

ADAM: Genomics Formats and Processing Patterns for Cloud Scale Computing

Matt Massie¹, Frank Austin Nothaft¹, Chris Hartl^{1,2}, Christos Kozanitis¹, and David Patterson¹

¹Department of Electrical Engineering and Computer Science, University of California, Berkeley

²The Broad Institute of MIT and Harvard

Abstract

Current genomics applications are dominated by the movement of data to and from disk. This data movement pattern is a significant bottleneck that prevents these applications from scaling well to distributed computing clusters. In this report, we introduce a new set of data formats and programming patterns for genomics applications using in-memory MapReduce frameworks. These formats and frameworks improve application performance, data storage efficiency, and programmer productivity.

1 Introduction

Although the cost of data processing has not historically been an issue for genomic studies, the falling cost of genetic sequencing will soon turn computational tasks into a dominant cost [24]. The process of transforming reads from alignment to variant-calling ready reads involves several processing stages including duplicate marking, base score quality recalibration, and local realignment. Currently, these stages have involved reading a Sequence/Binary Alignment Map (SAM/BAM) file, performing transformations on the data, and writing this data back out to disk as a new SAM/BAM file [22].

Because of the amount of time these transformations spend shuffling data to and from disk, these transformations have become the bottleneck in modern variant calling pipelines. It currently takes three days to perform these transformations on a high coverage BAM file¹. By rewriting these applications to make use of modern in-memory MapReduce frame-

works like Apache Spark [39], these transformations can now be completed in FIXME² hours.

In this paper, we introduce ADAM, which is a programming framework and a set of data formats for cloud scale genomic processing. These frameworks and formats scale efficiently to modern cloud computing performance, which allows us to parallelize read translation steps that are required between alignment and variant calling. In this paper, we provide an introduction to the data formats (§2) and pipelines used to process genomes (§3), introduce the ADAM formats and data access application programming interfaces (APIs) (§5) and programming model (§6). Finally, we review the performance and compression gains we achieve (§7), and outline future enhancements to ADAM that we are working on (§8).

Performance teasers.

2 Current Genomics Storage Standards

The current de facto standard for the storage and processing of genomics data in read format is BAM. BAM is a binary file format that implements the SAM standard for read storage [22]. BAM files can be operated on in several languages, using either the SAMtools ([22], C++), Picard ([5], Java), and Hadoop-BAM ([25], Hadoop MapReduce through Java) APIs. BAM provides efficient access and compression over the SAM file format, as its binary encoding reduces several fields into a single byte, and eliminates text processing on load. However, the file format has been criticized as it is difficult to process — the three main APIs that implement the format each note that they do not fully implement the format due to its com-

¹A 250GB BAM file at 30× coverage. See §7.2 for a longer discussion

²FIXME

plexity. Additionally, the file format is difficult to use in multiprocessing environments due to its use of a centralized header; the Hadoop-BAM implementation notes that its scalability is limited to distributed processing environments of less than 8 machines.

In response to the growing size of sequencing files³, a variety of compression methods have come to light [20, 14, 32, 26, 7, 12, 36, 18, 17]. SlimGene [20], cSRA [33], and CRAM [14] use reference based compression techniques to losslessly represent reads while they advocate in favor of lossy quality value representations. The former two use lower quantization levels to represent quality values and CRAM uses user defined budgets to store only fractions of a quality string. In the same spirit, Illumina presented recently a systematic approach of quality score removal in [17] which safely ignores quality scores from predictable nucleotides; these are bases that always appear after certain words. It is also worth mentioning that the standard configurations of cSRA and CRAM discard the optional fields of the BAM format and also simplify the QNAME field.

3 Genomic Processing Pipelines

After sequencing and alignment, there are a few common steps in modern genomic processing pipelines for producing variant call ready reads. These steps minimize the amount of erroneous data in the input read set by eliminating duplicate data, verifying the alignment of short inserts/deletions (indels), and calibrating the quality scores assigned to bases (base quality score recalibration, BQSR). The typical pipeline for variant calling is illustrated in figure 1.

Traditionally, bioinformaticians have focused on improving the accuracy of the algorithms used for alignment and variant calling. There is obvious merit to this approach — these two steps are the dominant contributors to the accuracy of the variant calling pipeline. However, the intermediate read processing stages are responsible for the majority of execution time. A breakdown of stage execution time is shown in table 1 for the version 2.7 of the Genome Analysis Toolkit (GATK), a popular variant calling pipeline [23]. The numbers in the table are derived from running on the NA12878 high-coverage human genome.

To provide the readers with a background about how the stages of this pipeline work, we will discuss

³High coverage BAM files can be approximately 250 GB for a human genome.

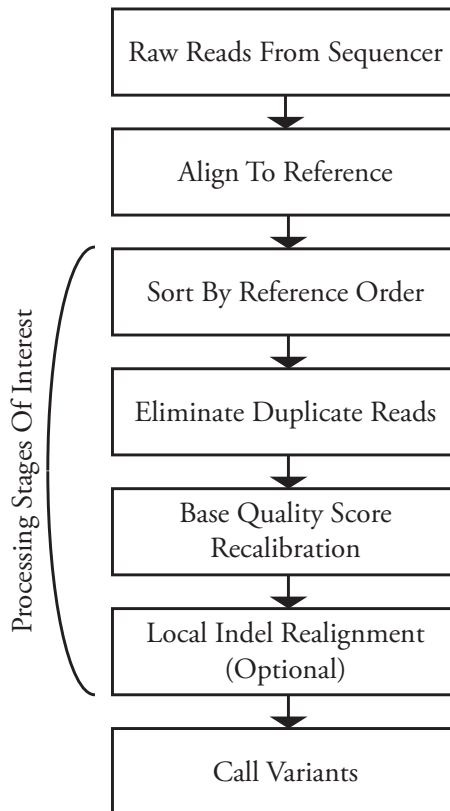


Figure 1: Variant Calling Pipeline

the algorithms that implement the intermediate read processing stages. For a detailed discussion of these algorithms, we refer readers to DePristo et al [13].

Sorting: This phase performs a pass over the reads and sorts them by the reference position at the start of their alignment.

Duplicate Removal: An insufficient number of sample molecules immediately prior to PCR can cause the generation of duplicate DNA sequencing reads. Detection of duplicate reads requires matching all reads by their 5' position and orientation after read alignment. Reads with identical position and orientation are assumed to be duplicates. When a group of duplicate reads is found, each read is scored and all, but the top-scoring read, are marked as duplicates. Approximately 66% of duplicates marked in this way are true duplicates caused by PCR-induced duplication while the remaining 33% are caused by the random distribution of read ends[9]. There is currently no way to separate "true" duplicates from

Table 1: Processing Stage Times for GATK Pipeline

Stage	GATK 2.7/NA12878
Alignment	52 hr
Sort	
Mark Duplicates	13 hr
BQSR	6 hr
Realignment	33 hr
Call Variants	4 hr
Total	67 hr

randomly occurring duplicates.

