

NetworKit: A tool suite for large-scale complex network analysis

CHRISTIAN L. STAUDT

*Institute of Theoretical Informatics, Karlsruhe Institute of Technology (KIT), 76131 Karlsruhe, Germany
(e-mail: christian.staudt@kit.edu)*

ALEKSEJS SAZONOV¹

*Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK
(e-mail: as45@sanger.ac.uk)*

HENNING MEYERHENKE

*Institute of Theoretical Informatics, Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany
(e-mail: meyerhenke@kit.edu)*

Abstract

We introduce NetworKit, an open-source software package for analyzing the structure of large complex networks. Appropriate algorithmic solutions are required to handle increasingly common large graph data sets containing up to billions of connections. We describe the methodology applied to develop scalable solutions to network analysis problems, including techniques like parallelization, heuristics for computationally expensive problems, efficient data structures, and modular software architecture. Our goal for the software is to package results of our algorithm engineering efforts and put them into the hands of domain experts. NetworKit is implemented as a hybrid combining the kernels written in C++ with a Python frontend, enabling integration into the Python ecosystem of tested tools for data analysis and scientific computing. The package provides a wide range of functionality (including common and novel analytics algorithms and graph generators) and does so via a convenient interface. In an experimental comparison with related software, NetworKit shows the best performance on a range of typical analysis tasks.

Keywords: *complex networks, network analysis, network science, parallel graph algorithms, data analysis software*

1 Motivation

Network science methodology is increasingly applied to study a variety of real-world phenomena (Costa et al., 2011; Boccaletti et al., 2006). Consequently, large network data sets comprising millions of edges are more and more common, and it is an active current research project to develop scalable methods for the analysis of large networks. In order to process such massive graphs, we need algorithms whose running time is essentially linear in the number of edges. Many analysis methods have been pioneered on small networks (e. g. for the study of social networks prior

¹ Parts of the work were performed while Aleksejs Sazonovs was with KIT as part of the RISE program of the German Academic Exchange Service DAAD.

to the arrival of massive online social networking services), so that underlying algorithms with higher complexity were viable. As we shall see in the following, developing a scalable analysis tool suite often entails replacing them with suitable linear- or nearly linear-time variants. Furthermore, solutions should employ parallel processing: While sequential performance is stalling, multicore machines become pervasive, and algorithms and software need to follow this development. Within the NetworKit project, scalable network analysis methods are developed, tested and packaged as ready-to-use software. In this process, we frequently apply the following algorithm and software engineering patterns: *parallelization*; *heuristics*; or *approximation algorithms* for computationally intensive problems; efficient *data structures*; and *modular* software architecture. With NetworKit, we intend to push the boundaries on the size of networks whose structure can be characterized on a shared-memory parallel computer.

In this work, we give an introduction to the tool suite and describe the methodology applied during development in terms of algorithm and software engineering aspects. We discuss methods to arrive at highly scalable solutions to common network analysis problems (Sections 2 and 3), describe the set of functionality (Sections 4 and 5), present example use cases (Section 6), compare with related software (Section 7), and evaluate the performance of analysis kernels experimentally (Section 8). Our experiments show that NetworKit is capable of quickly processing large-scale networks for a variety of analytics kernels, and does so faster and with a lower memory footprint than closely related software. We recommend NetworKit for the comprehensive structural analysis of massive complex networks (their size is primarily limited by the available memory). To this end, a new front end supports exploratory data analysis with fast graphical reports on structural features of the network (Section 6.2).

2 Methodology

2.1 Design goals

There is a variety of software packages which provide graph algorithms in general and network analysis capabilities in particular (see Section 7 for a comparison to related packages). However, NetworKit aims to balance a specific combination of strengths:

Performance. Algorithms and data structures are selected and implemented with high performance and parallelism in mind. Some implementations are among the fastest in published research. For example, community detection in a 3.3 billion edge web graph can be performed on a 16-core server with hyperthreading in less than 3 minutes (Staudt and Meyerhenke, 2016).

Usability and Integration. Networks are as diverse as the series of questions we might ask of them—e. g., what is the largest connected component, what are the most central nodes in it and how do they connect to each other? A practical tool for network analysis should therefore provide modular functions which do not restrict the user to predefined workflows. An interactive shell, which the Python language provides, is one prerequisite for that. While NetworKit works with the standard Python 3 interpreter, calling the module from the IPython shell and Jupyter Notebook

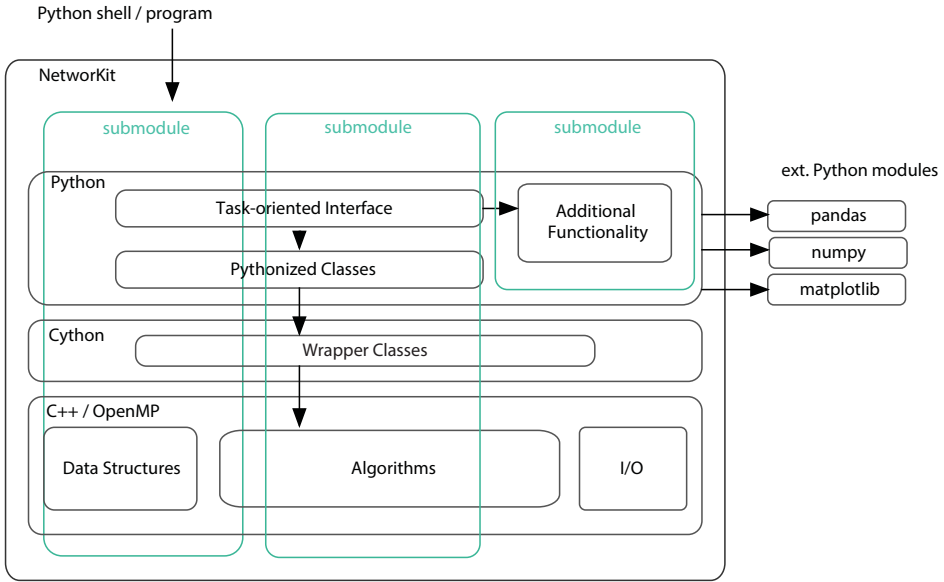


Fig. 1. NetworkKit architecture overview (→ represents call from/to). (Color online)

HTML interface (Perez et al., 2013) allows us to integrate it into a fully fledged computing environment for scientific workflows, from data preparation to creating figures. It is also easy to set up and control a remote compute server. As a Python module, NetworkKit can be seamlessly integrated with Python libraries for scientific computing and data analysis, e.g. *pandas* for data frame processing and analytics, *matplotlib* for plotting or *numpy* and *scipy* for numerical and scientific computing. For certain tasks, we provide interfaces to external tools, e.g. *Gephi* (Bastian et al., 2009) for graph visualization.

2.2 Architecture

In order to achieve the design goals described above, we implement NetworkKit as a two-layer hybrid of performance-aware code written in C++ with an interface and additional functionality written in Python. NetworkKit is distributed as a Python package, ready to be used interactively from a Python shell, which is the main usage scenario we envision for domain scientists. The code can be used as a library for application programming as well, either at the Python or C++ level. Throughout the project, we use object-oriented and functional concepts. Shared-memory parallelism is realized with OpenMP, providing loop parallelization and synchronization constructs while abstracting away the details of thread creation and handling. As illustrated in Figure 1, connecting these native implementations to the Python world is enabled by the Cython toolchain (Behnel et al., 2011), which is used to integrate native code by compiling it into a Python extension module. The resulting Python module *networkkit* is organized into several submodules for different areas of functionality, such as community detection or node centrality.

2.3 Framework foundations

As the central data structure, the `Graph` class implements a directed or undirected, optionally weighted graph using an adjacency array data structure with $O(n + m)$ memory requirement for a graph with n nodes and m edges. Nodes are represented by 64-bit integer indices from a consecutive range, and an edge is identified by a pair of nodes. Edges can be indexed as well, associating an integer index from a consecutive range with each edge, although this is optional, requiring $O(m)$ additional memory. We conclude that edge indexing does not impact analysis algorithm performance in general, on the basis of benchmarks showing no significant running time differences between enabled and disabled edge indices. The design described above enables a lean graph data structure, while also allowing arbitrary node and edge attributes to be stored in any container addressable by indices. In particular, it supports dynamic modifications to the graph in a flexible manner, unlike the *compressed sparse row* format common in high-performance scientific computing. Our graph API facilitates the concise formulation of graph algorithms on both the C++ and Python layer (see Figure 3 for an example).

