



Engineering and
Physical Sciences
Research Council

Spectral Toolkit of Algorithms for Graphs (STAG)

Technical Report (1)

Peter Macgregor and He Sun

School of Informatics
University of Edinburgh
United Kingdom
{peter.macgregor, h.sun}@ed.ac.uk

Current Version	STAG 1.2
Main Update	Local graph clustering
Date	7th April 2023
Website	https://staglibrary.io
Source Code	https://github.com/staglibrary
Funding Information	The development of STAG is supported by UK Research and Innovation under an EPSRC Early Career Fellowship scheme (EP/T00729X/1).

arXiv:2304.03170v1 [cs.SI] 5 Apr 2023



1 Summary

Spectral Toolkit of Algorithms for Graphs (STAG) is an open-source C++ and Python library of efficient spectral algorithms for graphs. Our objective is to implement advanced graph algorithms developed through algorithmic spectral graph theory, while making it practical to end users. This series of technical reports is to document our progress on STAG, including implementation details, engineering considerations, and the data sets against which our implementation is tested. The report is structured as follows:

- Section 2 describes the local clustering algorithm, which is the main update in this STAG release. The discussion is at a high level such that domain knowledge beyond basic algorithms is not needed.
- Section 3 provides a user guide to the essential features of STAG which allow a user to apply local clustering.
- Section 4 includes experiments and demonstrations of the functionality of STAG.
- Finally, Section 5 discusses several technical details; these include our choice of implemented algorithms, the default setup of parameters, and other technical choices. We leave these details to the final section, as it's not necessary for the reader to understand this when using STAG.

1.1 Implemented Algorithms

STAG 1.2 provides an implementation of the following key algorithms.

Local Graph Clustering. Given a large graph and some starting vertex v in the graph, the goal of local graph clustering is to find some cluster containing v . Moreover, the running time of the algorithm should depend only on the size of the returned cluster and should be independent of the total size of the graph [11].

STAG provides the first open-source local clustering algorithm which does not require the entire graph to be loaded into memory. This allows users to apply local clustering on massive graphs stored on disk or even in a cloud database, such as Neo4j¹. Section 4 demonstrates these applications.

Spectral Clustering. One of the most fundamental algorithms from Spectral Graph Theory is the spectral clustering algorithm [8, 10, 13]. Spectral clustering is “global” in the sense that it returns a partition of the entire vertex set of the graph.

Generating Graphs from Random Models. The Stochastic Block Model (SBM) and Erdős-Rényi model are popular random graph models which are frequently used to evaluate and analyse graph algorithms. STAG provides several convenient methods to generate graphs from these models.

2 Local Graph Clustering

Graph clustering algorithms are designed to partition an input graph into two or more clusters. As a basic technique in data science and machine learning, graph clustering has many applications in numerous areas of computer science and beyond. Most graph clustering algorithms need to read an entire input graph for the clustering task, which is computationally expensive if the graph is massive. If one is interested only in some “local” cluster information, then local graph clustering provides a more efficient method.

Typically, the objective of local graph clustering is to find some highly-connected vertex set (cluster) in an input graph. Let's assume that S is a highly-connected vertex set of an underlying undirected graph G , i.e., S forms a cluster. Then, a local clustering algorithm is given some vertex $v \in S$ as input, and returns some set S' such that S' is a reasonable approximation of the target set S . Moreover, the running time

¹<https://neo4j.com/>



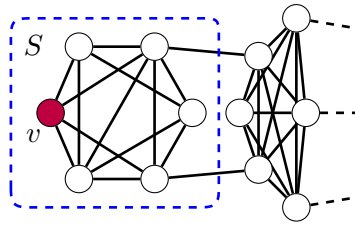


Figure 1: Given a massive graph containing some small cluster S , a local clustering algorithm takes as input a vertex $v \in S$ and returns an approximation of S without exploring the whole graph.

of the algorithm is proportional to the size of S and independent of the size of G . In applications, local clustering can be viewed as a search for related objects: given a query vertex, a local clustering algorithm returns a set of closely related vertices. Figure 1 illustrates local graph clustering.