Base Quality Score Recalibration: During the sequencing process, systemic errors occur that lead to the incorrect assignment of base quality scores. In this step, a statistical model of the quality scores is built and is then used to revise the measured quality scores.

Local Realignment: For performance reasons, all modern aligners use algorithms that provide approximate alignments. This can cause reads with evidence of indels to have slight misalignments with the reference. In this stage, we use fully accurate sequence alignment methods (Smith-Waterman algorithm [28]) to locally realign reads that contain evidence of short indels. This pipeline step is omitted in some variant calling pipelines, if the variant calling algorithm that is implemented is not susceptible to these local alignment errors.

For current implementations of these read processing steps, performance is limited by disk bandwidth. This is because the operations read in a SAM/BAM file, perform a bit of processing, and write the data to disk as a new SAM/BAM file. We address this by performing our processing iteratively in memory. The four read processing stages can then be chained together, eliminating three dumps to disk and an additional three long reads from disk.

4 Design Philosophy

Modern bioinformatics pipelines have been designed without a model for how the system should grow or for how components of the analytical system should connect. We seek to provide a more principled model for system composition. Our system architecture was inspired heavily by the Open Systems Interconnection (OSI) model for networking services [40]. This

conceptual model standardized the internal functionality of a networked computer system, and its "narrow waisted" design was critical to the development of the services that would become the modern internet. We present a similar decomposition of services for genomics data in figure 2.

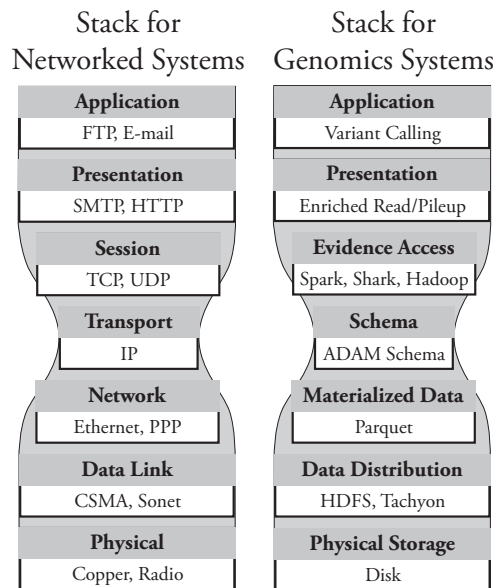


Figure 2: Stack Models for Networking and Genomics

REVIEW STACK

The seven levels of our stack model are decomposed as follows:

1. **Physical Storage:** This layer coordinates data writes to physical media.
2. **Storage:** This layer manages access, replication, and distribution of the genomics files that have been written to disk.
3. **Materialized Data:** This layer encodes the patterns for how data is encoded and stored. This provides read/write efficiency and compression.
4. **Data Schema:** This level specifies the representation of data when it is accessed, and forms the narrow waist of the pipeline.
5. **Evidence Access:** This layer implements efficient methods for performing common access patterns such as random database queries, or sequentially reading records from a flat file.

6. **Presentation:** The presentation layer sits on top of the data access layer and provides the application developer with efficient and straightforward methods for querying the characteristics of individual portions of genetic data.
7. **Application:** Applications like variant calling and alignment are implemented in this layer.

A well defined software stack has several significant advantages. By limiting application interactions with layers lower than the API, application developers are given a clear and consistent view of the data they are processing. By divorcing the API from the data access layer, we unlock significant flexibility. With careful design in the data format and data access layers, we can seamlessly support conventional flat file access patterns, while also allowing easy access to data with database methods (see §8.1). By treating the compute substrate and storage as separate layers, we also drastically increase the portability of the APIs that we implement. This approach is a significant improvement over current approaches which intertwine all layers together — if all layers are intertwined, new APIs must be developed to support new compute substrates [25], and new access patterns must be carefully retrofitted to the data format that is in use [19].

We can distill our philosophy into three observations that drive our design decisions:

1. **Scalability is a primary concern for modern genomics systems:** Processing pipelines operate on data that can range in size from tens of gigabytes to petabytes. We cannot remove processing stages, as this has an unacceptable impact on accuracy. To increase pipeline latency and throughput, our systems must be able to scale efficiently to tens to hundreds of nodes.
2. **Bioinformaticians should be concerned with data, not formats:** The SAM/BAM formats do not provide a schema for the data they store. This limits visibility into the structure of data being processed. Additionally, it significantly increases implementation difficulty: by making the format a primary concern, significant work has been required to introduce programming interfaces that allow these formats to be read from different programming languages and environments [22, 5, 25].
3. **Openness and flexibility are critical to adoption:** As the cost of sequencing continues

to drop [24], the processing of genomics data will become more widespread. By maintaining an open source standard that can be flexibly modified, we allow all the members of the community to contribute to improving our technologies. Developing a flexible stack is also important; genetic data will be processed in a host of heterogeneous computing environments, and with a variety of access patterns. By making our format easy to adopt to different techniques, we maximize its use.

In the next few sections, we discuss the implementation of ADAM, and how it was designed to satisfy these goals. We then analyze the performance of our system and discuss the greater impact of our design philosophy.

5 Data Format and API

ADAM contains formats for storing read and reference oriented sequence information, and variant/genotype data. The read oriented sequence format is forwards and backwards compatible with BAM/SAM, and the variant/genotype data is forwards and backwards compatible with VCF. In this section, we discuss the frameworks used to store this data. We then introduce the representations, and discuss the content of each representation.

The data representation for the ADAM format is described using the open source Apache Avro data serialization system [1]. The Avro system also provides a human readable schema description language that can auto-generate the schema implementation in several common languages including Scala, Java, C/C++/C#, Python, Ruby, and php. This provides a significant cross-compatibility advantage over the BAM/SAM format, where the data format has only been implemented for C/C++ through Samtools and for Java through Picard [22, 5].

We layer this data format inside of the Parquet column store [4]. Parquet is an open source columnar storage format that was designed by Cloudera and Twitter, which can be used as a single flat file, a distributed file, or as a database inside of Hive, Shark, or Impala. Columnar stores like Parquet provide several significant advantages when storing genomic data:

- Column stores enable predicate pushdown [21], which minimizes the amount of data read from disk. When using predicate pushdown, we deserialize specific fields in a record from disk, and apply them to a predicate function. We then

only read the full record if it passes the predicate. This is useful when implementing filters that check read quality.

