








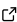

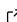
1 CM++ - A Meta-method for Well-Connected 2 Community Detection

3 **Vikram Ramavarapu** ¹, **Fábio Jose Ayres** ², **Minhyuk Park** ¹, **Vidya**
4 **Kamath Pailodi** ¹, **João Alfredo Cardoso Lamy** ², **Tandy Warnow** ¹, and
5 **George Chacko** ¹ ¶

6 ¹ Department of Computer Science, University of Illinois Urbana-Champaign, IL 61801, USA ² Insper
7 Institute, Sao Paulo, Brazil ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#)).

8 Introduction

9 Community detection methods help uncover the meso-scale structure of networks and have
10 broad applications ([Dey et al., 2022](#); [Haggerty et al., 2013](#); [Karatas & Sahin, 2018](#); [Waltman
& van Eck, 2012](#)). While communities can be defined in different ways ([Coscia et al., 2011](#)), a
11 common expectation is one of greater edge-density within and lesser edge density between
12 communities ([Fortunato & Newman, 2022](#)). A related expectation is that communities also
13 should be well-connected ([Bonchi et al., 2021](#); [Traag et al., 2019](#)).

14
15 In [Park et al. \(2023\)](#), we describe Connectivity Modifier (CM), a meta-method that enforces
16 well-connectedness in communities. As input, CM takes a network, a clustering of the network
17 generated by an algorithm, and a user-specified connectivity threshold. For each community
18 (cluster), CM uses VieCut ([Henzinger et al., 2018](#)) to find a small edge cut, and if the edge cut
19 size is below the specified connectivity threshold, then the cut is removed, and the partitions
20 are reclustered using the clustering algorithm. This process repeats until all clusters are well
21 connected.

22 We now present CM++ ([Ramavarapu et al., 2023](#)), which uses parallelism for scalability and
23 has new features. CM++ provides support for additional clustering paradigms and is designed
24 to be extensible by other developers. Concurrently, we present the CM++ Pipeline, a modular
25 and extensible workflow that automates CM++ operations. The pipeline consists of clustering,
26 pre-processing, connectivity modifier (CM++), and post-processing stages with generation of
27 cluster statistics.

28 Statement of Need

29 We have demonstrated that several widely-used codes do not generate well-connected clusters
30 ([Park et al., 2023](#)). A tool to enforce user-specified levels of well-connectedness to clusterings
31 from multiple community detection methods is not presently available.

CM Pipeline:

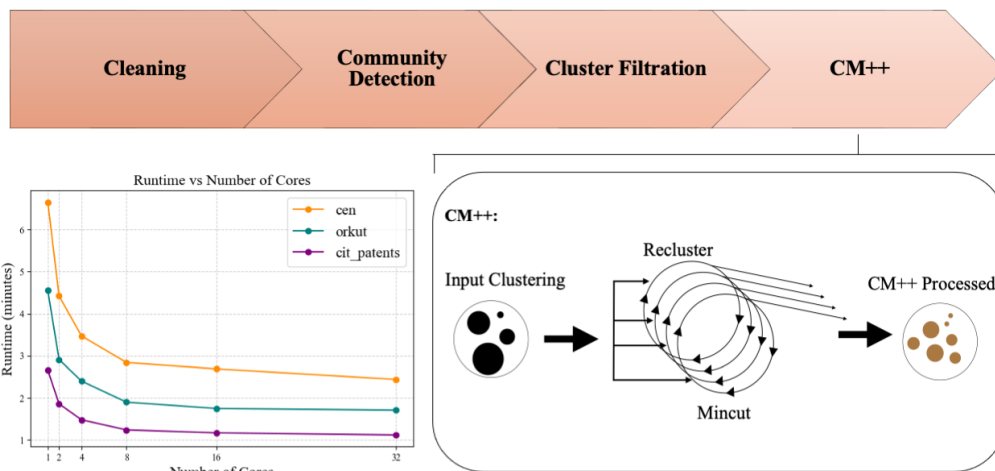


Figure 1: (Top) A visualization of a workflow created from the CM++ Pipeline. (Bottom right) Algorithmic schematic of CM++. CM++ splits the queue of clusters evenly between the spawned processes. Each process runs an instance of CM++ on its share of clusters (recursive mincutting and reclustering until well connected). (Bottom left) Runtime curve with respect to the number of parallel cores running CM++. CEN (14 million nodes, 1.3 billion edges) is the Curated Exosome Network. orkut (3.1 million nodes, 117 million edges) and cit_patents (3.7 million nodes, 16 million edges) are both from the SNAP database (Leskovec & Krevl, 2014) and processed through the removal of parallel edges and self-loops (Park et al., 2023).

CM++: Enforcing Well-Connectedness

Key Features:

- **Flexibility** For users to accompany their definition of a good community with well-connectedness, CM++ is designed to work with any clustering algorithm and presently provides built-in support for the Leiden algorithm (optimizing either the Constant Potts Model or modularity) (Traag et al., 2019), Iterative K-core Clustering (IKC) (Wedell et al., 2022), and Infomap (Rosvall & Bergstrom, 2008).
- **Dynamic Thresholding:** In order to allow the enforcement of connectivity to be flexible, connectivity thresholds can be constants, or functions of the number of nodes in the cluster, or the minimum node degree of the cluster.
- **Multi-processing:** For better performance, users can specify a larger number of cores to process clusters concurrently.