Andersen, Chung, and Lang [1] introduced a key local clustering algorithm which we refer to as the ACL algorithm. At a high level, the algorithm finds a local cluster by analysing the behaviour of random walks on the graph, beginning at the starting vertex. The ACL algorithm has proved extremely useful and has inspired extensive further research and applications [2, 7, 12, 14]. The `local_cluster` method of STAG provides an implementation of the ACL local clustering algorithm.

Existing open-source local clustering methods require that the entire graph is loaded into RAM in order to apply the local algorithm. In this sense, they are not truly “local” since they cannot be applied to graphs larger than the available memory and the total running time depends on the size of the graph. STAG provides the first open-source local clustering algorithm which can be applied to massive graphs without loading them into RAM. Moreover, the provided interface is simple and the algorithm can be applied to graphs stored in memory, on disk, or in a Neo4j database.

3 User’s Guide to Local Graph Clustering with STAG

This section provides a guide to the essential features of STAG which allow a user to apply local clustering. Section 3.1 describes how to install the STAG C++ and STAG Python libraries. Then, Section 3.2 introduces the graph file formats supported by STAG and demonstrates the methods for reading and writing graphs to disk. Finally, Section 3.3 explains the graph classes provided by STAG and Section 3.4 documents the `local_cluster` method for local clustering.

Although most of the examples in this section use C++, the functionality of STAG C++ is also available in STAG Python. Appendix B includes example code demonstrating how to perform local clustering with STAG Python. The full documentation of STAG C++ and STAG Python is available on the STAG library website.

3.1 Installation of STAG

This section describes how to install STAG for use with C++ and Python.

Installing STAG for C++. STAG is built on the Eigen and Spectra C++ libraries, and these must be installed before STAG. For information on installing Eigen and Spectra, please refer to their documentation. For convenience, Appendix A provides a bash script which, at the time of writing, will install Eigen and Spectra on a standard Linux system. Then, the latest version of STAG should be downloaded from

<https://github.com/staglibrary/stag/releases>.

After downloading and extracting the source code, STAG can be compiled and installed with `cmake`.

```

1  mkdir build_dir
2  cd build_dir
3  cmake ..
4  sudo make install

```



Once STAG has been installed, it is available for use in C++ projects built with the `cmake` build tools. The following `cmake` code will link a C++ project with STAG.

```
1 find_package(stag REQUIRED)
2 include_directories(${STAG_INCLUDE_DIRS})
3 target_link_libraries(YOUR_PROJECT stag)
```

An example STAG project demonstrating the full `cmake` configuration is available at

<https://github.com/staglibrary/example-stag-project>.

Installing STAG for Python. STAG Python can be installed from the Python Package Index with the `pip` tool.

```
1 python -m pip install stag
```

Then, the modules of STAG can be directly imported into any Python script.

3.2 File Formats

STAG supports two simple file formats for storing graphs on disk: `EdgeList` and `AdjacencyList`. Many graph datasets are provided in `EdgeList` format [6], and will work directly with STAG.

EdgeList File Format. In an `EdgeList` file, each line corresponds to one edge in the graph. A line consists of two integer node IDs and an optional edge weight, all separated with spaces. Here is an example of a simple `EdgeList` file.

```
1 # This is a comment
2 0 1 0.5
3 1 2 1
4 2 0 3
```

In this example, line 2 defines an edge between nodes 0 and 1 with weight 0.5.

AdjacencyList File Format. In an `AdjacencyList` file, each line corresponds to one node in the graph. A line consists of the node ID, followed by a list of adjacent nodes. The node IDs at the beginning of each line must be sorted in increasing order. Here is an example of a simple `AdjacencyList` file.

```
1 # This is a comment
2 0: 1 2
3 1: 0 3 2
4 2: 0 1
5 3: 1
```

In this example, node 1 has edges to nodes 0, 2, and 3.

Working with Files. STAG provides several methods for reading, writing, and converting between `EdgeList` and `AdjacencyList` files, as demonstrated in the following example.

```
1 #include <stag/graphio.h>
2 ...
3 // Read an AdjacencyList graph
4 std::string filename = "mygraph.adjacencylist";
5 stag::Graph myGraph = stag::load_adjacencylist(filename);
6
7 // Save an EdgeList graph
8 std::string new_filename = "mygraph.edgelist";
9 stag::save_edgelist(myGraph, new_filename);
10
11 // Convert an AdjacencyList to EdgeList directly
12 stag::adjacencylist_to_edgelist(filename, new_filename);
13 ...
```



3.3 Graph Classes

STAG provides several graph classes which can be applied for a wide variety of applications. Figure 2 summarises the available graph classes.

