

Performance comparison of Dask and Apache Spark on HPC systems

Mathieu Dugré, Valérie Hayot-Sasson, Tristan Glatard
Department of Computer Science and Software Engineering
Concordia University, Montréal, Québec, Canada
{mathieu.dugre, valerie.hayot-sasson, tristan.glatard}@concordia.ca

Abstract—

*Index Terms—*Performance, Big Data, Dask, Spark, Neuroimaging

From Mathieu: Discuss which keyword to use for the paper.

I. INTRODUCTION

II. BACKGROUND

A. Lustre

B. Dask

Dask is a Python-based Big Data engine with growing popularity in the scientific Python ecosystem. Dask was designed with data locality and in-memory computing in mind, to mitigate the data transfer bottleneck in Big Data workflows. Data locality, popularized by Map-Reduce [1], schedules tasks where the data reside. In-memory computing minimizes the overhead of transferring data to disk by keeping data in memory when possible. Dask uses lazy evaluation to reduce unnecessary communication and computation. The engine builds a dynamic graph before execution, allowing it to determine which task to compute. Dask workflows can further reduce data transfer by leveraging multithreading whenever Python's GIL does not restrict it. Fault-tolerance is achieved by recording data lineage: the sequence of operations used to modify the initial data.

Dask offers five data structures: [Array](#), [Bag](#), [DataFrame](#), [Delayed](#), and [Futures](#). Arrays offer a clone of NumPy API for distributed processing of large arrays. Bags are a distributed collection of Python object that offers a programming abstraction similar to [PyToolz](#). Dataframes are a parallel composition of [Pandas](#) Dataframes used to process a large amount of tabular data. Dask Delayed offers an API for distributing arbitrary functions that do not fit in the above frameworks. Lastly, Dask Futures can also execute arbitrary functions; however, it launches computation immediately rather than lazily. Dask modularity allows users to install only required components making it lightweight.

In Dask, a scheduler decides where and when to execute tasks using the Dask graph. API operations generate multiple fine-coarse tasks in the computation graph, allowing a more straightforward representation of complex algorithms.

The Dask engine is compatible with multiple distributed schedulers, including YARN and Mesos. Dask also provides its own *Dask Distributed scheduler*. We chose to use Dask Distributed scheduler to keep the environment balanced between the engines.

In the Dask Distributed scheduler, a *dask-scheduler* process administrates the resource provided by *dask-workers* in the cluster. The scheduler receives jobs from clients and assigns tasks to available workers. Task scheduler uses a LIFO (Last-In-First-Out) job scheduling policy. That is an utter process branch of the Dask graph before proceeding with the next one.

Dask offers multiple ways to deploy a cluster, including, but not limited to, SSH configs, Kubernetes, SLURM, PBS. For our experiments, we used the [Dask SLURM cluster](#) API.

C. Apache Spark

Apache Spark is a widely-used general-purpose Big Data engine. Like Dask, it aims at reducing data transfer costs by incorporating data locality, in-memory computing, and lazy evaluation.

Spark offers three options to schedule jobs: Spark Standalone, Mesos, and YARN. Spark Standalone is a simple built-in scheduler. YARN is mainly used to schedule Hadoop-based workflows, while Mesos can be used for various workflows. We limit our focus to Spark Standalone scheduler, as researchers are likely to execute their workflows in an HPC environment where, usually, neither YARN nor Mesos is available.

In the Spark Standalone scheduler, a *leader* [From Tristan: I am all in favor of inclusive terminology, but if Spark didn't update their wording we shouldn't do it for them](#) coordinates the resource provisioned by *workers* in the cluster. A *driver* process receives jobs from clients and requests workers from the leader. Jobs are divided into stages to be executed onto workers. Each operation in a stage is represented by a high-level task in the computation graph. Like Dask, Spark Standalone scheduler uses a LIFO policy to schedule tasks. Spark Standalone has two execution modes: (1) the client mode, where the driver process runs in a dedicated process, and (2) the cluster mode, where the driver runs within a worker

process. Our experiments use the client mode since cluster mode is not available in PySpark.

