</sup><http://wifo5-03.informatik.uni-mannheim.de/bizer/berlinsparqlbenchmark/>

two representations, i.e. an RDF representation along with a relational representation. Thus, BSBM also defines an SQL [sql 2008] representation of the queries. This allows comparison of SPARQL [Prud'hommeaux and Seaborne 2008] results to be compared against the performance of traditional RDBMSs.

*SP<sup>2</sup>Bench.* The SP<sup>2</sup>Bench<sup>12</sup> is a language-specific benchmark [Schmidt et al. 2010] which is based on the DBLP dataset. The generated datasets follow the key characteristics of the original DBLP dataset. In particular, the data mimics the correlations between entities. All random functions of the generator use a fixed seed that ensures that the data generation process is deterministic. SP<sup>2</sup>Bench is accompanied by 12 queries covering the various types of operators such as RDF access paths in addition to typical RDF constructs.

*Data Coherence.* A comparison of 4 RDF benchmarks (namely TPC-H [TPC 2016] data expressed in RDF, LUBM, BSBM, and SP<sup>2</sup>Bench) and 6 real-world data sets (such as, e.g., DBpedia, the Barton Libraries Dataset [Abadi et al. 2007] or WordNet [Miller 1995]) has been reported by [Duan et al. 2011]. The authors focus mainly on the *structuredness (coherence)* of each benchmark dataset claiming that a primitive metric (e.g., the number of triples or the average in/outdegree) quantifies only some target characteristics of each dataset. The degree of structuredness of a dataset  $D$  with respect to a type  $T$  is based on the regularity of instance data in  $D$  in conforming to type  $T$ . The type system is extracted from the data set by finding the RDF triples that have property `http://www.w3.org/1999/02/22-rdf-syntax-ns#type` and extract type  $T$  from their object. Properties of  $T$  are determined as the union of all the properties that the instances of type  $T$  have. The structuredness is then expressed as a weighted sum of share of set properties of each type, whereas higher weights are assigned to types with more instances. The authors show that the structuredness of the chosen benchmarks is fixed, whereas real-world RDF datasets are belonging to the non-tested area of the spectrum. As a consequence, they introduce a new benchmark that receives as input any dataset associated with a required level of structuredness and size (smaller than the size of the original data), and exploits the input documents as a seed to produce a subset of the original data with the target structuredness and size. In addition, they show that structuredness and size mutually influence each other.

*DBPSB.* DBpedia SPARQL Benchmark (DBPSB)<sup>13</sup> proposed at the University of Leipzig has been designed using workloads that have been issued by humans and applications over existing RDF data [Morsey et al. 2011; Morsey et al. 2012]. In addition, the authors argue that benchmarks like LUBM, BSBM, or SP<sup>2</sup>Bench resemble relational database benchmarks involving relational data structures with few and homogeneously structured classes, whereas, in reality, RDF datasets are increasingly heterogeneous. For example, DBpedia version 3.6 consists of 289,016 classes of which 275 classes are defined based on the DBpedia ontology. In addition, different data types and object references of the various types are used in property values. Hence, they presented a universal SPARQL benchmark generation approach which uses a flexible data production mechanism that mimics the input data source. This dataset generation process begins using an input dataset; then multiple datasets with different sizes are then generated by duplicating all the RDF triples with changing their namespaces. For generating smaller datasets, an adequate selection of all triples is selected randomly or using a sampling mechanism over the various classes in the dataset. The methodology

<sup>12</sup><http://dbis.informatik.uni-freiburg.de/forschung/projekte/SP2B/>

<sup>13</sup><http://aksw.org/Projects/DBPSB.html>

is applied on the DBpedia SPARQL endpoint and a set of 25 SPARQL query templates is derived to cover the frequently used SPARQL features.

*LODIB.* The Linked Open Data Integration Benchmark (LODIB)<sup>14</sup> has been designed with the aim of reflecting the real-world heterogeneities that exist on the Web of Data in order to enable testing of Linked Data translation systems [Rivero et al. 2012]. It provides a catalogue of 15 data translation patterns (e.g., rename class, remove language tag etc.), each of which is a common data translation problem in the context of Linked Data. The benchmark provides a data generator that produces three different synthetic data sets that need to be translated by the system under test into a single target vocabulary. They reflect the pattern distribution in analyzed 84 data translation examples from the LOD Cloud. The data sets reflect the same e-commerce scenario used for BSBM.

*SIB.* The developers of the Social Network Intelligence BenchMark (SIB)<sup>15</sup> based the design of their benchmark on the claim that existing benchmarks are limited in reflecting the characteristics of the real RDF databases and are mostly focusing on the relational style aspects. Hence, they proposed a benchmark for query evaluation using real graphs [Boncz et al. 2013]. The proposed benchmark mimics using an RDF store for a social network. The distribution of the generated data for each type follows the distribution of the associated type inferred from real-world social networks. Additionally, association rules are exploited for representing the real-world data correlation in the generated synthetic data. The generated data is linked with the RDF datasets from DBpedia. The benchmark specification contains 3 query mixes – interactive, update, and analysis – expressed in SPARQL 1.1 Working Draft.

*Geographica.* The Geographica benchmark<sup>16</sup> has been designed to target the area of geospatial data [Garbis et al. 2013] and respective SPARQL extensions GeoSPARQL [Battle and Kolas 2012] and stSPARQL [Koubarakis and Kyzirakos 2010]. The benchmark involves a real-world workload that uses openly available datasets that cover a range of geometry types (e.g., points, lines, polygons) and a synthetic workload. In the former case there is a (1) a micro benchmark that evaluates primitive spatial functions (involving 29 queries) and (2) macro benchmark that tests the performance of RDF engines in various application scenarios such as map exploring and search (consisting of 11 queries). In the latter case of a synthetic workload the generator produces synthetic datasets of different sizes that corresponds to an ontology based on OpenStreetMap and instantiates query templates. The generated SPARQL query workload is corresponding to spatial selection and spatial joins by instantiating 2 query templates.

*WatDiv.* The Waterloo SPARQL Diversity Test Suite (WatDiv)<sup>17</sup> has been designed at the University of Waterloo. It implements stress testing tools that focus on addressing the observation that existing SPARQL benchmarks are not adequate for evaluating systems using a verity of queries and workloads [Aluç et al. 2014]. The benchmark focuses on two types of query aspects – structural and data-driven – and performs a detailed analysis on existing SPARQL benchmarks (LUBM, BSBM, DBPSB, and SP<sup>2</sup>Bench) using the two classes of query features. The structural features involve triple pattern count, join vertex count, join vertex degree, and join vertex count. The data-driven features involve result cardinality and several types of selectivity. The

<sup>14</sup><http://lodib.wbmg.de/>

<sup>15</sup>[https://www.w3.org/wiki/Social\\_Network\\_Intelligence\\_BenchMark](https://www.w3.org/wiki/Social_Network_Intelligence_BenchMark)

<sup>16</sup><http://geographica.di.uoa.gr/>

<sup>17</sup><http://dsg.uwaterloo.ca/watdiv/>

analysis of the four benchmarks reveals that they are not sufficiently diverse to evaluate the strengths and weaknesses of the various physical design alternatives that have been implemented by the different RDF systems. In particular, WatDiv, provides (1) a data generator which generates scalable datasets according to the WatDiv schema, (2) a query template generator which follows the WatDiv schema and produces a set of query templates, (3) a query generator that uses the generated templates and instantiates them with real RDF values from the dataset, and (4) a feature extractor which extracts the structural features of the generated data and workload.

*RBench.* RBench [Qiao and Özsoyoğlu 2015] is an application-specific benchmark which receives any RDF dataset as an input and produces as an output a set of datasets, that have similar characteristics of the input dataset, using size scaling factor  $s$  and (node) degree scaling factor  $d$ . A query generation process has been implemented to produce 5 different types of queries (edge-based queries, node-based queries, path queries, star queries, subgraph queries) for any generated data. The benchmark project FEASIBLE [Saleem et al. 2015] is also an application-specific benchmark; however, contrary to RBench, it is designed to produce benchmarks from a set of queries (in particular from query logs) by relying on sample queries of a user-defined size from the input set of queries.

*LDBC.* The Linked Data Benchmark Council<sup>18</sup> (LDBC) [Angles et al. 2014] is a result of a (closed) EU project that brought together a community of academic researchers and industry that had the goal of developing an open source, yet industrial grade benchmarks for RDF and graph databases. In the Semantic Web domain, the project released the Semantic Publishing Benchmark (SPB) [spb 2015] that has been inspired by the Media/Publishing industry (namely BBC<sup>19</sup>). The application scenario of this benchmark simulates a media or a publishing organization that handles large amount of streaming content (e.g., news, articles). The data generator mimics three types of relations in the generated synthetic data: clustering of data, correlations of entities, and random tagging of entities. Two workloads are provided: (1) basic, involving an interactive query-mix querying the relations between entities in reference data, and (2) advanced, focusing on interactive and analytical query-mixes. The LDBC has designed two other benchmarks: the Social Network Benchmark (SNB) [Erling et al. 2015] for the social network domain (see Section 3.4) and Graphalytics [Iosup et al. 2016] for the analytics domain.

*LinkGen.* LinkGen is a synthetic linked data generator that has been designed to generate RDF datasets for a given vocabulary [Joshi et al. 2016]. The generator is designed to receive a vocabulary as an input and supports two statistical distributions for generating entities: Gaussian distribution and Zipf's power-law distribution. LinkGen can augment the generated data with inconsistent and noisy data such as updating a given datatype property with two conflicting values or adding triples with syntactic errors. The generator also provides a feature to inter-link the generated objects with real ones given that the end-user provides entities from real datasets. The datasets can be generated in any of of two modes: on-disk and streaming.

### 3.3. Graph Databases

Currently there exists a number of papers which compare the efficiency of graph databases with regards to distinct use cases, such as the community detection problem [Beis et al. 2015], social tagging systems [Giatsoglou et al. 2011], graph travers-

<sup>18</sup><http://ldbpcouncil.org/industry/organization/origins>

<sup>19</sup><http://www.bbc.com/>

sal [Ciglan et al. 2012], graph pattern matching [Pobiedina et al. 2014], data provenance [Vicknair et al. 2010], or even several distinct use cases [Grossniklaus et al. 2013]. However, the number of graph data generators and benchmarks that have been designed specifically for graph data management systems (Graph DBMS) is relatively small. Either a general graph generator is used for benchmarking graph databases, such as, e.g., the HPC Scalable Graph Analysis Benchmark [Dominguez-Sal et al. 2010] or the graph DBMS benchmarking tools are designed while having in mind a more general scope. Hence it is questionable whether a benchmark that is targeted specifically for graph databases is necessary. [Dominguez-Sal et al. 2011] discussed this question and related topics. On the basis of a review of applications of graph databases (namely, social network analysis, genetic interactions, recommendation systems, and travel planning and routing), the authors analyzed and discussed the features of the graphs for these types of applications and how they could affect the benchmarking process, different types of operations used in these applications and the characteristics of the evaluation setup of the benchmark. In this section, we focus on graph data generators and benchmarks that have been primarily targeting graph DBMSs.

