

# ADAM Enables Distributed Analyses Across Large Scale Genomic Datasets

Frank Austin Nothaft<sup>1,2,\*</sup>, Arun Ahuja<sup>3</sup>, Timothy Danford<sup>1,4</sup>, Michael Heuer<sup>1</sup>, Jey Kottalam<sup>1</sup>, Matt Massie<sup>1</sup>, Audrey Musselman-Brown<sup>5</sup>, Beau Norgeot<sup>5,6</sup>, Ravi Pandya<sup>7</sup>, Justin Paschall<sup>1</sup>, Hannes Schmidt<sup>5</sup>, Eric Tu<sup>1</sup>, Ryan Williams<sup>3</sup>, Carl Yeksigian<sup>4</sup>, Michael Linderman<sup>3</sup>, Jeff Hammerbacher<sup>3</sup>, **Benedict Paten**<sup>5,\*</sup>, Uri Laserson<sup>3,9</sup>, Gaddy Getz<sup>10</sup>, David Haussler<sup>5</sup>, Anthony D. Joseph<sup>1</sup>, David A. Patterson<sup>1,2</sup>

<sup>1</sup>AMPLab, University of California, Berkeley, CA, <sup>2</sup>ASPIRE Lab, University of California, Berkeley, CA, <sup>3</sup>Icahn School of Medicine at Mount Sinai, New York, NY, <sup>4</sup>Tamr, Inc., Cambridge, MA, <sup>5</sup>Genome Informatics Lab, University of California, Santa Cruz, CA, <sup>6</sup>Pharmaceutical Science and Pharmacogenomics, University of California, San Francisco, CA, <sup>7</sup>Microsoft Research, Redmond, WA, <sup>8</sup>GenomeBridge, Cambridge, MA, <sup>9</sup>Cloudera, Inc., San Francisco, CA, <sup>10</sup>The Broad Institute of MIT and Harvard, Cambridge, MA