- Column stores achieve extremely high compression [6]. This reduces storage space on disk, and also improves serialization performance inside of MapReduce frameworks. We are able to achieve a $0.75\times$ lossless compression ratio when compared to BAM, and have especially impressive results on quality scores. This is discussed in more detail in §7.4.
- Varied data projections can be achieved with column stores. This means that we can choose to only read several fields from a record. This improves performance for applications that do not read all the fields of a record. Additionally, we do not pay for null fields in records that we store. This allows us to easily implement lossy compression on top of the ADAM format, and also allows us to eliminate the FASTQ standard.

By supporting varied data projections with low cost for field nullification, we can implement lossy compression on top of the ADAM format. We discuss this further in §7.4.

A visualization of how ADAM in Parquet compares with BAM can be seen in figure 3. We remove the file header and instead distribute these values across all of the records stored. This eliminates global information and makes the file much simpler to distribute across machines. This is effectively free in a columnar store, as the store just notes that the information is replicated across many reads. Additionally, Parquet writes data to disk in regularly sized row groups (from [4], see Parquet format readme) — this allows parallelism across rows by enabling the columnar store to be split into independent groups of reads.

BAM File Format	ADAM File Format
Header: n = 500	chr: c20 * 5, ... ; seq:
Reference 1	TCGA, GAAT, CCGAT,
Sample 1	TTGCAC, CCGT, ... ;
1: c20, TCGA, 4M; 2: c20,	cigar: 4M, 4M1D, 5M,
GAAT, 4M1D; 3: c20, CCGAT,	6M, 3M1D1M, ... ; ref:
5M; 4: c20, TTGCAC, 6M; 5:	ref1 * 500; sample:
c20, CCGT, 3M1D1M; ...	sample1 * 500;

Figure 3: Comparative Visualization of BAM and ADAM File Formats

As discussed in §4, the internal data types presented in this section make up the narrow waist of our proposed genomics stack. In §8.1 and §9.1, we

review how this abstraction allows us to support different data access frameworks and storage techniques to tailor the stack for the needs of a specific application.

5.1 Read Oriented Storage

Our default read oriented data format is defined in table 2. This format provides all of the fields supported by the SAM format. To make it easier to split the file between multiple machines for distributed processing, we have eliminated the file header. Instead, the data encoded in the file header can be reproduced from every single record. Because our columnar store supports dictionary encoding and these fields are replicated across all records, replicating this data across all records has a negligible cost in terms of file size on disk.

This format is used to store all data on disk. In addition to this format, we provide an enhanced API for accessing read data.

5.2 Reference Oriented Storage

In addition to storing sequences as reads, we provide a storage format for reference oriented (pileup) data. This format is documented in table 3. This pileup format is also used to implement a data storage format similar to the GATK’s ReducedReads format [13].

We treat inserts as an inserted sequence range at the locus position. For bases that are an alignment match against the reference⁴, we store the read base and set the range offset and length to 0. For deletions, we perform the same operation for each reference position in the deletion, but set the read base to null. For inserts, we set the range length to the length of the insert. We step through the insert from the start of the read to the end of the read, and increment the range offset at each position.

As noted earlier, this datatype supports an operation similar to the GATK’s ReducedRead datatype. We refer to these datatypes as aggregated pileups. We perform aggregation by grouping together all bases that share a common reference position and read base. Within an insert, we group further by the position in the insert. Once the bases have been grouped together, we average the base and mapping quality scores. We also count the number of

⁴Following the conventions for CIGAR strings, an alignment match does not necessarily correspond to a sequence match. An alignment match simply means that the base is not part of an indel. The base may not match the reference base at the loci.

Table 2: Read Oriented Format

Group	Field	Type
General	Reference Name	String
	Reference ID	Int
	Start	Long
	Mapping Quality	Int
	Read Name	String
	Sequence	String
	Mate Reference	String
	Mate Alignment Start	Long
	Mate Reference ID	Int
	Cigar	String
	Base Quality	String
Flags	Read Paired	Boolean
	Proper Pair	Boolean
	Read Mapped	Boolean
	Mate Mapped	Boolean
	Read Negative Strand	Boolean
	Mate Negative Strand	Boolean
	First Of Pair	Boolean
	Second Of Pair	Boolean
	Primary Alignment	Boolean
Attributes	Mismatching Positions	String
	Other Attributes	String
Read Group	Sequencing Center	String
	Description	String
	Run Date	Long
	Flow Order	String
	Key Sequence	String
	Library	String
	Median Insert	Int
	Platform	String
	Platform Unit	String
	Sample	String

Table 3: Reference Oriented Format

Group	Field	Type
General	Reference Name	String
	Reference ID	Int
	Position	Long
	Range Offset	Int
	Range Length	Int
	Reference Base	Base
	Read Base	Base
	Base Quality	Int
	Mapping Quality	Int
	Number Soft Clipped	Int
	Number Reverse Strand	Int
	Count At Position	Int
Read Group	Sequencing Center	String
	Description	String
	Run Date	Long
	Flow Order	String
	Key Sequence	String
	Library	String
	Median Insert	Int
	Platform	String
	Platform Unit	String
	Sample	String

Table 4: Genotype Format

Field	Type
Sample ID	String
Genotype	Array[Int]
Likelihood (Phred)	Array[Int]
Format	String

6 In-Memory Programming Model

soft clipped bases that show up at this location, and the amount of bases that are mapped to the reverse strand.

5.3 Variant and Genotype Storage

Chris H.

Variant type is one of SNP, multiple nucleotide polymorphism (MNP), insertion, deletion, structural variant (SV), or complex.

As discussed in §3, the main bottleneck in current genomics processing pipelines is packaging data up on disk in a BAM file after each processing step. While this process is useful as it maintains data lineage⁵, it significantly increases the latency of the processing pipeline.

For the processing pipeline we have built on top

⁵Since we store the intermediate data from each processing step, we can later redo all pipeline processing after a certain stage without needing to rerun the earlier stages. This comes at the obvious cost of space on disk. This tradeoff makes sense if the cost of keeping data on disk and reloading it later is lower than the cost of recomputing the data. This does not necessarily hold for modern MapReduce frameworks [38].

Table 5: Variant Format

Field	Type
Reference Name	String
Position	Long
Variant ID	String
Reference Allele	String
Alternate Alleles	Array[String]
Allele Count	Array[Int]
Chromosome Count	Int
Type	Variant Type
Info	String
Filter	String
Genotypes	Array[AdamGenotype]

of the ADAM format, we have minimized disk accesses (read one file at the start of the pipeline, and write one file after all transformations have been completed). Instead, we cache the reads that we are transforming in memory, and chain multiple transformations together. Our pipeline is implemented on top of the Spark MapReduce framework, where reads are stored as a map using the Resilient Distributed Dataset (RDD) primitive.