Spark’s primary data structure is Resilient Distributed Dataset (RDD) [2], a fault-tolerant, parallel collection of data elements. RDDs are the basis of the other Spark data structure: Datasets and DataFrames. Datasets are similar to RDD but benefit additional performance by leveraging the Spark SQL’s optimized execution engine. The DataFrames are Datasets organized into named-columns and are used to process tabular data. While the DataFrame API is available in all supported languages, Datasets are limited to Scala and Java.

Python is a standard programming language in the scientific community, offering numerous data processing libraries. While serialization from Python to Java, an operation required when using Spark’s Python API, creates overhead, we found it minimal [3]. We focus on PySpark API to have a more balanced environment between the different engines and for its suitability to neuroimaging.

III. METHODS

A. Infrastructure

For our experiments we used the SlashBin cluster at Concordia University. Each compute node has 2 16-cores CPUs, 256 GB 2666 MHz memory in 8-channel, and 2.88 TB SSDs of mounted storage. The nodes are interconnected with a dedicated 10 Gbit/s Ethernet connection. Centos 8 with Linux kernel *4.18.0-240.1.1.el8_lustre.x86_64* installed on the compute nodes.

Both Spark and Dask are configured to have 8 worker processes each with 8 threads. Each worker is allocated 31.5 GB of memory. A new cluster is spinned up and teared down for each experiment.

Dask ?? and Spark ?? **From Mathieu: Add version when after locking it for experiments.** was used for our experiments.

B. Dataset

We used BigBrain [4], a 3-D image of the human brain with voxel intensity ranging between 0 and 65,535. We converted the blocks into the NIfTI format, a popular format in neuroimaging. We left the NIfTI blocks uncompressed, resulting in a total data size of 648 GB. To evaluate the effect of block size, we resplit these blocks into 1000, 2500, and 5000 blocks of 648 MB, 259.2 MB, and 129.6 MB, respectively. [sam](#) library was used to resplit the image.

We also used the dataset provided by the Consortium for Reliability and Reproducibility (CoRR) [5], freely available on [datalad](#). The entire dataset is 408.4 GB, containing anatomical, diffusion and functional images of 1,397 subjects acquired in 29 sites. We used all 3,491 anatomical images, representing 39 GB overall (11.17 MB per image on average).

C. Applications

1) *Increment*: We adapted the increment application used in [6]. This synthetic application reads blocks of the BigBrain from Lustre and simulates computation by sleeping for a specified period. To simulate intermediate results, we repeat

the sleep process for a configurable amount of time. We prevent data caching of the blocks by incrementing their voxels value by one after each sleep operation. Finally, we write the resulting NIfTI image back to Lustre. This application allows us to study the engines when their inputs are processed independently. The map-only scenario of this application mimics the processing of multiple independent subjects in parallel.

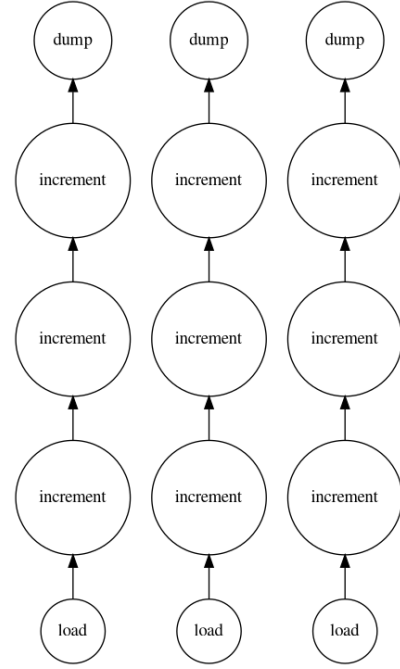


Fig. 1: Task graph for Incrementation with 3 iterations and 3 BigBrain blocks.