\* = {fnothaft@berkeley.edu, benedict@soe.ucsc.edu}



## Background

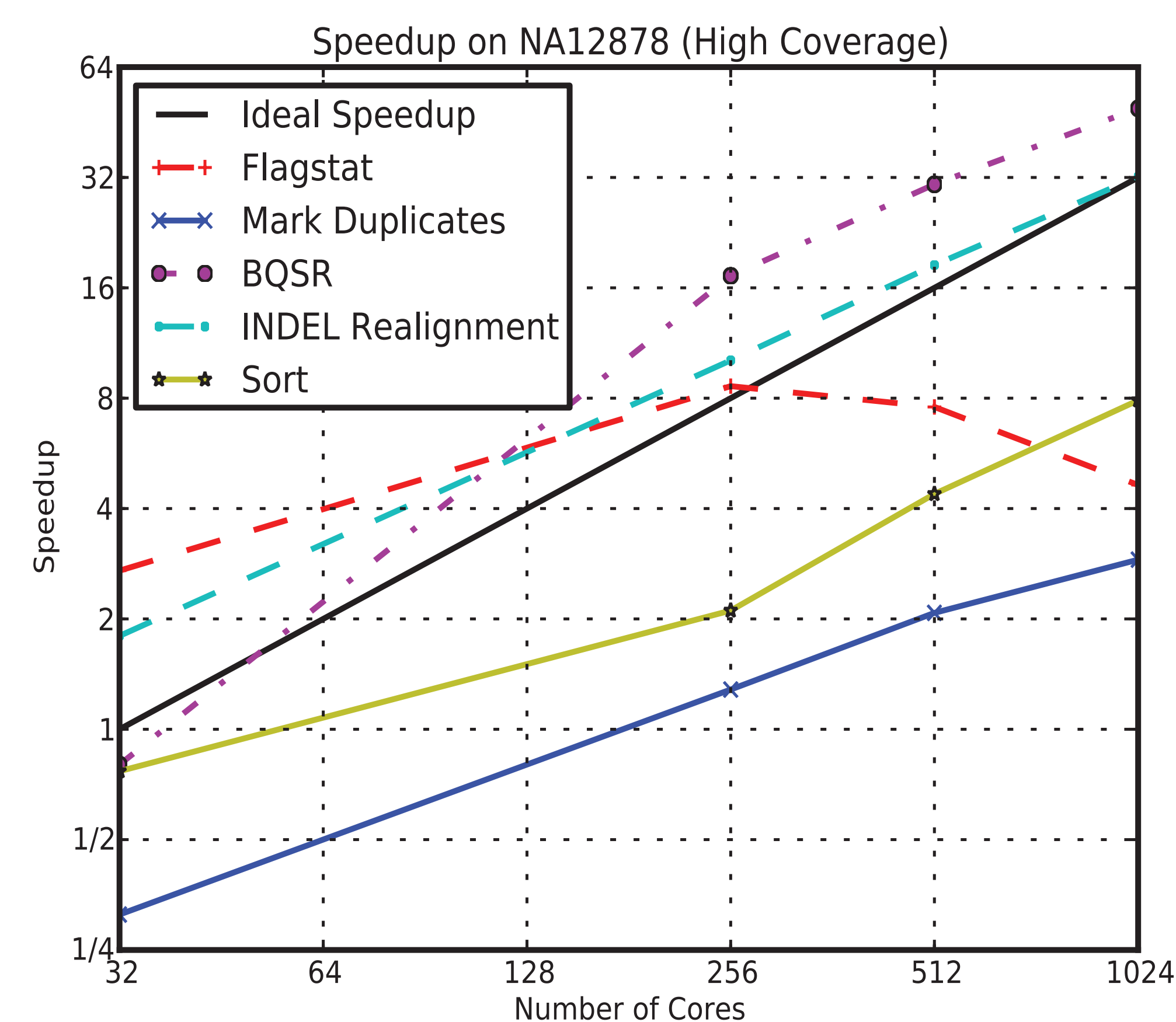
Currently, it is hard to write analyses that scan across large genomic datasets:

- High performance computing systems have poor I/O perf.
- Users frequently struggle with inconsistent file formats
- Current computational model is too low level

ADAM is a framework that allows for the efficient parallelism of genomic queries using Apache Spark. ADAM outperforms traditional tools on a single node, and can scale to hundreds of nodes.

## Performance

- Compared to GATK, Picard, samtools, and Sambamba
- Evaluated core processing steps on 234GB NA12878 dataset
- Evaluated using 1 i2.8xlarge and 32–128 r3.2xlarge instances on EC2

