

OVERCOMING THE MEMORY BOUND OF BIG DATA ANALYTICS TO IMPROVE SERVER THROUGHPUT USING FAST STORAGE

Emmanouil Anagnostakis

Graduate Research Assistant

Computer Architecture and VLSI Systems Laboratory, ICS-FORTH

Heraklion, Greece

manosanag@ics.forth.gr

Abstract

Managing big data analytics i.e. Apache Spark poses challenges due to limited memory resources in data centers. The memory pressure that arises during data processing can result in low server throughput, causing delays and inefficiencies. Memory is wasted in long GC (Garbage Collection) cycles leaving no room for useful work. In this paper, we propose a novel approach to improve server throughput for managed big data analytics using smart heap offloading to fast storage devices and reducing memory pressure. Our approach involves offloading data from heap memory to fast storage devices in a smart and efficient manner, thereby freeing up heap memory and reducing memory pressure without suffering from storage latencies. We present a detailed methodology for running Apache Spark using our proposed mechanism of smart heap offloading, which significantly improves server throughput for managed big data analytics. We implement our proposed approach in Oracle's OpenJDK8 and evaluate its performance using various workloads of the Spark Bench suite on a real-world cluster. Our experimental results show that our approach significantly improves server throughput while reducing memory usage against native Spark, making it a promising solution for managed big data analytics in data centers. We also include results to show that our implementation can save money for someone if deployed in a world cluster like Amazon's EC2 that is available to everyone.

1. Introduction

With the exponential growth of data in various fields such as finance, healthcare, social media, and e-commerce, there is

a significant need for scalable and efficient big data processing frameworks. Apache Spark [?] is one such framework that has gained popularity due to its ability to handle large-scale data processing and analytics. Spark provides a distributed computing platform that can process data in parallel across multiple nodes in a cluster. However, with the increasing size and complexity of big data workloads, Spark clusters are facing significant challenges in meeting the performance and throughput requirements.

One of the main challenges in Spark clusters is the high computational and memory requirements of big data analytics workloads. These requirements can result in excessive CPU and memory usage on Spark workers, leading to performance bottlenecks and slow job completion times. To address these challenges, researchers have proposed various techniques to optimize the performance of Spark clusters, including data partitioning, caching, and resource allocation.

In this paper, we focus on the memory limit problem of servers becoming an obstacle for further throughput increase and we propose a new technique for improving the performance and job throughput of Spark clusters by moving parts of the main managed Java Heap to a fast storage device such as NVMe, thereby saving memory for other more useful tasks. Our approach leverages the capabilities of the underlying machine to create less memory-consuming computation tasks, thereby reducing the workload on the Spark workers to improve job throughput, while maintaining effective single instance performance under the colocation of multiple instances required to achieve max throughput.

Specifically, in order to achieve higher throughput and better performance for Spark, we use TeraHeap, a secondary managed

memory-mapped heap over an NVMe storage device, which is used to hold the Resilient Distributed Datasets (Spark RDDs) instead of the main managed Java Heap and completely remove any Serialization/Deserialization and Garbage Collection (GC) cost over them.

TeraHeap 1) eliminates Serialization/Deserialization overheads posed by this kind of frameworks when moving data off-heap to/from fast storage devices 2) eliminates GC pauses over the secondary heap, therefore significantly minimizing overall GC overhead. By offloading the managed Java Heap and relaxing computation-intensive tasks, we aim to reduce the workload on Spark workers, thereby improving their performance and job throughput. We also explore the trade-offs between the cost of offloading and the performance gains achieved.

The main contribution of this paper is a comprehensive evaluation of the performance and cost trade-offs of creating lightweight computation tasks in Spark clusters. We demonstrate the effectiveness of our approach using various big data analytics workloads on a real-world Spark cluster. We also compare our approach with the native Spark distribution and show that our approach can be used instead of this distribution to improve performance and server throughput.

The rest of the paper is organized as follows. In section 2 we discuss related work on Spark optimization techniques and offloading techniques. In section 3, we describe our experimental methodology in order for someone to achieve the desired performance using TeraHeap. In section 4, we present our experimental results and evaluate the performance and cost trade-offs of our approach. In section 5, we discuss future research directions. Finally we conclude the paper in section 6 with an outline of our work.

