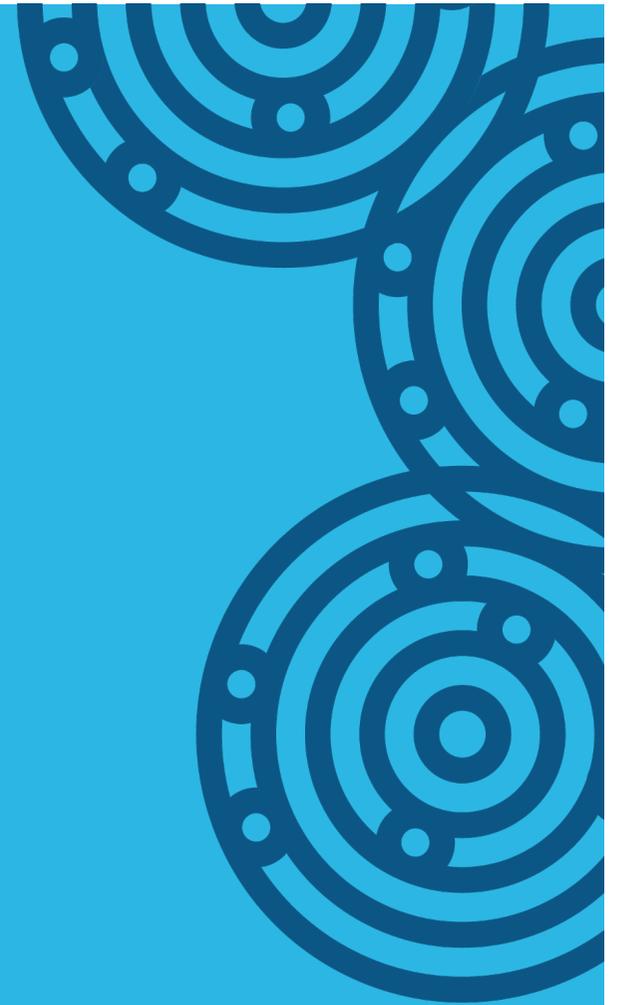


cloudera®

Petascale Analytics in Genomics with Hadoop

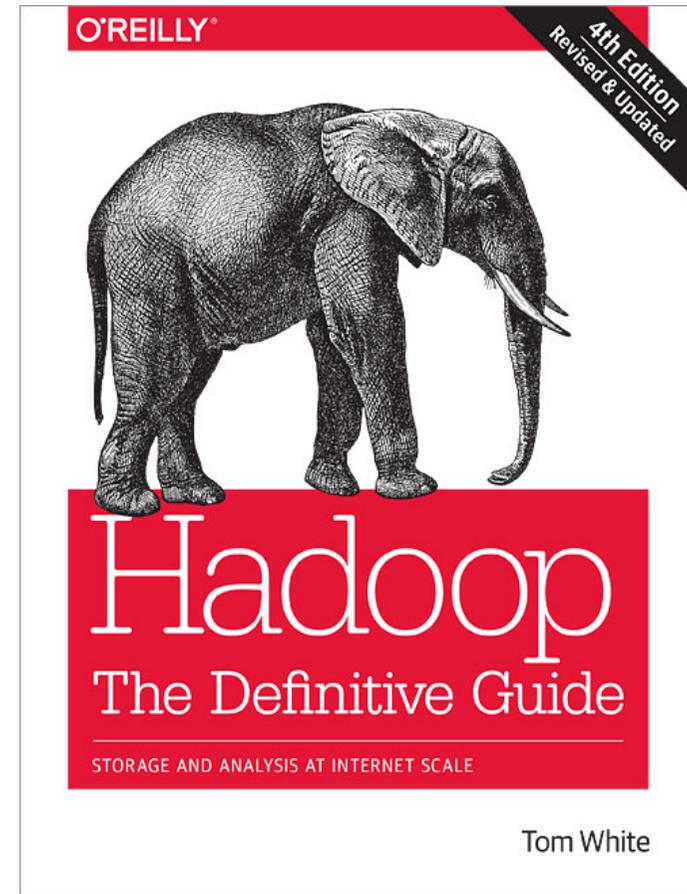
2 June 2016, Strata+Hadoop World, London

Tom White | @tom_e_white

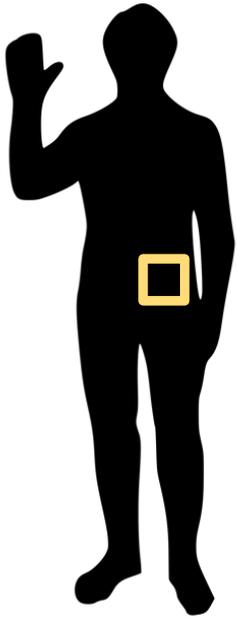


About Me

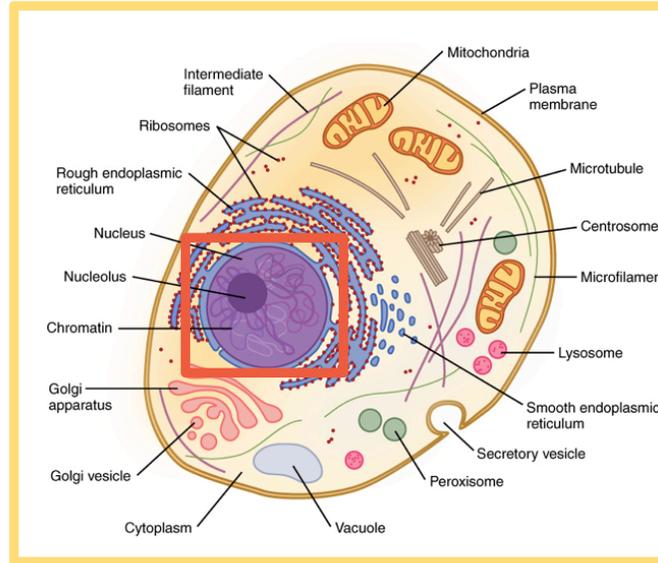
- Data Science Team at Cloudera
- Apache Hadoop Committer, PMC Member, Apache Member
- Author of “Hadoop: The Definitive Guide”



