

Design Principles for Modular and Scalable Scientific Analysis Systems

ABSTRACT

Revolutions in data acquisition are drastically changing how science conducts experiments. For example, “next-generation” sequencing technologies have decreased the cost of sequencing a genome by 10,000×, which has driven exponential growth in the total volume of genome sequence data available. However, many traditional “scientific computing” systems are a poor fit for these analyses, as they either provide poor programming abstractions, or require too much effort to program. As a result, there has been an inefficient duplication of work in the scientific community.

In this paper, we introduce a set of principles for decomposing scientific analysis systems so that they can be implemented efficiently on top of existing systems, while providing productive programming interfaces. We motivate these principles with an example genomics pipeline which leverages open-source MapReduce and columnar storage techniques to achieve a $> 50\times$ speedup over traditional genomics systems, at half the cost.

Categories and Subject Descriptors

L.4.1 [Applied Computing]: Life and medical sciences—*Computational biology*; H.1.3.2 [Information Systems]: Data management systems—*Database management system engines, parallel and distributed DBMSs*; E.3.2 [Software and its Engineering]: Software creation and management—*Software Development Process Management*

General Terms

Design

Keywords

Analytics, MapReduce, Genomics, Scientific Computing

1. INTRODUCTION

[9]

1. Data science is a growing trend in both academia and industry
 - (a) Driven by dramatic improvements in acquisition systems (e.g., sequencing, mass spectrometry, MRI systems)
 - (b) Also driven by rise of statistical systems which are easy to use for non-experts (e.g., Scikit-learn [11], MLI [14])
 - (c) Few queries in modern data science resemble traditional scientific computing patterns; not as communication heavy, but very UDF heavy
2. Computing is becoming a dominant cost for science
 - (a) Both in terms of literal costs \rightarrow paying for compute, storage, machines [15, 12]
 - (b) And in terms of NRE effort [17]
3. The design of these scientific processing systems must confront the sources of inefficiency:
 - (a) Efficiency is both computational cost and development cost
 - (b) An efficient system should be fast *enough*
 - (c) An efficient system should be able to support the common queries we need
 - (d) An efficient system should minimize the number of wheels that are reinvented
 - (e) An efficient system should have a simple programming interface and layered design [1]
4. Network stack achieves a similar goal:
 - (a) Can swap out layers to tailor implementation
 - (b) Abstract lower levels of the stack to provide a simple programming interface
5. Also, current computer systems fail to address important characteristics of scientific workloads:
 - (a) *Huge* data sizes, e.g., TB for neuroscience [3, 4], PB for genomics; may be too large to stage locally, or have small “hot” set
 - (b) Different join patterns; need to join objects in a coordinate system
 - (c) Spatial/temporal analysis: esp. for neuroscience, and other “signal processing” sciences \rightarrow similar to stream processing, but subtly different; may need “window sweeping” function

- (d) Programming models! Need to enable:
 - Scientists to write UDFs → SQL is a bad interface
 - Scientists to do *interactive data analysis*

6. Contributions of this work:

- Provide principles for the design of scientific analysis systems
- Implemented fast genomic system
- Implemented coordinate plane joins
- Efficient lookup from block stores

2. BACKGROUND

This section will compare and contrast the various “big data” analysis systems with existing scientific systems.

1. MapReduce-based workflows

- (a) In CS, development of MapReduce → Hadoop → Spark
- (b) Equivalent systems in bioinformatics → GATK [10]
- (c) Hadoop-based genomics tools [13, 7]
- (d) Use of Spark for neuroscience

2. Database driven systems

- (a) SciDB [2]
- (b) GQL [6, 1]

3. Storage layers

- (a) CRAM [5]
- (b) YT [16]

Takeaways:

- There are significant computer system design problems in science:
 1. Compression → column stores (CRAM [5], YT [16])
 2. Performance/parallelism → MapReduce (GATK [10])
- Traditional SC/DB systems provide poor abstractions for most scientists
- As a result, there are lots of “roll your own” systems in science

3. PRINCIPLES FOR SCIENTIFIC ANALYSIS SYSTEMS

3.1 Workloads

1. Characteristics of data

- (a) Scientific data tends to be sparse
- (b) Different users want to look at different subsets of both rows and columns
- (c) Data may not always be in a single site, or stored locally

- (d) *Experimental data* is immutable.