2. Related Work

Several studies have been conducted to improve the performance of big data processing systems. One approach is to utilize memory-aware task co-location to improve Spark application throughput, which has been investigated by Marco et al. in [3]. Meanwhile, in [4], Kirisame et al. proposed optimal heap limits to reduce browser memory use. Another research direction is to leverage far memory to improve job throughput, as studied by Amaro et al. in [5]. To facilitate memory offloading in datacenters, Weiner et al. presented TMO, a transparent memory offloading system in [6]. In cloud computing platforms, Sharma et al. proposed per-VM page cache partitioning to improve performance in [7]. Chen and Wang introduced Spark on Entropy, a reliable and efficient scheduler for low-latency parallel jobs in heterogeneous clouds, in [8]. Thamsen et al. developed Mary, Hugo, and Hugo*, three learning-based schedulers for distributed data-parallel processing jobs on shared clusters in [9]. Additionally, Bhimani et al. proposed a lightweight virtualization framework for accelerating big data applications on

enterprise cloud in [10], while Zhang et al. focused on understanding and improving disk-based intermediate data caching in Spark in [11]. Finally, Intasorn et al. investigated using compression tables to improve HiveQL performance with Spark in a case study on NVMe storage devices in [12].

These studies demonstrate a variety of approaches for optimizing big data processing systems, ranging from memory-aware task co-location and memory offloading to scheduler design and virtualization frameworks. The findings from these studies can provide insights and guidance for future research in the field of big data processing.

3. Experimental Methodology

To evaluate the effectiveness of our proposed approach, we conduct a set of experiments using various workloads from the Spark Bench suite on a real-world cluster. Our experimental methodology consists of several steps. First, we set up our research server, using various configurations of datacenter machines, including CPU, memory, and storage. Next, we install and configure Spark and OpenJDK8 on the cluster. We use Spark's default configuration settings for our experiments, except for the garbage collector settings, which we tune according to our proposed approach.

We select two workloads from the Spark Bench suite that represent different types of data processing, such as machine learning and graph processing, (and SQL queries ???). We run each workload using our proposed approach and compare it with the performance of the same workload using the default configuration, garbage collector tuning, and heap offloading approaches.

We use all the DRAM provided by our server leaving 8-10 GB for Operating System, while increasing the number of Spark instances that are executed. We show that by using TeraHeap, each individual instance requires less memory therefore memory becomes available for more instances to be deployed, achieving more total throughput than Native Spark in similar time windows.

To measure the performance of our approach, we use several metrics, including server throughput, memory usage, and execution time. We measure server throughput as the megabytes of the dataset processed per second, memory usage as the amount of memory used during data processing, and execution time as the time taken to complete the workload.

We repeat each experiment several times to ensure statistical significance and calculate the mean and standard deviation of the metrics for each approach. We also perform a significance test to determine if the difference in performance between our proposed approach and the other approaches is statistically significant.

In summary, our experimental methodology involves setting up a real-world cluster, selecting appropriate workloads, measuring performance using several metrics, and repeating each experiment several times to ensure statistical significance. By follow-

ing this methodology, we can evaluate the effectiveness of our proposed approach for improving server throughput for managed big data analytics.

3.1 Server Characteristics

The server used in our experiments is a high-performance machine with hardware specifications found in real-world clusters like Amazon EC2. It is equipped with 8x DDR4 32-GB 2.4 GHz 64-bit DIMMs, providing a total of 256 GB of memory. The DDR4 memory technology is known for its high bandwidth and low power consumption, making it ideal for data-intensive applications like big data analytics. The server also features 32x Intel Xeon E5-2630 2.4 GHz 64-bit CPUs, each with 512 KB L1, 2 MB L2, and 20 MB L3 (LLC) cache. The Xeon E5-2630 CPU is a high-performance processor designed for data centers, offering a high core count, high clock speed, and advanced features like hyper-threading and Turbo Boost. The large L3 cache helps reduce memory latency, enabling faster data access for CPU-bound workloads. In addition to the powerful CPUs and memory, the server also has 2x KVS NVMe storage devices. NVMe is a high-performance storage technology that uses PCIe to connect directly to the CPU, providing low latency and high throughput. The KVS (Key-Value Store) storage devices are designed for fast, random access to data, making them ideal for storing and retrieving large amounts of data in big data applications. Overall, the server's hardware specifications make it a powerful platform for conducting experiments on managed big data analytics and evaluating the performance of our proposed approach.