3 Algorithm and implementation patterns

As explained in Section 1, our main focus are scalable algorithms in order to support network analysis on massive networks. We identify several algorithm and implementation patterns that help to achieve this goal and present them below by means of case studies. For experimental results, we express processing speed in “edges per second,” aggregating real running time over a set of graphs and normalizing by graph size.

3.1 Parallelism

Our first case study concerns the *core decomposition* of a graph. The sequential kernel implemented in NetworKit runs in $O(m)$ time, matching other implementations (Batagelj and Zaveršnik, 2011). The main algorithmic idea we reuse for computing the core numbers is to start with $k = 0$ and increase k iteratively. Within each iteration phase, all nodes with degree k are successively removed (thus, also nodes whose degree was larger at the beginning of the phase can become affected by a removal of a neighbor). Using a bucket priority queue, we can extract and update nodes accordingly in amortized constant time, resulting in $O(m)$ in total. While this already scales to large inputs, we can achieve further speedup through parallelization. The sequential algorithm cannot be made parallel easily due to its sequential access to the bucket priority queue. For achieving a higher degree of parallelism, we follow Dasari et al. (2014). Their ParK algorithm replaces the extract-min operation in the above algorithm by identifying the node set V' with nodes of minimum residual degree while iterating in parallel over all (active) nodes. V' is then further processed similarly to the node retrieved by extract-min in the above algorithm, only in parallel again. ParK thus performs more sequential work, but with thread-local buffers it relies on a minimal amount of synchronization. Moreover, its data access pattern is more cache-friendly, which additionally contributes to better performance.

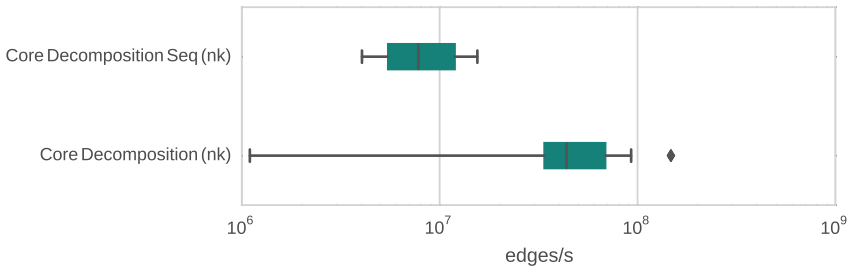


Fig. 2. Core decomposition: sequential versus parallel performance. (Color online)

Figure 2 is the result of running time measurements on a test set of networks (see Section 8 for the setup). On average, processing speed is increased by almost an order of magnitude through parallelization. Some overhead of the parallel algorithm implies that speedup is only noticeable on large graphs, hence the large variance. For example, processing time for the 260 million edge uk-2002 web graph is reduced from 22 to 2 seconds.

3.2 Heuristics and approximation algorithms

In this example, we illustrate how inexact methods deliver appropriate solutions for an otherwise computationally impractical problem. *Betweenness centrality* is a well-known node centrality measure that has an intuitive interpretation in transport networks: Assuming that the transport processes taking place in the network are efficient, they follow shortest paths through the network, and therefore preferably pass through nodes with high betweenness. For instance, their removal would interfere strongly with the function of the network. It is clear that network analysts would like to be able to identify such nodes in networks of any size. NetworKit comes with an implementation of the currently fastest known algorithm for betweenness (Brandes, 2001), which has $O(nm)$ running time in unweighted graphs.

In order to parallelize the algorithm, several single-source shortest path searches can be run in parallel to compute the intermediate *dependency* values whose sum yields a node's betweenness. Figure 3 shows C++ code for the parallel version, which is simplified to focus on the core algorithm. To avoid race conditions, each thread works on its own dependency array, which need to be aggregated into one betweenness array in the end (lines 35–39).

We now evaluate the performance of the implementations experimentally (see Section 8 for settings). Figure 4 shows aggregated running speed over a set of smaller networks (from Table 3). In practice, this means that the sequential version of Brandes' algorithm (*BetweennessSeq*) takes almost 8 hours to process the 600k edge graph *caidaRouterLevel* (representing internet router-level topology (CAIDA, 2003)). Parallelism with 32 (hyper)threads (*Betweenness*) reduces the running time to ca. 90 minutes. Still, parallelization does not change the algorithm's inherent complexity. This means that running times rise so steeply with the size of the input graph that computing an exact solution to betweenness is not viable on the large networks we want to target. In typical use cases, obtaining exact values for betweenness is not necessary, though. An approximate result is likely good enough to appreciate the structure of the network for exploratory analysis, and to

```

1 // thread-local scores for efficient parallelism
2 count maxThreads = omp_get_max_threads();
3 count z = G.upperNodeIdBound();
4 std::vector<std::vector<double>> scorePerThread(maxThreads, std::vector<double>(z));
5
6 auto computeDependencies = [&](node s) {
7     std::vector<double> dependency(z, 0.0);
8     // run SSSP algorithm and count paths
9     std::unique_ptr<SSSP> sssp;
10    if (G.isWeighted()) {
11        sssp.reset(new Dijkstra(G, s, true, true));
12    } else {
13        sssp.reset(new BFS(G, s, true, true));
14    }
15    sssp->run();
16    // compute dependencies for nodes in order of decreasing distance from s
17    std::vector<node> stack = sssp->getStack();
18    while (!stack.empty()) {
19        node t = stack.back();
20        stack.pop_back();
21        for (node p : sssp->getPredecessors(t)) {
22            double w = sssp->numberOfPaths(p)/sssp->numberOfPaths(t);
23            double c = w * (1 + dependency[t]);
24            dependency[p] += c;
25        }
26        if (t != s) {
27            scorePerThread[omp_get_thread_num()][t] += dependency[t];
28        }
29    }
30 };
31
32 // iterate over nodes in parallel and apply
33 G.balancedParallelForNodes(computeDependencies);
34
35 // add up all thread-local values
36 for (auto& localScore : scorePerThread) {
37     G.parallelForNodes([&](node v){
38         scoreData[v] += localScore[v];
39     });
40 }

```

Fig. 3. Code example: Parallel calculation of betweenness centrality. (Color online)

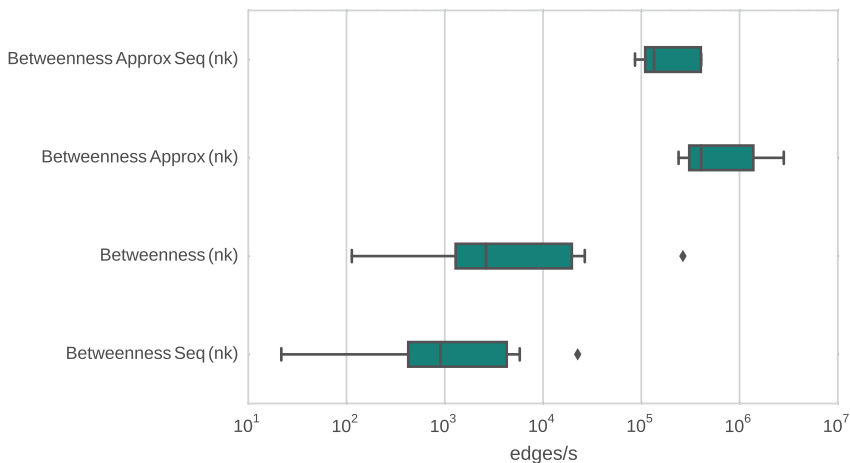


Fig. 4. Processing speed of exact and inexact algorithms for betweenness centrality. (Color online)

identify a set of top betweenness nodes. Therefore, we use a heuristic approach based on computing a relatively small number of randomly chosen shortest path trees (Geisberger et al., 2008). In contrast to the exact algorithm, running the approximative algorithm with 42 samples takes 6 seconds sequentially. Applying this algorithm cuts running time by orders of magnitude, but still yields a ranking of nodes that is highly similar to a ranking by exact betweenness values. We observe that the distribution of relative rank errors (exact rank divided by approximated rank) has little variance around 1.0. Nodes on average maintain the rank they would have according to exact betweenness even with such a small number of samples. Experiments of this type (see Geisberger et al. 2008) confirm that in typical cases betweenness can be closely approximated with a relatively small number s of shortest-path searches. Therefore, we can replace an $O(nm)$ algorithm with one of time complexity $O(sm)$ in many use cases. The inexact algorithm offers the same opportunities for parallelization, yielding additional speedups: In the example above, parallel running time is down to 1.5 seconds on 32 (hyper)threads.