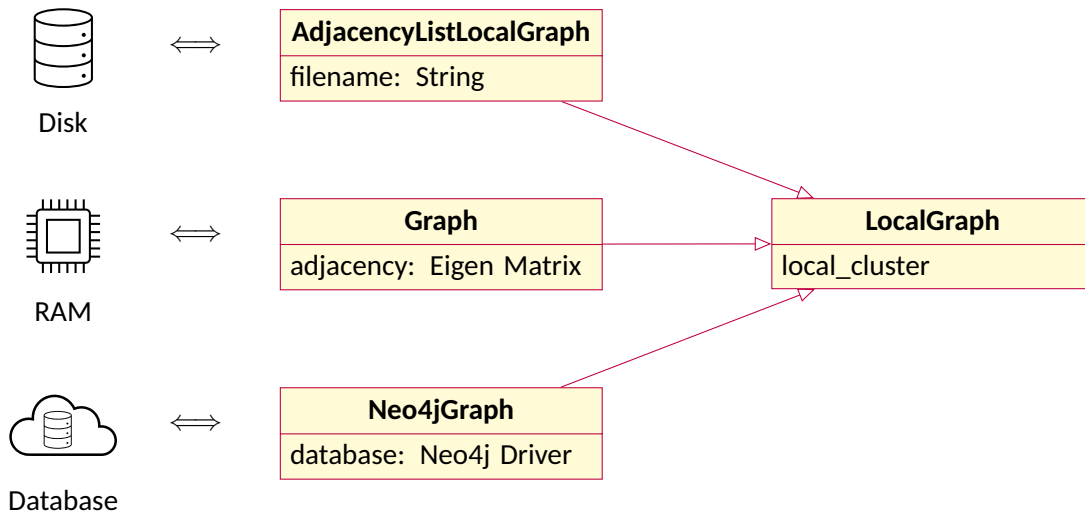


Figure 2: The graph classes provided by STAG. Every class inherits from the abstract `LocalGraph` class, which provides the local clustering method. The key difference between the different classes is the location of the graph data. The `AdjacencyListLocalGraph` class reads node adjacency data from disk, the `Graph` class stores the entire graph in RAM, and the `Neo4jGraph` class queries node adjacency data from a Neo4j database.

The LocalGraph Class. STAG provides an abstract `stag::LocalGraph` class which defines the data structure necessary to apply local clustering. The only required method on the data structure is `neighbors(v)` which returns a list of the neighbors of node v . Every graph class provided by STAG inherits from `LocalGraph`.

The Graph Class. The `stag::Graph` class is the basic graph object within the STAG library. The class stores the adjacency matrix of the graph in memory as a sparse matrix.

The AdjacencyListLocalGraph Class. The `stag::AdjacencyListLocalGraph` class provides an implementation of the `stag::LocalGraph` interface for a graph stored on disk as an `AdjacencyList`. The graph is loaded into memory in a local way only. This allows for local algorithms to be executed on very large graphs stored on disk without loading the whole graph into memory. The following example demonstrates how to create an `AdjacencyListLocalGraph` with STAG C++.

```
1 #include <stag/graph.h>
2 ...
3 // Create an AdjacencyListLocalGraph
4 std::string filename = "mygraph.adjacencylist";
5 stag::AdjacencyListLocalGraph myGraph(filename);
6
7 // Get the neighbours of node 0
8 std::vector<long long> neighbors = myGraph.neighbors_unweighted(0);
9 ...
```

The Neo4jGraph Class. STAG Python additionally provides the `Neo4jGraph` class, which provides an implementation of the `LocalGraph` interface for a graph stored in a Neo4j database. The following example shows how to create the `Neo4jGraph` using the database connection information.



```

1  import stag.neo4j
2
3  # Connect to the database
4  uri = "<Database URI>"
5  username = "neo4j"
6  password = "<password>"
7  my_graph = stag.neo4j.Neo4jGraph(uri, username, password)
8
9  # Print the neighbors of node 0
10 print(my_graph.neighbors_unweighted(0))

```

The `Neo4jGraph` class provides additional methods for querying the properties of the nodes in the database. The details are available in the full STAG documentation.