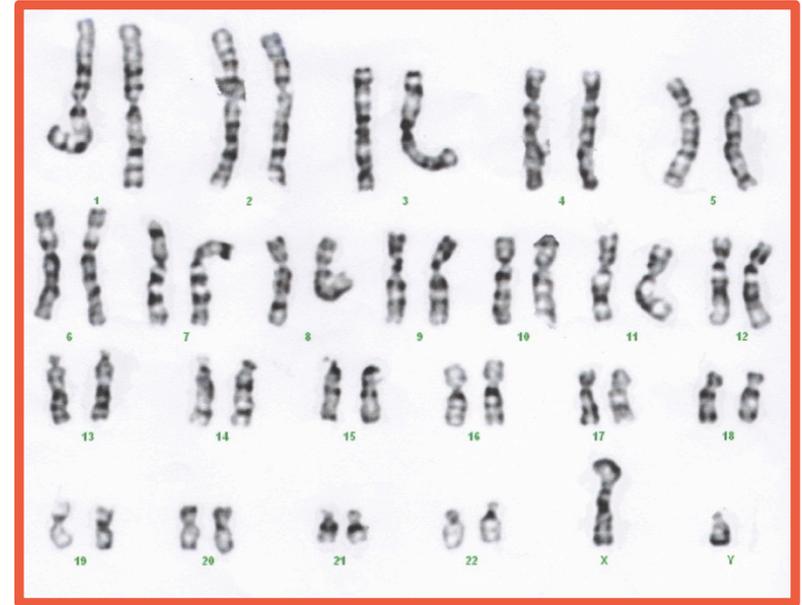
What is **genomics**?