- (e) What are access patterns?

2. Characteristics of a ideal storage system:

- (a) Efficient support for projection of different columns
- (b) Efficient support for per-record predicates
- (c) Should not relegate user to a single execution environment

3. Processing:

- (a) Workloads are highly variable by field
- (b) For genomics, workloads are trivially data-parallel
- (c) Similar for fields with heavy image processing workloads
- (d) Simulation based fields are tougher; have all-to-all computation pattern, run on supercomputer
- (e) Defer discussion to §3.3
- (f) Ideally, cross-platform.

3.2 Layering

Discussion of Figure 1.

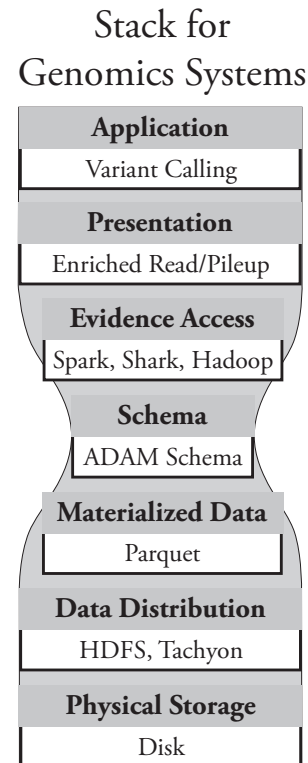


Figure 1: A Stack Model for Scientific Analysis

Specifically, we need to:

1. Show how current systems fit into the stack model, and how our proposed stack is different

2. Elucidate why it is more efficient to build systems that are decomposed as per our stack above (reference network stack and protocol interchange), see Bafna et al [1], talk about costs of programming without good stack model

3.3 Execution Platforms

TL;DW; need to have a good discussion of what applications are good for MapReduce, what are good on top of a database, what are good on an HPC farm, what should be done with an abacus, etc. Needs to be written carefully to show the virtues of the Figure 1 stack, while being frank about weaknesses.

4. IMPLEMENTATION

4.1 Genomics Pipeline

Compare/contrast to current pipelines; talk about what the stages do in reasonable but not excessive detail. Make reference to WHAM [8] to show that this is an application domain that SIGMOD has determined to be important.

4.2 Coordinate System Joins

This will be a compare/contrast discussion of the multiple join algorithms we’ve created. TBD.

4.3 Loading Remote Data

1. Data may not be kept locally
 - (a) Too much data to keep locally
 - (b) Not all data is hot
2. May push data off local disks into block store
3. Manually re-staging data has high latency cost → impacts throughput
4. What do we need to do to accommodate this?
 - (a) Efficient indexing
 - (b) Remote push-down predicate
5. Discuss S3/Parquet interaction

5. PERFORMANCE

This section will address:

- Performance of ADAM on real datasets
- Compression achieved by Parquet
- Examples extending the proposed stack to Astronomy

Experiments to run:

- General demonstration of scaling for genomics pipeline; updated experiments from TR
- Experiments on coordinate system joins; broadcast vs. partition join strategies
- Experiments showing benefit from performing remote data access without staging

6. DISCUSSION

6.1 Scientific Processing on MapReduce

Big critique from SciDB camp is that MR is an inappropriate platform for scientific computing due to lack of support for linear algebra. We need to counter this point, by allusion to performance on algorithms above, and by alluding to specialized libraries for ML & graph processing [14, 18].

Also, note that we don’t argue that MR is the correct platform for particle simulations and other traditional MPI workloads. However, MPI is the *wrong* platform for most analytics.

Also, note that most scientific workloads require applying a UDF across a large set of data. This is not inefficient to *run* on a database, but it is inefficient to *write*; SQL is a poor language for scientific/statistical computing.

6.2 Cost of Non-Commodity Systems

The advantage of the stack model we propose is that it enables the use and reuse of commodity systems, instead of re-inventing the wheel (or, inventing a *slightly* different wheel).