3.4 Local Clustering

STAG provides the following `local_cluster` method.

```

1  std::vector<long long> stag::local_cluster(stag::LocalGraph* graph,
2                                          long long      seed_vertex,
3                                          double         target_volume)

```

Given a graph and a starting vertex, the `local_cluster` method finds a cluster close to the starting vertex. The running time of the algorithm is proportional to the size of the returned cluster and independent of the size of the entire graph. The parameters of the method are described as follows:

- **graph** - a `LocalGraph` object. This could be a `Graph`, an `AdjacencyListLocalGraph`, Or a `Neo4jGraph`.
- **seed_vertex** - the starting vertex in the graph.
- **target_volume** - an estimate of the volume of the target cluster. This parameter does *not* impose a hard constraint on the algorithm and so an approximate volume is sufficient.

When working with very large graphs, it is recommended to use the `AdjacencyListLocalGraph` object for local clustering in order to avoid the overhead of reading the entire graph into memory. Section 4.1 demonstrates the advantage of using the `AdjacencyListLocalGraph` for local clustering. The following code demonstrates a complete program which uses STAG C++ to find a local cluster in a graph stored in an `AdjacencyList` file on disk.

```

1  #include <iostream>
2  #include <stag/graph.h>
3  #include <stag/cluster.h>
4
5  int main() {
6      // Create the graph backed by a large file on disk
7      std::string filename = "mygraph.adjacencylist";
8      stag::AdjacencyListLocalGraph mygraph(filename);
9
10     // Perform local clustering
11     int start_vertex = 1;
12     double target_volume = 100;
13     auto cluster = stag::local_cluster(&mygraph, start_vertex, target_volume);
14
15     // Print the returned cluster
16     for (auto v : cluster) std::cout << v << ", ";
17     std::cout << std::endl;
18 }

```



4 Showcase studies

STAG makes local clustering straightforward for a variety of applications, and this section presents some examples of local clustering with STAG. The code used to produce all experimental results is available at

<https://github.com/staglibrary/local-clustering-case-study>.

All experiments are performed on an HP ZBook laptop with an 11th Gen Intel(R) Core(TM) i7-11800H @ 2.30GHz processor and 32 GB RAM.

4.1 Example 1

The advantage of local clustering over other clustering algorithms is that the running time of local clustering is proportional to the size of the returned cluster and independent of the total size of the graph. If we first load the entire graph into memory before applying local clustering, then we lose the advantage of the sub-linear running time. For this reason, STAG provides the `AdjacencyListLocalGraph` class which provides local access to a graph stored on disk without reading the entire graph. In this example, we compare the running time of local clustering on an `AdjacencyListLocalGraph` object, which accesses the graph locally on disk, and a `Graph` object, which loads the entire graph into memory.

We generate graphs of various sizes from the stochastic block model as follows. Given parameters k , p , and q , we create a graph with k clusters C_1, \dots, C_k , each containing 1,000 vertices. For every pair of vertices $(u, v) \in V \times V$, we add the edge (u, v) with probability p if u and v are in the same cluster and with probability q otherwise. We always set $p = 0.01$ and $q = 0.001/k$. This ensures that the conductance of the constructed clusters is always close to 0.1.

We perform local clustering on the constructed graphs for a random starting node and target volume 20,000, and compare the following two methods:

- **In memory:** the entire graph is loaded into memory as a `Graph` object before applying local clustering.
- **On disk:** the graph is read locally from a file on disk with an `AdjacencyListLocalGraph` object.

Figure 3 shows the running time of the local clustering algorithm for each method across a range of graph sizes. These results demonstrate that for large graphs, the overhead of reading the entire graph into memory dominates the running time of the algorithm and reading the graph directly from disk is significantly more efficient.