Organism



Cell



Genome

cloudera

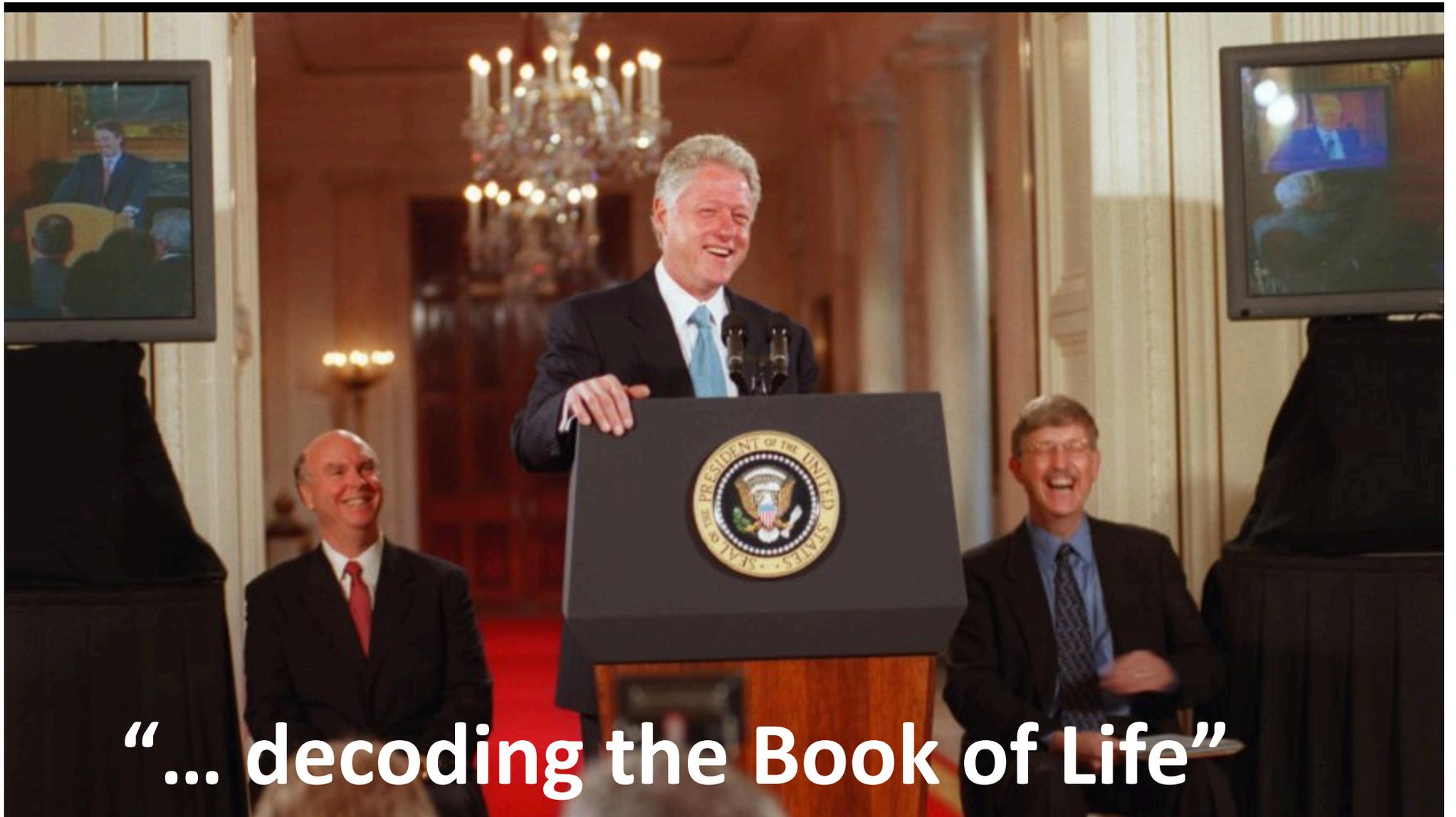




Reference chromosome

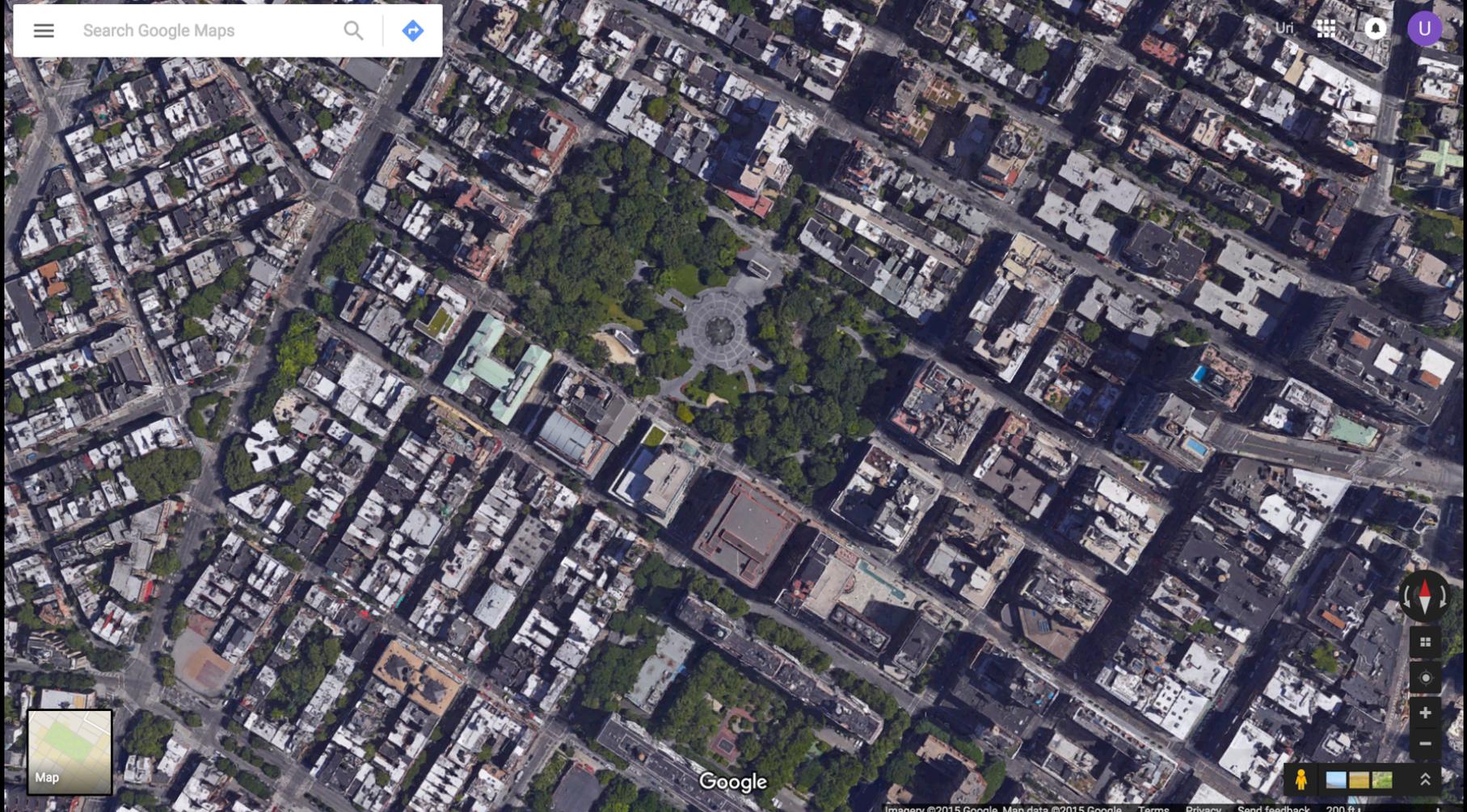