If a true approximation with a guaranteed error bound is desired, NetworkKit users can apply another inexact algorithm (Riondato and Kornaropoulos, 2016) which accepts an error bound parameter ϵ . It sacrifices some computational efficiency but allows a probabilistic guarantee that the resulting betweenness scores have at most $\pm\epsilon$ difference from the exact scores.

3.3 Efficient data structures

The case study on data structures deals with a *generative network model*, which simplify complex network research in several respects (see Section 5). Theoretical analyses have shown that *Random hyperbolic graphs* (RHGs) (Krioukov et al., 2010) have many features also found in real complex networks (Bode et al., 2015; Gugelmann et al., 2012; Kiwi and Mitsche, 2015). During the geometric generation process, nodes are distributed randomly on a hyperbolic disk of radius R and edges are inserted for every node pair whose distance is below R . The straightforward RHG generation process would probe the distance of all pairs, yielding a quadratic time complexity. This impedes the creation of massive networks. NetworkKit provides the first generation algorithm for RHGs with subquadratic running time ($O((n^{3/2} + m) \log n)$ with high probability) (von Looz et al., 2015). The acceleration stems primarily from the reduction of distance computations through a polar quadtree adapted to hyperbolic space. Instead of probing each pair of nodes, the generator performs for each node one efficient range query supported by the quadtree. In practice, this leads to an acceleration of at least two orders of magnitude. With the quadtree-based approach, networks with billions of edges can be generated in parallel in a few minutes (von Looz et al., 2015).

3.4 Modular design

In terms of software design, we aim at a modular architecture with encapsulation of algorithms into software components (classes and modules). Among the benefits are extensibility and code reuse: For example, new centrality measures can be easily added by implementing a subclass with the code specific to the centrality

Table 1. Selection of analysis algorithms contained in NetworKit. Complexity expressed in terms of n nodes, m edges, s samples, and maximum node degree d .

Category	Task	Algorithm	Time	Space
Centrality	Degree	–	$O(n)$	$O(n)$
	Betweenness	(Brandes, 2001)	$O(nm)$	$O(n + m)$
	ap. betweenness	(Geisberger et al., 2008),(Riondato and Kornaropoulos, 2016)	$O(sm)$	$O(n + m)$
	Closeness	Shortest-path search from each node	$O(mn)$	$O(n)$
	ap. closeness	(Eppstein and Wang, 2004)	$O(sm)$	$O(n)$
	PageRank	Power iteration	$O(m)$ Typical (Section 4.2)	$O(n)$
	Eigenvector centrality	Power iteration	$O(m)$ Typical	$O(n)$
	Katz centrality	(Katz, 1953)	$O(m)$ Typical	$O(n)$
	k -path centrality	(Alahakoon et al., 2011)	See (Alahakoon et al., 2011)	$O(n)$
	Local clustering coefficient	Parallel iterator	$O(nd^2)$	$O(n)$
Partitions	k -core decomposition	(Dasari et al., 2014)	$O(m)$	
	Connected components	BFS	$O(m)$	$O(n)$
	Community detection	PLM, PLP (Staudt and Meyerhenke, 2016)	$O(m)$	$O(m), O(n)$
Global	Diameter	iFub (Crescenzi et al., 2013)	$O(m)$ Typical	$O(n)$

computation, while code applicable to all centrality measures and a common interface remains in the base class. Through these and other modularizations, developers can add a new centrality measure and get derived measures almost “for free.” These include for instance the *centralization index* (Freeman, 1979) and the *assortativity coefficient* (Freeman, 1979), which can be defined with respect to any node centrality measure and may in each case be a key feature of the network. Modular design also allows for optimizations on one algorithm to benefit other client algorithms. For instance, betweenness and other centrality measures (such as closeness) require the computation of shortest paths, which is done via breadth-first search in unweighted graphs and Dijkstra’s algorithm in weighted graphs, decoupled to avoid code redundancy (see lines 10–14 in Figure 3).

4 Analytics

The following describes the core set of network analysis algorithms implemented in NetworKit as of version 4.0. Table 1 summarizes the core set of algorithms for typical problems. Where applicable, implementations process both weighted and unweighted graphs appropriately.

4.1 Global network properties

Global properties include simple statistics such as the number of nodes and edges and the graph's density, as well as properties related to distances: The *diameter* of a graph is the maximum length of a shortest path between any two nodes. We use the iFUB algorithm (Crescenzi et al., 2013) both for the exact computation as well as an estimation of a lower and upper bound on the diameter. iFub has a worst case complexity of $O(nm)$ but has shown excellent typical-case performance on complex networks, where it often converges on the exact value in linear time.

4.2 Node centrality

Node centrality measures quantify the structural importance of a node within a network. More precisely, we consider a node centrality measure as any function which assigns to each node an attribute value of (at least) ordinal scale of measurement. The assigned value depends on the position of a node within the network as defined by a set of edges.

The simplest measure that falls under this definition is the *degree*, i.e. the number of connections of a node, to which the graph data structure provides constant-time access. *Eigenvector centrality* and its variant *PageRank* (Brin and Page, 2012) are implemented in NetworKit based on parallel power iteration, whose convergence time depends on a numerical error tolerance parameter and spectral properties of the network, but is among the fast linear-time algorithms for typical inputs. For *betweenness centrality*, we provide the solutions discussed in Section 3.2. Similar techniques are applied for computing *closeness centrality* exactly and approximately (Eppstein and Wang, 2004). Our current research extends the former approach to dynamic graph processing (Bergamini and Meyerhenke, 2015; Bergamini et al., 2015). The *local clustering coefficient* expresses how many of the possible connections between neighbors of a node exist, which can be treated as a node centrality measure according to the definition above (Newman, 2010). In addition to a parallel algorithm for clustering coefficients, NetworKit also implements a sampling approximation algorithm (Schank and Wagner, 2005), whose constant time complexity is independent of graph size.

4.3 Edge centrality, sparsification and link prediction

The concept of centrality can be extended to edges: Not all edges are equally important for the structure of the network, and scores can be assigned to edges depending on the graph structure such that they can be ranked (e.g. *edge betweenness*, which depends on the number of shortest paths passing through an edge).

A ranking of this kind can also be used to filter edges and thereby reduce the size of data. NetworKit includes a wide set of edge ranking methods, with a focus on sparsification techniques meant to preserve certain properties of the network. For instance, we show that a method that ranks edges leading to high-degree nodes (hubs) closely preserves many properties of social networks, including diameter, degree distribution, and centrality measures. Other methods, including a family of *Simmelian backbones*, assign higher importance to edges within dense

regions of the graph and hence preserve or emphasize communities. Details are reported in our recent experimental study (Hamann et al., 2016). While currently experimental and focused on one application, namely structure-preserving sparsification, the design is extensible so that general edge centrality indices can be easily implemented.

A somewhat related problem, conceptually and through common methods, is the task of *link prediction*. Link prediction algorithms examine the edge structure of a graph to derive similarity scores for unconnected pairs of nodes. Depending on the score, the existence of a future or missing edge is inferred. NetworKit includes implementations for a wide variety of methods (Esders, 2015).

4.4 Partitioning the network

Another class of analysis methods partitions the set of nodes into subsets depending on the graph structure. For instance, all nodes in a *connected component* are reachable from each other. A network's connected components can be computed in linear time using breadth-first search. *Community detection* is the task of identifying groups of nodes in the network which are significantly more densely connected among each other than to the rest of nodes. The task can be turned into a well-defined though NP-hard optimization problem by using community quality measures, first and foremost *modularity* (Girvan and Newman, 2002). We approach community detection from the perspective of modularity maximization and engineer parallel heuristics which deliver a good trade-off between solution quality and running time (Staudt and Meyerhenke, 2016). The PLP (Parallel Label Propagation) algorithm implements community detection by label propagation (Raghavan et al., 2007), which extracts communities from a labeling of the node set. The *Louvain method* for community detection (Blondel et al., 2008) can be classified as a locally greedy, bottom-up multilevel algorithm. We recommend the PLM (Parallel Louvain Method) algorithm with optional refinement step as the default choice for modularity-driven community detection in large networks. PLM can be parametrized to use the multi-resolution modularity objective function (Lambiotte, 2010) as a remedy for modularity's resolution limit, i.e. the behavior of detected communities growing with the input network size. For very large networks in the range of billions of edges, PLP delivers a better time to solution, albeit with a qualitatively different solution and worse modularity.