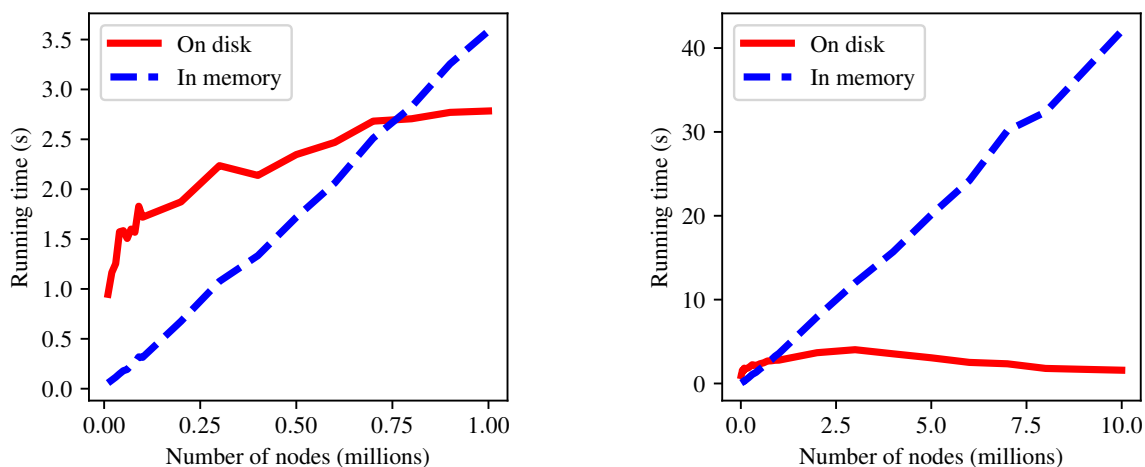


Figure 3: The comparison of the running time of local clustering on graphs in memory and on disk. For massive graphs, reading the graph locally from disk is significantly faster.



4.2 Example 2

In the second example, we demonstrate the applicability of local clustering for finding sets of related nodes in a real-world graph. We use the `wiki-topcats` dataset [14] which is a graph of Wikipedia hyperlinks constructed in 2011. The graph includes the pages in the top 100 Wikipedia categories and includes 1,791,489 vertices and 28,511,807 edges. The dataset is available on the SNAP datasets page [6] as an EdgeList file. We convert the EdgeList to an AdjacencyList with the `edgelist_to_adjacencylist` method, and use the `AdjacencyListLocalGraph` object for local clustering.

With a few lines of code, STAG allows us to create a “related pages” search using local clustering. By providing a search page and setting the target volume to be 100, the local clustering returns a set of pages which are closely connected to the search page. Figure 4 shows some example of local clustering results from the Wikipedia graph.

Search page: Emacs	Search page: Stag Hound	Search page: Stagflation
Emacs	Stag Hound	Stagflation
Robert J. Chassell	Over-canvassed sailing	Agflation
Zmacs	Sea Serpent (clipper)	Biflation
Climacs	Medium Clipper	Differential accumulation
Aquamacs	Donald McKay	Embedded liberalism
Dunnet (game)	Extreme clipper	Supply shock
Bernard Greenberg	Memnon (clipper)	Shimshon Bichler
Macro recorder	Flying Cloud (clipper)	Douglas Harper
GNU Manifesto	Eleanor Creesy	Online savings account
TNT (instant messenger)	Weigh anchor	Jonathan Nitzan
Agda (theorem prover)	Stowage	

Figure 4: Example of local clusters found in the Wikipedia dataset.

4.3 Example 3

In this example, we demonstrate local clustering on a Neo4j database in the cloud. We first follow the Neo4j documentation to create a cloud database using the AuraDB service [9]. We use the “Movies” example dataset provided by Neo4j. Then, by creating a `Neo4jGraph` object with STAG Python, we are able to search for related movies using local clustering. Listing 1 shows the complete Python script used to perform this search, and Figure 5 shows some of the search results.

```
1 import stag.neo4j
2 import stag.cluster
3
4 # Initialise the graph object with the Neo4j database credentials
5 uri = "<Database URI>"
6 username = "neo4j"
7 password = "<password>"
8 g = stag.neo4j.Neo4jGraph(uri, user, password)
9
10 # Perform local clustering
11 seed_id = g.query_id("title", "The Matrix")
12 cluster = stag.cluster.local_cluster(g, seed_id, 5)
13
14 # Print the names of the returned movies.
15 for node_id in cluster:
16     title = g.query_property(node_id, 'title')
17     print(title)
```

Listing 1: Python code for local clustering with a Neo4j database.



