

Partition Balancing in Distributed Nested Data-Parallel Systems

Frank Austin Nothaft

Department of Electrical Engineering and Computer
Science, University of California, Berkeley
fnothaft@eecs.berkeley.edu

Michael Linderman

Carl Icahn School of Medicine at Mount Sinai
michael.linderman@mssm.edu

Abstract

Map-reduce frameworks such as Apache Hadoop and Spark provide the abstraction of a large, flat array that is processed in parallel across many machines. While this simple programming model has enabled the broad adoption of data-parallel distributed frameworks, these systems cannot express irregular parallel computation, and their performance is impacted by data imbalance across nodes. In this paper, we present a distributed framework for implementing *nested data parallel* (NDP) computation. Unlike previous NDP systems described by Blelloch and Bergstrom et al. that relied on the use of a *flattening* transformation at compile-time, we present a cost model that is used to select between multiple partitioning strategies at run-time. Through this, we provide a user-tunable means for trading node-to-node imbalance versus communication when executing distributed NDP programs.

1. Introduction

The use of map-reduce as a flat data-parallel (FDP) programming model for distributed systems has grown rapidly since it was introduced in Google’s seminal 2004 paper [7]. The development of the open-source Apache Hadoop system enabled the use of this programming model outside of Google. Modern map-reduce systems such as Apache Spark [11] have refined the programming model further by reducing the dependency of the framework on disk via in-memory caching. While this refinement has enabled the use of map-reduce for iterative workloads, the programming model remains confined to computation that can be expressed via flat data-parallel operations. Specifically, Spark presents users with the abstraction of a resilient distributed dataset (RDD), which appears as a flat array that is distributed across compute nodes in a cluster [10].

To expand the algorithms that could be expressed as a data-parallel computation, Blelloch introduced the nested data-parallel (NDP) model [5]. In this programming model, users are provided the abstraction of an array whose elements are a nested level of arrays and data-parallel operators are applied on the nested arrays. A primary complication in the implementation of the NDP model is that the nested arrays frequently do not have uniform size. Several approaches have been suggested for balancing work in NDP programs, including the compile-time vectorization of NDP

programs for execution on single-instruction multiple-data (SIMD) machines [5, 6] and flattening for multiple-instruction multiple-data (MIMD) machines [3]. Additionally, dynamic work-stealing approaches have been implemented [4].

In this paper, we introduce a distributed programming framework for NDP algorithms. Our implementation is built on top of Apache Spark. We track the structure of the nested arrays at run-time and choose between two strategies (*uniform* and *segmented*) for partitioning data across machines based on the estimated cost of each strategy. In the segmented partitioning strategy, all values in a single nested segment are co-located on a single compute node, and most computation can proceed without communication. The uniform strategy provides perfect load balance across all nodes, but many operations will need to communicate to execute. To choose between these strategies at runtime, we provide a cost model that evaluates the performance tradeoff of communication overhead versus node-to-node imbalance. Since the partitioning strategy is chosen at runtime, we can leverage knowledge of the nested array structure.

2. Background

2.1 Data-Parallel Programming Models

2.2 Data-Parallel Distributed Computing

This work builds upon the infrastructure provided by the Apache Spark distributed data-parallel computing framework [2, 11]. Spark was designed for “cloud computing” platforms, where machines may be unreliable, and where network performance may preclude the use of traditional distributed message passing systems such as the Message Passing Interface (MPI). Unlike Apache Hadoop, where data is shuffled to/from disk between all processing stages [1, 9], Spark has an in-memory processing model. The in-memory processing model leads to large (100×) performance for iterative jobs implemented using Spark instead of Hadoop [11].

Spark’s programming model is implemented on top of the *resilient distributed dataset* (RDD) abstraction [10]. The RDD abstraction presents a view of an array which is chopped into *partitions* that are distributed over the computers within the Spark cluster. Programs enqueue data-parallel transformations on the Spark master, which are then interpreted into a directed-acyclic graph (DAG) for execution; this approach is similar to the DAG scheduling approach pioneered in Dryad [8]. Spark executes the DAG whenever a disk shuffle is required, or if a newly enqueued operation would create a cycle in the DAG (i.e., an iterative computation is scheduled).

To execute a stage, the Spark master serializes the user defined function (UDF) which is being applied, transmits this function to all worker nodes, and then applies the computation. While the general application programming interface (API) that Spark provides is data-parallel across all elements, the API transformations are implemented by applying the transformation iteratively across all ele-

ments of a partition—Spark also exposes this parallel-by-partition interface via the `mapPartitions` call. If data must be moved between partitions after an execution stage, a disk shuffle will occur. Disk shuffles are analogous to an execution stage with interleaved computation and all-to-all communication followed by a barrier before the next stage. This process is similar to how the results of map phases are moved to reducers in both MapReduce and Hadoop [7]—the main distinction is that Spark allows successive map phases, and that shuffles do not occur after all stages.

3. Implementation

3.1 Characterizing and Modeling Imbalance

3.2 Partitioning Strategies

4. Performance

5. Discussion

6. Conclusion

References

- [1] Apache Hadoop. <http://hadoop.apache.org>.
- [2] Apache Spark. <http://spark.apache.org>.
- [3] BERGSTROM, L., FLUET, M., RAINEY, M., REPPY, J., ROSEN, S., AND SHAW, A. Data-only flattening for nested data parallelism. In *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '13)* (2013), ACM, pp. 81–92.
- [4] BERGSTROM, L., RAINEY, M., REPPY, J., SHAW, A., AND FLUET, M. Lazy tree splitting. In *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming (ICFP '10)* (2010), ACM, pp. 93–104.
- [5] BLELLOCH, G. E. *Vector models for data-parallel computing*, vol. 356. MIT Press Cambridge, 1990.
- [6] BLELLOCH, G. E., AND SABOT, G. W. Compiling collection-oriented languages onto massively parallel computers. *Journal of Parallel and Distributed Computing (JPDC)* 8, 2 (1990), 119–134.
- [7] DEAN, J., AND GHEMAWAT, S. MapReduce: Simplified data processing on large clusters. In *Proceedings of the 6th Conference on Symposium on Operating Systems Design and Implementation (OSDI '04)* (2004), USENIX Association, pp. 10–10.
- [8] ISARD, M., BUDIU, M., YU, Y., BIRRELL, A., AND FETTERLY, D. Dryad: distributed data-parallel programs from sequential building blocks. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007 (EuroSys '07)* (2007), vol. 41, ACM, pp. 59–72.
- [9] SHVACHKO, K., KUANG, H., RADIA, S., AND CHANSLER, R. The Hadoop distributed file system. In *Proceedings of the IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST '10)* (2010), IEEE, pp. 1–10.
- [10] ZAHARIA, M., CHOWDHURY, M., DAS, T., DAVE, A., MA, J., MCCAULEY, M., FRANKLIN, M. J., SHENKER, S., AND STOICA, I. 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 (NSDI '12)* (2012), USENIX Association, pp. 2–2.
- [11] ZAHARIA, M., CHOWDHURY, M., FRANKLIN, M. J., SHENKER, S., AND STOICA, I. Spark: cluster computing with working sets. In *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing (HotCloud '10)* (2010), p. 10.