*XGDBench.* XGDBench [Dayarathna and Suzumura 2014] is an extensible benchmarking platform for graph databases used in cloud-based systems. Its intent is to automate the process of graph database benchmarking in the cloud by focusing on the domain social networking services. It extends the Yahoo! Cloud Multiplicative Attribute (MAG) Graph Serving Benchmark (YCSB) [Cooper et al. 2010] and provides a set of standard workloads representing various performance issues. In particular, the workload of XGDBench involves basic operations such as read / insert / update / delete an attribute, loading of the list of neighbours and BFS traversal. Using the generators, 7 workloads are created, such as update heavy, read mostly, short range scan, traverse heavy etc. The data model of XGDBench is a simplified version of the Multiplicative Attribute Graph (MAG) [Kim and Leskovec 2010] model, a synthetic graph model which models the interactions between node attributes and graph structure. The generated graphs are thus in MAG format, with power-law degree distribution closely simulating real-world social networks. The simplified MAG algorithm accepts the required number of nodes, the number of attributes per each node, a threshold value for random attribute initialization, an edge affinity threshold determining existence of an edge between two nodes, and an affinity matrix. Large graphs can be generated on multi-core systems by a multi-threaded version of the generator.

*gMark.* gMark [Bagan et al. 2016] is a schema-driven and domain-agnostic generator of both graph instances and graph query workloads. It can generate instances under the form of N-triples and queries in various concrete query languages, including OpenCypher<sup>20</sup>, recursive SQL, SPARQL and LogicQL. In gMark, it is possible to specify a *graph configuration* involving the admitted edge predicates and node labels occurring in the graph instance along with additional parameters such as degree distribution, occurrence constraints, etc. The *Query workload configuration* describes parameters of the query workload to be generated, by including the number of queries, arity, shape and selectivity of the queries. The problem of deciding whether there exists a graph that satisfies a defined graph specification  $G$  is NP-complete. The same applies to the problem of deciding whether there exists a query workload compliant with a given query workload configuration  $Q$ . In view of this, gMark adopts a best effort approach in which the parameters specified in the configuration files are attained in a relaxed fashion in order to achieve linear running time whenever possible.

<sup>20</sup><https://neo4j.com/developer/cypher-query-language/>

*GraphGen.* GraphAware GraphGen<sup>21</sup> is a graph generation engine based on Neo4j's<sup>22</sup> query language OpenCypher [Gra 2015]. It creates nodes and relationships based on a schema definition expressed in Cypher, and it can also generate property values on both nodes and edges. As such, GraphGen is a precursor of property graphs generators. The resulting graph can be exported to several formats (namely GraphJSON<sup>23</sup> and CypherQueries) or loaded directly to a DBMS. However, it is very likely that it is not maintained anymore due to the lack of available recent commits.

*Strengths and Weaknesses of Graph Database Generators.* The graph DBMS generators discussed in this section have in common the fact that they can generate semantically rich labeled graphs with properties (ranging from properties values in GraphGen to MAG structures in XGDBench). They are also capable of generating query workloads in concrete syntaxes (among which OpenCypher in GraphGen and gMark) and one of them (XGDBench) can also handle update operations on both graph structure and content. However, more comprehensive graph DBMS generators that also produce data manipulation operations (such as updates for graph databases) are urgently needed.

Additionally, none of these generators is enabled to work on property graphs and on corresponding languages for property graphs, such as the newly emerging standard GQL [Chafi et al. 2018] and G-Core [Angles et al. 2018]. However, a full-fledged graph DBMS generator for property graphs and property graph query workloads [Bonifati et al. 2018b] is still missing and could be an interesting opportunity to build such a generator in the next future. Another weakness is represented by the fact that explicit correlations among graph elements cannot be encoded for instance in gMark or GraphGen, whereas they could be fruitful in order to reproduce the behavior of real-world graphs in which attribute values are correlated one with another. On the other hand, social network and Linked Data generators that support correlations (as highlighted in Section 3.4 and Section 3.2) typically exhibits a fixed schema and are not multi-domain as the graph DBMS generators discussed in this Section.

### 3.4. Social Networks

On-line social networks, like Facebook, Twitter, or LinkedIn, have become a phenomenon used by billions of people every day and thus providing extremely useful information for various domains. However, an analysis of such type of graph has to cope with two problems: (1) availability of the data and (2) privacy of the data. Hence, data generators which provide realistic synthetic social network graphs are in a great demand.

In general, any analysis of social networks identifies their various specific features [Chakrabarti and Faloutsos 2006]. For example, a social graph often has high *clustering coefficient*, i.e. the degree of transitivity of a graph. Or, its *diameter*, i.e. the longest shortest path amongst some fraction (e.g. 90%) of all connected nodes, is usually low due to weak ties joining faraway cliques.

Another important aspect of social networks is the community effect. A detailed study of structure of communities in 70 real-world networks is provided, e.g., in [Leskovec et al. 2008]. [Prat-Pérez and Dominguez-Sal 2014] analyzed the structure of communities (clustering coefficient, triangle participation ratio, bridges, diameter, conductance and size) in both real-world graphs and outputs of existing graph generators LFR [Lancichinetti et al. 2008] and the LDBC-SNB [Erling et al. 2015]. They

<sup>21</sup><http://graphgen.graphaware.com/>

<sup>22</sup><https://neo4j.com/>

<sup>23</sup><https://github.com/GraphAlchemist/GraphJSON/wiki/GraphJSON>

found out that discovered communities in different graphs have common distributions and that communities in a single graph have diverse nature and are difficult to fit with a single model.

The existing social network generators try to reproduce different aspects of the generated network. They can be categorized into statistical and agent-based. *Statistical approaches* [Lancichinetti et al. 2008; Yao et al. 2011; Armstrong et al. 2013; Pham et al. 2013; Ali et al. 2014; Erling et al. 2015; Nettleton 2016] focused on reproducing aspects of the network. In *agent-based approaches* [Barrett et al. 2009; Bernstein and O'Brien 2013] the networks are constructed by directly simulating the agents' social choices.

*Realistic Social Network.* [Barrett et al. 2009] focused on the construction of realistic social networks. For this purpose the authors combine both public and private data sets with large-scale agent-based techniques. The process works as follows: In the first step it generates a synthetic data by combining public and commercial databases. In the second step, it determines a set of activity templates. A 24-hour activity sequence including geolocations is assigned to each synthetic individual. To demonstrate the approach, the authors create a synthetic US population consisting of people and households together with respective geolocations. For this purpose the authors combine simulation and data fusion techniques utilizing various real-world data sources such as U.S. Census data, responses to an activity survey or a time-use survey. The result is captured by a dynamic social contact network. Similar methods for agent-based strategies have been reported in [Bernstein and O'Brien 2013].

*Linkage vs. Activity Graphs.* [Yao et al. 2011] distinguished between two types of social network graphs – the *linkage graph*, where nodes represent people and edges correspond to their friendships, and the *activity graph*, where nodes also represent people but correspond to their interactions. On the basis of the analysis of Flickr<sup>24</sup> social links and Epinions<sup>25</sup> network of user interactions, the authors discover that they both exhibit high clustering coefficient (community structure), small diameter, and power-law degree distribution. Considering the dynamic properties they both have relatively stable clustering coefficient over time and follow the densification law. On the other hand, diameter shrinking is not observed in Epinions activity graph and there is a difference in degree correlation (i.e., frequency of mutual connections of similar nodes) – linkage graphs have positive degree correlation, activity graphs have neutral degree correlation. With regards to the findings, the proposed generator focusses on linkage graphs with positive degree correlation. For this purpose it extends the forest fire spreading process algorithm [Leskovec et al. 2005b] with link symmetry. It has two parameters: the *burning probability*  $P_b$  and the *symmetry probability*  $P_s$ .  $P_b$  ensures a BFS-based forward burning process in which fire burns strongly with  $P_b$  approaching 1.  $P_s$  ensures backward linking from old nodes to new nodes and “adds fuel to the fire as it brings more links”.

*LinkBench.* The LinkBench benchmark [Armstrong et al. 2013] has been designed for the purpose of analysis of efficiency of a database storing Facebook's production data. The benchmark considers true Big Data and related problems with sharding, replication etc. The social graph at Facebook comprises objects (nodes with IDs, version, timestamp and data) and associations (directed edges, pairs of node IDs, with visibility, timestamp and data). The size of the target graph is the number of nodes. Graph edges and nodes are generated concurrently during bulk loading. The space of

<sup>24</sup><https://www.flickr.com/>

<sup>25</sup><http://www.epinions.com/>



node IDs is divided into chunks which enable parallel processing. The edges of the graph are generated in accordance with the results of analysing real-world Facebook data (such as outdegree distribution). A workload corresponding to 10 graph operations (such as insert object, count the number of associations etc.) and their respective characteristics over the real-world data is generated for the synthetic data.

**S3G2.** The Scalable Structure-correlated Social Graph Generator (S3G2) [Pham et al. 2013] is a general framework which produces a directed labeled graph whose vertices represent objects having property values. The respective classes determine the structure of the properties. S3G2 does not aim at generating near real-world data, but at generating synthetic graphs with a correlated structure. Hence, the existing data values influence the probability of choosing a certain property value from a pre-defined dictionary, or the probability of connecting two nodes. For example, the degree distribution can be correlated with the properties of a node and thus, e.g., people having many friend relationships typically post more comments and pictures. The data generation process starts with generating a number of nodes with property values generated according to specified property value correlations and then adding respective edges according to specified correlation dimensions. It has multiple phases, each focusing on one correlation dimension. Data is generated in a Map phase corresponding to a pass along one correlation dimension. Then the data are sorted along the correlation dimension in the following Reduce phase. A heuristic observation that “the probability that two nodes are connected is typically skewed with respect to some similarity between the nodes” enables to focus only on sliding window of most probable candidates. The core idea of the framework is demonstrated using an example of a social network (consisting of persons and social activities). The dictionaries for property values are inspired by DBpedia and provided with 20 property value correlations. The edges are generated according to 3 correlation dimensions.

*Cloning of Social Networks.* [Ali et al. 2014] introduced a synthetic network generator designed for cloning social network statistics of an existing dataset. The network starts with a small number of nodes, and new nodes are added until the network reaches the required number. It has two basic parameters: homophily and link density. A high *homophily* value reflects that links have higher chances to be established among the nodes with the same label; these labels can be considered as being equivalent to community membership.

Attribute Synthetic Generator (ASG) is a network generator for reproducing the node feature distribution of standard networks and rewiring the network to preferentially connect nodes that exhibit a high feature similarity. The network is initialized with a group of three nodes, and new nodes and links are added to the network based on link density, homophily, and feature similarity. As new nodes are created, their labels are assigned based on the prior label distribution. After the network has reached the same number of nodes as the original social media dataset, each node initially receives a random attribute assignment. Then a stochastic optimization process is used to move the initial assignments closer to the target distribution extracted from social media dataset using the Particle Swarm Optimization algorithm. The tuned attributes are then used to add additional links to the network based on the feature similarity parameter – a source node is selected randomly and connected to the most similar node. Multi-Link Generator (MLG) further uses link co-occurrence statistics from the original dataset to create a multiplex network. MLG uses the same network growth process as ASG. Based on the link density parameter, either a new node is generated with a label based on the label distribution of the target dataset or a new link is created between two existing nodes.