5 Network generators

Generative network models aim to explain how networks form and evolve specific structural features. Such models and their implementations as generators have at least two important uses: On the one hand, algorithm or software engineers want generators for synthetic datasets which can be arbitrarily scaled and parametrized and produce graphs which resemble the real application data. On the other hand, network scientists employ models to increase their understanding of network phenomena. NetworKit provides a versatile collection of graph generators for this purpose, summarized in Table 2.

Table 2. Overview of network generators.

Model [and algorithm]	Description
<i>Erdős-Rényi</i> (P. Erdős, 1960) [(Batagelj and Brandes, 2005)]	Random edges with uniform probability
<i>Planted partition/stochastic blockmodel</i>	Dense areas with sparse connections
<i>Barabasi-Albert</i> (Albert and Barabási, 2002)	Preferential attachment process resulting in power-law degree distribution
<i>Recursive Matrix (R-MAT)</i> (Chakrabarti et al., 2004)	Power-law degree distribution, small-world property, self-similarity
<i>Chung-Lu</i> (Aiello et al., 2000)	Replicate a given degree distribution
<i>Havel-Hakimi</i> (Hakimi, 1962)	Replicate a given degree distribution
<i>hyperbolic unit-disk model</i> (Krioukov et al., 2010) [(von Looz et al., 2015)]	Large networks, power-law degree distribution and High clustering
<i>LFR</i> (Lancichinetti and Fortunato, 2009)	Complex networks containing communities

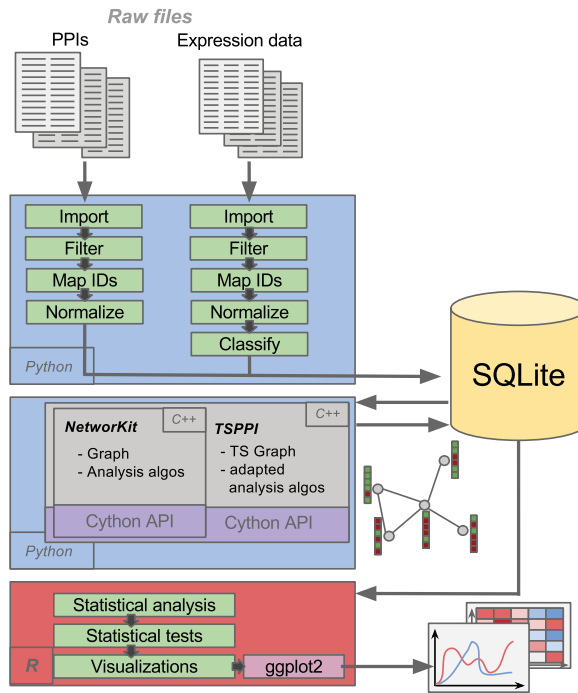


Fig. 5. PPI network analysis pipeline with NetworkKit as central component. (Color online)

6 Example use cases

In the following, we present possible workflows and use cases, highlighting the capabilities of NetworkKit as a data analysis tool and a library.

6.1 As a library in an analysis pipeline

A master’s thesis (Flick, 2014) provides an early example of NetworkKit as a component in an application-specific data mining pipeline (Figure 5). This pipeline performs

analysis of *protein-interaction (PPI) networks* and implements a preprocessing stage in Python, in which networks are compiled from heterogeneous data sets containing interaction data as well as *expression data* about the occurrence of proteins in different cell types. During the network analysis stage, preprocessed networks are retrieved from a database, and NetworkKit is called via the Python frontend. The C++ core has been extended to enable more efficient analysis of *tissue-specific* PPI networks, by implementing in-place filtering of the network to the subgraphs of proteins that occur in given cell types. Finally, statistical analysis and visualization is applied to the network analysis data. The system is close to how we envision NetworkKit as a high-performance algorithmic component in a real-world data analysis scenario, and we therefore place emphasis on the toolkit being easily scriptable and extensible.

6.2 Exploratory network analysis with network profiles

Utilizing NetworkKit as a library requires writing some custom code and some expertise in selecting algorithms and their parameters. This is one reason why we also provide an interface that requires just a few lines of code for an exploratory analysis of large network, and returns an extensive overview. We treat networks as statistical data sets whose properties should be determined via graph algorithms and the results summarized via statistical graphics. The underlying module assembles many algorithms into one analysis pipeline, automates analysis tasks and produces a graphical report to be displayed in the Jupyter Notebook or exported to an HTML or \LaTeX report document. Such a *network profile* gives a statistical overview over the properties of the network. It consists of the following parts: First global properties such as size and density are reported. The report then focuses on a variety of node centrality measures, showing an overview of their distributions in the network (see Figure 6). Detailed views for centrality measures (see Figure 7) follow: Their distributions are plotted in histograms and characterized with standard statistics, and network-specific measures such as centralization and assortativity are shown. Furthermore, correlations between centralities are treated as noteworthy empirical features of the network. For instance, betweenness may or may not be positively correlated with increasing node degree. The prevalence of low-degree, high-betweenness nodes may influence the resilience of a transport network, as only few links then need to be severed in order to significantly disrupt transport processes following shortest paths. The report displays correlations in the form of a matrix of Spearman's correlation coefficients, showing how node ranks derived from the centrality measures correlate with each other (see Figure 8(b)). Furthermore, scatter plots for each combination of centrality measure are shown, suggesting the type of correlation (see Figure 9(a)). The report continues with different ways of partitioning the network, showing histograms and pie charts for the size distributions of connected components, modularity-based communities (see Figure 9(b)) and k -shells, respectively. The default configuration of the module is such that even networks with hundreds of millions of edges can be characterized in minutes on a parallel workstation. Furthermore, it can be configured by the user depending on the desired choice of analytics and level of detail, so that custom reports can be generated.

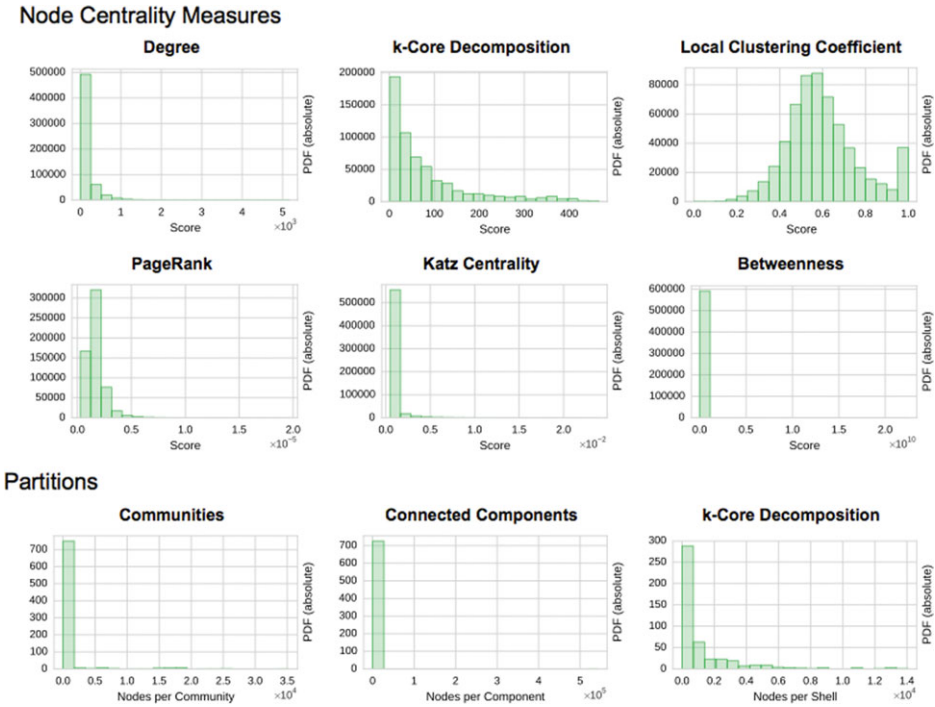


Fig. 6. Overview on the distributions of node centrality measures and size distributions of different network partitions—here, a Facebook social network. (Color online)

To pick an example from a scientific domain, the human connectome network *con-fiber_big* maps brain regions and their anatomical connections at a relatively high resolution, yielding a graph with ca. 46 million edges. As the resolution of brain imaging technology improves, connectome analysis is likely to yield ever more massive network data sets, considering that the brain at the neuronal scale is a complex network on the order of 10^{10} nodes and 10^{14} edges. On a first look, the network has properties similar to a social network, with a skewed degree distribution and high clustering. The pattern of correlations (Figure 8(b)) differs from that of a typical friendship network (Figure 8(a)), with weaker positive correlations across the spectrum of centrality measures. As one observation to focus on, we may pick the strong negative correlation between the local clustering coefficient on the one hand and the PageRank and betweenness centrality on the other. High betweenness nodes are located on many shortest paths, and high PageRank results from being connected to neighbors which are themselves highly connected. Thus, the correlations point to the presence of structural hub nodes that connect different brain regions which are not directly linked. Also, a look at a scatter plot generated (Figure 9(a)) reveals more details on the correlations: The local clustering coefficient steadily falls with node degree, a majority of nodes having high clustering and low degree, a few nodes having low clustering and high degree. Both observations are consistent with the finding of *connector hub* regions situated along the midline of the brain, which are highly connected and link otherwise separated brain modules organized around smaller *provincial hubs* (Sporns and Betzel, 2016).