3.2 Native Spark Configuration

We use Spark v3.3.0 with Kryo Serializer, a state-of-the-art highly optimized S/D Library for Java that Spark recommends. We run Spark with Native OpenJDK8 as a baseline. We use the Parallel Scavenge garbage collector which is the one TeraHeap is implemented for. Parallel Scavenge is also the go-to collector for applications that need high throughput like Spark. We use an executor with eight mutator threads for each instance of Spark we deploy on our server. For Parallel Scavenge, we use 8 GC Threads for minor GC and the default single-threaded old generation GC. Table XXX summarizes the Spark configuration we use as baseline. Spark uses the MEMORY-AND-DISK storage level to place executor memory (heap) in DRAM and cache RDDs in the on-heap cache, up to 50% of the total heap size. Any remaining RDDs are serialized in the off-heap cache over an NVMe SSD. This is device is also used by Spark for shuffling. We run each instance of Spark in a cgroup (***) containing two JVM instances, one for Spark driver and one for Spark executor and all the processes needed to measure performance for this instance. Each cgroup has a limited DRAM Budget. A part of the budget is the Java Heap which, for the rest of the paper, we call H1. We do this in order to be sure that every instance of Spark

running on our server has a fair amount of total DRAM available for it to use. We choose to try two different amounts for H1, 40% and 80% of total DRAM budget. 80% is the go-to percentage of total DRAM RedHat uses in its datacenters (***). What remains is going to be used for JVM Native memory (i.e. CodeCache) and for the operating system's Page Cache.

3.3 TeraHeap

3.3.1 What is TeraHeap? TeraHeap is a high-capacity managed heap that is memory-mapped over a fast storage device (preferably block-addressable NVMe or byte-addressable NVM). The high speeds these kind of devices operate on erase any overhead caused by the use of MMIO. TeraHeap is designed as an extension of the main Java Heap. It holds specific long-lived objects that have the same lifetime span. This makes TeraHeap a GC-free heap as it can delete entire regions of objects at once without a need to scan the heap over and over again for dead objects, which would be a performance kill as it would require scans over the storage device. The two main contributions of TeraHeap are the following: 1) MMIO keeps the objects that reside in the storage device deserialized, thus eliminating the need for Serialization/Deserialization, which is the no 1 overhead when running MEMORY-AND-DISK Spark 2) As discussed, TeraHeap reduces GC overheads without wasting DRAM by avoiding scans on long-lived objects, specifically Spark RDDs.

3.3.2 Spark Configuration The configuration for TeraHeap is pretty much the same as with Native, with some necessary differences to achieve our goal. TeraHeap is mapped to a different storage device (NVMe) than that Spark is using for shuffling. We do this in order for TeraHeap to utilize its device to its fullest. MMIO allows TeraHeap Spark to run in MEMORY-ONLY storage level as Spark is unaware of using any device and the OS takes control of the I/O. We also make the same decisions for the DRAM budget trying 40% and 80% H1 of total DRAM budget. That way we have a configuration where H1 dominates PageCache and the reverse, which shows what the needs of Spark applications are, running alone or colocated with other Spark applications.

3.4 What workloads did we choose to use for our experiments and why?

For our experiments, we selected two specific workloads from the Spark Bench suite: PageRank and LinearRegression. The primary reason for selecting these workloads is that they represent different types of big data analytics tasks: PageRank is a graph-based workload, while LinearRegression is a machine learning workload. By selecting workloads from both categories, we can investigate the performance of our proposed approach across a range of big data analytics tasks. Furthermore, both

PageRank and LinearRegression are well-established workloads that are commonly used for benchmarking big data analytics systems, making them a suitable choice for our experiments. Overall, the selection of these workloads allows us to evaluate the performance of our approach in a variety of contexts and provide insights into the effectiveness of our approach for improving server throughput in managed big data analytics systems.

3.4.1 PageRank PageRank is a widely used graph-based algorithm that measures the importance of nodes in a network. It has become a popular benchmark for evaluating the performance of distributed systems, including big data analytics systems like Apache Spark. PageRank is computationally intensive and requires significant memory and I/O resources, making it a suitable workload for evaluating the performance of our proposed approach for improving server throughput. Additionally, PageRank is a common algorithm in real-world applications, such as search engines and social networks, making it relevant for practical use cases.