7. CONCLUSION

In the end, we conclude.

APPENDIX

A. REFERENCES

- [1] V. Bafna, A. Deutsch, A. Heiberg, C. Kozanitis, L. Ohno-Machado, and G. Varghese. Abstractions for genomics. *Communications of the ACM*, 56(1):83–93, 2013.
- [2] P. G. Brown. Overview of SciDB: large scale array storage, processing and analysis. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 963–968. ACM, 2010.
- [3] J. P. Cunningham. Analyzing neural data at huge scale. *Nature methods*, 11(9):911–912, 2014.
- [4] J. Freeman, N. Vladimirov, T. Kawashima, Y. Mu, N. J. Sofroniew, D. V. Bennett, J. Rosen, C.-T. Yang, L. L. Looger, and M. B. Ahrens. Mapping brain activity at scale with cluster computing. *Nature methods*, 11(9):941–950, 2014.
- [5] M. H.-Y. Fritz, R. Leinonen, G. Cochrane, and E. Birney. Efficient storage of high throughput DNA sequencing data using reference-based compression. *Genome research*, 21(5):734–740, 2011.
- [6] C. Kozanitis, A. Heiberg, G. Varghese, and V. Bafna. Using Genome Query Language to uncover genetic variation. *Bioinformatics*, 30(1):1–8, 2014.
- [7] B. Langmead, M. C. Schatz, J. Lin, M. Pop, and S. L. Salzberg. Searching for SNPs with cloud computing. *Genome Biology*, 10(11):R134, 2009.
- [8] Y. Li, A. Terrell, and J. M. Patel. WHAM: A high-throughput sequence alignment method. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data (SIGMOD ’11)*, SIGMOD ’11, pages 445–456, New York, NY, USA, 2011. ACM.
- [9] M. Massie, F. Nothaft, C. Hartl, C. Kozanitis, A. Schumacher, A. D. Joseph, and D. A. Patterson.

- ADAM: Genomics formats and processing patterns for cloud scale computing. Technical report, UCB/EECS-2013-207, EECS Department, University of California, Berkeley, 2013.
- [10] A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytzsky, K. Garimella, D. Altshuler, S. Gabriel, M. Daly, et al. The Genome Analysis Toolkit: a mapreduce framework for analyzing next-generation DNA sequencing data. *Genome research*, 20(9):1297–1303, 2010.
 - [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011.
 - [12] E. E. Schadt, M. D. Linderman, J. Sorenson, L. Lee, and G. P. Nolan. Computational solutions to large-scale data management and analysis. *Nature Reviews Genetics*, 11(9):647–657, 2010.
 - [13] M. C. Schatz. CloudBurst: highly sensitive read mapping with mapreduce. *Bioinformatics*, 25(11):1363–1369, 2009.
 - [14] E. R. Sparks, A. Talwalkar, V. Smith, J. Kottalam, X. Pan, J. Gonzalez, M. J. Franklin, M. I. Jordan, and T. Kraska. MLI: An API for distributed machine learning. In *2013 IEEE 13th International Conference on Data Mining (ICDM’ 13)*, pages 1187–1192. IEEE, 2013.
 - [15] L. D. Stein et al. The case for cloud computing in genome informatics. *Genome Biology*, 11(5):207, 2010.
 - [16] M. J. Turk, B. D. Smith, J. S. Oishi, S. Skory, S. W. Skillman, T. Abel, and M. L. Norman. yt: A multi-code analysis toolkit for astrophysical simulation data. *The Astrophysical Journal Supplement Series*, 192(1):9, 2011.
 - [17] G. Wilson, D. Aruliah, C. T. Brown, N. P. C. Hong, M. Davis, R. T. Guy, S. H. Haddock, K. D. Huff, I. M. Mitchell, M. D. Plumbley, et al. Best practices for scientific computing. *PLoS biology*, 12(1):e1001745, 2014.
 - [18] R. S. Xin, J. E. Gonzalez, M. J. Franklin, and I. Stoica. GraphX: A resilient distributed graph system on Spark. In *First International Workshop on Graph Data Management Experiences and Systems*, page 2. ACM, 2013.