**LDBC SNB.** Despite having a common Facebook-like dataset, thanks to three distinct workloads the Social Network Benchmark (SNB) [Erling et al. 2015] provided by LDBC represents three distinct benchmarks. The network nodes correspond to people and the edges represent their friendship and messages they post in discussion trees on their forums. The three query workloads involve: (1) SNB-Interactive, i.e., complex read-only queries accessing a significant portion of data, (2) SNB-BI, i.e., queries accessing a high percentage of entities and grouping them in various dimensions, and (3) SNB-Algorithms, involving graph analysis algorithms, such as community detection, PageRank, BFS, and clustering. The graph generator realizes power laws, uses skewed value distributions, and ensures reasonable correlations between graph structures and property values. To provide scalability it is implemented on top of Hadoop.

**Towards More Realistic Data.** [Nettleton 2016] argued that the majority of existing works focuses on topology generation which approximates the characteristics of a real social network (such as a small graph diameter, skew degree distribution, small average path length, and community structures), however, this is usually done without any data. Hence, they introduced a general stochastic modeling approach that enables the users to populate a graph topology with data. The approach has three steps: (1) topology generation (using R-MAT) plus community identification using the Louvain method [Blondel et al. 2008] or usage of a real-world topology from SNAP<sup>26</sup>, (2) data definition that describes distribution profiles, attribute value definitions, using a parameterizable set of data propagation rules and affinities, and (3) data population.

### 3.5. Graph Analytics

Graph analytics, especially in the context of Big Data, is a popular area of studying interesting structural specifics of graphs, usually various types of networks. Hence, the respective benchmarks and graph data generators developed for the purpose of testing graph analytics tools (such as, e.g., PowerGraph [Gonzalez et al. 2012] or Parallel Graph analytiX [Sevenich et al. 2016]) have to focus mainly on graphs with complex structures.

**HPC Scalable Graph Analysis Benchmark.** The HPC Scalable Graph Analysis Benchmark [HPC 2009; Bader and Madduri 2005] consists of a weighted, directed graph that has a power-law distribution and four related analysis techniques (namely graph construction, graph extraction with BFS, classification of large vertex sets, and graph analysis with betweenness centrality). The generator has the following parameters: the number of nodes, the number of edges, and maximum weight of an edge. It outputs a list of tuples containing identifiers of vertices of an edge (with the direction from the first one to the second one) and weights (positive integers with a uniform random distribution) assigned to the edges of the multigraph. The algorithm of the generator is based on R-MAT. To reduce the data locality, the vertex ids are randomly permuted and the output tuples are shuffled. A related project from the same authors developed for the 9th DIMACS Shortest Paths Challenge is GTgraph [Bader and Madduri 2006]. It involves three types of graphs: (1) Erdős-Rényi random graphs, (2) different graphs from the DARPA HPCS SSCA#2 graph theory benchmark (version 1.0), and (3) small-world graphs generated using the R-MAT graph generator.

### 3.6. Community Detection

Community detection is one of the many graph analytics algorithms typically used on domains such as social networks or bioinformatics. Communities are groups of nodes that are highly connected among them, while being scarcely connected to the rest of

<sup>26</sup><https://snap.stanford.edu/data/>

the graph. Such communities emerge from the fact that real graphs are not random, but follow real-world dynamics that make similar entities to have a larger probability to be connected. As a consequence, detected communities are used to reveal vertices with similar characteristics, for instance to discover functionally equivalent proteins in protein-protein interaction networks, or persons with similar interests in social networks. Such applications have made community detection a hot topic during the last fifteen years with tens of developed algorithms and detection strategies [Zhao 2017; Kim and Lee 2015]. For comparing the quality of the different proposed techniques, one needs graphs with *reference communities*, that is, communities known beforehand. Since it is very difficult to have large real graphs with reference communities (mainly because these would require a manual labeling), graphs for benchmarking community detection algorithms are typically generated synthetically.

*Danon et al.*. The first attempts to compare community detection algorithms using synthetic graphs proposed the use of random graphs composed by several Erdős-Rényi subgraphs, connected more internally than externally [Danon et al. 2005]. Each of these subgraphs has the same size and the same internal/external density of edges. However, such graphs miss the realism observed in real graphs, where communities are of different sizes and densities, thus several proposals exist to overcome such an issue.

*LFR*. Lancichinetti, Fortunato and Radicchi (hence LFR) [Lancichinetti et al. 2008] propose a class of benchmark graphs for community detection where communities are of diverse sizes and densities. The generated communities follow a power-law distribution whose parameters can be configured. The degree of the nodes is also sampled from a power-law distribution. Additionally, the generator introduces the concept of the “mixing factor”, which consists of the percentage of edges in the graph connecting nodes that belong to different communities. Such parameter allows the degree of modularity of the generated graph to be tuned, thus testing the robustness of the algorithms under different conditions. The generation process is implemented as an optimization process starting with an empty graph and progressively filling it with nodes and edges guided by the specified constraints.

*LFR-Overlapping*. Lancichinetti, Fortunato and Radicchi [Lancichinetti and Fortunato 2009] extended LFR to support the notion of directed graphs and overlapping communities. Overlapping communities extend the notion of communities by allowing the sharing of vertices, thus a vertex can belong to more than one community. This extended generator allows controlling the same parameters of LFR, as well as the amount of overlap of the generated communities.

Besides synthetic graph generators, Yang and Leskovec [Yang and Leskovec 2015] proposed the use of real-world graphs with explicit group annotations (e.g., forums in a social network, categories of products, etc.) to infer what they call “meta-communities”, and use them to evaluate overlapping community detection algorithms. However, a recent study from Hric, Darst and Fortunato [Hric et al. 2014] reveal a loose correspondence between communities (the authors refer to them as structural communities) and meta-communities. This result reveals that algorithms working for structural communities do not work well for finding meta-communities and vice versa, suggesting significantly different underlying characteristics between the two types of communities, which are yet to be identified.

In this regard and to the best of our knowledge, there does not exist a set of generators to generate graphs with meta-communities for community detection algorithm benchmarking. The closest one is the LDBC-SNB data generator [Erling et al. 2015] which has been provided by the generation of groups of users in the social network.

Even though the generation process does not specifically enforce the generation of groups (meta-communities) for benchmarking community detection algorithms, the study from Prat-Pérez and Dominguez-Sal reveals that these groups are more similar to the real meta-communities than those structural communities generated by the LFR benchmark [Prat-Pérez and Dominguez-Sal 2014].

The differences observed between structural and meta-communities reveal the need of more accurate community definitions tight more specifically to the domain or use case. Current community detection algorithms and graph generators for community detection are stuck to the traditional (and vague) definition of community, assuming that there exists a single algorithm that would fit all the use cases. Thus, future work requires the study of domain-specific community characteristics that can be used to generate graphs with a community structure that accurately resembles that of specific use cases, and thus revealing which are the best algorithms for each particular scenario.

### 3.7. Graph Data Streaming

One way for dealing with big graphs is to process them using the *streaming* mode where the data stream could consist of the edges of the graph. In this mode, the graph processing algorithms can process the input stream in the order it arrives while using only a limited amount of memory [McGregor 2014].

*S2Gen*. The streaming mode has mainly attracted the attention of the RDF and Semantic Web community. Thus, Phuoc et al. [Le-Phuoc et al. 2012] presented an evaluation framework for linked stream data processing engines. The framework uses a datasets generated with the Stream Social network data Generator (S2Gen), which simulates streams of user interactions (or events) in social networks (e.g., posts) in addition to the user metadata such as users' profile information, social network relationships, posts, photos and GPS information. The data generator of this framework provides the users the flexibility to control the characteristics of the generated stream by tanning a range of parameters, which includes the frequency at which interactions are generated, limits such as the maximum number of messages per user and week, and the correlation probabilities between the different objects (e.g., users) in the social network.

*RSPLab*. Tommasini et al. [Tommasini et al. 2017] introduced another framework for benchmarking RDF Stream Processing systems, RSPLab. The Streamer component of this framework is designed to publish RDF streams from the various existing RDF benchmarks (e.g., BSBM, LUBM) (see Section 3.2). In particular, the Streamer component uses TripleWave<sup>27</sup>, an open-source framework for publishing and sharing RDF streams on the Web [Mauri et al. 2016]. TripleWave acts as a means for plugging-in and combining streams from multiple Web data sources, in both push and pull mode.

## 4. GRAPH DATASETS

Besides graph data generators, many graph algorithms and systems are benchmarked using real-world graph data sets, including new graph data generation techniques. In this section, we review those graph datasets most widely used in the literature and their purpose.

*Large Scale Graph Analytics*. Large scale graph analytics systems are usually benchmarked using a combination of real-world and synthetically generated graphs. Because of the rapid evolution of such systems and their increasing capability to pro-

<sup>27</sup><http://streamreasoning.github.io/TripleWave/>

cess larger graphs, the datasets typically used in the literature quickly change over-time. The most commonly used datasets (which can be downloaded from several graph dataset repositories [Leskovec and Krevl 2014; degli studio di Milano 2018]) are:

- **DBLP** [Yang and Leskovec 2015]: Represents a co-authorship network where researchers from DBLP are represented as vertices and there exists an edge between them if they have published a paper together.
- **Amazon** [Yang and Leskovec 2015]: Represents a product co-purchasing network, such that each vertex represents a product, while edges between two products exists if a person bought both of them.
- **Road networks** [Leskovec et al. 2009]: Consist of a set of road networks from the US.
- **LiveJournal** [Yang and Leskovec 2015]: A social network where vertices represent the users and the edges the acquaintances between them.
- **Orkut** [Yang and Leskovec 2015]: Like LiveJournal, this is a social network where vertices represent persons and the edges represent the friendship relationship between them.
- **Twitter** [Kwak et al. 2010]: This is a directed network representing the follower-followee interaction of Twitter.
- **Friendster** [Yang and Leskovec 2015]: A social network like LiveJournal and Orkut, where nodes represent persons and edges their relationships.
- **WebUK** [Boldi et al. 2008]: A web graph of the UK subdomain, where nodes represent websites and edges represent the hyperlinks between them.
- **ClueWeb2012** [Project 2012]. This is a crawl of the web from 2012, where vertices represent websites and the edges between them represent the hyperlinks.

*Recommender Systems.* Many recommender systems are based on graphs, either bipartite graphs or many-to-many graphs. Such systems are usually tested on several real-world datasets, many of them containing information about the rating of products or other items such as movies, books or songs. The following are the most widely used datasets for testing recommender systems:

- **MovieLens** [GroupLens 2018]: Consists of a bipartite graph between users and movies, where edges represent the ratings the different users have made to the movies they watched.
- **Book-Crossings** [Ziegler et al. 2005]: Is a dataset with book ratings from users.
- **Jester** [Goldberg et al. 2001]: This dataset contains jokes (with their text) and ratings from users.
- **Last.fm** [GroupLens 2011]: This dataset contains the list of the top most listened artists per user, including the number of times the songs from those artists were played. Additionally, it contains connections between users (the social network they form), and a set of tags attached to artists that can be used to create content vectors.
- **Netflix** [Zhou et al. 2008]: This dataset was used during the famous Netflix prize. It consists of a bipartite graph with movie ratings from users.

*Reputation Algorithms.* Another application using social network data is that of assessing the degree of reputation of a user based on their past interactions and rankings [Kamvar et al. 2003; Katz and Golbeck 2006; Kumar et al. 2016]. Such algorithms are typically evaluated on networks with explicit user rankings or voting from other users, being the most widely used datasets the following examples:

- **Bitcoin** [Moore and Christin 2013]: This dataset contains information from several Bitcoin exchanges, where users are able to rate other users after their transactions.

Here, the vertices represent the users and the edges, which are labelled, the ratings between them.

