

15618: Spring 2023 - Final Project Proposal

Wei-Lun Chiu (weilunc) (weilunc@andrew.cmu.edu)
Jhao-Ting Chen (jhaoting) (jhaoting@andrew.cmu.edu)

March 29, 2023

1 Title

Supercomputer or Cloud Computing: Large-scale, High-computation Data Clustering Tasks with PSC/GHC+PyOMP and AWS+PySpark

2 URL

<https://jtchen0528.github.io/cmu-15618-project/>

3 Summary

We aim to analyze the performances and costs of executing data clustering algorithms, between executing on a supercomputer (PSC/GHC with PyOMP) and on a set of workers with comparable resource constraints (AWS instances with PySpark). We aim to analyze the serial and parallelized implementation, on performance, data imbalances, and different algorithms.

4 Background

4.1 Data Clustering

The dataset size for machine learning tasks grows larger by days, and self-supervised tasks become more popular than before. The performance and effectiveness for data clustering tasks such as dataset reduction[1] or pseudo-labeling[2] hence raise more research awareness. For our project, we choose to analyze the k-means, DBSCAN(Density-Based Spatial Clustering of Applications with Noise)[3], and different parallel approaches (grid-based, partition-based, and task-based). We will implement serially and in parallel, on multi-core machines and multi-node clusters.

4.2 Cloud Computing

Cloud computing provides a cost-efficient alternative to maintaining expensive high-performance computing infrastructure. With on-demand computing resources, scalability, and pay-as-you-go services, one may access computing power as needed without the need for investing in specialized hardware. Cloud computing is especially beneficial for those who don't have access to resources like PSC or multi-core machines at home. By executing large-scale, high-computing tasks on several cloud service instances in parallel with distributed technologies, organizations can achieve comparable performance and accurate results while also significantly reducing costs.

4.3 PyOMP

PyOMP [4] caught our eyes when we're researching through our project. PyOMP, released by *Intel Corp.*, provides an easy-to-use API that allows developers to parallelize their Python code with minimal modifications, by adding OpenMP directives to their code. Instead of embedded parallelism inside Numpy [5], which has overhead when creating and destroying threads, PyOMP library was

implemented in Numba. Numba utilized JIT compiler to generate optimized machine code at run-time, provides significant speedups for computationally intensive tasks. PyOMP then offers users an OpenMP-similar syntax for accessibility. Programs written in C with OpenMP performs only 2.8% faster than programs in PyOMP.

4.4 PySpark

PySpark [6] is a Python API for Apache Spark [7], a powerful open-source distributed computing system used for big data processing and analytics. Spark is designed to work with large-scale data processing tasks that require high-speed data processing and distributed computing capabilities. PySpark provides even more easy-to-use interface for users to handle data analytic tasks in Python. For our project, we aim to implement our program in PySpark, so that data can be clustered in parallel on multiple cloud instances.

5 The Challenge

The challenges can be described in the following bullet-points:

- (a) **New parallelization libraries** Familiarize with PyOMP and PySpark, implement parallelized clustering algorithms such as k-means, DBSCAN.
- (b) **Several parallelization algorithms** Implement different parallelized methods of the clustering algorithms, such as grid-based, partition-based, and task-based. Inspect workload imbalances for each processor or cloud instance. With PyOMP and PySpark.
- (c) **Explore communications in the program** The data were divided for each processors/instances. The communication per computation for the same program may differ when the number of workers scales up. Experiments are required for proper explanation.
- (d) **Cloud service set-up** Configure instances on cloud services with similar hardware that could be compared with PSC or GHC. Configure PySpark environments on master instance and slaves.
- (e) **Scalability on processors** Implement and analyze the scalability of the processors executed, within or across machines.
- (f) **Constraints in supercomputer and cloud instances** Discover execution constraints in single high-compute machine or several cloud instances, in memory, disk, network bandwidth, etc.

6 Resources

We will derive our starter serial k-means clustering program from a naive implementation [8], and implementation from scikit-learn [9]. We then plan to implement PyOMP version leveraging parallelized k-means [10]. We will refer to [11] and [11] as our DBSCAN implementation. For PyOMP implementation, we will refer to their official page [12] and their paper [4]. For PySpark implementation, we will refer to their document [13].