3.4.2 LinearRegression LinearRegression is a machine learning algorithm that is used to predict numerical values based on input data. It is a well-known and widely used algorithm in machine learning, and is commonly used for regression analysis in fields such as economics, finance, and engineering. LinearRegression is computationally intensive and requires significant memory and I/O resources, making it a suitable workload for evaluating the performance of our proposed approach for improving server throughput. Furthermore, the inclusion of a machine learning workload like LinearRegression allows us to investigate the performance of our approach across different types of big data analytics tasks and gain insights into the effectiveness of our approach for improving server throughput in a range of contexts.

3.5 Is Spark in need of Java Heap or more cache for I/O?

3.6 What kind of metrics should someone use to be accurate when measuring performance?

When measuring performance, it's important to choose metrics that provide a comprehensive view of the system's behavior. In the case of measuring the performance of Spark instances, there are several key metrics that one should consider. These include heap capacity, which is the amount of memory allocated to the Java Virtual Machine (JVM) running Spark, and total memory used by the instance, which is the actual amount of memory consumed by the Spark instance, as measured by the cgroup budget. GC time is also an important metric, as it measures the amount of time spent by the JVM garbage collector in freeing up memory. Serialization/deserialization time, measured using

a Java async profiler, is important for understanding how much time is spent in this operation, which can be a bottleneck for some workloads. Other time, which is simply the difference between total time and GC and serialization/deserialization time, can provide insight into other factors that may be affecting performance. Device traffic, measured using iostat, is important for understanding how much data is being read from and written to storage devices. CPU idle and IO wait, measured using mpstat, can help identify how much of the CPU and IO resources are being utilized. Finally, average throughput, measured using Spark Bench, is a good indicator of the overall performance of the system. Other metrics, such as the total amount of data processed and the number of minor and major garbage collections, as measured using jstat, can also provide valuable insights into system behavior. By considering a range of metrics, one can get a more accurate and comprehensive view of the performance of Spark instances.

3.7 Is cost a contributing factor to pursuing higher throughput for a server?

4. Evaluation

5. Future Work

While our proposed offloading technique shows promising results in improving job throughput for big data analytics workloads on Spark clusters, there are several avenues for future work to further improve the performance and scalability of Spark clusters.

Firstly, one potential direction for future work is to investigate the use of other types of storage mediums such as the hybrid NVM. This medium could improve the performance of Big data analytics further by combining the advantages of memory and storage.

Secondly, another area for future work is to develop techniques for dynamically adjusting the heap offloading decisions based on workload characteristics and resource availability. For example, the offloading decision can be based on the size of the input data or the availability of DRAM capacity in the cluster. Such techniques can help maximize the performance gains achieved by offloading while minimizing the cost of offloading.

Thirdly, an interesting direction for future work is to explore the use of heap offloading in environments where Spark clusters are deployed across multiple machines using RDMA to achieve communication between the different machines. This can help utilize the DRAM, CPU and storage availability in more than one machine and provide a more cost-effective solution for big data processing.

Finally, another potential area for future work is to investigate the use of heap offloading for other big data processing frameworks beyond Spark. Many other big data processing frameworks such as Apache Giraph can potentially benefit from offloading techniques to improve their performance and scalability.

ity.

Overall, there are many exciting avenues for future work in improving the performance and scalability of big data processing frameworks such as Spark. Our proposed offloading technique provides a solid foundation for future work and offers a promising approach for addressing the challenges of big data processing.

6. Conclusion

In this paper, we proposed a new technique for improving the performance and job throughput of Spark clusters by moving parts of the managed Java Heap to a secondary memory-mapped heap over fast storage devices such as NVMe. Our approach leverages the capabilities of the underlying running machine to free computation-intensive tasks running on the Spark workers from memory pressure, thereby reducing the workload on the workers and improving their performance and job throughput.

Our experimental results demonstrate the effectiveness of our approach using various big data analytics workloads on a Spark cluster. We also compare our approach with the native Spark distribution and showed that our approach can be used instead of this distribution to further improve performance.

Our work contributes to the growing body of research on improving the performance and scalability of Spark clusters for big data analytics workloads. Our approach offers a scalable solution for processing increasingly large and complex big data workloads and can be easily integrated into existing Spark clusters.

Overall, our offloading technique offers a promising approach to improving job throughput for big data analytics workloads on Spark clusters, particularly for computation-intensive tasks. With the increasing demand for efficient and scalable big data processing frameworks, our approach provides a valuable contribution to the field of big data analytics and memory management.