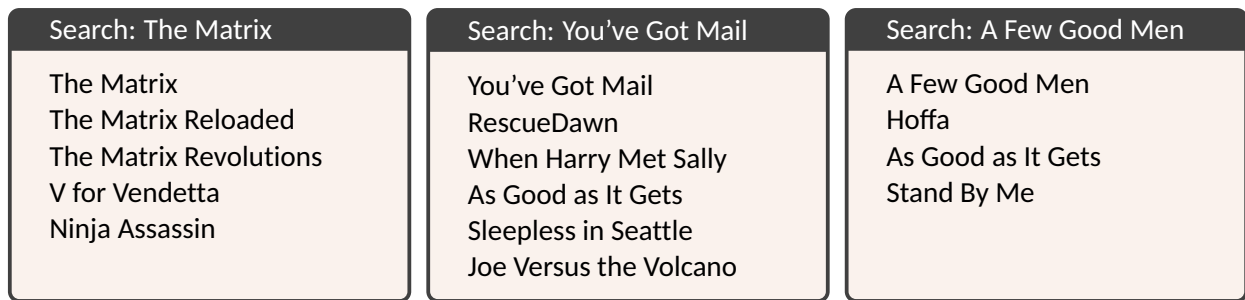


Figure 5: Example of local clusters found in the Neo4j movies dataset.

5 Technical Considerations

In this section we discuss the technical choices made in the design and implementation of STAG. We discuss our choice of implemented algorithm, the setting of the default parameters, and the implementation of the `AdjacencyListLocalGraph` class for reading a graph locally from a file on disk.

5.1 Implemented Algorithm

Many algorithms have been proposed for local clustering, including those based on PageRank [1, 2], the evolving set process [3, 4], and network flows [5]. We chose to implement the algorithm based on PageRank presented by Andersen, Chung, and Lang [1], and we refer to this as the ACL algorithm. We chose this algorithm because it is relatively simple, easy to understand, and effective in practice. Furthermore, the theoretical guarantees for the ACL algorithm are optimal up to constant factors.² The ACL algorithm requires two parameters:

- the α parameter controls the “teleport probability” of the personalised PageRank; and
- the ϵ parameter controls the approximation error of the approximate PageRank calculation.

STAG provides the `local_cluster_acl` method which allows the user to specify the parameters α and ϵ directly.

```

1  std::vector<long long> stag::local_cluster_acl(stag::LocalGraph* graph,
2                                             long long      seed_vertex,
3                                             double         alpha,
4                                             double         epsilon)

```

For convenience, STAG also provides the `local_cluster` method which requires only an estimate of the volume of the target cluster. Given a volume γ , the `local_cluster` method uses the parameters $\alpha = 1/2000$ and $\epsilon = 1/(20\gamma)$ for the ACL algorithm.

5.2 Reading Graphs Locally From Disk

A key feature of STAG is the `AdjacencyListLocalGraph` class which reads the neighbourhood information of a graph in a local way from an `AdjacencyList` file on disk. Since the data in an `AdjacencyList` file is sorted according to the node ID, we can query the neighbors of any node in $O(\log(n))$ time by binary search of the `AdjacencyList` file. As demonstrated in Section 4.1, this additional logarithmic factor in the running time is much preferable to the cost of reading the entire graph into memory when applying local algorithms to massive graphs.

²The original analysis by Andersen et al. [1] has an extra factor of $\log(n)$ in the approximation guarantee. This factor is not necessary and has been removed in later analysis using the same technique [7, 12].