- **Wikipedia RFA** [West et al. 2014]: This network is extracted from Wikipedia and the vertices represent users, while the edges, which are labelled, the votes emitted by administrators for the user to become an administrator. Each vote is accompanied with a text explaining the vote's sign.
- **WikiSigned** [Maniu et al. 2011]: This is a network of Wikipedia editors where the vertices are the editors and the edges, which are labelled, represent the trust level between two editors.
- **Extended Epinions** [Massa and Avesani 2007]: This is an extended version of the Epinions network which also contains the levels of distrust between the users, expressed by means of edges between them.
- **Twitter Indian Elections** [Kagan et al. 2015]: This network represents a Twitter network where vertices represent users and there is an edge between two users if a user mentions another one in a tweet. The edges are labelled with the average sentiment of one user towards another.

*Graph partitioning, Clustering and Community Detection.* Many of the already discussed datasets are also typically used to test and compare graph partitioning, clustering and community detection algorithms. One can find a comprehensive list of datasets for such algorithms in [Bader et al. 2012; Yang and Leskovec 2015], some of which we have already discussed for other applications. Here, we summarize the most widely used in the literature:

- **DBLP, Amazon, LiveJournal, Orkut and Friendster** [Yang and Leskovec 2015] are widely used for evaluating community detection since they provide information about the meta-communities they contain. For instance, in Livejournal, users can join groups of users talking about given topics. Such groups are exported as meta-communities, and similar approaches are used for other graphs. Community detection algorithms are then evaluated by trying to infer such meta-communities without prior knowledge of the user to group assignment.
- **Zachary Karate Club** [Zachary 1977]: This graph consists of a small network of members of a karate club that was dismissed and split into two new karate clubs. Vertices represent persons and edges friendship relationships. Information about what people joined each of the two new karate clubs is provided.
- **PolBlogs** [Adamic and Glance 2005]: This is a network consisting of blogs talking about the US 2005 political elections, where a vertex represents a blog and the edges the hyperlinks between the blogs. Each blog has an associated label, whether it is left or right oriented.
- **PolBooks** [Bader et al. 2012]: Is a network of books of Amazon that talk about politics, and edges between two books exist if they are co-purchased together. The books have labels of whether they are left or right oriented.
- **Football** [Girvan and Newman 2002]: This network represents American university football teams. Each node is a football team and the edges are games between teams during the season. Each team is associated to a division.

*Information Diffusion.* Information networks are studied using graph algorithms in order to understand how information propagates. As such, there several datasets used to understand such networks.

- **Higgs-Twitter** [De Domenico et al. 2013]: This dataset contains the Twitter network before, during and after the discovery of the Higgs boson. Specifically, it contains the tweets, the mentions, retweets, followers/followees, etc.

- **Memetracker** [Leskovec et al. 2009]: Memetracker is a dataset that contains the quotes and phrases that appear more frequently on the entire news spectrum. It consists of the links of the news, the time and memes, and is used to understand how information spawns, evolves and dies.

*Semantic Web and Knowledge bases.* Semantic web and RDF engines have a well accepted set of real-world datasets, used to test reasoning engines as well as to use as baselines for the creation of new synthetic generators for the semantic web.

- **DbPedia** [Bizer et al. 2009]: The DBPedia project aims at creating a structured version of Wikipedia. It allows users to semantically query the relationships between the different resources at Wikipedia.
- **Freebase** [Bollacker et al. 2008]: This is a collaborative knowledge base made of metadata maintained by community members. It contains structured data from multiple sources in an attempt to create a global resource of information accessible both from persons and machines. The project was discontinued in 2014 and replaced by Wikidata.
- **Yago** [Suchanek et al. 2007], **Yago2** [Hoffart et al. 2013] & **Yago3** [Mahdisoltani et al. 2013]: These are open source knowledge bases developed at the Max Planck Institute of Computer Science in Saarbrücken. They contain structured data that has been automatically extracted from Wikipedia and other sources.
- **Wikidata** [Vrandečić and Krötzsch 2014]: This is a knowledge base collaboratively edited by the community and hosted by the Wikimedia Foundation. It contains structured data from other related projects such as Wikipedia, Wikivoyage or Wikisource.
- **Billion Triple Challenge** [Käfer and Harth 2014]: This is a semantic dataset crawled from different sources, including Freebase, DBPedia or the BBC. A detailed description of the crawling process can be found in [Käfer et al. 2012].

## 5. CHALLENGES AND OPEN PROBLEMS

To conclude the overview of the state-of-the-art of graph data generation, in this section we discuss several of the open challenges.

### 5.1. Simple Usage, Simple Parameters

The proposal of a data generator (not necessarily for graph data) has to face an important schism. On one hand, it must provide the user with as many parameters as possible in order to enable him/her to generate arbitrary data. This approach seems to be reasonable, but it entails a shortcoming due to the fact that ordinary users are unwilling to use complex benchmarking tools. This observation can be seen, for example, in the case of XML benchmarks – even though there exist robust and complex data generators (such as ToXGene [Barbosa et al. 2002], which supports the specification of structural aspects, value distributions, references etc.), the most popular benchmarking tool is XMark [Schmidt et al. 2002], which models a single use case and enables its users to specify just the size of the data. Hence, the other extreme is to provide a simple data generator which does not require any complex settings and thus guarantees a simple and fast benchmarking process.

Considering the complex structure of graph data and the variety of applications requiring highly specific types of graphs, the latter solution is difficult to implement. A reasonable compromise can be found in a data generator which is provided with sample graph data and is capable of automatic analysis of its structural and value features in order to learn the complex parameters. We could see this type approach in some cases, such as Semantic Web generators LBBM or DBPSB.

## 5.2. Large Scale Graphs with Realistic Structure

Most of existing graph generators are focused on generating large graphs with realistic structural characteristics and focus principally on reproducing the degree distribution and the clustering coefficient [Kolda et al. 2014; Edunov et al. 2016]. However, there are other structural characteristics that one might be interested in reproducing for a large graph, such as the diameter, the size of the largest connected component, or the hierarchical community structure. Graph practitioners are highly interested in knowing how other high-level structural characteristics affect the performance of graph queries and graph algorithms. Hence, a compelling open challenge consists of creating graph generators that allow one to reproduce diverse structural characteristics of the graphs along with large scale sizes.

## 5.3. Single- vs. Multi-domain

Most of existing graph generators also generate graphs that are either not labeled or are specific to a given domain (e.g. social networks). Graphs from different domains have different schemas, structural characteristics, property distributions, etc. which might have an impact on the performance of the application under test. Thus, graph processing engine developers are asking for generators or tools to generate multi-domain graphs in a flexible and holistic manner, allowing to configure aspects such as size, schemata, data distributions and other structural characteristics such as degree distributions, clustering coefficients, and so on.

## 5.4. Generating Noisy Graphs and Graphs with Anomalies

Injecting noise and/or anomalies and errors into graphs is crucial for testing both machine learning algorithms working on this complex data and data quality techniques aiming at detecting anomalies and repairing graph data.

Concerning the former, analyzing and labeling structural networks is deemed to be more difficult for graph datasets in the presence of noise. Since de-noising graph data is difficult to achieve, various machine learning-based approaches have been adapted to work with noise (i.e. mislabeled samples) or outliers, such as imbalanced graph classification [Pan and Zhu 2013] and binary graph classification with positive and negative weights [Cheung et al. 2016]. Synthetic graph generators that take into account noisy and missing data have been studied in [Jr. and Getoor 2010], where graph identification is presented in order to model the inference of a cleaned output network from a noisy input graph. Concerning the latter, data quality techniques handling graph data are recently considering ad-hoc generation of graph data and graph quality rules in order to evaluate the effectiveness of error detection and data repairing algorithms [Fan et al. 2016; Arioua and Bonifati 2018]. The corresponding graph quality rules are typically handcrafted by domain experts, whereas an automatic generation of such rules along with the graph data generation in tandem would be an interesting future challenge for the community.

## 5.5. Streaming Graph Generators

Stream computing is a new paradigm that is necessitated by various modern data generation scenarios such as the ubiquity of mobile devices, location services, sensor pervasiveness and emerging IoT applications. These applications generate the data with high Velocity, one of the main 3V characteristics of Big Data applications [Sakr 2016]. In most of these high speed data generation scenarios, various objects are connected together with different relations and data exchanges in a graph-structured manner. The Semantic Web community has been considering the aspect of implementing streaming RDF generator and benchmarks, however, there is still a clear lack on considering this



aspect in other important and timely domains such as IoT. In addition, graph streaming generators should consider some specific aspects for the stream processing domains such as the out-of-order handling (late arrivals) [Li et al. 2008] and the variety in the schemas and formats of the different data streaming sources. It is also recommended for the streaming graph generators to support the distributed environment as this is the most common scenario for such type of applications.

### 5.6. Evolving Graph Data

As user requirements as well as environments change, most of the existing applications naturally evolve over time, at least to some extent. This evolution usually influences the structure of the data and consequently all the related parts of the application (i.e., storage strategies, operations, indexes etc.). In the world of graph data such graphs that change with time are denoted as *evolving*, *temporal*, *dynamic*, or *time-varying*. They can be modeled as labeled graphs, where the labels capture some measure of time [Michail 2015].

The evolution of graphs can be considered from multiple perspectives. We can assume a static set of nodes and a varying set of edges. Or, there are applications where the graphs only “grow”, i.e., the set of nodes and/or edges is only extended with new items. In the most general case we can assume any changes in both set of nodes and set of edges. Anyway with the evolution aspect the complexity of classical graph problems increases significantly [Michail 2015; Wu et al. 2014a]. In some graph applications, such as, e.g., social networks, the evolution of the data is a significant aspect, especially in the activity graphs [Doreian and Stokman 1997; Kumar et al. 2006; Hellmann and Staudigl 2014; Wang et al. 2013; Kossinets and Watts 2006; Viswanath et al. 2009]. However, as shown in [Leskovec et al. 2005a; Leskovec et al. 2005b], evolving graphs have further specific features. For example, some graphs grow over time according to a *densification power law* which means that in real graphs, edges tend to appear at a higher pace than vertices, meaning that these graphs densify as they grow. Also the way the new edges are distributed has the effect of a shrinking diameter that ends up stabilizing as the graph grows with time.

A related problem is *data versioning* and its respective ability to query across multiple versions of data or to carry out general analysis. This problem has been investigated for instance within the domain of Linked Open Data [Papakonstantinou et al. 2016; Meimaris and Papastefanatos 2016; Fernández et al. 2015; Fernández et al. 2015].

The respective data generator should hence be able to simulate a natural growth and/or changes in the structure of the graph with regards to the various features of distinct use cases. However, even though the area of dynamic graphs is intensively studied, surprisingly there seem to exist only very few proposals of a generator for dealing with this area. In [Goerke et al. 2012] the authors focus on *clustering dynamic graphs*, i.e. graphs where the clustering corresponds to the partition of nodes into natural groups based on the concept of intra-cluster versus inter-cluster density of edges. The generator generates a time series of random graphs  $G_0, G_1, \dots, G_n$ , where  $G_t$  emerges from  $G_{t-1}$  via successive atomic updates like for instance, the insertion of a vertex or the removal of an edge. The generator dynamically monitors the ground truth clustering, and the probability of the updates is chosen in such a way that the ground truth is maintained while the randomness of the generated graph is kept.