Local Clustering Coefficient

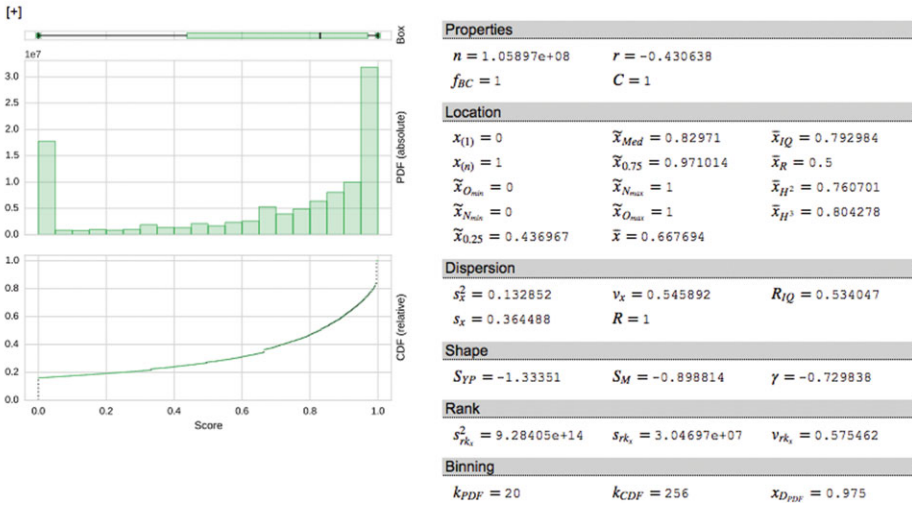


Fig. 7. Detailed view on the distribution of node centrality scores—here, local clustering coefficients of the 3 billion edge web graph uk-2007. (Color online)

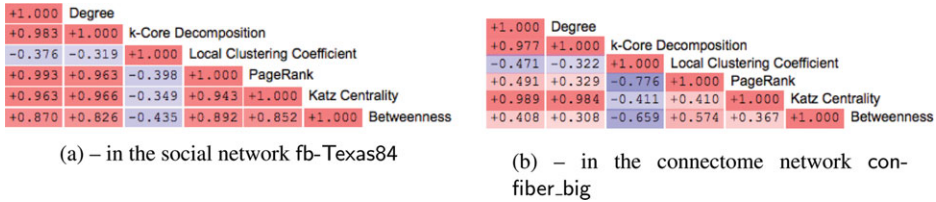


Fig. 8. Correlation between node centrality measures. (a) in the social network fb-Texas84. (b) in the connectome network con-fiber_big. (Color online)

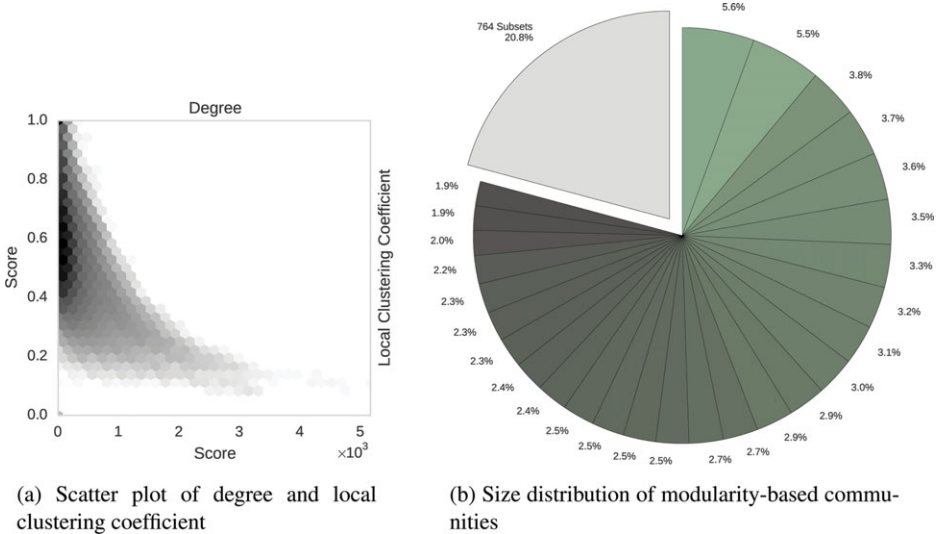


Fig. 9. Statistical graphics from the profile of the connectome graph con-fiber_big. (a) Scatter plot of degree and local clustering coefficient. (b) Size distribution of modularity-based communities. (Color online)

Another aspect we can focus on is community structure. There has been extensive research on the modular structure of brain networks, indicating that communities in the connectivity network can be interpreted as functional modules of the brain (Sporns and Betzel, 2016). The communities found by the PLM modularity-maximizing heuristic in the `con-fiber.big` graph can be interpreted accordingly. Their size distribution (Figure 9(b), in which a green pie slice represents the size of a community) shows that a large part of the network consists of about 30 communities of roughly equal size, in addition to a large number of very small communities (grey). While a thorough analysis of these findings would certainly need to incorporate domain-specific knowledge, these examples illustrate how `NetworKit`'s capability to quickly generate an overview of structural properties can be used to generate hypotheses about the network data.

7 Comparison to related software

Recent years have seen a proliferation of graph processing and network analysis software which vary widely in terms of target platform, user interface, scalability and feature set. We therefore locate `NetworKit` relative to these efforts. Although the boundaries are not sharp, we would like to separate network analysis toolkits from general purpose graph frameworks (e.g. Boost Graph Library and JUNG (O'Madadhain et al., 2003)), which are less focused on data analysis workflows.

As closest in terms of architecture, functionality and target use cases, we see `igraph` (Csardi and Nepusz, 2006) and `graph-tool` (Peixoto, 2006). They are packaged as Python modules, provide a broad feature set for network analysis workflows, and have active user communities. Like `NetworKit`, `igraph`, and `graph-tool` address the scalability issue by implementing core data structures and algorithms in C or C++. `graph-tool` builds on the Boost Graph Library and parallelizes some kernels using OpenMP. These similarities make those packages ideal candidates for an experimental comparison with `NetworKit` (see Section 8.2). `NetworkX` (Hagberg et al., 2008) is also a mature toolkit and the de-facto standard for the analysis of small to medium networks in a Python environment, but not suitable for massive networks due to its pure Python implementations, and often orders of magnitude slower than `NetworKit`. The aforementioned packages have a substantial core feature set in common with `NetworKit` and a broader feature set or different feature choices in some areas. Python as a lingua franca makes it easy to build workflows combining features from multiple packages.

Other projects are geared toward network science but differ in important aspects from `NetworKit`. `Gephi` (Bastian et al., 2009), a GUI application for the Java platform, has a strong focus on visual network exploration. `Pajek` (Batagelj and Mrvar, 2004), a proprietary GUI application for the Windows operating system, also offers analysis capabilities similar to `NetworKit`, as well as visualization features. The variant `PajekXXL` uses less memory and thus focuses on large datasets.

The `SNAP` (Leskovec and Sosič, 2014) network analysis package has also recently adopted the hybrid approach of C++ core and Python interface. Related efforts from the algorithm engineering community are `KDT` (Lugowski et al., 2012) (built on an algebraic, distributed parallel backend), `GraphCT` (Ediger et al., 2013) (focused on massive multithreading architectures such as the Cray XMT), `STINGER` (a dynamic

Table 3. Networks used in this paper.