CM++ Pipeline: A Flexible and User-Friendly Community Detection Pipeline

Key Features:

- **Graph Cleaning:** Removal of parallel and duplicate edges as well as self loops.
- **Community Detection:** Clusters an input network with one of Leiden, IKC, and InfoMap.
- **Cluster Filtration:** A pre-processing stage that allows users to filter out clusters that are trees or have size below a given threshold.
- **Community Statistics Reporting:** Generates node and edge count, modularity score, Constant Potts Model score, conductance, and edge-connectivity at multiple stages.
- **Extensibility:** Developers can remove stages and design new ones.

54 Limitations

55 The current version of CM++ offers a limited range of built-in clustering options, but is
56 designed to simplify extension by other developers. With IKC, CM++ has failed to complete
57 on clusters on the order of a million nodes due to very high memory usage. CM++ is limited
58 to community detection algorithms that yield disjoint communities, so algorithms that yield
59 overlapping communities are not supported by CM++.

60 Parallel Strategy

61 In its current form, CM++ distributes the input clusters roughly equally across available cores.
62 Each core runs an instance of CM++, and outputs are aggregated at the end. This strategy
63 may suffer from load balancing issues if there are large outliers in cluster size. A version is
64 being developed that uses a shared memory queue of clusters that each core can fetch from.

65 Conclusions

66 CM++ offers performant improvements over its predecessor CM. The accompanying pipeline
67 provides additional functionality and is customizable, allowing users to re-order modules and
68 add custom modules.

69 Acknowledgements

70 This work was supported in part by the Insper-Illinois Collaboration and by Oracle Research
71 Awards to Tandy Warnow and George Chacko. The authors thank Nathan Bryans and Christine
72 Ballard from Oracle Research for their assistance with the Oracle Cloud Infrastructure.

73 References

- 74 Bonchi, F., García-Soriano, D., Miyauchi, A., & Tsourakakis, C. E. (2021). Finding densest
75 k-connected subgraphs. *Discrete Applied Mathematics*, 305, 34–47. <https://doi.org/10.1016/j.dam.2021.08.032>
- 76
- 77 Coscia, M., Giannotti, F., & Pedreschi, D. (2011). A classification for community discovery
78 methods in complex networks. *Statistical Analysis and Data Mining: The ASA Data
79 Science Journal*, 4(5), 512–546. <https://doi.org/10.1002/sam.10133>
- 80 Dey, A. K., Tian, Y., & Gel, Y. R. (2022). Community detection in complex networks: From
81 statistical foundations to data science applications. *WIREs Computational Statistics*, 14(2),
82 e1566. <https://doi.org/10.1002/wics.1566>
- 83 Fortunato, S., & Newman, M. E. J. (2022). 20 years of network community detection. *Nature
84 Physics*, 18(8), 848–850. <https://doi.org/10.1038/s41567-022-01716-7>
- 85 Haggerty, L. S., Jachiet, P.-A., Hanage, W. P., Fitzpatrick, D. A., Lopez, P., O'Connell, M. J.,
86 Pisani, D., Wilkinson, M., Baptiste, E., & McInerney, J. O. (2013). A pluralistic account
87 of homology: Adapting the models to the data. *Molecular Biology and Evolution*, 31(3),
88 501–516. <https://doi.org/10.1093/molbev/mst228>
- 89 Henzinger, M., Noe, A., Schulz, C., & Strash, D. (2018). Practical minimum cut algorithms.
90 *ACM Journal of Experimental Algorithmics*, 23.
- 91 Karatas, A., & Sahin, S. (2018, December). Application areas of community detection: A
92 review. *2018 International Congress on Big Data, Deep Learning and Fighting Cyber
93 Terrorism (IBIGDELFT)*. <https://doi.org/10.1109/ibigdelft.2018.8625349>

- 94 Leskovec, J., & Krevl, A. (2014). *SNAP Datasets: Stanford large network dataset collection*.
95 <http://snap.stanford.edu/data>.
- 96 Park, M., Tabatabaee, Y., Ramavarapu, V., Liu, B., Pailodi, V. K., Ramachandran, R.,
97 Korobskiy, D., Ayres, F., Chacko, G., & Warnow, T. (2023). Well-connected communities
98 in real-world and synthetic networks. *To Appear, Proceedings of COMPLEX Networks*
99 *2023*. <https://arxiv.org/abs/2303.02813>
- 100 Ramavarapu, V., Kamath, V., Park, M., Ayres, F., & Chacko, G. (2023). *Connectivity modifier*
101 *pipeline*. https://github.com/illinois-or-research-analytics/cm_pipeline.
- 102 Rosvall, M., & Bergstrom, C. T. (2008). Maps of random walks on complex networks reveal
103 community structure. *Proceedings of the National Academy of Sciences*, *105*(4), 1118–1123.
104 <https://doi.org/10.1073/pnas.0706851105>
- 105 Traag, V. A., Waltman, L., & van Eck, N. J. (2019). From Louvain to Leiden: Guaranteeing
106 well-connected communities. *Scientific Reports*, *9*(1), 5233. <https://doi.org/10.1038/s41598-019-41695-z>
107
- 108 Waltman, L., & van Eck, N. J. (2012). A new methodology for constructing a publication-level
109 classification system of science. *Journal of the American Society for Information Science*
110 *and Technology*, *63*(12), 2378–2392. <https://doi.org/10.1002/asi.22748>
- 111 Wedell, E., Park, M., Korobskiy, D., Warnow, T., & Chacko, G. (2022). Center–periphery
112 structure in research communities. *Quantitative Science Studies*, *3*(1), 289–314. https://doi.org/10.1162/qss_a_00184
113

DRAFT