7 Goals and Deliverables

Our step-by-step goals are listed below:

- (a) Implement k-means and DBSCAN serial algorithm in python, and parallelized version with PyOMP. Sample datasets would be MNIST [14], CIFAR-10 [15], and CIFAR-100.
- (b) Conduct experiments on GHC and PSC with different number of processors, analyze the speedups. Expecting speedups to be in proportion to the number of processors executed on.
- (c) (Optional) Implement different k-means strategies such as grid-based, partition-based, or task-based for performance comparison and workload imbalances inspection.

- (d) (Optional) Implement serial and parallelized program in C. Analyze the performance between OpenMP and PyOMP.
- (e) Implement k-means and DBSCAN distributed/parallelized algorithm with PySpark, on set of instances with similar hardware constraints added up as GHC/PSC.
- (f) Conduct experiments on performance difference, data communication overheads, and results between PySpark and PyOMP, with different number of workers.
- (g) (Optional) Scale up the problem with larger datasets. Repeat the experiments.
- (h) Conclude the performance, costs, complexity trade-offs between renting/maintaining a super-computer/cloud services.

For our final presentation, we will demonstrate the clustered results, the performance comparison, workloads across processors/instances, costs, ... etc for every implementations we have.

8 Platform Choice

We are running our code in Unix environment. We code in python for consistency when comparing with PySpark and PyOMP, also python is less complex for data analysis tasks. We will run our code with the following processors:

- (a) Intel Core i7-9700 CPU on GHC machines.
- (b) Two AMD EPYC 7742 CPUs on PSC regular memory machines.
- (c) Sets of AWS spot instances with single-node performance similar to a core in the above two CPUs.

9 Schedule

For the following 4 weeks in April, we will finish (a) in the first week (due to final exam on April 5th),

Week 1

- Establish development environment and tools (PyOMP and AWS).
- Create serial implementations of k-means and DBSCAN.
- Develop partition-based parallel implementations of k-means and DBSCAN using PyOMP.

Week 2

- Create different parallel approaches (grid-based and task-based) using PyOMP.
- Analyze the PyOMP version and compare it with the serial version.
- Implement clustering algorithms on PySpark.
- Prepare a Milestone report.

Week 3

- Review feedback and modify the project plan as necessary.
- Scale up the clustering algorithm to handle large datasets with PySpark.
- Analyze the experiment results and begin writing the final report.
- Draft the final report.

Week 4

- Finish any remaining experiments.
- Allow for extra time to make any final edits and revisions.
- Finalize the project report and poster.

References

- [1] Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S. Morcos, “Semdedup: Data-efficient learning at web-scale through semantic deduplication,” 2023.
- [2] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” 2021.
- [3] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*. 1996, KDD’96, p. 226–231, AAAI Press.
- [4] Todd Anderson and Tim Mattson, “Multithreaded parallel Python through OpenMP support in Numba,” in *Proceedings of the 20th Python in Science Conference*, Meghann Agarwal, Chris Calloway, Dillon Niederhut, and David Shupe, Eds., 2021, pp. 140 – 147.
- [5] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant, “Array programming with NumPy,” *Nature*, vol. 585, no. 7825, pp. 357–362, Sept. 2020.
- [6] “Pyspark,” <https://github.com/apache/spark/tree/master/python>, 2017.
- [7] Matei Zaharia, Reynold S Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shivaram Venkataraman, Michael J Franklin, et al., “Apache spark: a unified engine for big data processing,” *Communications of the ACM*, vol. 59, no. 11, pp. 56–65, 2016.
- [8] “corvasto/simple-k-means-clustering-python,” <https://github.com/corvasto/Simple-k-Means-Clustering-Python>.
- [9] “scikit-learn/_kmeans.py,” https://github.com/scikit-learn/scikit-learn/blob/9aaed4987/sklearn/cluster/_kmeans.py#L1161.
- [10] “Christineharris/parallel-k-means-clustering,” <https://github.com/ChristineHarris/Parallel-K-Means-Clustering>.
- [11] “wangyiqu/dbscan-python,” <https://github.com/wangyiqu/dbscan-python>.
- [12] “Python-for-hpc/pyomp,” <https://github.com/Python-for-HPC/pyomp>.
- [13] “Spark core - pyspark 3.3.2 documentation,” <https://spark.apache.org/docs/latest/api/python/reference/pyspark.html>.
- [14] Li Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [15] Alex Krizhevsky, “Learning multiple layers of features from tiny images,” 2009.