7 Performance

This section reviews performance metrics pertinent to the ADAM file format, and applications written on top of ADAM. Our performance analysis is partitioned into three sections:

- **7.1. Microbenchmarks:** In this section, we review several testcases that have little computation, and review ADAM against Hadoop-BAM. This section looks to evaluate the disk access performance of ADAM. In this section, we are most interested in validating scalability and building intuition for the efficiencies provided by variable projection and predicate pushdown.
- **7.2. Applications:** The main goal of ADAM is to accelerate the read processing pipelines that were introduced in §3. In this section, we review the throughput of our processing pipeline.
- **7.4. Compression:** Here, we review the reductions in disk space that we are able to achieve when compared against BAM. We also review how we can achieve further compression through aggregating pileup data, and how to implement both lossy and reference-based compression top of ADAM.

7.1 Microbenchmarks

To validate the performance of ADAM, we created several microbenchmarks. These microbenchmarks are meant to demonstrate the pure read/write throughput of ADAM, and to show how the system scales in a cluster. To validate these benchmarks, we implemented the relevant tests in both ADAM and Hadoop-BAM [25], running on Spark on an Amazon Elastic Compute 2 (EC2) cluster. We used the *m2.4xlarge* instance type for all of our tests. This instance type is an 8 core memory optimized machine with 68.4 Gibibytes (GiB) of memory per machine. We ran these benchmarks against two genomes from the 1000 Genomes project: HG00096, a low-coverage human genome (15 GB BAM), and NA12878, a high-coverage human genome (246 GB BAM).

7.1.1 Read/Write Performance

The decision to use a columnar store was based off of the observation that most genomics applications are read-heavy. In this section, we aim to quantify the improvements in read performance that we achieve by using a column store, and how far these gains can be maximized through the use of projections. We also seek to observe how significant of a penalty we pay when writing data back to disk.

Data Materialization and Write Performance:

Read Length and Quality: This microbenchmark scans all the reads in the file, and collects statistics about the length of the read, and the mapping quality score. The results of this benchmark can be found in Figure 4.

7.1.2 Predicate Pushdown

Predicate pushdown is one of the significant advantages afforded to us through the use of Parquet [4]. In traditional pipelines, the full record is deserialized from disk, and the filter is implemented as a Map/Reduce processing stage. Parquet provides us with predicate pushdown: we deserialize only the columns needed to identify whether the record should be fully deserialized. We then process those fields; if they pass the predicate, we then deserialize the full record.

We can prove that this process requires no additional reads. Intuitively this is true as the initial columns read would need to be read if they were to

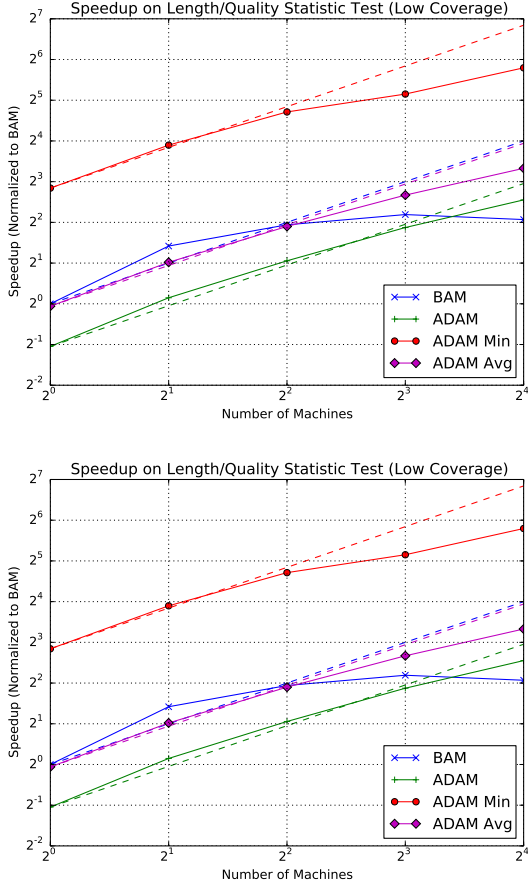


Figure 4: Read Length and Quality Speedup for HG00096 and NA12878

be used in the filter later. However, this can be formalized.

Theorem 1. *Predicate pushdown requires no additional disk accesses.*

Proof. If R is the set of records to read, $Proj$ is the projection we are using, and $size(i, j)$ returns the amount of data recorded in column j of record i , the total data read without predicate pushdown is:

$$\sum_m^R \sum_i^{Proj} size(m, i) \quad (1)$$

By using predicate pushdown, we reduce this total to:

$$\sum_m^R \sum_i^{Pred} size(m, i) + \sum_n^{R_{p=True}} \sum_i^{Proj \setminus Pred} size(n, i) \quad (2)$$

$Pred$ represents the set of columns used in the predicate. The set $R_{p=True}$ represents the subset of records in R that pass the predicate function. We can show that if no records fail the predicate (i.e. $R = R_{p=True}$), then equations 1 and 2 become equal. This is because by definition $Proj = (Pred \cup (Proj \setminus Pred))$. \square

In this section, we seek to provide an intuition for how predicate pushdown can improve the performance of actual genomics workloads. We apply predicate pushdown to implement read quality predication, mapping score predication, and predication by chromosome.

Quality Flags Predicate: This microbenchmark scans the reads in a file, and filters out reads that do not meet a specified quality threshold. We use the following quality threshold:

- The read is mapped.
- The read is at its primary alignment position.
- The read did not fail vendor quality checks.
- The read has not been marked as a duplicate.

This threshold is typically applied as one of the first operations in any read processing pipeline.

7.2 Applications

ADAM relies on Apache Spark for in-memory, map-reduce processing. At the heart of Spark is a programming construct, called a Resilient Distributed Dataset (RDD), which represents a fault-tolerant collection of elements that are operated on in parallel. Programming with an RDD is very similar to programming with any standard Scala collection like `Seq` or `List`. This polymorphism allows developers to focus on the algorithm they’re developing without concerning themselves with the details of distributed computation. The example code presented in this section highlight the succinctness and expressiveness of Spark Scala code. While ADAM is currently written exclusively in Scala, Apache Spark also has Java and Python support.

For these applications to be run in a distributed fashion, the data needs to be loaded onto the Hadoop Distributed File System (HDFS). HDFS breaks the file into blocks, distributes them across the cluster and maintains information about block (and replica) locations. Breaking large files, like BAM files, into blocks

allows systems, like Spark, to create tasks which run across the entire cluster processing the records inside of each block. The wall clock times presented in the applications below do not consider the time it takes to load data to HDFS.