2) *Multi-Increment*: Our second application is an adaptation of the increment application. A significant difference is that, at each iteration, it uses a random BigBrain block as the increment value. This change allows the multi-increment application to have inter-worker communication while remaining simple.

3) *Histogram*: As our third application, we calculate the histogram of the BigBrain image. The application reads the BigBrain blocks from Lustre, calculates each intensity’s frequency, and then writes the aggregated result back on Lustre. This map-reduce application has a very high read overwrite ratio. Moreover, this application requires shuffling, albeit of a limited amount of data. The amount of inter-worker communication is in-between the increment and multi-increment applications.

4) *Kmeans*: For our fourth application, we apply Kmeans clustering to the voxel intensities of the BigBrain image. We set the number of clusters to 3, to segment the white and grey matter and the noise. The application starts by reading the image blocks, combining all voxels in a 1-D array, and choosing initial centroids using the min, max, and intermediate values. It assigns each voxel to the centroids it is the closest and updates each centroid by computing the average of the

voxels associated with it. It repeats the assignment and update steps for a configurable amount of time. Finally, the voxels of the image blocks are classified and written back to the file system. Updating the centroids involves substantial data communication between the workers.

For this application, the Spark and Dask implementations

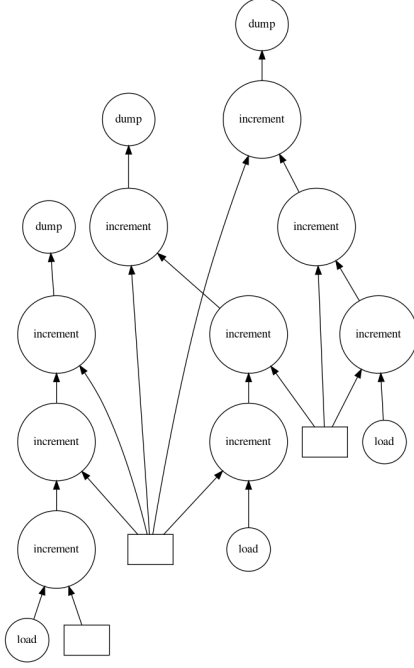


Fig. 2: Task graph for Multi-Incrementation with 3 iterations and 3 BigBrain blocks.

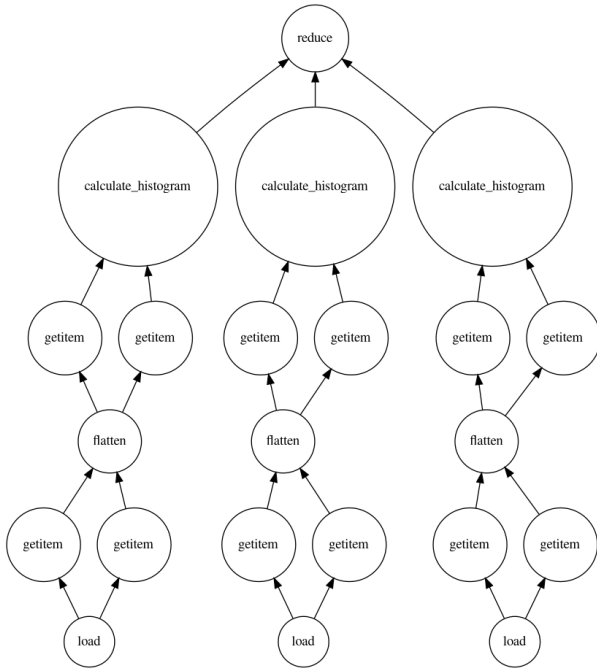


Fig. 3: Task graph for Histogram with 3 BigBrain blocks.

differ slightly, to take advantage of the best-suited API from both engines. The Spark implementation uses the Map-Reduce paradigm, while the Dask one uses array programming.