Another recent proposal of a generator [Purohit et al. 2018] of temporal graphs results from an observation that small subgraph patterns in networks, called *network motifs* or *graphlets*, are crucial indicators of the structure and the evolution of the graphs [Paranjape et al. 2017]. For a given graph and a predefined ordered list of structural atomic motifs the generator first computes the distribution of the motifs in

the graph. The distribution is then used to generate a synthetic graph with the same features.

### 5.7. Multi-model Data

With the dawn of Big Data and especially its Variety, another key 3V characteristic, new types of database management systems have emerged. One of the most interesting ones are the so-called *multi-model databases* that enable to store and thus query across structurally different data. There exist various types of multi-model systems combining distinct subsets of Big Data structures including graph data. For example, OrientDB<sup>28</sup> which was implemented on the basis of an object DBMS currently supports graph, document, key/value, and object models. Such type of DBMSs also needs a specific data generator that would enable to test new features and analyze efficiency of operations. However, since the multi-model systems are in the context of Big Data rather new, there exist only a few benchmarks targeting multi-model DBMSs (such as Bigframe [Kunjir et al. 2014] or UniBench [Lu 2017]) with limited capabilities.

Another interesting approach to multi-model data is to adopt a unifying expressive graph data model, known as property graph data model [Bonifati et al. 2018b]. Such a model allows to specify multi-edges and list of properties for the nodes. Synthetic graph generators for property graphs and its companion standard graph query language [Angles 2018; Angles et al. 2018] are also needed in order to boost their availability and adoption for different communities.

### 5.8. Machine Learning Based Graph Generation

With the advent of neural networks and specially generative adversarial networks (GANs) [Goodfellow et al. 2014], several researchers have started to explore their application to generate graphs. This is the case of [Kipf and Welling 2016; Grover et al. 2018; Simonovsky and Komodakis 2018; Li et al. 2018; You et al. 2018], which present several generative models to generate realistic graphs. Such techniques still suffer from several problems. For instance, some of them are limited to learn from a single graph [Kipf and Welling 2016; Grover et al. 2018] or generate small graphs [Simonovsky and Komodakis 2018; Li et al. 2018; You et al. 2018]. The technique proposed in [You et al. 2018] is capable of generating graphs with complex edge dependencies (e.g. community structure) and is not restricted to graphs of a fixed size. However, there are still in general several open challenges, including the capability of learning from and generating large graphs comparable in size to those typically used for benchmarking, and robust generation techniques with structural guarantees (e.g. degree distribution, clustering coefficient, etc.).

### 5.9. Privacy-preserving Graph Generation

A lot of work has been conducted on techniques for publishing social network graphs with privacy guarantees [Wu et al. 2010]. However, the topic of generating social graphs with a realistic structure yet private has been barely explored.

Most of the existing work falls within the topic of graph generation with “differential privacy” [Dwork and Lei 2009] guarantees. More specifically, in [Wang and Wu 2013] the authors develop a differential privacy graph generation approach based on the dK-graph generation model [Mahadevan et al. 2006] that outperforms the Stochastic Kronecker Graph Model [Leskovec et al. 2005a] in terms of the produced structural properties, even though the results show that there is still room for improvement.

Following this line of research, recent work [Qin et al. 2017] extends the notion of differential privacy and propose an “edge local differential privacy” based graph gen-

<sup>28</sup><http://orientdb.com/orientdb/>

eration method. The proposed method allows generating privacy preserving synthetic social graphs without the need of a centralized data curator, while preserving structural properties more accurately than straw-hat methods such as Randomized Neighbor Lists (based on randomized response [Dwork et al. 2014]) and Degree-based Graph Generation (which perturbs the original graph degrees using the Laplace mechanism [Dwork and Lei 2009]). Again, even though the proposed technique outperforms the baselines, the results show that there is still room for improving the structural properties of the generated graph.

## 6. CONCLUSION

Graph data occur in a vast amount of distinct applications, such as biology, chemistry, physics, computer science, or social sciences, to name just a few. Graphs form one of the most complex data structures requiring specific and usually sophisticated approaches for processing and analysis. The history of graph theory, that started from when these structures and their respective algorithms were studied, can be traced back to the 18th century.

With the recent dawn of Big Data there have been more occurrences of large scale graphs where the efficiency of processing methods is critical. Approaches that work for smaller scale graphs often cannot be used, the data need to be processed in a distributed way and hence the efficiency is influenced by other aspects, such as limits of data transport. In addition, distribution of graphs, especially for highly connected cases, is a difficult task. Thus extensive testing of these methods for graphs of various sizes and structural complexity is extremely important.

The aim of this survey was to provide a thorough overview and comparison of graph data generators. We do not limit ourselves to a single application domain, but we cover the currently most popular areas of graph data processing. We believe that this wide scope provides a uniquely useful insight into state-of-the-art tools as well as open issues for both researchers and practitioners.

## ACKNOWLEDGMENTS

The authors would like to thank Dr. Kamesh Madduri for consultations and suggestions on the covered areas.

## REFERENCES

- 2004. *OWL Web Ontology Language Overview*. W3C. <http://www.w3.org/TR/owl-features/>
- 2008. *ISO/IEC 9075-1:2008 Information technology – Database languages – SQL – Part 1: Framework (SQL/Framework)*. ISO. [http://www.iso.org/iso/catalogue\\_detail.htm?csnumber=45498](http://www.iso.org/iso/catalogue_detail.htm?csnumber=45498)
- 2009. *HPC Scalable Graph Analysis Benchmark*. GraphAnalysis.org. <http://www.graphanalysis.org/benchmark/>
- 2015. *GraphGen: Graph generator for Neo4j*. Graph Aware Ltd. <http://graphgen.graphaware.com/>
- 2015. *Semantic Publishing Benchmark v2.0*. LDBC. <http://ldbouncil.org/developer/spb>
- 2016. *DBLP Computer Science Bibliography*. DBLP team. <http://dblp.uni-trier.de/>
- 2016. *The TPC Benchmark – H*. <http://www.tpc.org/tpch/>
- Daniel J. Abadi, Adam Marcus, Samuel R. Madden, and Kate Hollenbach. 2007. *Using The Barton Libraries Dataset As An RDF Benchmark*. Technical Report MIT-CSAIL-TR-2007-036. MIT.
- Lada A Adamic and Natalie Glance. 2005. The political blogosphere and the 2004 US election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery*. ACM, 36–43.
- Charu C. Aggarwal and Karthik Subbian. 2014. Evolutionary Network Analysis: A Survey. *ACM Comput. Surv.* 47, 1 (2014), 10:1–10:36.
- Leman Akoglu, Mary McGlohon, and Christos Faloutsos. 2008. RTM: Laws and a Recursive Generator for Weighted Time-Evolving Graphs. In *ICDM*. IEEE Computer Society, 701–706.

- Awrad Mohammed Ali, Hamidreza Alvani, Alireza Hajibagheri, Kiran Lakkaraju, and Gita Sukthankar. 2014. Synthetic Generators for Cloning Social Network Data. In *Proceedings of the ASE International Conference on Social Informatics*. Cambridge, MA.
- Güneş Aluç, Olaf Hartig, M. Tamer Özsu, and Khuzaima Daudjee. 2014. Diversified Stress Testing of RDF Data Management Systems. In *Proceedings of the 13th International Semantic Web Conference - Part I (ISWC '14)*. Springer-Verlag New York, Inc., New York, NY, USA, 197–212. DOI: [http://dx.doi.org/10.1007/978-3-319-11964-9\\_13](http://dx.doi.org/10.1007/978-3-319-11964-9_13)
- Renzo Angles. 2018. The Property Graph Database Model. In *AMW (CEUR Workshop Proceedings)*, Vol. 2100. CEUR-WS.org.
- Renzo Angles, Marcelo Arenas, Pablo Barceló, Peter A. Boncz, George H. L. Fletcher, Claudio Gutierrez, Tobias Lindaaaker, Marcus Paradies, Stefan Plantikow, Juan F. Sequeda, Oskar van Rest, and Hannes Voigt. 2018. G-CORE: A Core for Future Graph Query Languages. In *SIGMOD Conference*. ACM, 1421–1432.
- Renzo Angles, Peter Boncz, Josep Larriba-Pey, Irini Fundulaki, Thomas Neumann, Orri Erling, Peter Neubauer, Norbert Martinez-Bazan, Venelin Kotsev, and Ioan Toma. 2014. The Linked Data Benchmark Council: A Graph and RDF Industry Benchmarking Effort. *SIGMOD Rec.* 43, 1 (May 2014), 27–31. DOI: <http://dx.doi.org/10.1145/2627692.2627697>
- Abdallah Arioua and Angela Bonifati. 2018. User-guided Repairing of Inconsistent Knowledge Bases. In *EDBT*. OpenProceedings.org, 133–144.
- Timothy G. Armstrong, Vamsi Ponnkanti, Dhruva Borthakur, and Mark Callaghan. 2013. LinkBench: A Database Benchmark Based on the Facebook Social Graph. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data (SIGMOD '13)*. ACM, New York, NY, USA, 1185–1196. DOI: <http://dx.doi.org/10.1145/2463676.2465296>
- David A. Bader and Kamesh Madduri. 2005. Design and Implementation of the HPCS Graph Analysis Benchmark on Symmetric Multiprocessors. In *Proceedings of the 12th International Conference on High Performance Computing (HiPC'05)*. Springer-Verlag, Berlin, Heidelberg, 465–476. DOI: [http://dx.doi.org/10.1007/11602569\\_48](http://dx.doi.org/10.1007/11602569_48)
- David A. Bader and Kamesh Madduri. 2006. *GTgraph – A Suite of Synthetic Random Graph Generators*. <http://www.cse.psu.edu/~kxm85/software/GTgraph/index.html>
- David A. Bader, Henning Meyerhenke, Peter Sanders, and Dorothea Wagner. 2012. *10th DIMACS Implementation Challenge Workshop*. Georgia Institute of Technology, Atlanta, GA, USA. <http://www.cc.gatech.edu/dimacs10/>
- Guillaume Bagan, Angela Bonifati, Radu Ciucanu, George Fletcher, Aurélien Lemay, and Nicky Advokaat. 2016. gMark: Schema-Driven Generation of Graphs and Queries. *IEEE Transactions on Knowledge and Data Engineering* (Nov. 2016). <https://hal.inria.fr/hal-01402575>
- Albert-Laszlo Barabasi and Reka Albert. 1999. Emergence of Scaling in Random Networks. *Science* 286, 5439 (1999), 509–512. DOI: <http://dx.doi.org/10.1126/science.286.5439.509>
- Denilson Barbosa, Alberto O. Mendelzon, John Keenleyside, and Kelly A. Lyons. 2002. ToXgene: An extensible template-based data generator for XML.. In *WebDB* (2003-03-06). 49–54. <http://dblp.uni-trier.de/db/conf/webdb/webdb2002.html#BarbosaMKL02>
- Christopher L. Barrett, Richard J. Beckman, Maleq Khan, V. S. Anil Kumar, Madhav V. Marathe, Paula E. Stretz, Tridib Dutta, and Bryan Lewis. 2009. Generation and Analysis of Large Synthetic Social Contact Networks. In *Winter Simulation Conference (WSC '09)*. Winter Simulation Conference, 1003–1014. <http://dl.acm.org/citation.cfm?id=1995456.1995598>
- Robert Battle and Dave Kolas. 2012. Enabling the geospatial semantic web with parliament and geosparql. *Semantic Web* 3, 4 (2012), 355–370.
- Sotirios Beis, Symeon Papadopoulos, and Yiannis Kompatsiaris. 2015. *Benchmarking Graph Databases on the Problem of Community Detection*. Springer International Publishing, Cham, 3–14. DOI: [http://dx.doi.org/10.1007/978-3-319-10518-5\\_1](http://dx.doi.org/10.1007/978-3-319-10518-5_1)
- Garrett Bernstein and Kyle O'Brien. 2013. Stochastic Agent-based Simulations of Social Networks. In *Proceedings of the 46th Annual Simulation Symposium (ANSS 13)*. Society for Computer Simulation International, San Diego, CA, USA, Article 5, 8 pages. <http://dl.acm.org/citation.cfm?id=2499604.2499609>
- Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. 2009. DBpedia - A Crystallization Point for the Web of Data. *Web Semant.* 7, 3 (Sept. 2009), 154–165. DOI: <http://dx.doi.org/10.1016/j.websem.2009.07.002>
- Christian Bizer and Andreas Schultz. 2009. The Berlin SPARQL Benchmark. *International Journal On Semantic Web and Information Systems* (2009).

- Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. 2008. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* 2008, 10 (2008), P10008. <http://stacks.iop.org/1742-5468/2008/i=10/a=P10008>
- Paolo Boldi, Massimo Santini, and Sebastiano Vigna. 2008. A Large Time-Aware Graph. *SIGIR Forum* 42, 2 (2008), 33–38.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. AcM, 1247–1250.
- Peter Boncz, Minh-Duc Pham, Orri Erling, Ivan Mikhailov, and Yrjana Rankka. 2013. *Social Network Intelligence BenchMark*. [https://www.w3.org/wiki/Social\\_Network\\_Intelligence\\_BenchMark](https://www.w3.org/wiki/Social_Network_Intelligence_BenchMark)
- Angela Bonifati, George Fletcher, Jan Hidders, and Alexandre Iosup. 2018a. A Survey of Benchmarks for Graph-Processing Systems. In *Graph Data Management*. Springer.
- Angela Bonifati, George Fletcher, Hannes Voigt, and Nikolay Yakovets. 2018b. *Querying Graphs*. Morgan & Claypool Publishers.
- J. M. Carlson and John Doyle. 2000. Highly Optimized Tolerance: Robustness and Design in Complex Systems. *Phys. Rev. Lett.* 84 (Mar 2000), 2529–2532. Issue 11. DOI: <http://dx.doi.org/10.1103/PhysRevLett.84.2529>
- Hassan Chafi, Jason Crawford, Alastair Green, and Keith Hare. 2018. Graph Query Language GQL. <https://www.gqlstandards.org/>. (2018).
- Deepayan Chakrabarti and Christos Faloutsos. 2006. Graph Mining: Laws, Generators, and Algorithms. *ACM Comput. Surv.* 38, 1, Article 2 (June 2006). DOI: <http://dx.doi.org/10.1145/1132952.1132954>
- Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. 2004. R-MAT: A Recursive Model for Graph Mining. In *Proceedings of the Fourth SIAM International Conference on Data Mining, Lake Buena Vista, Florida, USA, April 22-24, 2004*. 442–446. DOI: <http://dx.doi.org/10.1137/1.9781611972740.43>
- James Cheng, Yiping Ke, Wilfred Ng, and An Lu. 2007. Fg-index: Towards Verification-free Query Processing on Graph Databases. In *Proceedings of the 2007 ACM SIGMOD International Conference on Management of Data (SIGMOD '07)*. ACM, New York, NY, USA, 857–872. DOI: <http://dx.doi.org/10.1145/1247480.1247574>
- Gene Cheung, Weng-Tai Su, Yu Mao, and Chia-Wen Lin. 2016. Robust Semi-Supervised Graph Classifier Learning with Negative Edge Weights. *CoRR* abs/1611.04924 (2016).
- Avery Ching, Sergey Edunov, Maja Kabiljo, Dionysios Logothetis, and Sambavi Muthukrishnan. 2015. One trillion edges: Graph processing at facebook-scale. *Proceedings of the VLDB Endowment* 8, 12 (2015), 1804–1815.
- Marek Ciglan, Alex Averbuch, and Ladi Lav Hluchy. 2012. Benchmarking Traversal Operations over Graph Databases. In *Proceedings of the 2012 IEEE 28th International Conference on Data Engineering Workshops (ICDEW '12)*. IEEE Computer Society, Washington, DC, USA, 186–189. DOI: <http://dx.doi.org/10.1109/ICDEW.2012.47>
- Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking Cloud Serving Systems with YCSB. In *Proceedings of the 1st ACM Symposium on Cloud Computing (SoCC '10)*. ACM, New York, NY, USA, 143–154. DOI: <http://dx.doi.org/10.1145/1807128.1807152>
- Leon Danon, Albert Diaz-Guilera, Jordi Duch, and Alex Arenas. 2005. Comparing community structure identification. *Journal of Statistical Mechanics: Theory and Experiment* 2005, 09 (2005), P09008.
- Miyuru Dayarathna and Toyotaro Suzumura. 2014. Graph Database Benchmarking on Cloud Environments with XGDBench. *Automated Software Engg.* 21, 4 (Dec. 2014), 509–533. DOI: <http://dx.doi.org/10.1007/s10515-013-0138-7>
- Manlio De Domenico, Antonio Lima, Paul Mougél, and Mirco Musolesi. 2013. The anatomy of a scientific rumor. *Scientific reports* 3 (2013), 2980.
- Universita degli studio di Milano. 2018. Laboratory of Web Algorithmics. <http://law.di.unimi.it/>. (2018). [Online; accessed 03-July-2018].
- David Dominguez-Sal, Norbert Martinez-Bazan, Victor Mentes-Mulero, Pere Baleta, and Josep Lluís Larriba-Pay. 2011. A Discussion on the Design of Graph Database Benchmarks. In *Proceedings of the Second TPC Technology Conference on Performance Evaluation, Measurement and Characterization of Complex Systems (TPCTC'10)*. Springer-Verlag, Berlin, Heidelberg, 25–40. <http://dl.acm.org/citation.cfm?id=1946050.1946053>
- D. Dominguez-Sal, P. Urbón-Bayes, A. Giménez-Vañó, S. Gómez-Villamor, N. Martínez-Bazán, and J. L. Larriba-Pey. 2010. 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)*. Springer-Verlag, Berlin, Heidelberg, 37–48. <http://dl.acm.org/citation.cfm?id=1927585.1927590>

- P. Doreian and F.N. Stokman. 1997. *Evolution of Social Networks*. Number sv. 1 in The journal of mathematical sociology. Gordon and Breach Publishers. <https://books.google.com/books?id=ZL4zCCgfmOkC>
- Songyun Duan, Anastasios Kementsietsidis, Kavitha Srinivas, and Octavian Udrea. 2011. Apples and Oranges: A Comparison of RDF Benchmarks and Real RDF Datasets. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data (SIGMOD '11)*. ACM, New York, NY, USA, 145–156. DOI: <http://dx.doi.org/10.1145/1989323.1989340>
- Cynthia Dwork and Jing Lei. 2009. Differential privacy and robust statistics. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*. ACM, 371–380.
- Cynthia Dwork, Aaron Roth, and others. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.
- Sergey Edunov, Dionysios Logothetis, Cheng Wang, Avery Ching, and Maja Kabiljo. 2016. Darwini: Generating realistic large-scale social graphs. *arXiv preprint arXiv:1610.00664* (2016).
- Paul Erdos and Alfred Renyi. 1960. On the evolution of random graphs. *Publ. Math. Inst. Hungary. Acad. Sci.* 5 (1960), 17–61.
- Orri Erling, Alex Averbuch, Josep Larriba-Pey, Hassan Chafi, Andrey Gubichev, Arnau Prat, Minh-Duc Pham, and Peter Boncz. 2015. The LDBC Social Network Benchmark: Interactive Workload. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM, New York, NY, USA, 619–630. DOI: <http://dx.doi.org/10.1145/2723372.2742786>
- Wenfei Fan, Yinghui Wu, and Jingbo Xu. 2016. Functional Dependencies for Graphs. In *SIGMOD Conference*. ACM, 1843–1857.
- Javier D Fernández, Axel Polleres, and Jürgen Umbrich. 2015. Towards Efficient Archiving of Dynamic Linked Open Data. *Proc. of DIACHRON* (2015), 34–49.
- Javier D. Fernández, Jürgen Umbrich, and Axel Polleres. 2015. BEAR: Benchmarking the Efficiency of RDF Archiving. (2015).
- Alfio Ferrara, Davide Lorusso, Stefano Montanelli, and Gaia Varese. 2008. Towards a Benchmark for Instance Matching. In *Ontology Matching (OM 2008) (CEUR Workshop Proceedings)*, Pavel Shvaiko, Jerome Euzenat, Fausto Giunchiglia, and Heiner Stuckenschmidt (Eds.), Vol. 431. CEUR-WS.org.
- George Garbis, Kostis Kyzirakos, and Manolis Koubarakis. 2013. Geographica: A Benchmark for Geospatial RDF Stores (Long Version). In *The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part II*. 343–359. DOI: [http://dx.doi.org/10.1007/978-3-642-41338-4\\_22](http://dx.doi.org/10.1007/978-3-642-41338-4_22)
- Maria Giatsoglou, Symeon Papadopoulos, and Athena Vakali. 2011. *Massive Graph Management for the Web and Web 2.0*. Springer Berlin Heidelberg, Berlin, Heidelberg, 19–58. DOI: [http://dx.doi.org/10.1007/978-3-642-17551-0\\_2](http://dx.doi.org/10.1007/978-3-642-17551-0_2)
- M Girvan and MEJ Newman. 2002. Network of american football games between division ia colleges during regular season fall 2000. *Proc. Natl. Acad. Sci. USA* 99 (2002), 7821–7826.
- Robert Goerke, Roland Kluge, Andrea Schumm, Christian Staudt, and Dorothea Wagner. 2012. *An Efficient Generator for Clustered Dynamic Random Networks*. Technical Report 17. Karlsruhe.
- Ken Goldberg, Theresa Roeder, Dhruv Gupta, and Chris Perkins. 2001. Eigentaste: A constant time collaborative filtering algorithm. *information retrieval* 4, 2 (2001), 133–151.
- Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. 2012. PowerGraph: Distributed Graph-parallel Computation on Natural Graphs. In *Proceedings of the 10th USENIX Conference on Operating Systems Design and Implementation (OSDI'12)*. USENIX Association, Berkeley, CA, USA, 17–30. <http://dl.acm.org/citation.cfm?id=2387880.2387883>
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.
- Michael Grossniklaus, Stefania Leone, and Tilmann Zschke. 2013. *Towards a benchmark for graph data management and processing*. Technical Report.
- GroupLens. 2011. HetRec. <https://grouplens.org/datasets/hetrec-2011/>. (2011). [Online; accessed 03-July-2018].
- GroupLens. 2018. Movielens. <https://grouplens.org/datasets/movielens/>. (2018). [Online; accessed 03-July-2018].
- Aditya Grover, Aaron Zweig, and Stefano Ermon. 2018. Graphite: Iterative generative modeling of graphs. *arXiv preprint arXiv:1803.10459* (2018).
- Yuanbo Guo, Zhengxiang Pan, and Jeff Heflin. 2005. LUBM: A benchmark for {OWL} knowledge base systems. *Web Semantics: Science, Services and Agents on the World Wide Web* 3, 2–3 (2005), 158 –