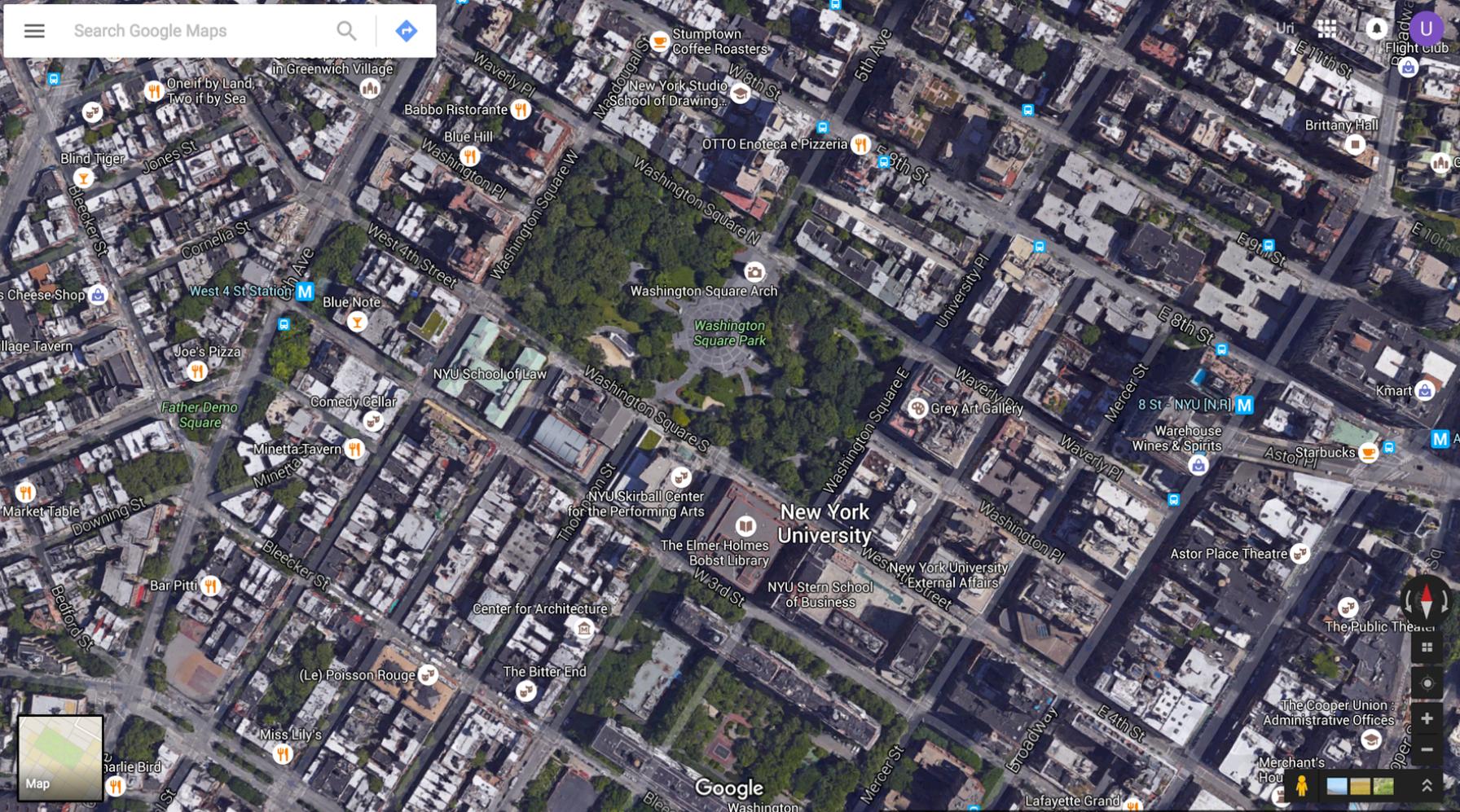
Location



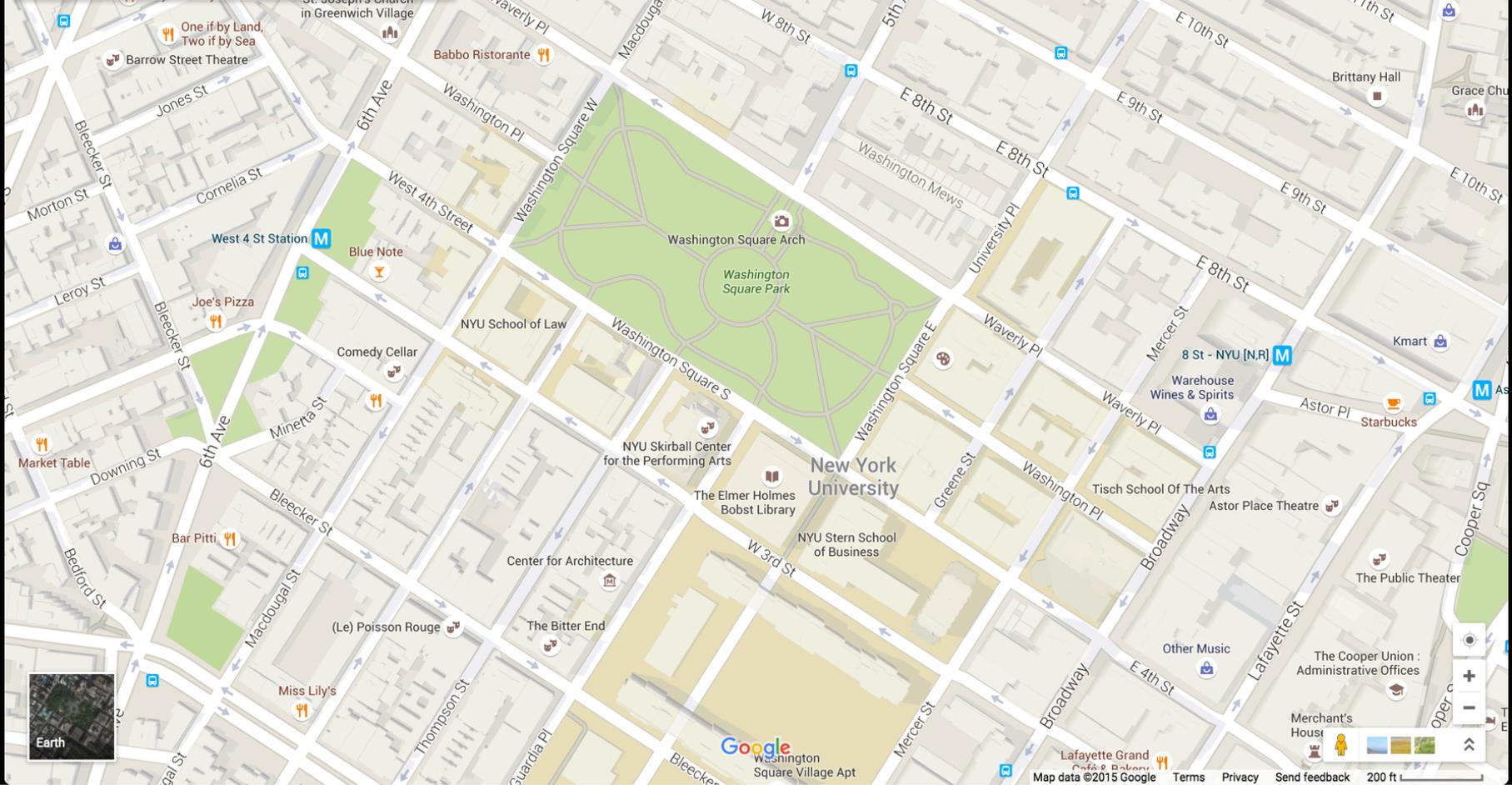
“... decoding the Book of Life”