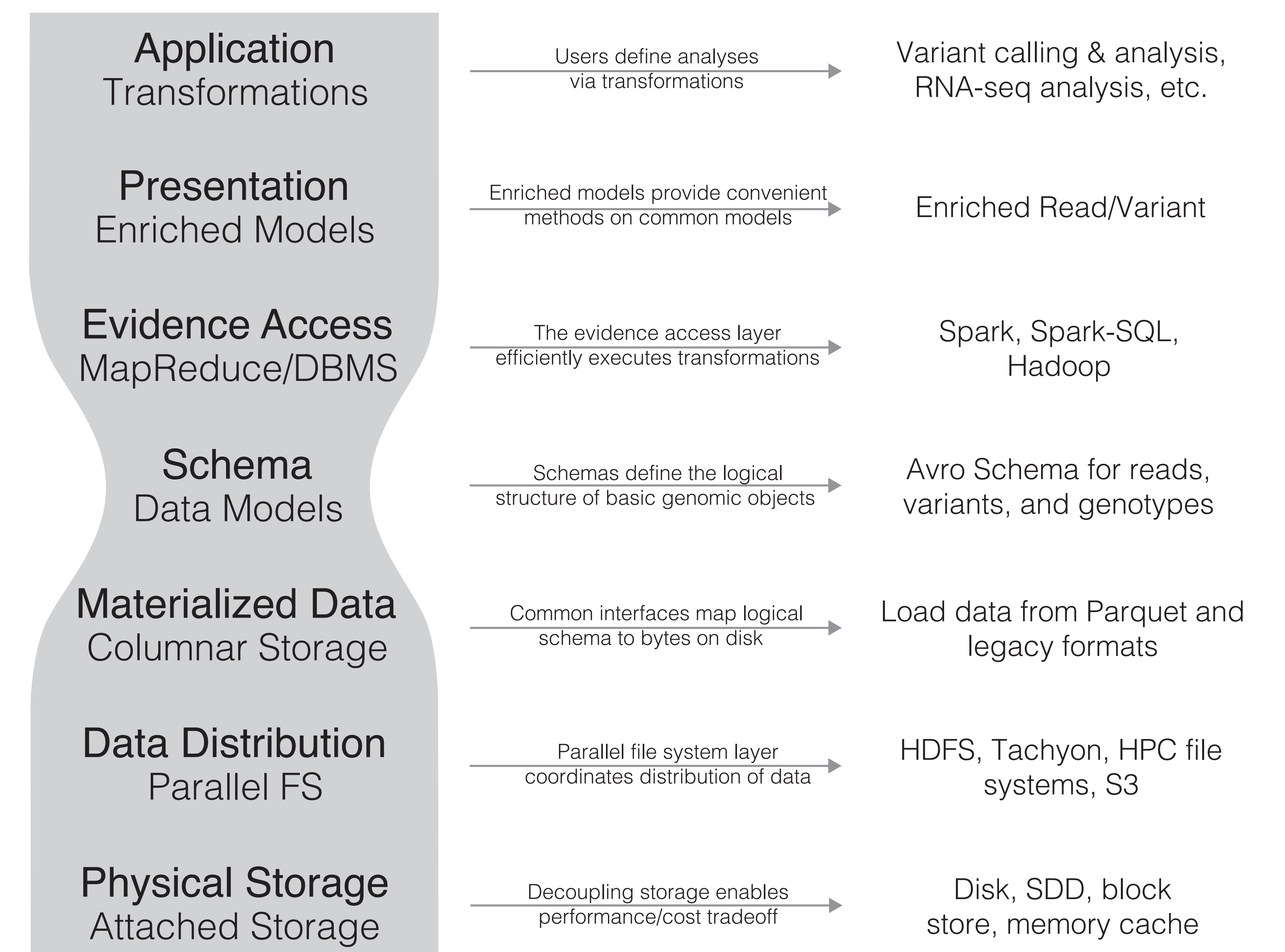


See detailed numbers in Nothaft et al., "Rethinking data-intensive science using scalable analytics systems." In Proceedings of the International Conference on Management of Data, May 2015 (SIGMOD '15).

## Architecture

ADAM uses a decomposed stack model. This has important benefits:

- Queries are programmed against a schema. The user doesn't need to know the format of data on disk, or where data is physically stored.
- ADAM builds upon Apache Spark's RDD model. RDDs are parallel arrays, and all transformations to an RDD run in parallel.
- Most systems use lower level abstractions, like an iterator over the genome. ADAM queries are written with higher level primitives: duplicate marking maps to a groupBy, finding overlapping genomic objects is implemented as an optimized parallel join.



## Accuracy Against GATK Best Practices

- We evaluated ADAM by replacing the GATK "Best Practices" pre-processing stages with an ADAM based reimplementation
  - GATK was run on a single i2.8xlarge node, ADAM was run on 16 r3.4xlarge nodes
  - The ADAM-based pipeline is  $3.55\times$  faster, and  $2\times$  cheaper
  - The two pipelines generate statistically equivalent variant calls
- During this process, we identified two bugs in the GATK/Picard. Both of these issues are caused by sort order invariants necessitated by programming at a lower level of abstraction.