182. DOI: <http://dx.doi.org/10.1016/j.websem.2005.06.005> Selected Papers from the International Semantic Web Conference, 2004 ISWC, 43rd. International Semantic Web Conference, 2004.
- Tim Hellmann and Mathias Staudigl. 2014. Evolution of social networks. *European Journal of Operational Research* 234, 3 (2014), 583 – 596. DOI: <http://dx.doi.org/10.1016/j.ejor.2013.08.022>
- Johannes Hoffart, Fabian M Suchanek, Klaus Berberich, and Gerhard Weikum. 2013. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *Artificial Intelligence* 194 (2013), 28–61.
- Darko Hric, Richard K Darst, and Santo Fortunato. 2014. Community detection in networks: Structural communities versus ground truth. *Physical Review E* 90, 6 (2014), 062805.
- Alexandru Iosup, Tim Hegeman, Wing Lung Ngai, Stijn Heldens, Arnau Prat-Pérez, Thomas Manhardt, Hassan Chafio, Mihai Capotă, Narayanan Sundaram, Michael Anderson, Ilie Gabriel Tănase, Ying-long Xia, Lifeng Nai, and Peter Boncz. 2016. LDBC Graphalytics: A Benchmark for Large-scale Graph Analysis on Parallel and Distributed Platforms. *Proc. VLDB Endow.* 9, 13 (Sept. 2016), 1317–1328. DOI: <http://dx.doi.org/10.14778/3007263.3007270>
- Amit Krishna Joshi, Pascal Hitzler, and Guozhu Dong. 2016. LinkGen: Multipurpose Linked Data Generator. In *The Semantic Web – ISWC 2016*, Paul Groth, Elena Simperl, Alasdair Gray, Marta Sabou, Markus Krötzsch, Freddy Lecue, Fabian Flöck, and Yolanda Gil (Eds.). Springer International Publishing, Cham, 113–121.
- Galileo Mark S. Namata Jr. and Lise Getoor. 2010. Identifying graphs from noisy and incomplete data. *SIGKDD Explorations* 12, 1 (2010), 33–39.
- Tobias Käfer and Andreas Harth. 2014. Billion Triples Challenge data set. Downloaded from <http://km.aifb.kit.edu/projects/btc-2014/>. (2014).
- Tobias Käfer, Jürgen Umbrich, Aidan Hogan, and Axel Polleres. 2012. Towards a dynamic linked data observatory. *LDOW at WWW* (2012).
- Vadim Kagan, Andrew Stevens, and VS Subrahmanian. 2015. Using twitter sentiment to forecast the 2013 pakistani election and the 2014 indian election. *IEEE Intelligent Systems* 1 (2015), 2–5.
- Sepandar D Kamvar, Mario T Schlosser, and Hector Garcia-Molina. 2003. The eigentrust algorithm for reputation management in p2p networks. In *Proceedings of the 12th international conference on World Wide Web*. ACM, 640–651.
- Yarden Katz and Jennifer Golbeck. 2006. Social network-based trust in prioritized default logic. In *AAAI*, Vol. 6. 1345–1350.
- Jungeun Kim and Jae-Gil Lee. 2015. Community Detection in Multi-Layer Graphs: A Survey. *SIGMOD Rec.* 44, 3 (Dec. 2015), 37–48. DOI: <http://dx.doi.org/10.1145/2854006.2854013>
- Myunghwan Kim and Jure Leskovec. 2010. *Multiplicative Attribute Graph Model of Real-World Networks*. Springer Berlin Heidelberg, Berlin, Heidelberg, 62–73. DOI: [http://dx.doi.org/10.1007/978-3-642-18009-5\\_7](http://dx.doi.org/10.1007/978-3-642-18009-5_7)
- Thomas N Kipf and Max Welling. 2016. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308* (2016).
- Donald Knuth. 2008. *The Stanford GraphBase*. <http://www3.cs.stonybrook.edu/~algorithm/implement/graphbase/implement.shtml>
- Tamara G Kolda, Ali Pinar, Todd Plantenga, and Comandur Seshadhri. 2014. A scalable generative graph model with community structure. *SIAM Journal on Scientific Computing* 36, 5 (2014), C424–C452.
- Gueorgi Kossinets and Duncan J. Watts. 2006. Empirical Analysis of an Evolving Social Network. *Science* 311, 5757 (2006), 88–90. DOI: <http://dx.doi.org/10.1126/science.1116869>
- Manolis Koubarakis and Kostis Kyzirakos. 2010. Modeling and querying metadata in the semantic sensor web: The model stRDF and the query language stSPARQL. In *Extended Semantic Web Conference*. Springer, 425–439.
- Ravi Kumar, Jasmine Novak, and Andrew Tomkins. 2006. Structure and Evolution of On-line Social Networks. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '06)*. ACM, New York, NY, USA, 611–617. DOI: <http://dx.doi.org/10.1145/1150402.1150476>
- Srijan Kumar, Francesca Spezzano, VS Subrahmanian, and Christos Faloutsos. 2016. Edge Weight Prediction in Weighted Signed Networks.. In *ICDM*. 221–230.
- Mayuresh Kunjir, Prajakta Kalmegh, and Shivnath Babu. 2014. Thoth: Towards Managing a Multi-System Cluster. *PVLDB* 7, 13 (2014), 1689–1692.
- Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. 2010. What is Twitter, a social network or a news media?. In *Proceedings of the 19th international conference on World wide web*. ACM, 591–600.

- Andrea Lancichinetti and Santo Fortunato. 2009. Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities. *Phys. Rev. E* 80, 1 (July 2009), 016118. DOI: <http://dx.doi.org/10.1103/PhysRevE.80.016118>
- Andrea Lancichinetti, Santo Fortunato, and Filippo Radicchi. 2008. Benchmark graphs for testing community detection algorithms. *Phys. Rev. E* 78, 4 (Oct. 2008), 046110. DOI: <http://dx.doi.org/10.1103/PhysRevE.78.046110>
- Danh Le-Phuoc, Minh Dao-Tran, Minh-Duc Pham, Peter Boncz, Thomas Eiter, and Michael Fink. 2012. Linked stream data processing engines: Facts and figures. In *International Semantic Web Conference*. Springer, 300–312.
- Jure Leskovec, Lars Backstrom, and Jon Kleinberg. 2009. Meme-tracking and the dynamics of the news cycle. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 497–506.
- Jurij Leskovec, Deepayan Chakrabarti, Jon Kleinberg, and Christos Faloutsos. 2005a. Realistic, Mathematically Tractable Graph Generation and Evolution, Using Kronecker Multiplication. In *Proceedings of the 9th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD'05)*. Springer-Verlag, Berlin, Heidelberg, 133–145. DOI: [http://dx.doi.org/10.1007/11564126\\_17](http://dx.doi.org/10.1007/11564126_17)
- Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. 2005b. Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining (KDD '05)*. ACM, New York, NY, USA, 177–187. DOI: <http://dx.doi.org/10.1145/1081870.1081893>
- Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>. (June 2014).
- Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney. 2008. Statistical Properties of Community Structure in Large Social and Information Networks. In *Proceedings of the 17th International Conference on World Wide Web (WWW '08)*. ACM, New York, NY, USA, 695–704. DOI: <http://dx.doi.org/10.1145/1367497.1367591>
- Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney. 2009. Community structure in large networks: Natural cluster sizes and the absence of large well-defined clusters. *Internet Mathematics* 6, 1 (2009), 29–123.
- Jin Li, Kristin Tufte, Vladislav Shkapenyuk, Vassilis Papadimos, Theodore Johnson, and David Maier. 2008. Out-of-order processing: a new architecture for high-performance stream systems. *Proceedings of the VLDB Endowment* 1, 1 (2008), 274–288.
- Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, and Peter Battaglia. 2018. Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324* (2018).
- Jiaheng Lu. 2017. Towards Benchmarking Multi-Model Databases. In *CIDR 2017, 8th Biennial Conference on Innovative Data Systems Research, Chaminade, CA, USA, January 8-11, 2017, Online Proceedings*. [http://cidrdb.org/cidr2017/gongshow/abstracts/cidr2017\\\_20.pdf](http://cidrdb.org/cidr2017/gongshow/abstracts/cidr2017\_20.pdf)
- Li Ma, Yang Yang, Zhaoming Qiu, Guotong Xie, Yue Pan, and Shengping Liu. 2006. Towards a Complete OWL Ontology Benchmark. In *Proceedings of the 3rd European Conference on The Semantic Web: Research and Applications (ESWC'06)*. Springer-Verlag, Berlin, Heidelberg, 125–139. DOI: [http://dx.doi.org/10.1007/11762256\\_12](http://dx.doi.org/10.1007/11762256_12)
- Priya Mahadevan, Dmitri Krioukov, Kevin Fall, and Amin Vahdat. 2006. Systematic topology analysis and generation using degree correlations. In *ACM SIGCOMM Computer Communication Review*, Vol. 36. ACM, 135–146.
- Farzaneh Mahdisoltani, Joanna Biega, and Fabian M Suchanek. 2013. Yago3: A knowledge base from multilingual wikipedias. In *CIDR*.
- Grzegorz Malewicz, Matthew H Austern, Aart JC Bik, James C Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM, 135–146.
- Silviu Maniu, Bogdan Cautis, and Talel Abdesslem. 2011. Building a signed network from interactions in Wikipedia. In *Databases and Social Networks*. ACM, 19–24.
- Paolo Massa and Paolo Avesani. 2007. Trust-aware recommender systems. In *Proceedings of the 2007 ACM conference on Recommender systems*. ACM, 17–24.
- Andrea Mauri, Jean-Paul Calbimonte, Daniele Dell'Aglio, Marco Balduini, Marco Brambilla, Emanuele Della Valle, and Karl Aberer. 2016. Triplewave: Spreading RDF streams on the web. In *International Semantic Web Conference*. Springer, 140–149.
- Andrew McGregor. 2014. Graph stream algorithms: a survey. *ACM SIGMOD Record* 43, 1 (2014), 9–20.
- Brendan McKay and Adolfo Piperno. 2014. *nauty and Traces*. <http://pallini.di.uniroma1.it/>