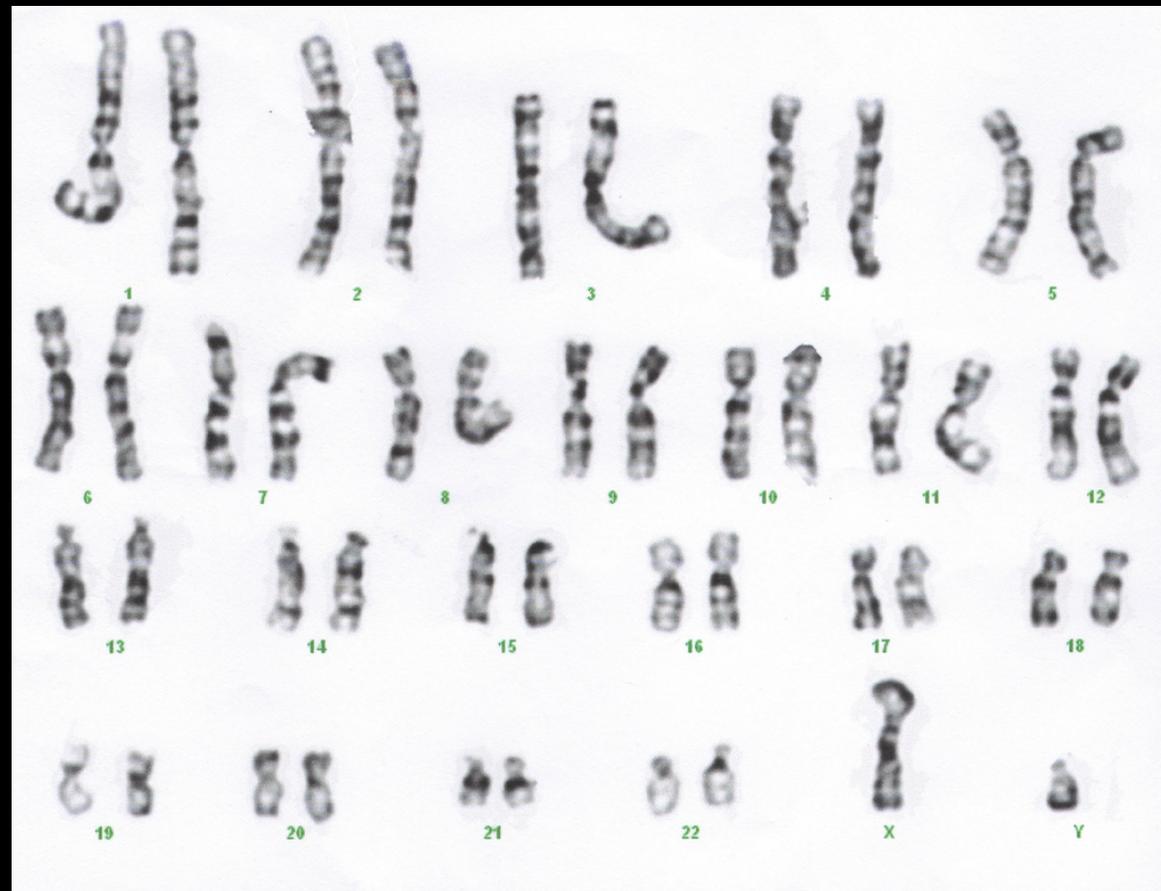
Search Google Maps





Search Google Maps

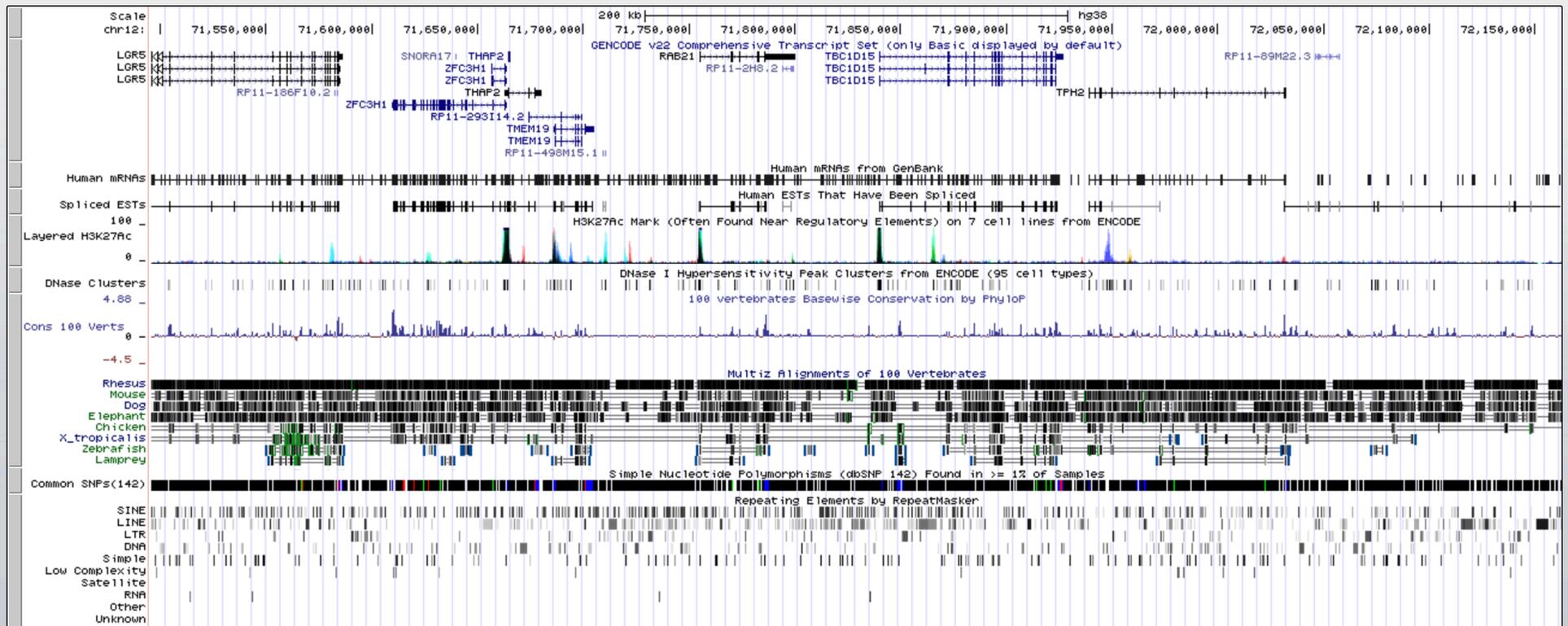




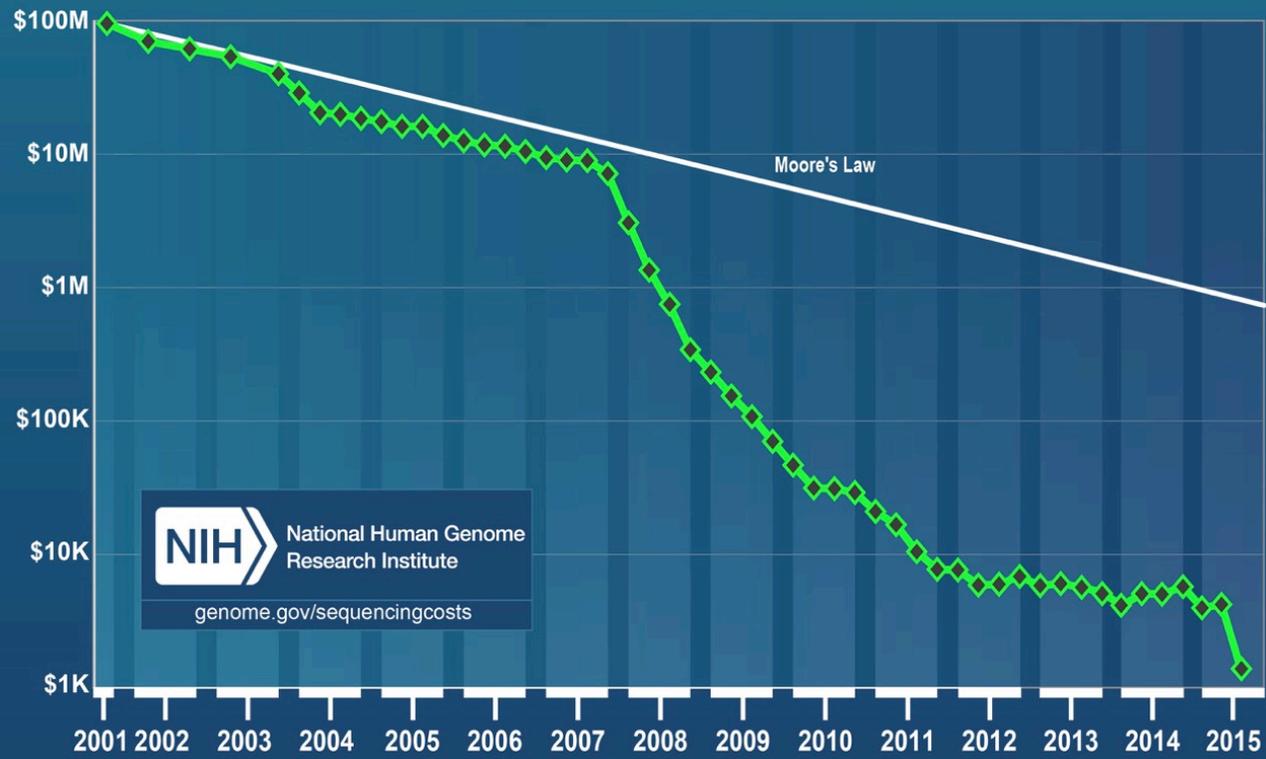
UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

move <<< << < > >> >>> ZOOM in 1.5x 3x 10x base ZOOM out 1.5x 3x 10x 100x

chr12:71,495,864-72,162,530 666,667 bp. enter position, gene symbol or search terms go



Cost per Genome



What is **bioinformatics**?



```
>read1  
TTGGACATTTTCGGGGTCTCAGATT  
>read2  
AATGTTGTTAGAGATCCGGGATTT  
>read3  
GGATTCCCCGCCGTTTGAGAGCCT  
>read4  
AGGTTGGTACCGCGAAAAGCGCAT
```



Bioinformatics!

Pipelines!



[Skip to content](#)

Font size: A⁻ A⁺ Contrast: C C C C

[Home](#) in v y



[About Us](#) | [100,000 Genomes Project](#) | [Research](#) | [Industry Partnerships](#) | [Library & resources](#) | [News & Events](#)

[Home](#) > [The 100,000 Genomes Project](#)

The 100,000 Genomes Project

The project will sequence 100,000 genomes from around 70,000 people. Participants are NHS patients with a rare disease, plus their families, and patients with cancer.

The aim is to create a new genomic medicine service for the NHS – transforming the way people are

Understanding genomics

Our Head of Engagement, Vivienne Parry, explains more about genomics in this film courtesy of our partners at Health Education England.

This website uses cookies to improve your experience. By continuing to use the site, you agree to the use of cookies.

[Accept](#)

[Read More](#)

It's pipelines all the way down!