Table 6: Time to place a file on HDFS

File Size	Cluster Size
TODO	TODO
TODO	TODO

TODO: provide some numbers on HDFS load times (including the speedup found using distcp – e.g. a 250GB on S3 can be written to HDFS in just over 20 minutes – ganglia reported sustained 450Mb/s). Note: we did not load Hadoop-native for these tests which would have helped with performance.

The NA12878 high-coverage BAM file used for these experiments can be found on the 1000 Genomes ftp site, [ftp.1000genomes.ebi.ac.uk](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/NA12878/high_coverage_alignment/) in directory /vol1/ftp/data/NA12878/high_coverage_alignment/.

TODO: talk about conversion from BAM to ADAM

TODO: Explain an ADAMRecord

7.2.1 Flagstat

Even the simplest application can benefit from the ADAM format and execution environment. For example, the samtools **flagstat** command prints read statistics on a BAM file, e.g.

```
$ samtools flagstat NA12878.bam
1685229894 + 0 in total (QC-passed reads + QC-failed reads)
0 + 0 duplicates
1622164506 + 0 mapped (96.26%:-nan%)
1685229894 + 0 paired in sequencing
842614875 + 0 read1
842615019 + 0 read2
1602206371 + 0 properly paired (95.07%:-nan%)
1611286771 + 0 with itself and mate mapped
10877735 + 0 singletons (0.65%:-nan%)
3847907 + 0 with mate mapped to a different chr
2174068 + 0 with mate mapped to a different chr (mapQ>=5)
```

ADAM has a **flagstat** command which provides identical information in the same text format.

Table 7: Time to Run Flagstat on High-Coverage NA12878

Software	Wall clock time
Samtools 0.1.19	25m24s (todo: verify)
ADAM 0.5.0	0m46s

The ADAM **flagstat** command was run on EC2 with 32 **cr1.8xlarge** machines which is, not suprisingly, why ADAM ran 34 times faster. While the samtools **flagstat** command is multi-threaded, it can only be run on a single machine.

TODO: Collect numbers on

7.2.2 Sort

When coordinate sorting a BAM file, reads are re-ordered by reference id and alignment start position. This sorting makes it easier for downstream tools to find reads for a specific area of the reference and makes finding mates in paired-ends reads more efficient. Many utilities will only except BAM files that have been sorted, e.g. Picard MarkDuplicates.

The ADAM Scala code for sorting a BAM file is succinct and relatively easy to follow, e.g.

```
1: def adamSortReadsByReferencePosition(): RDD[ADAMRecord] = {
2:   rdd.map(p => {
3:     val referencePos = ReferencePosition(p) match {
4:       case None =>
5:         // Move unmapped reads to the end
6:         ReferencePosition(Int.MaxValue, Long.MaxValue)
7:       case Some(pos) => pos
8:     }
9:     (referencePos, p)
10:  }).sortByKey().map(p => p._2)
11: }
```

Each read is mapped to a tuple with the two values: position and read. Unmapped reads that a value of **None** for **ReferencePosition** are sent to the end of the file.

Table 8: Sort NA12878

Software	EC2 profile	Wall clock time
Picard Tools 1.102	TODO	13h (todo: verify this)
ADAM 0.5.0	32 cr1.8xlarge	32m38s
ADAM 0.5.0	100 m2.4xlarge	20m57s

ADAM was run on an EC2 cluster with 32 **cr1.8xlarge** machines and a cluster with 100 **m2.4xlarge** instances.

Monitoring the Spark console during sorting revealed that the **sortByKey** operation (on line 10 above) took a total of 1 min 37 secs. The **map** operation (at line 2 above) took 1 min 48 secs, and the remainder of the time, 29 min 25 secs, was spent writing the results to HDFS. There were a total of 1586 tasks with a majority finishing in under 10 minutes. A sole traggling task took 29 min 25 secs to write all unmapped reads. This indicates that using a custom partitioner could dramatically improve performance.

7.2.3 Mark Duplicates

The ADAM algorithm for marking duplicate is nearly identical to Picard’s MarkDuplicates command. The following ADAM code, reformatted for this report, expresses the algorithm succinctly in 42 lines of Scala code, e.g.

```
1:for (((leftPos, library), readsByLeftPos) <-
2:  rdd.adamSingleReadBuckets()
3:  .keyBy(ReferencePositionPair(_))
4:  .groupBy(leftPositionAndLibrary);
5:
6:buckets <- {
7:
8:leftPos match {
9:  // These are all unmapped reads.
10:  // There is no way to determine if they are duplicates
11:  case None =>
12:    markReads(readsByLeftPos.unzip._2, areDups = false)
13:  // These reads have their left position mapped
14:  case Some(leftPosWithOrientation) =>
15:    // Group the reads by their right position
16:    val readsByRightPos = readsByLeftPos.groupBy(rightPosition)
17:    // Find any reads with no right position
18:    val fragments = readsByRightPos.get(None)
19:    // Check if we have any pairs (reads with a right position)
20:    val hasPairs = readsByRightPos.keys.exists(_.isDefined)
21:    if (hasPairs) {
22:      // Since we have pairs, mark all fragments as duplicates
23:      val processedFrgs = if (fragments.isDefined) {
24:        markReads(fragments.get.unzip._2, areDups = true)
25:      } else {
26:        Seq.empty
27:      }
28:      val processedPairs =
29:        for (buckets <- (readsByRightPos - None).values;
30:            processedPair <-
31:              scoreAndMarkReads(buckets.unzip._2))
32:        yield processedPair
33:      processedPairs ++ processedFrgs
34:    } else if (fragments.isDefined) {
35:      // No pairs. Score the fragments.
36:      scoreAndMarkReads(fragments.get.unzip._2)
37:    } else {
38:      Seq.empty
39:    }
40:}
41:};
42:read <- buckets.allReads() yield read
```

For lines 1-4, all reads with the same record group name and read name are collected into buckets. These buckets contain the read and optional mate or secondary alignments. Each read bucket is then keyed by 5' position and orientation and grouped together by left (lowest coordinate) position, orientation and library name.

For lines 8-41, we processed each read bucket with a common left 5' position. Unmapped reads are never marked duplicate as their position is not known. Mapped reads with a common left position are separated into paired reads and fragments. Fragments, in this context, are reads that have no mate or should have a mate but it doesn’t exist.

If there are pairs in a group, all fragments are marked duplicates and the paired reads are grouped by their right 5' position. All paired reads that have

a common right and left 5' position are scored and all but the highest scoring read is marked a duplicate.

If there are no pairs in a group, all fragments are scored and all but the highest scoring fragment are marked duplicates.