- Marios Meimaris and George Papastefanatos. 2016. The EvoGen Benchmark Suite for Evolving RDF Data. In *Joint Proceedings of the 2nd Workshop on Managing the Evolution and Preservation of the Data Web (MEPDaW 2016) and the 3rd Workshop on Linked Data Quality (LDQ 2016) co-located with 13th European Semantic Web Conference (ESWC 2016), Heraklion, Crete, Greece, May 30th, 2016*. 20–35. [http://ceur-ws.org/Vol-1585/mepdaw2016\\_paper.03.pdf](http://ceur-ws.org/Vol-1585/mepdaw2016_paper.03.pdf)
- Othon Michail. 2015. *An Introduction to Temporal Graphs: An Algorithmic Perspective*. Springer International Publishing, Cham, 308–343. DOI: [http://dx.doi.org/10.1007/978-3-319-24024-4\\_18](http://dx.doi.org/10.1007/978-3-319-24024-4_18)
- George A. Miller. 1995. WordNet: A Lexical Database for English. *Commun. ACM* 38, 11 (Nov. 1995), 39–41. DOI: <http://dx.doi.org/10.1145/219717.219748>
- Tyler Moore and Nicolas Christin. 2013. Beware the middleman: Empirical analysis of Bitcoin-exchange risk. In *International Conference on Financial Cryptography and Data Security*. Springer, 25–33.
- Mohamed Morsey, Jens Lehmann, Sören Auer, and Axel-Cyrille Ngonga Ngomo. 2012. Usage-centric Benchmarking of RDF Triple Stores. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence (AAAI'12)*. AAAI Press, 2134–2140. <http://dl.acm.org/citation.cfm?id=2900929.2901031>
- Mohamed Morsey, Jens Lehmann, Sören Auer, and Axel-Cyrille Ngonga Ngomo. 2011. *DBpedia SPARQL Benchmark – Performance Assessment with Real Queries on Real Data*. Springer Berlin Heidelberg, Berlin, Heidelberg, 454–469. DOI: [http://dx.doi.org/10.1007/978-3-642-25073-6\\_29](http://dx.doi.org/10.1007/978-3-642-25073-6_29)
- David F. Nettleton. 2016. A synthetic data generator for online social network graphs. *Social Network Analysis and Mining* 6, 1 (2016), 44. DOI: <http://dx.doi.org/10.1007/s13278-016-0352-y>
- Shirui Pan and Xingquan Zhu. 2013. Graph Classification with Imbalanced Class Distributions and Noise. In *IJCAI 2013, Proceedings of the 23rd International Joint Conference on Artificial Intelligence, Beijing, China, August 3-9, 2013*. 1586–1592.
- Vassilis Papakonstantinou, Giorgos Flouris, Irimi Fundulaki, Kostas Stefanidis, and Giannis Roussakis. 2016. Versioning for Linked Data: Archiving Systems and Benchmarks. In *Proceedings of the Workshop on Benchmarking Linked Data (BLINK 2016) co-located with the 15th International Semantic Web Conference (ISWC), Kobe, Japan, October 18, 2016*. <http://ceur-ws.org/Vol-1700/paper-05.pdf>
- Ashwin Paranjape, Austin R. Benson, and Jure Leskovec. 2017. Motifs in Temporal Networks. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining (WSDM '17)*. ACM, New York, NY, USA, 601–610. DOI: <http://dx.doi.org/10.1145/3018661.3018731>
- Minh-Duc Pham, Peter Boncz, and Orri Erling. 2013. *S3G2: A Scalable Structure-Correlated Social Graph Generator*. Springer Berlin Heidelberg, Berlin, Heidelberg, 156–172. DOI: [http://dx.doi.org/10.1007/978-3-642-36727-4\\_11](http://dx.doi.org/10.1007/978-3-642-36727-4_11)
- Natalia Pobiedina, Stefan Rümmele, Sebastian Skritek, and Hannes Werthner. 2014. *Benchmarking Database Systems for Graph Pattern Matching*. Springer International Publishing, Cham, 226–241. DOI: [http://dx.doi.org/10.1007/978-3-319-10073-9\\_18](http://dx.doi.org/10.1007/978-3-319-10073-9_18)
- Arnau Prat-Pérez and David Dominguez-Sal. 2014. How Community-like is the Structure of Synthetically Generated Graphs?. In *Proceedings of Workshop on GRaph Data Management Experiences and Systems (GRADES'14)*. ACM, New York, NY, USA, Article 7, 9 pages. DOI: <http://dx.doi.org/10.1145/2621934.2621942>
- The Lemur Project. 2012. ClueWeb2012. <http://law.di.unimi.it/webdata/clueweb12/>. (2012). [Online; accessed 02-July-2018].
- Eric Prud'hommeaux and Andy Seaborne. 2008. *SPARQL Query Language for RDF*. <http://www.w3.org/TR/rdf-sparql-query/>
- Sumit Purohit, Lawrence B Holder, and George Chin. 2018. Temporal Graph Generation Based on a Distribution of Temporal Motifs. In *Proceedings of the 14th International Workshop on Mining and Learning with Graphs (MLG)*.
- Shi Qiao and Z. Meral Özsoyoglu. 2015. RBench: Application-Specific RDF Benchmarking. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM, New York, NY, USA, 1825–1838. DOI: <http://dx.doi.org/10.1145/2723372.2746479>
- Zhan Qin, Ting Yu, Yin Yang, Issa Khalil, Xiaokui Xiao, and Kui Ren. 2017. Generating synthetic decentralized social graphs with local differential privacy. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 425–438.
- Carlos R. Rivero, Andreas Schultz, Christian Bizer, and David Ruiz. 2012. Benchmarking the Performance of Linked Data Translation Systems. In *WWW2012 Workshop on Linked Data on the Web, Lyon, France, 16 April, 2012*. <http://ceur-ws.org/Vol-937/ldow2012-paper-09.pdf>
- Sherif Sakr. 2016. *Big data 2.0 processing systems: a survey*. Springer.
- Sherif Sakr and Eric Pardede (Eds.). 2011. *Graph Data Management: Techniques and Applications*. IGI Global. DOI: <http://dx.doi.org/10.4018/978-1-61350-053-8>

- Muhammad Saleem, Qaiser Mehmood, and Axel-Cyrille Ngonga Ngomo. 2015. *FEASIBLE: A Feature-Based SPARQL Benchmark Generation Framework*. Springer International Publishing, Cham, 52–69. DOI: [http://dx.doi.org/10.1007/978-3-319-25007-6\\_4](http://dx.doi.org/10.1007/978-3-319-25007-6_4)
- Albrecht Schmidt, Florian Waas, Martin Kersten, Michael J. Carey, Ioana Manolescu, and Ralph Busse. 2002. XMark: A Benchmark for XML Data Management. In *Proceedings of the 28th International Conference on Very Large Data Bases (VLDB '02)*. VLDB Endowment, 974–985. <http://dl.acm.org/citation.cfm?id=1287369.1287455>
- Michael Schmidt, Thomas Hornung, Michael Meier, Christoph Pinkel, and Georg Lausen. 2010. *SP2Bench: A SPARQL Performance Benchmark*. Springer Berlin Heidelberg, Berlin, Heidelberg, 371–393. DOI: [http://dx.doi.org/10.1007/978-3-642-04329-1\\_16](http://dx.doi.org/10.1007/978-3-642-04329-1_16)
- Martin Sevenich, Sungpack Hong, Oskar van Rest, Zhe Wu, Jayanta Banerjee, and Hassan Chafi. 2016. Using Domain-specific Languages for Analytic Graph Databases. *Proc. VLDB Endow.* 9, 13 (Sept. 2016), 1257–1268. DOI: <http://dx.doi.org/10.14778/3007263.3007265>
- Martin Simonovsky and Nikos Komodakis. 2018. GraphVAE: Towards Generation of Small Graphs Using Variational Autoencoders. *arXiv preprint arXiv:1802.03480* (2018).
- Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2007. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web*. ACM, 697–706.
- Riccardo Tommasini, Emanuele Della Valle, Andrea Mauri, and Marco Brambilla. 2017. RSPLab: RDF stream processing benchmarking made easy. In *International Semantic Web Conference*. Springer, 202–209.
- Chad Vicknair, Michael Macias, Zhendong Zhao, Xiaofei Nan, Yixin Chen, and Dawn Wilkins. 2010. 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)*. ACM, New York, NY, USA, Article 42, 6 pages. DOI: <http://dx.doi.org/10.1145/1900008.1900067>
- Bimal Viswanath, Alan Mislove, Meeyoung Cha, and Krishna P. Gummadi. 2009. On the Evolution of User Interaction in Facebook. In *Proceedings of the 2Nd ACM Workshop on Online Social Networks (WOSN '09)*. ACM, New York, NY, USA, 37–42. DOI: <http://dx.doi.org/10.1145/1592665.1592675>
- Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.
- Chong Jun Wang, Gang Wang, Yu Li Lei, and Shao Jie Qiao. 2013. Community Evolution in Dynamic Social Networks. In *Information Technology Applications in Industry, Computer Engineering and Materials Science (Advanced Materials Research)*, Vol. 756. Trans Tech Publications, 2634–2638. DOI: <http://dx.doi.org/10.4028/www.scientific.net/AMR.756-759.2634>
- Sui-Yu Wang, Yuanbo Guo, Abir Qasem, and Jeff Heflin. 2005. *Rapid Benchmarking for Semantic Web Knowledge Base Systems*. Springer Berlin Heidelberg, Berlin, Heidelberg, 758–772. DOI: [http://dx.doi.org/10.1007/11574620\\_54](http://dx.doi.org/10.1007/11574620_54)
- Yue Wang and Xintao Wu. 2013. Preserving differential privacy in degree-correlation based graph generation. *Transactions on data privacy* 6, 2 (2013), 127.
- Robert West, Hristo S Paskov, Jure Leskovec, and Christopher Potts. 2014. Exploiting social network structure for person-to-person sentiment analysis. *arXiv preprint arXiv:1409.2450* (2014).
- Huanhuan Wu, James Cheng, Silu Huang, Yiping Ke, Yi Lu, and Yanyan Xu. 2014a. Path Problems in Temporal Graphs. *Proc. VLDB Endow.* 7, 9 (May 2014), 721–732. DOI: <http://dx.doi.org/10.14778/2732939.2732945>
- Hongyan Wu, Toyofumi Fujiwara, Yasunori Yamamoto, Jerven Bolleman, and Atsuko Yamaguchi. 2014b. BioBenchmark Toyama 2012: an evaluation of the performance of triple stores on biological data. *Journal of Biomedical Semantics* 5, 1 (2014), 32. DOI: <http://dx.doi.org/10.1186/2041-1480-5-32>
- Xintao Wu, Xiaowei Ying, Kun Liu, and Lei Chen. 2010. A survey of privacy-preservation of graphs and social networks. In *Managing and mining graph data*. Springer, 421–453.
- Jaewon Yang and Jure Leskovec. 2015. Defining and evaluating network communities based on ground-truth. *Knowledge and Information Systems* 42, 1 (2015), 181–213.
- Yuan Yao, Jiufeng Zhou, Lixin Han, Feng Xu, and Jian Lü. 2011. *Comparing Linkage Graph and Activity Graph of Online Social Networks*. Springer Berlin Heidelberg, Berlin, Heidelberg, 84–97. DOI: [http://dx.doi.org/10.1007/978-3-642-24704-0\\_14](http://dx.doi.org/10.1007/978-3-642-24704-0_14)
- Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. 2018. GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models. In *International Conference on Machine Learning*. 5694–5703.
- Wayne W Zachary. 1977. An information flow model for conflict and fission in small groups. *Journal of anthropological research* 33, 4 (1977), 452–473.

- Yunpeng Zhao. 2017. A survey on theoretical advances of community detection in networks. *Wiley Interdisciplinary Reviews: Computational Statistics* 9, 5 (2017). DOI: <http://dx.doi.org/10.1002/wics.1403>
- Yunhong Zhou, Dennis Wilkinson, Robert Schreiber, and Rong Pan. 2008. Large-scale parallel collaborative filtering for the netflix prize. In *International Conference on Algorithmic Applications in Management*. Springer, 337–348.
- Cai-Nicolas Ziegler, Sean M McNee, Joseph A Konstan, and Georg Lausen. 2005. Improving recommendation lists through topic diversification. In *Proceedings of the 14th international conference on World Wide Web*. ACM, 22–32.

Received October 2018; revised October 2018; accepted October 2018