5) *BIDS App example*: Our fifth application is BIDS App example: a neuroimaging pipeline to measure the brain volume from MRIs. For this application, we use the CoRR dataset. The application extracts the brain volume of each participant, then computes the average for each group of participants. Unlike the other applications, BIDS App example is a command-line executed in a Docker image (bids/example on DockerHub). We converted the Docker image to a Singularity image for use in HPC environments, [From Mathieu: Cite paper on reason why this is done](#), using [docker2singularity](#)

6) *BIDS App MRIQC*: Our last application is BIDS App MRIQC: a neuroimaging pipeline to perform quality control of a brain image. For this application, we use the CoRR dataset. [From Mathieu: Validate if this one was used after running the experiments. Maybe another dataset will be required; e.g. ADHD200](#) This application verifies the quality of images on a per-subject basis. Like the BIDS App example, this application is run using a command-line tool. We use the same method to convert the Docker image to a Singularity image.

D. Experiments

Table I the four parameters that are varied throughout the experiments. We varied (1) the number of workers to assess the scalability of the scheduler for the engine, (2) the BigBrain block size in Increment, Multi-Increment, and Histogram to measure the effect of the different I/O pattern and parallelization degrees, (3) the number of iterations to evaluate the effect of number of task, and (4) the sleep delay to study the effect of task duration. It should note that increasing the number of iterations for a given sleep delay also increases the total compute time of an application.

To avoid potential external bias such as caching, background process and network load, we ran the applications in randomized order and cleared the page cache in between every experiments. Each benchmark was run ten times.

For each run, we measure the makespan of the application as well as the cumulative time spent in the different functions for read, processing, and writing data. The overhead calculation for each CPU thread is the end time of the last processed task minus the total runtime of the tasks ran for this thread. Summing those results gives the total overhead for the application.

[From Mathieu: Decide which row to keep for \(# fo nodes / # of workers\) and \(Block resplit / Block size\)](#)

IV. RESULTS

V. DISCUSSION

VI. CONCLUSION

ACKNOWLEDGMENT

REFERENCES

- [1] J. Dean and S. Ghemawat, "MapReduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.

TABLE I: Parameters for the experiments

	Increment	Multi-Increment	Kmeans	Histogram	BIDS Example	BIDS MRIQC
# of Nodes	2, 4, 8					
# of Workers	16, 32 , 64					
Block resplit	1000, 2500, 5000				n/a	
Block Size [MB]	648, 259.2, 129.6				n/a	
# of Iterations	1, 8, 64			n/a		
Sleep Delay [s]	0.25, 1, 4, 16		n/a			

- [2] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, “Resilient Distributed Datasets: A Fault-tolerant Abstraction for In-memory Cluster Computing,” in *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation*, ser. NSDI’12. Berkeley, CA, USA: USENIX Association, 2012, pp. 2–2.
- [3] M. Dugr, V. Hayot-Sasson, and T. Glatard, “A performance comparison of dask and apache spark for data-intensive neuroimaging pipelines,” in *2019 IEEE/ACM Workflows in Support of Large-Scale Science (WORKS)*, 2019, pp. 40–49.
- [4] K. Amunts, C. Lepage, L. Borgeat, H. Mohlberg, T. Dickscheid, M.-É. Rousseau, S. Bludau, P.-L. Bazin, L. B. Lewis, A.-M. Oros-Peusquens, N. J. Shah, T. Lippert, K. Zilles, and A. C. Evans, “BigBrain: An Ultrahigh-Resolution 3D Human Brain Model,” *Science*, vol. 340, no. 6139, pp. 1472–1475, 2013.
- [5] X.-N. Zuo, J. S. Anderson, P. Bellec, R. M. Birn, B. B. Biswal, J. Blautzik, J. C. Breitner, R. L. Buckner, V. D. Calhoun, F. X. Castellanos *et al.*, “An open science resource for establishing reliability and reproducibility in functional connectomics,” *Scientific data*, vol. 1, p. 140049, 2014.
- [6] V. Hayot-Sasson, S. T. Brown, and T. Glatard, “Performance Evaluation of Big Data Processing Strategies for Neuroimaging,” in *19th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CC-Grid)*, Larnaca, Cyprus, 2019.