Table 9: Mark Duplicates for NA12878

Software	EC2 profile	Wall clock time
Samtools rmdup	TODO?	TODO?
Picard MarkDuplicates 1.102	TODO	13h (todo: verify this)
ADAM 0.5.0	100 m2.4xlarge	28m17

TODO: talk about differences between ADAM and Picard, e.g. ADAM does not require sorted input,

7.2.4 Mark Duplicates and Sort In-Memory

Table 10: Sort and Mark Duplicates for NA12878

Software	EC2 profile	Wall clock time
Picard Tools 1.102	TODO	26h (todo: verify this)
ADAM 0.5.0	100 m2.4xlarge	50m45s

7.3 Tuning Spark on EC2

All applications were run on Amazon EC2 with minimal tuning. The performance numbers in this report should not be considered best case performance numbers.

An important Apache Spark property is `spark.default.parallelism` which controls the number of tasks to use for distributed shuffle operations. This value needs to be tuned to match the underlying cluster that Spark is running on.

Table 11: Time to Sort NA12878

<code>spark.default.parallelism</code>	Wall clock time
500	47 min 10 secs
48	32 min 38 secs

The following graph shows CPU time series for two sort operations on a 32-node `cr1.8xlarge` EC2 cluster running back to back with different values for `spark.default.parallelism`. The early CPU profile is when `spark.default.parallelism` is set to 500 and the latter profile is when it’s set to 48. You can see from the graph that with the higher level

of parallelism there is a significantly higher amount of system CPU (in red) as 32 CPUs contend for access to the two disks in each `cr1.8xlarge` instance. While `cr1.8xlarge` instances provide plenty of memory, they do not have enough disks to help spread the I/O load around.

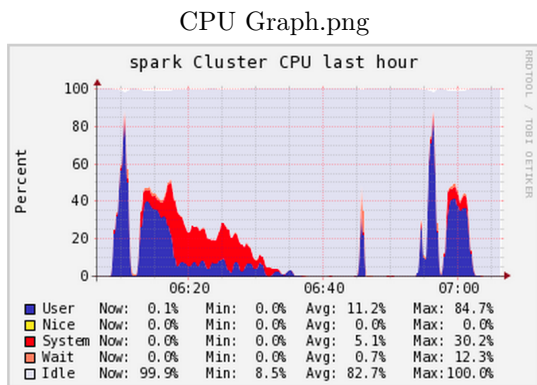


Figure 5: Comparison of 500 to 48 parallelism

7.4 Compression

By using dictionary coding, run length encoding with bit-packing, and GZIP compression, ADAM files can be stored using less disk space than an equivalent BAM file.

8 Future Work

Beyond the initial release of the ADAM data formats, API, and read transformations, we are working on several other extensions to ADAM. These extensions strive to make ADAM accessible to more users and to more programming patterns.

8.1 Database Integration

ADAM format plays a central role in a library that we develop which will allow sequencing data to be queried from the popular SQL based warehouse software such as Shark [35], Hive [2] and Impala [3]. In fact, the Avro/Parquet storage backend enables the direct usage of ADAM files as SQL tables; this integration is facilitated through libraries that integrate Parquet with the aforementioned software [4].

Of course there is a lot of work remaining before SHARK provides biologically relevant functionality, such as the one that is outlined in Bafna et al [8];

Shark needs to be enhanced with libraries that provide efficient indexing, storage management and interval query handling. The indexing scheme that we develop enables range based read retrieval, such as the indexing mechanism of samtools, quick access to mate pair reads such as the indexing mechanism of GQL [19] and also text based search features on fields such as the read sequence and the CIGAR string. To keep storage under control we implement a set of heuristics that prevent the materialization of queries with reproducible output. Finally, given that most biological data is modelled after intervals on the human genome, our library includes user defined functions that implement interval operators such as the IntervalJoin operator of GQL.

8.2 API Support Across Languages

Currently, the ADAM API and example read processing pipeline are implemented in the Scala language. However, long term we plan to make the API accessible to users of other languages. ADAM is built on top of Avro, Parquet, and Spark. Spark has bindings for Scala, Java, and Python, and Parquet implementations exist for Java and C++. As the data model is implemented on top of Avro, the data model can also be autogenerated for C/C++/C#, Python, Ruby, and php. We hope to make API implementations available for some of these languages at a later date.

9 Discussion

Table to go here.

9.1 Columnar vs. Row Storage

In our current implementation, we use the Parquet columnar store to implement the data access layer. We chose to use columnar storage as we believe that columnar storage is a better fit for bioinformatics workloads than flat row storage; this discussion is introduced in §5. Typically, column based storage is preferable unless our workload is write heavy; write heavy workloads perform better with row based storage formats [30]. Genomics pipelines tend to skew in favor of reading data — pipelines tend to read a large dataset, prune/reduce data, perform an analysis, and then write out a significantly smaller file that contains the results of the analysis.

However, it is conceivable that emerging workloads could change to be write heavy. We note that although our current implementation would not be op-

timized for these workloads, the layered model proposed in §4 allows us to easily change our implementation. Specifically, we can swap out columnar storage in the data access layer with row oriented storage⁶. The rest of the implementation in the system would be isolated from this change: algorithms in the pipeline would be implemented on top of the higher level API, and the compute substrate would interact with the data access layer to implement the process of reading and writing data.

Impact on number of seeks.

9.2 Dancehall vs. Datacenter

Matt

9.3 New Applications Enabled By ADAM

The ADAM format and environment provides a significant performance improvement over current BAM based solutions. By reducing processing latency, ADAM will make genomic data available for use in latency-critical applications, such as genomic driven treatment of acute diseases. Additionally, the genomics stack design we proposed in §4 unlocks new computational patterns that can be applied to genomic data. In this section, we introduce a few of these novel applications, discuss why they are important, and explain how these applications are enabled by ADAM.

9.3.1 Personalized Treatment of Acute Diseases

There is significant interest in using genomic data as another layer of evidence for the treatment of diseases that have genetic causes [15, 31]. While chronic hereditary diseases like Huntington’s Disease may not be sensitive to the latency of genomic processing pipelines, reduced latency is critical for neonatal diagnosis and for the treatment of some cancers. In acute myeloid leukemia (AML), there exists a well defined set of mutations that can be used for guiding patient treatment [27]. However, due to the rapid progression of AML, treatment decisions must be made in under one week. As illustrated in §3, this target is not within the reach of current processing pipelines.

It is worth noting that the numbers presented in §3 would be optimistic in reality. To improve accuracy, the current state of the art in mutation calling uses

⁶It is worth noting that the Avro serialization system that we use to define ADAM is a performant row oriented system.

very high coverage sequencing data ($\sim 100 - 150\times$) for the tumor cell line, along with high coverage sequencing data ($\sim 30\times$) for the normal cell line [11]. We are developing a variant and mutation calling pipeline that uses the ADAM processing pipeline and the high-performance SNAP aligner [37] to achieve performance that will satisfy the latency requirements for use in acute disease diagnosis.