Name	Type	n	m	Source
fb-Caltech36	Social (friendship)	769	16,656	(Traud et al., 2012)
PGPgiantcompo	Social (trust)	10,680	24,316	(Boguña et al. 2014)
CoAuthorsDBLP	Coauthorship (science)	299,067	977,676	(Bader et al., 2014)
fb-Texas84	Social (friendship)	36,371	1,590,655	(Traud et al., 2012)
Foursquare	Social (friendship)	639,014	3,214,986	(Zafarani and Liu, 2009)
Lastfm	Social (friendship)	1,193,699	4,519,020	(Zafarani and Liu, 2009)
wiki-Talk	Social	2,394,385	4,659,565	(Leskovec and Krevl, 2014)
Flickr	Social (friendship)	639,014	55,899,882	(Zafarani and Liu, 2009)
in-2004	Web	1,382,908	13,591,473	(Boldi et al., 2004)
Actor-collaboration	Collaboration (film)	382,219	15,038,083	(Kunegis, 2013)
eu-2005	Web	862,664	16,138,468	(Boldi et al., 2004)
Flickr-growth-u	Social (friendship)	2,302,925	33,140,017	(Kunegis, 2013)
Con-fiber.big	Brain (connectivity)	591,428	46,374,120	http://openconnecto.me
Twitter	Social (followership)	15,395,404	85,331,845	(Zafarani and Liu, 2009)
uk-2002	Web	18,520,486	261,787,258	(Boldi et al., 2004)
uk-2007-05	Web	105,896,555	3,301,876,564	(Boldi et al., 2004)

graph data structure with some analysis capabilities) (Ediger et al., 2012), and Ligra (Shun and Blelloch, 2013) (a recent shared-memory parallel library). They offer high performance through native, parallel implementations of certain kernels. However, to characterize a complex network in practice, we need a substantial set of analytics which those frameworks currently do not provide.

Among solutions for large-scale graph analytics, distributed computing frameworks (for instance, GraphLab (Low et al., 2012)) are often prominently named. However, graphs arising in many data analysis scenarios are not bigger than the billions of edges that fit into a conventional main memory and can therefore be processed far more efficiently in a shared-memory parallel model (Shun and Blelloch, 2013), which we confirm experimentally in a recent study (Koch et al., 2015).

8 Performance evaluation

This section presents an experimental evaluation of the performance of NetworKit's algorithms. Our platform is a shared-memory server with 256 GB RAM and 2×8 Intel(R) Xeon(R) E5-2680 cores (32 threads due to hyperthreading) at 2.7 GHz, using the GCC 4.8 compiler and the openSUSE 13.1 OS. We use NetworKit 4.0, igraph 0.7.0, and graph-tool 2.9.

8.1 Benchmark

Figure 10 shows results of a benchmark of the most important analytics kernels in NetworKit. The algorithms were applied to a diverse set of 15 real-world networks in the size range from 16 k to 260 M edges, including web graphs, social networks, connectome data, and internet topology networks (see Table 3 for a description). Kernels with quadratic running time (like Betweenness) were restricted to the subset of the four smallest networks. The box plots illustrate the range of processing rates achieved (dots are outliers). The benchmark illustrates that a set of efficient linear-time kernels, including ConnectedComponents, the community detectors,

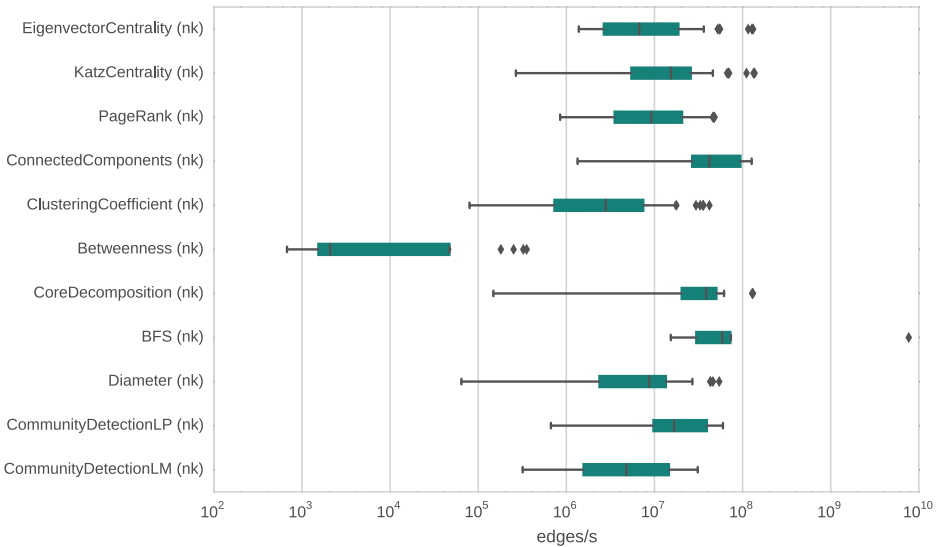


Fig. 10. Processing rates of NetworkKit analytics kernels. (Color online)

PageRank, CoreDecomposition, and ClusteringCoefficient, scales well to networks in the order of 10^8 edges. The iFub (Crescenzi et al., 2013) algorithm demonstrates good performance on complex networks, moving diameter calculation effectively into the class of linear-time kernels, illustrating that performance is often strongly dependent on the specific structure of complex networks. Algorithms like BFS and ConnectedComponents actually scan every edge at a rate of 10^7 to 10^8 edges per second.

8.2 Comparative benchmark

NetworkKit, igraph, and graph-tool rely on the same hybrid architecture of C/C++ implementations with a Python interface. igraph uses non-parallel C code while graph-tool also features parallelism. We benchmarked typical analysis kernels for the three packages in comparison on the aforementioned parallel platform and present the measured performance in Figure 11. Where applicable, algorithm parameters were selected to ensure a fair comparison. graph-tool's implementation of Brandes' betweenness algorithm does more work as it also calculates edge betweenness scores during the run. graph-tool also takes a different approach to community detection, hence the comparison is between igraph and NetworkKit only. We summarize the benchmark results as follows: In our benchmark, NetworkKit was the only framework that could consistently run the set of kernels (excluding the quadratic-time betweenness) on the full set of networks in the timeframe of an overnight run. For some of igraph's and graph-tool's implementations, the test set had to be restricted to a subset of smaller networks to make it possible to run the complete benchmark overnight. NetworkKit has the fastest average processing rate on all of these typical analytics kernels. Our implementations have a slight edge over the others for breadth-first search, connected components, clustering coefficients, and betweenness. Considering that the algorithms are very similar, this is likely due

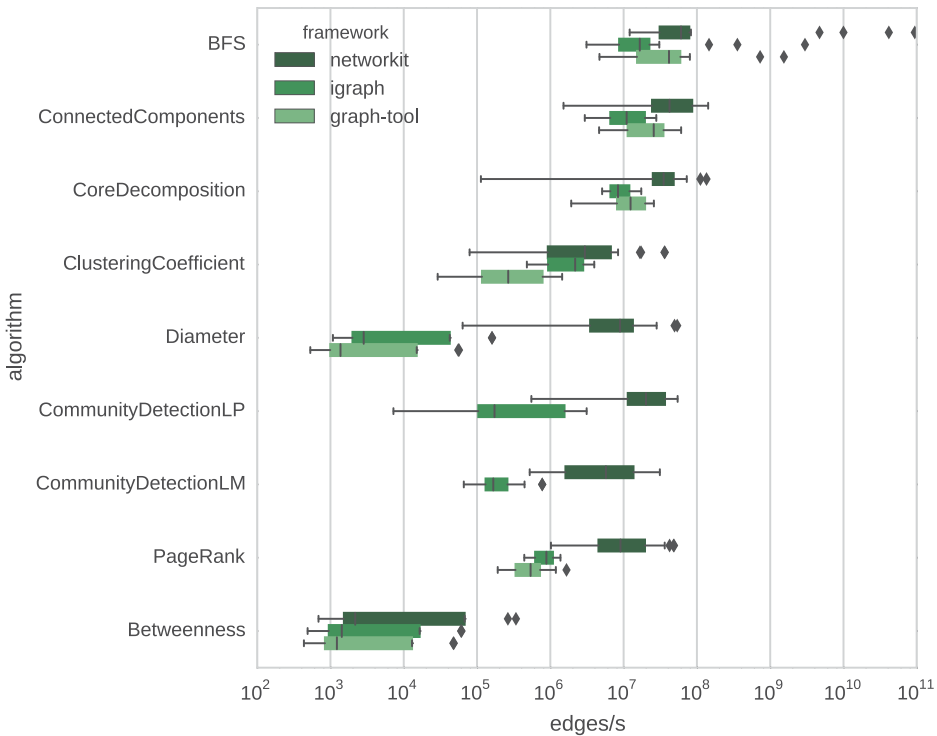


Fig. 11. Processing rates of typical analytics tasks: NetworKit in comparison with igraph and graph-tool. (Color online)

to subtle differences and optimizations in the implementation. For PageRank, core decomposition and the two community detection algorithms, our parallel methods lead to a larger speed advantage. The massive difference for the diameter calculation is due to our choice of the iFub algorithm (Crescenzi et al., 2013), which has better running time in the typical case (i.e. complex networks with hub structure) and enables the processing of inputs that are orders of magnitudes larger.