Node 1

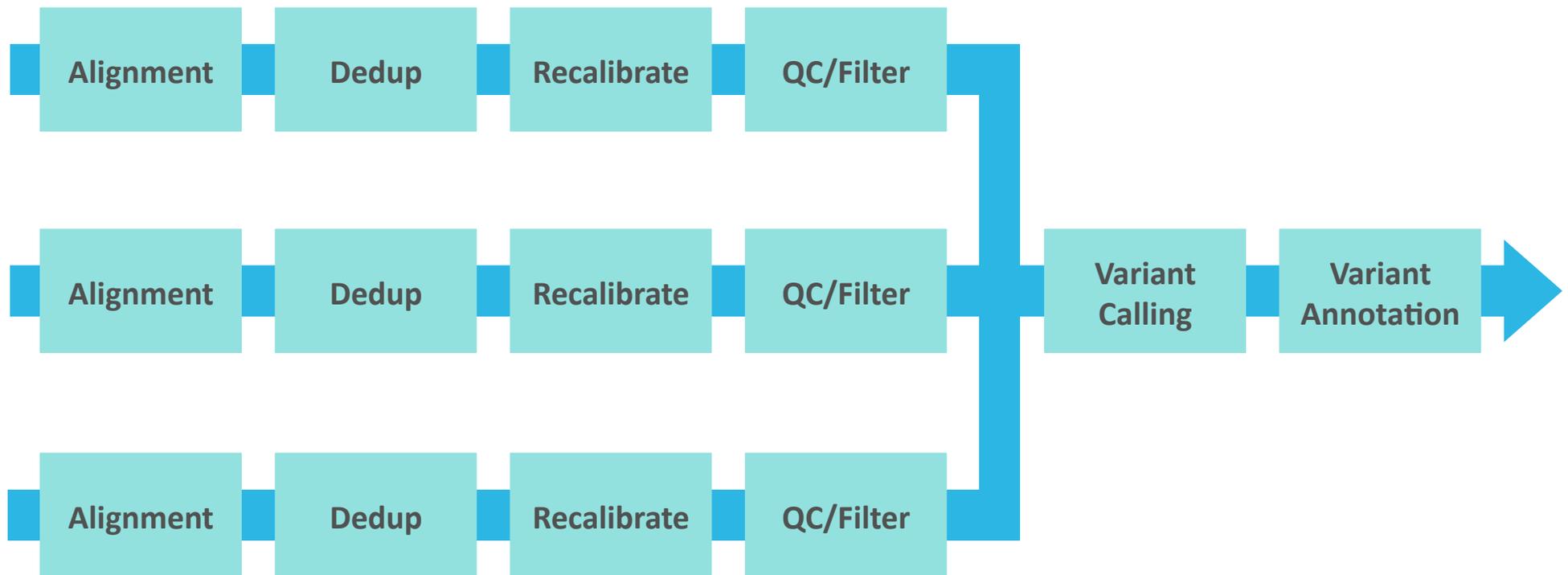


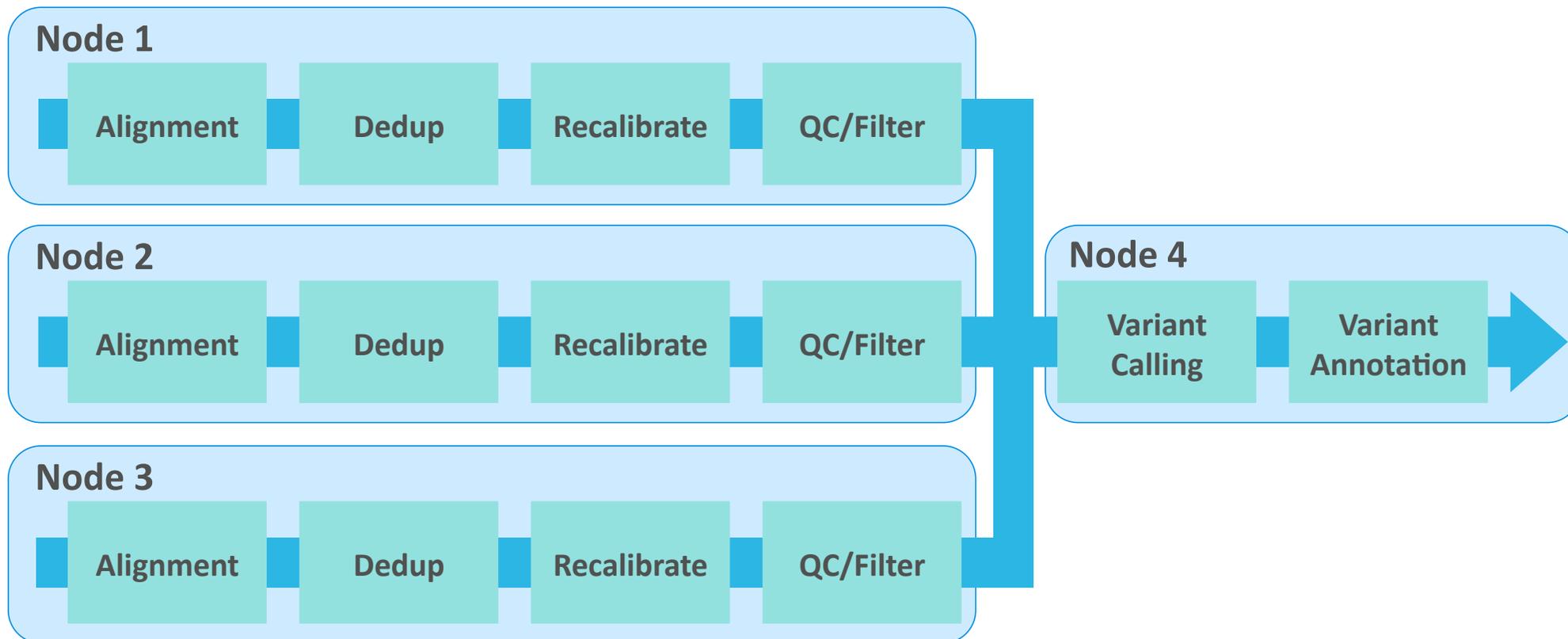
Node 2

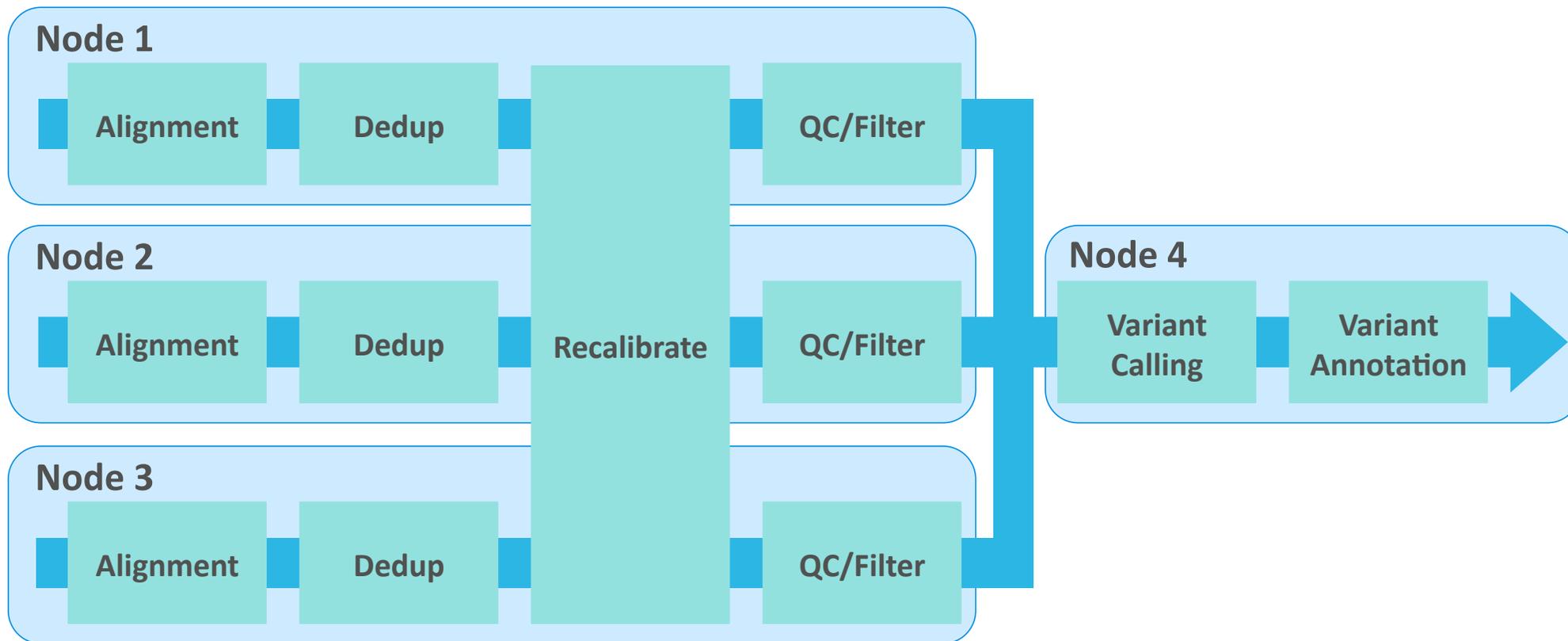


Node 3









How can Hadoop be used in bioinformatics?