9.3.2 Viral Outbreak Monitoring

There is significant interest in tools for the analysis of viral outbreaks in both the medical and defense epidemiology communities [29, 10]. Genomic techniques can be used for identifying infection vectors in epidemics that feature rapidly mutating diseases, as well as for the rapid identification of unknown infectious agents in an attack using biological weapons. Similar to the acute genetic driven diseases discussed in §9.3.1, both of these applications are latency sensitive; the United States Department of Defense (DoD) has set a goal of developing a system which can identify an unknown infectious agent within one hour with the use of genomic data.

In addition to providing significantly lower latency than current pipelines, there are other advantages to ADAM. Current approaches to tracking viral outbreaks perform a graph based analysis that uses variant data and temporal data to identify possible modes of disease transmission [29]. As ADAM is implemented on top of Spark, data in ADAM format can be processed as a graph by using the GraphX distributed graph analysis framework [34].

9.3.3 Ad Hoc Genomic Analysis

ADAM is implemented on top of the Spark MapReduce framework. A major contribution of the Spark framework was the introduction of a command line interface that allows for the real-time ad hoc analysis of data kept in memory [38]. This interactive framework can be coupled with the ability to perform interactive queries in Shark to allow for rapid data aggregation and analysis. Several applications for this approach are presented by Kozanitis et al. [19] such as allowing direct evidence access to physicians who are working in a clinical setting with patients with hereditary diseases. However, GQL performance limited the performance of these queries to ~ 45 minutes. We are still working on the implementation of our database interface (see §8.1), but expect to achieve improved performance.

In addition to clinical applications, interactive data access could be applied to pathway analysis. Specifi-

cally, if a biologist has genetic data for organisms who demonstrate a specific phenotype, interactive data processing could be used for correlating lab observations with genotypes. Instead of genotyping the full genome/exome of all reference organisms, the raw reads for all organisms could be aligned in stored in a database. Once lab measurements and cytological observations have allowed the biologist to hypothesize that a specific pathway had failed, the database of raw reads could be queried for all genes of interest (genes that influence the pathway under study). This would lead to a significant reduction in the amount of genomic data that needed to be processed, which would allow for interactive read processing and variant calling.

9.4 Programming Model Improvements

Beyond improving performance by processing data in memory, our programming model improves programmer efficiency. We build our operations as a set of transforms that extend the RDD processing available in Spark. These transformations allow a very straightforward programming model that has been demonstrated to significantly reduce code size [38]. Additionally, the use of Scala, which is an object functional statically typed language that supports type inference couples provides further gains.

Additional gains come from our data model. We expose a clear schema, and build an API that expands upon this schema to implement commonly needed functions. This API represents level 6 in the stack we proposed in section §4. To reduce the cost of this API, we make any additional transformations from our internal schema to the data types presented in our API lazy. This eliminates the cost of providing these functions if they are not used. Additionally, if an additional transformation is required, we only pay the cost the first time the transformation is performed.

In total, these improvements can provide a 2× improvement in programmer productivity. This is demonstrated through GAParquet, which demonstrates some of the functionality in the GATK on top of the ADAM stack [16]. In the DiagnoseTargets stage of GATK, the number of lines of code (LOC) needed to implement the stage dropped from 400 to 200. We also see a reduction in lines of code needed to implement the read transformation pipeline described in §3. A LOC comparison of GATK/GAParquet and the read transformation pipeline is found in table 12.

Table 12: Lines of Code for ADAM and Original Implementation

Application	Original	ADAM
GATK Diagnose Targets		
Walker	326	134
Subclass	93	66
Total	419	200
Translation Pipeline		
BQSR		
Mark Duplicates		
Realign Indels		
Sort		
Total		

9.5 Benefits of Explicit Schema

Matt

10 Summary

In this technical report, we have presented ADAM which is a new data storage format and processing pipeline for genomics data. ADAM makes use of efficient columnar storage systems to improve the lossless compression available for storing read data, and uses in-memory processing techniques to eliminate the read processing bottleneck faced by modern variant calling pipelines. On top of the file formats that we have implemented, we also present APIs that enhance developer access to read, pileup, genotype, and variant data. We are currently in the process of extending ADAM to support SQL querying of genomic data, and extending our API to more programming languages. ADAM promises to improve the development of applications that process genomic data, by removing current difficulties with the extraction and loading of data and by providing simple and performant programming abstractions for processing this data.

A Availability

The ADAM source code is available at Github at <http://www.github.com/massie/adam>, and the ADAM project website is at <http://adam.cs.berkeley.edu>. Additionally, ADAM is deployed through Maven with the following dependency:

TBD

At publication time, the current version of ADAM is 0.5.0. ADAM is open source and is released under the Apache 2 license.