Another scalability factor is the memory footprint of the graph data structure. NetworKit provides a lean implementation in which the 260 M edges of the uk-2002 web graph occupy only 9 GB, compared with igraph (93 GB) and graph-tool (14 GB). After indexing the edges, e.g. in order to compute edge centrality scores, NetworKit requires 11 GB for the graph.

A third factor that should not be ignored for real workflows is I/O. Getting a large graph from hard disk to memory often takes far longer than the actual analysis. For our benchmark, we chose the GML graph file format for the input files, because it is supported by all three frameworks. We observed that the NetworKit parser is significantly faster for these non-attributed graphs (Figure 12).

9 Open-source development and distribution

Through open-source development, we would like to encourage usage and contributions by a diverse community, including data mining users and algorithm engineers. While the core developer team is located at KIT, NetworKit is conceived

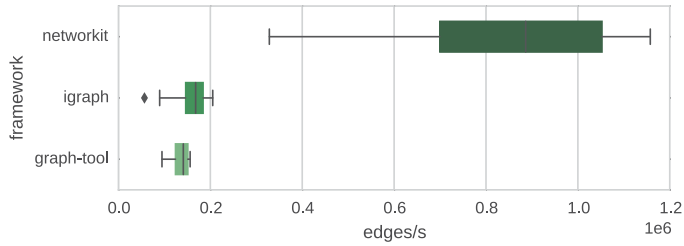


Fig. 12. I/O rates of reading a graph from a GML file: NetworkKit in comparison with igraph and graph-tool. (Color online)

as a community project with a growing number of external users and contributors. The code is free software licensed under the permissive MIT License. The package source, documentation, and additional resources can be obtained from <http://networkkit.iti.kit.edu>. The package networkkit is also installable via the Python package manager pip.

10 Conclusion

The NetworkKit project exists at the intersection of graph algorithm research and network science. Its contributors develop and collect state-of-the-art algorithms for network analysis tasks and incorporate them into ready-to-use software. The open-source package is under continuous development. The result is a tool suite of network analytics kernels, network generators and utility software to explore and characterize large network data sets on typical multicore computers. We detailed techniques that allow NetworkKit to scale to large networks, including appropriate algorithm patterns (parallelism, heuristics, data structures) and implementation patterns (e. g. modular design). The interface provided by our Python module allows domain experts to focus on data analysis workflows instead of the intricacies of programming. This is facilitated by a new frontend that generates comprehensive statistical reports on structural features of the network. Users familiar with the Python ecosystem of data analysis tools will appreciate the possibility to seamlessly integrate our toolkit.

Among similar software packages, NetworkKit yields the best performance for common analysis workflows. Our experimental study showed that NetworkKit is capable of quickly processing large-scale networks for a variety of analytics kernels in a reliable manner. This translates into faster workflows and extended analysis capabilities in practice. We recommend NetworkKit for the comprehensive structural analysis of large complex networks, as well as processing large batches of smaller networks. With fast parallel algorithms, scalability is in practice primarily limited by the size of the shared memory: A standard multicore workstation with 256 GB RAM can therefore process up to 10^{10} edge graphs.

Acknowledgments

This work was partially supported by the project *Parallel Analysis of Dynamic Networks—Algorithm Engineering of Efficient Combinatorial and Numerical Methods*, which is funded by the Ministry of Science, Research and the Arts Baden-Württemberg, and by DFG grant ME-3619/3-1 (FINCA) within the SPP 1736

Algorithms for Big Data. Aleksejs Sazonovs acknowledges support by the RISE program of the German Academic Exchange Service (DAAD). We thank Maximilian Vogel and Michael Hamann for continuous algorithm and software engineering work on the package. We also thank Lukas Barth, Miriam Beddig, Elisabetta Bergamini, Stefan Bertsch, Pratistha Bhattarai, Andreas Bilke, Simon Bischof, Guido Brückner, Mark Erb, Kolja Esders, Patrick Flick, Lukas Hartmann, Daniel Hoske, Gerd Lindner, Moritz v. Looz, Yassine Marrakchi, Mustafa Özdayi, Marcel Radermacher, Klara Reichard, Marvin Ritter, Arie Slobbe, Florian Weber, Michael Wegner, and Jörg Weisbarth for contributing to the project.

References

- Aiello, W., Chung, F., & Lu, L. (2000). A random graph model for massive graphs. In *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, Portland, OR, USA. ACM, pp. 171–180.
- Alahakoon, T., Tripathi, R., Kourtellis, N., Simha, R., & Iamnitchi, A. (2011). K-path centrality: A new centrality measure in social networks. In *Proceedings of the 4th Workshop on Social Network Systems*, Salzburg, Austria. ACM, p. 1.
- Albert, R., & Barabási, A. (2002). Statistical mechanics of complex networks. *Reviews of Modern Physics*, **74**(1), 47.
- Bader, D. A., Meyerhenke, H., Sanders, P., Schulz, C., Kappes, A., & Wagner, D. (2014). Benchmarking for graph clustering and partitioning. In *Encyclopedia of Social Network Analysis and Mining* (pp. 73–82). Springer.
- Bastian, M., Heymann, S., & Jacomy, M. (2009). Gephi: An open source software for exploring and manipulating networks. In *Proceedings of the Third International Conference on Weblogs and Social Media, ICWSM 2009*, San Diego, CA, USA, May 17–20, 2009.
- Batagelj, V., & Brandes, U. (2005). Efficient generation of large random networks. *Physical Review E*, **71**(3), 036113.
- Batagelj, V., & Mrvar, A. (2004). *Pajek—analysis and visualization of large networks*, Lecture Notes in Computer Science, vol. 2265 (pp 477–478). Springer.
- Batagelj, V., & Zaveršnik, M. (2011). Fast algorithms for determining (generalized) core groups in social networks. *Advances in Data Analysis and Classification*, **5**(2), 129–145.
- Behnel, S., Bradshaw, R., Citro, C., Dalcin, L., Seljebotn, D. S., & Smith, K. (2011). Cython: The best of both worlds. *Computing in Science & Engineering*, **13**(2), 31–39.
- Bergamini, E., & Meyerhenke, H. (2015). Fully-dynamic approximation of betweenness centrality. In *Proceedings of the Algorithms - ESA 2015 - 23rd Annual European Symposium*, Patras, Greece, September 14–16. Springer, pp. 155–166.
- Bergamini, E., Meyerhenke, H., & Staudt, C. (2015). Approximating betweenness centrality in large evolving networks. In *Proceedings of the 17th Workshop on Algorithm Engineering and Experiments, ALENEX 2015*. San Diego, CA, USA, January 5, 2015. SIAM, pp. 133–146.
- Blondel, V. D., Guillaume, J.-L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, **2008**(10), P10008.
- Boccaletti, S., Latora, V., Moreno, Y., Chavez, M., & Hwang, D.-U. (2006). Complex networks: Structure and dynamics. *Physics Reports*, **424**(4), 175–308.
- Bode, M., Fountoulakis, N., & Müller, T. (2015). On the largest component of a hyperbolic model of complex networks. *Electronic Journal of Combinatorics*, **22**(3), P3.24.
- Boldi, P., Codenotti, B., Santini, M., & Vigna, S. (2004). Ubicrawler: A scalable fully distributed web crawler. *Software: Practice & Experience*, **34**(8), 711–726.