ACKNOWLEDGMENT

REFERENCES

- [1] Zaharia, M., Xin, R. S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M. J., Ghodsi, A., Gonzalez, J., Shenker, S., and Stoica, I., 2016. "Apache spark: A unified engine for big data processing". In Communications of the ACM, Association for Computing Machinery.
- [2] Kolokasis, I. G., Papagiannis, A., Pratikakis, P., Bilas, A., Zakkak, F., Evdrou, G., Akram, S., and Kozanitis, C., 2023. "Teraheap: Reducing memory pressure in managed big data frameworks". In ASPLOS '23, March 25–29, 2023, Vancouver, BC, Canada, Association for Computing Machinery.
- [3] Marco, V. S., Taylor, B., Porter, B., and Wang, Z., 2017. "Improving spark application throughput via memory aware task co-location: A mixture of experts approach". In Proceedings of Middleware '17, Las Vegas, NV, USA, Association for Computing Machinery.
- [4] Kirisame, M., Shenoy, P., and Panchekha, P., 2022. "Optimal heap limits for reducing browser memory use". In OOPSLA, Association for Computing Machinery.
- [5] Amaro, E., Branner-Augmon, C., Luo, Z., Ousterhout, A., Aguilera, M. K., Panda, A., Ratnasamy, S., and Shenker, S., 2020. "Can far memory improve job throughput?". In EuroSys '20, April 27–30, 2020, Heraklion, Greece, Association for Computing Machinery.
- [6] Weiner, J., Agarwal, N., Schatzberg, D., Yang, L., Wang, H., Sanouillet, B., Sharma, B., Heo, T., Jain, M., Tang, C., and Skarlatos, D., 2022. "Tmo: Transparent memory offloading in datacenters". In ASPLOS '22, February 28 – March 4, 2022, Lausanne, Switzerland, Association for Computing Machinery.
- [7] Sharma, P., Kulkarni, P., and Shenoy, P. "Per-vm page cache partitioning for cloud computing platforms". Association for Computing Machinery.
- [8] Chen, H., and Wang, F. Z., 2015. "Spark on entropy: A reliable and efficient scheduler for low-latency parallel jobs in heterogeneous cloud". In LCN 2015, Clearwater Beach, Florida, USA, Association for Computing Machinery.
- [9] Thamsen, L., Beilharz, J., Tran, V. T., Nedelkoski, S., and Kao, O., 2019. "Mary, hugo, and hugo*: Learning to schedule distributed data-parallel processing jobs on shared clusters". In Euro-Par 2019, Association for Computing Machinery.
- [10] Bhimani, J., Yang, Z., Leeser, M., and Mi, N., 2017. "Accelerating big data applications using lightweight virtualization framework on enterprise cloud". In 2017 IEEE High Performance Extreme Computing Conference (HPEC), IEEE.
- [11] Zhang, K., Tanimura, Y., Nakada, H., and Ogawa, H., 2017. "Understanding and improving disk-based intermediate data caching in spark". In 2017 IEEE International Conference on Big Data (Big Data), IEEE.
- [12] Intasorn, Y., Rattanaopas, K., and Chuchuen, Y., 2022. "Using compression tables to improve hiveql performance with spark a case study on nvme storage devices". In 2022 26th International Computer Science and Engineering Conference (ICSEC), IEEE.
- [13] Apache, 2023. "Rdd programming guide (spark 3.4.0 - 2023 update) - <https://spark.apache.org/docs/latest/rdd-programming-guide.html>".
- [14] Apache, 2023. "Building spark (spark 3.4.0 - 2023 update) - <https://spark.apache.org/docs/latest/building-spark.html>".

- [15] Apache, 2023. “Tuning spark (spark 3.4.0 - 2023 update) - <https://spark.apache.org/docs/latest/tuning.html>”.
- [16] Apache, 2023. “Spark configuration (spark 3.4.0 - 2023 update) - <https://spark.apache.org/docs/latest/configuration.html>”.
- [17] Apache, 2023. “Monitoring and instrumentation (spark 3.4.0 - 2023 update) - <https://spark.apache.org/docs/latest/monitoring.html>”.
- [18] chriswhocodes, 2023. “Vm options explorer - openjdk8 hotspot - <https://chriswhocodes.com>”.