B Algorithm Implementations

B.1 Sort Implementation

B.2 BQSR Implementation

B.3 Indel Realignment Implementation

B.4 Duplicate Marking Implementation

References

- [1] Apache Avro. <http://avro.apache.org>.
- [2] Apache Hive. <http://hive.apache.org/>.
- [3] Cloudera Impala. <http://www.cloudera.com/content/cloudera/en/products/cdh/impala.html>.
- [4] Parquet. <http://www.parquet.io>.
- [5] Picard. <http://picard.sourceforge.org>.
- [6] ABADI, D., MADDEN, S., AND FERREIRA, M. Integrating compression and execution in column-oriented database systems. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data* (2006), ACM, pp. 671–682.
- [7] ASNANI, H., BHARADIA, D., CHOWDHURY, M., OCHOA, I., SHARON, I., AND WEISSMAN, T. Lossy Compression of Quality Values via Rate Distortion Theory. *ArXiv e-prints* (July 2012).
- [8] BAFNA, V., DEUTSCH, A., HEIBERG, A., KOZANITIS, C., OHNO-MACHADO, L., AND VARGHESE, G. Abstractions for Genomics: Or, which way to the Genomic Information Age? *Communications of the ACM* (2013). (To appear).
- [9] BAINBRIDGE, M., WANG, M., BURGESS, D., KOVAR, C., RODESCH, M., D’ASCENZO, M., KITZMAN, J., WU, Y.-Q., NEWSHAM, I., RICHMOND, T., JEDDELOH, J., MUZYNY, D., ALBERT, T., AND GIBBS, R. Whole exome capture in solution with 3 gbp of data. *Genome Biology* 11, 6 (2010), R62.
- [10] BEARD, J. DARPA’s bio-revolution. *Biomedical services* (2008), 4–6.
- [11] CIBULSKIS, K., LAWRENCE, M. S., CARTER, S. L., SIVACHENKO, A., JAFFE, D., SOUGNEZ, C., GABRIEL, S., MEYERSON, M., LANDER, E. S., AND GETZ, G. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature biotechnology* 31, 3 (2013), 213–219.
- [12] COX, A. J., BAUER, M. J., JAKOBI, T., AND ROSONE, G. Large-scale compression of genomic sequence databases with the Burrows-Wheeler transform. *Bioinformatics* 28, 11 (Jun 2012), 1415–1419.
- [13] DEPRISTO, M. A., BANKS, E., POPLIN, R., GARIMELLA, K. V., MAGUIRE, J. R., HARTL, C., PHILIPPAKIS, A. A., DEL ANGEL, G., RIVAS, M. A., HANNA, M., ET AL. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature genetics* 43, 5 (2011), 491–498.
- [14] FRITZ, M. H.-Y., LEINONEN, R., COCHRANE, G., AND BIRNEY, E. Efficient storage of high throughput DNA sequencing data using reference-based compression. *Genome research* 21, 5 (2011), 734–740.
- [15] GHOSH, D., AND POISSON, L. M. “omics” data and levels of evidence for biomarker discovery. *Genomics* 93, 1 (2009), 13–16.
- [16] HARTL, C. GAParquet. <http://www.github.com/chart1/gaparquet>.
- [17] JANIN, L., ROSONE, G., AND COX, A. J. Adaptive reference-free compression of sequence quality scores. *Bioinformatics* (Jun 2013).
- [18] JONES, D. C., RUZZO, W. L., PENG, X., AND KATZE, M. G. Compression of next-generation sequencing reads aided by highly efficient de novo assembly. *Nucleic Acids Res.* (Aug 2012).
- [19] KOZANITIS, C., HEIBERG, A., VARGHESE, G., AND BAFNA, V. Using genome query language to uncover genetic variation. *Bioinformatics* (2013).
- [20] KOZANITIS, C., SAUNDERS, C., KRUGLYAK, S., BAFNA, V., AND VARGHESE, G. Compressing genomic sequence fragments using SlimGene. *J. Comput. Biol.* 18, 3 (Mar 2011), 401–413.

- [21] LAMB, A., FULLER, M., VARADARAJAN, R., TRAN, N., VANDIVER, B., DOSHI, L., AND BEAR, C. The vertica analytic database: C-store 7 years later. *Proceedings of the VLDB Endowment* 5, 12 (2012), 1790–1801.
- [22] LI, H., HANDSAKER, B., WYSOKER, A., FENNEL, T., RUAN, J., HOMER, N., MARTH, G., ABECASIS, G., DURBIN, R., ET AL. The sequence alignment/map format and SAMtools. *Bioinformatics* 25, 16 (2009), 2078–2079.
- [23] MCKENNA, A., HANNA, M., BANKS, E., SIVACHENKO, A., CIBULSKIS, K., KERNYTSKY, A., GARIMELLA, K., ALTSHULER, D., GABRIEL, S., DALY, M., ET AL. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome research* 20, 9 (2010), 1297–1303.
- [24] NHGRI. DNA sequencing costs. <http://www.genome.gov/sequencingcosts/>.
- [25] NIEMENMAA, M., KALLIO, A., SCHUMACHER, A., KLEMELÄ, P., KORPELAINEN, E., AND HELJANKO, K. Hadoop-BAM: directly manipulating next generation sequencing data in the cloud. *Bioinformatics* 28, 6 (2012), 876–877.
- [26] POPITSCH, N., AND VON HAESELER, A. NGC: lossless and lossy compression of aligned high-throughput sequencing data. *Nucleic Acids Res.* (Oct 2012).
- [27] SMITH, C. C., WANG, Q., CHIN, C.-S., SALERNO, S., DAMON, L. E., LEVIS, M. J., PERL, A. E., TRAVERS, K. J., WANG, S., HUNT, J. P., ET AL. Validation of ITD mutations in FLT3 as a therapeutic target in human acute myeloid leukaemia. *Nature* 485, 7397 (2012), 260–263.
- [28] SMITH, T. F., AND WATERMAN, M. S. Comparison of biosequences. *Advances in Applied Mathematics* 2, 4 (1981), 482–489.
- [29] SNITKIN, E. S., ZELAZNY, A. M., THOMAS, P. J., STOCK, F., HENDERSON, D. K., PALMORE, T. N., SEGRE, J. A., ET AL. Tracking a hospital outbreak of carbapenem-resistant klebsiella pneumoniae with whole-genome sequencing. *Science translational medicine* 4, 148 (2012), 148ra116–148ra116.
- [30] STONEBRAKER, M., ABADI, D. J., BATKIN, A., CHEN, X., CHERNIACK, M., FERREIRA, M., LAU, E., LIN, A., MADDEN, S., O’NEIL, E., ET AL. C-store: a column-oriented DBMS. In *Proceedings of the 31st international conference on Very large data bases* (2005), VLDB Endowment, pp. 553–564.
- [31] VAN ALLEN, E. M., WAGLE, N., AND LEVY, M. A. Clinical analysis and interpretation of cancer genome data. *Journal of Clinical Oncology* 31, 15 (2013), 1825–1833.
- [32] WAN, R., ANH, V. N., AND ASAI, K. Transformations for the compression of FASTQ quality scores of next-generation sequencing data. *Bioinformatics* 28, 5 (Mar 2012), 628–635.
- [33] WHEELER, D. L., ET AL. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* 36 (Jan 2008), 13–21.
- [34] XIN, R. S., GONZALEZ, J. E., FRANKLIN, M. J., STOICA, I., AND AMPLAB, E. Graphx: A resilient distributed graph system on spark.
- [35] XIN, R. S., ROSEN, J., ZAHARIA, M., FRANKLIN, M. J., SHENKER, S., AND STOICA, I. Shark: SQL and rich analytics at scale. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 2013), SIGMOD ’13, ACM, pp. 13–24.
- [36] YANOVSKY, V. ReCoil - an algorithm for compression of extremely large datasets of dna data. *Algorithms Mol Biol* 6 (2011), 23.
- [37] ZAHARIA, M., BOLOSKEY, W. J., CURTIS, K., FOX, A., PATTERSON, D., SHENKER, S., STOICA, I., KARP, R. M., AND SITTTLER, T. Faster and more accurate sequence alignment with snap. *arXiv preprint arXiv:1111.5572* (2011).
- [38] ZAHARIA, M., CHOWDHURY, M., DAS, T., DAVE, A., MA, J., MCCAULEY, M., FRANKLIN, M., 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* (2012), USENIX Association, p. 2.
- [39] 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* (2010), p. 10.
- [40] ZIMMERMANN, H. OSI reference model—the ISO model of architecture for open systems interconnection. *IEEE Transactions on Communications* 28, 4 (1980), 425–432.