References

- [1] Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using Pagerank vectors. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pages 475–486, 2006.
- [2] Reid Andersen, Fan Chung, and Kevin Lang. Local partitioning for directed graphs using Pagerank. In *5th International Workshop on Algorithms and Models for the Web-Graph (WAW'07)*, pages 166–178, 2007.
- [3] Reid Andersen, Shayan Oveis Gharan, Yuval Peres, and Luca Trevisan. Almost optimal local graph clustering using evolving sets. *Journal of the ACM*, 63(2):1–31, 2016.
- [4] Reid Andersen and Yuval Peres. Finding sparse cuts locally using evolving sets. In *41st Annual ACM Symposium on Theory of Computing (STOC'09)*, pages 235–244, 2009.
- [5] Kimon Fountoulakis, Di Wang, and Shenghao Yang. p -Norm flow diffusion for local graph clustering. In *37th International Conference on Machine Learning (ICML'20)*, page 3222–3232, 2020.
- [6] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [7] Peter Macgregor and He Sun. Local algorithms for finding densely connected clusters. In *38th International Conference on Machine Learning (ICML '21)*, pages 7268–7278, 2021.
- [8] Peter Macgregor and He Sun. A tighter analysis of spectral clustering, and beyond. In *39th International Conference on Machine Learning (ICML '22)*, pages 14717–14742, 2022.
- [9] Neo4j. Neo4j graph database. <http://neo4j.org/>, 2023.
- [10] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *15th Advances in Neural Information Processing Systems (NeurIPS'01)*, pages 849–856, 2001.
- [11] Daniel A. Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):1–26, 2013.
- [12] Yuuki Takai, Atsushi Miyauchi, Masahiro Ikeda, and Yuichi Yoshida. Hypergraph clustering based on Pagerank. In *26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'20)*, pages 1970–1978, 2020.
- [13] Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- [14] Hao Yin, Austin R Benson, Jure Leskovec, and David F Gleich. Local higher-order graph clustering. In *23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'17)*, pages 555–564, 2017.

A Installing STAG Dependencies

For convenience, we provide the following bash script for installing the STAG C++ dependencies. At the time of writing, this will download and install the Eigen and Spectra libraries.

```
1 # Create a directory to work in
2 mkdir libraries
3 cd libraries
4
5 # Install Eigen
6 wget https://gitlab.com/libeigen/eigen/-/archive/3.4.0/eigen-3.4.0.tar.gz
7 tar xzvf eigen-3.4.0.tar.gz
8 cd eigen-3.4.0
```



```

9  mkdir build_dir
10 cd build_dir
11 cmake ..
12 sudo make install
13 cd ../../
14
15 # Install Spectra
16 wget https://github.com/yixuan/spectra/archive/v1.0.1.tar.gz
17 tar xzvf v1.0.1.tar.gz
18 cd spectra-1.0.1
19 mkdir build_dir
20 cd build_dir
21 cmake ..
22 sudo make install
23 cd ../../

```

B User Guide Examples using STAG Python

This section includes example code omitted from Section 3 demonstrating how to use STAG Python for local clustering.

B.1 Working with Files

The following example demonstrates how to read and write AdjacencyList and Edgelist files with STAG Python.

```

1  import stag.graphio
2
3  # Read an AdjacencyList graph
4  filename = "mygraph.adjacencylist"
5  myGraph = stag.graphio.load_adjacencylist(filename)
6
7  # Save an Edgelist graph
8  new_filename = "mygraph.edgelist"
9  stag.graphio.save_edgelist(myGraph, new_filename)
10
11 # Convert an AdjacencyList to Edgelist directly
12 stag.graphio.adjacencylist_to_edgelist(filename, new_filename)

```

B.2 Graph Classes

The following example shows how to create an AdjacencyListLocalGraph object with STAG Python.

```

1  import stag.graph
2
3  # Create the graph object
4  filename = "mygraph.adjacencylist"
5  my_graph = stag.graph.AdjacencyListLocalGraph(filename)
6
7  # Show the neighbors of node 0
8  print(my_graph.neighbors_unweighted(0))

```

B.3 Local Clustering

The following example gives a complete program for finding a local cluster in a graph stored in an AdjacencyList file with STAG Python.



```
1 import stag.graph
2 import stag.cluster
3
4 # Create the graph object
5 filename = "mygraph.adjacencylist"
6 mygraph = stag.graph.AdjacencyListLocalGraph(filename)
7
8 # Find a local cluster
9 start_vertex = 1
10 target_volume = 100
11 cluster = stag.cluster.local_cluster(mygraph, start_vertex, target_volume)
12
13 # Display the result
14 print(cluster)
```