- Brandes, U. (2001). A faster algorithm for betweenness centrality. *Journal of Mathematical Sociology*, **25**(2), 163–177.
- Brin, S., & Page, L. (2012). Reprint of: The anatomy of a large-scale hypertextual web search engine. *Computer Networks*, **56**(18), 3825–3833.
- CAIDA (2003). Caida skitter router-level topology measurements. Retrieved from <http://www.caida.org/data/router-adjacencies/>.
- Chakrabarti, D., Zhan, Y., & Faloutsos, C. (2004). R-MAT: A recursive model for graph mining. In *SDM* (vol. 4, pp. 442–446). Orlando, FL, USA: SIAM.
- Costa, L. d. F., Oliveira Jr, O. N., Traverso, G., Rodrigues, F. A., Villas Boas, P. R., Antiqueira, L., . . . Correa Rocha, L. E. (2011). Analyzing and modeling real-world phenomena with complex networks: A survey of applications. *Advances in Physics*, **60**(3), 329–412.
- Crescenzi, P., Grossi, R., Habib, M., Lanzi, L., & Marino, A. (2013). On computing the diameter of real-world undirected graphs. *Theoretical Computer Science*, **514**, 84–95.
- Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network research. *InterJournal, Complex Systems*, **1695**(5), 1–9.
- Dasari, N. S., Ranjan, D., & Zubair, M. (2014). ParK: An efficient algorithm for k-core decomposition on multicore processors. In J. Lin, J. Pei, X. Hu, W. Chang, R. Nambiar, C. Aggarwal, N. Cercone, V. Honavar, J. Huan, B. Mobasher, & S. Pyne (Eds.), *IEEE International Conference on Big Data, Big Data 2014* (pp. 9–16). Washington, DC: IEEE, October 27–30, 2014.
- Ediger, D., Jiang, K., Riedy, E. J., & Bader, D. A. (2013). GraphCT: Multithreaded algorithms for massive graph analysis. *IEEE Transactions on Parallel and Distributed Systems*, **24**(11), 2220–2229.
- Ediger, D., McColl, R., Riedy, E. J., & Bader, D. A. (2012). STINGER: High performance data structure for streaming graphs. In *IEEE Conference on High Performance Extreme Computing HPEC 2012*, Waltham, MA, USA, September 10–12, 2012, pp. 1–5.
- Eppstein, D., & Wang, J. (2004). Fast approximation of centrality. *Journal of Graph Algorithms Applications*, **8**, 39–45.
- Esders, K. (2015). Link prediction in large-scale complex networks. Master's thesis, Karlsruhe Institute of Technology.
- Flick, P. (2014). Analysis of human tissue-specific protein-protein interaction networks. Master's thesis, Karlsruhe Institute of Technology.
- Freeman, L. C. (1979). Centrality in social networks conceptual clarification. *Social Networks*, **1**(3), 215–239.
- Geisberger, R., Sanders, P., & Schultes, D. (2008). Better approximation of betweenness centrality. In *Proceedings of the Tenth Workshop on Algorithm Engineering and Experiments (ALENEX)*, San Francisco, CA, USA. SIAM, pp. 90–100.
- Girvan, M., & Newman, M. (2002). Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, **99**(12), 7821.
- Gugelmann, L., Panagiotou, K., & Peter, U. (2012). Random hyperbolic graphs: Degree sequence and clustering - (extended abstract). In *Automata, Languages, and Programming - 39th International Colloquium, ICALP 2012, Proceedings, Part II*, pp. 573–585.
- Hagberg, A., Swart, P., & S Chult, D. (2008). Exploring network structure, dynamics, and function using NetworkX. Technical report, Los Alamos National Laboratory (LANL).
- Hakimi, S. L. (1962). On realizability of a set of integers as degrees of the vertices of a linear graph. I. *Journal of the Society for Industrial & Applied Mathematics*, **10**(3), 496–506.
- Hamann, M., Lindner, G., Meyerhenke, H., Staudt, C. L., & Wagner, D. (2016). Structure-preserving sparsification methods for social networks. *Social Network Analysis and Mining*, **6**(1), 22:1–22:22.
- Katz, L. (1953). A new status index derived from sociometric analysis. *Psychometrika*, **18**(1), 39–43.

- Kiwi, M. A., & Mitsche, D. (2015). A bound for the diameter of random hyperbolic graphs. In *Proceedings of the Twelfth Workshop on Analytic Algorithmics and Combinatorics, ANALCO 2015*, San Diego, CA, USA, January 4, 2015, pp. 26–39.
- Koch, J., Staudt, C. L., Vogel, M., & Meyerhenke, H. (2015). Complex network analysis on distributed systems: An empirical comparison. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, Paris, France. ACM, pp. 1169–1176.
- Krioukov, D., Papadopoulos, F., Kitsak, M., Vahdat, A., & Boguñá, M. (2010). Hyperbolic geometry of complex networks. *Physical Review E*, **82**, 036106.
- Kunegis, J. (2013). KONECT: The koblenz network collection. In *22nd International World Wide Web Conference, WWW '13*, Rio de Janeiro, Brazil, May 13–17, 2013, Companion Volume, pp. 1343–1350.
- Lambiotte, R. (2010). Multi-scale modularity in complex networks. In *Proceedings of the Eighth International Symposium on Modeling and Optimization in Mobile, Ad-Hoc and Wireless Networks (WiOpt 2010)*, May 31–June 4, 2010, University of Avignon, Avignon, France, pp. 546–553.
- Lancichinetti, A., & Fortunato, S. (2009). Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities. *Physical Review E*, **80**(1), 016118.
- Leskovec, J., & Krevl, A. (2014). SNAP Datasets: Stanford large network dataset collection. Retrieved from <http://snap.stanford.edu/data>.
- Leskovec, J., & Sosič, R. (2014). SNAP: A general purpose network analysis and graph mining library in C++. Retrieved from <http://snap.stanford.edu/snap>.
- Low, Y., Bickson, D., Gonzalez, J., Guestrin, C., Kyrola, A., & Hellerstein, J. M. (2012). Distributed graphlab: A framework for machine learning and data mining in the cloud. *Proceedings of the VLDB Endowment*, **5**(8), 716–727.
- Lugowski, A., Alber, D.M., Buluç, A., Gilbert, J. R., Reinhardt, S.P., Teng, Y., & Waranis, A. (2012). A flexible open-source toolbox for scalable complex graph analysis. In *Proceedings of the Twelfth SIAM International Conference on Data Mining*, Anaheim, CA, USA, April 26–28, 2012, pp. 930–941.
- Newman, M. (2010). *Networks: An introduction*. New York, NY, USA: Oxford University Press.
- O'Madadhain, J., Fisher, D., White, S., & Boey, Y. (2003). The JUNG (java universal network/graph) framework. Irvine, California: University of California.
- Erdős, P., & Rényi, A. (1960). On the evolution of random graphs. *Publication of the Mathematical Institute of the Hungarian Academy of Sciences*, **5**.
- Peixoto, T. P. (2006). graph-tool. Retrieved from <http://graph-tool.skewed.de>.
- Perez, F., Granger, B. E., & Obispo, C. (2013). An open source framework for interactive, collaborative and reproducible scientific computing and education.
- Raghavan, U. N., Albert, R., & Kumara, S. (2007). Near linear time algorithm to detect community structures in large-scale networks. *Physical Review E*, **76**(3), 036106.
- Riondato, M., & Kornaropoulos, E. M. (2016). Fast approximation of betweenness centrality through sampling. *Data Mining and Knowledge Discovery*, **30**(2), 438–475.
- Schank, T., & Wagner, D. (2005). Approximating clustering coefficient and transitivity. *Journal of Graph Algorithms and Applications*, **9**(2), 265–275.
- Shun, J., & Blelloch, G. E. (2013). Ligra: A lightweight graph processing framework for shared memory. In *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPoPP '13*, Shenzhen, China, February 23–27, 2013. ACM, pp. 135–146.
- Sporns, O., & Betzel, R. F. (2016). Modular brain networks. *Annual Review of Psychology*, **67**, 613–640.

- Staudt, C. L., & Meyerhenke, H. (2016). Engineering parallel algorithms for community detection in massive networks. *IEEE Transactions on Parallel and Distributed Systems*, **27**(1), 171–184.
- Traud, A. L., Mucha, P. J., & Porter, M. A. (2012). Social structure of facebook networks. *Physica A: Statistical Mechanics and its Applications*, **391**(16), 4165–4180.
- von Looz, M., Meyerhenke, H., & Prutkin, R. (2015). Generating random hyperbolic graphs in subquadratic time. In *Proceedings of 26th International Symposium on Algorithms and Computation (ISAAC 2015)*, LNCS, Nagoya, Japan. Springer.
- Zafarani, R., & Liu, H. (2009). Social computing data repository at ASU. Retrieved from <http://socialcomputing.asu.edu>.