Genomics on Hadoop – A Potted History

- 2010 - Hadoop-BAM - MR input/output formats for bio (BAM, VCF, etc)
- 2011 - Seal - MR tools for reads
- 2012 - SeqPig - Pig interface for Hadoop-BAM
- 2013 - ADAM - a genomics analysis platform on Spark, Avro, and Parquet
- 2013 - OpenCGA - a variant store built on HBase
- 2014 - Halvade - a tool to run the GATK best practices pipeline using MR
- 2014 - Guacamole - Spark variant caller for ADAM
- 2015 - GATK4 - a toolkit for running genomics pipelines on Spark
- 2016 - Hail - PLINK-like tool for whole genome association analysis

Spark + Genomics = ADAM

- Hosted at Berkeley and the AMPLab
- Apache 2 License
- Contributors from both research and commercial organizations
- Core spatial primitives, variant calling
- Avro and Parquet for data models and file formats

cloudera

The screenshot displays the GitHub profile for Big Data Genomics. At the top, the repository name "Big Data Genomics" is shown with a settings gear icon and the URL "http://bdgenomics.org/". Below this is a search bar with "Find a repository..." and a "+ New repository" button. The main content area lists several repositories:

- adam**: A genomics processing engine and specialized file format built using Apache Avro, Apache Spark and Parquet. Apache 2 licensed. Updated 2 hours ago. (Scala, 234 stars, 83 forks)
- PacMin**: Assembler for PacBio reads. Apache 2 licensed. Updated 3 days ago. (Scala, 1 star, 1 fork)
- eggo**: Ready-to-go Parquet-formatted public 'omics datasets. Updated 5 days ago. (Python, 5 stars, 3 forks)
- recipes**: Recipes using BDG projects. Apache 2 licensed. Updated 6 days ago. (Shell, 1 star, 3 forks)

On the right side, there are two panels:

- People**: A grid of 21 contributor avatars, including the "amp lab" logo. An "Invite someone" button is located below the grid.
- Teams**: A search bar "Jump to a team" and three team listings:
 - Owners**: 4 members · 15 repositories
 - ADAM Committers**: 8 members · 2 repositories
 - avocado committers**: 8 members · 1 repository

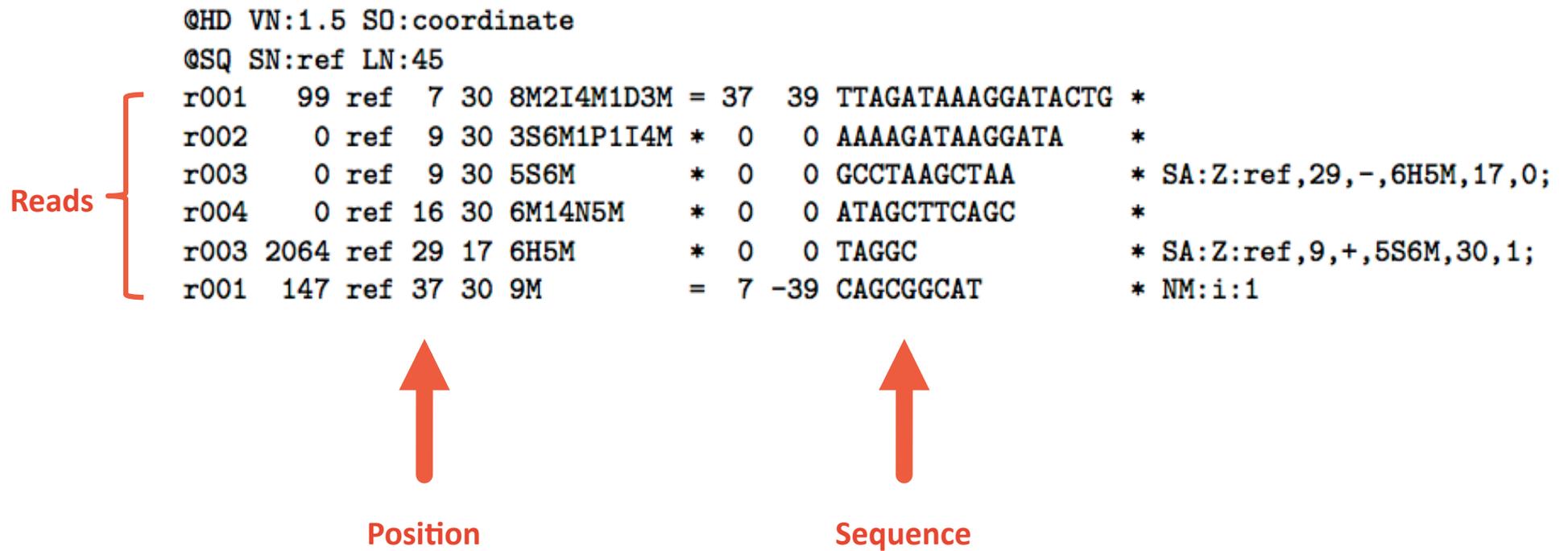
Genome Analysis Toolkit (GATK)

- Developed by the Broad Institute
- Core is MIT license, some proprietary tools on top
- Version 4 has been re-written to use Spark, now competitive with ADAM for speed
- Uses existing bio file formats for input and output, but Spark RDDs for intermediate data

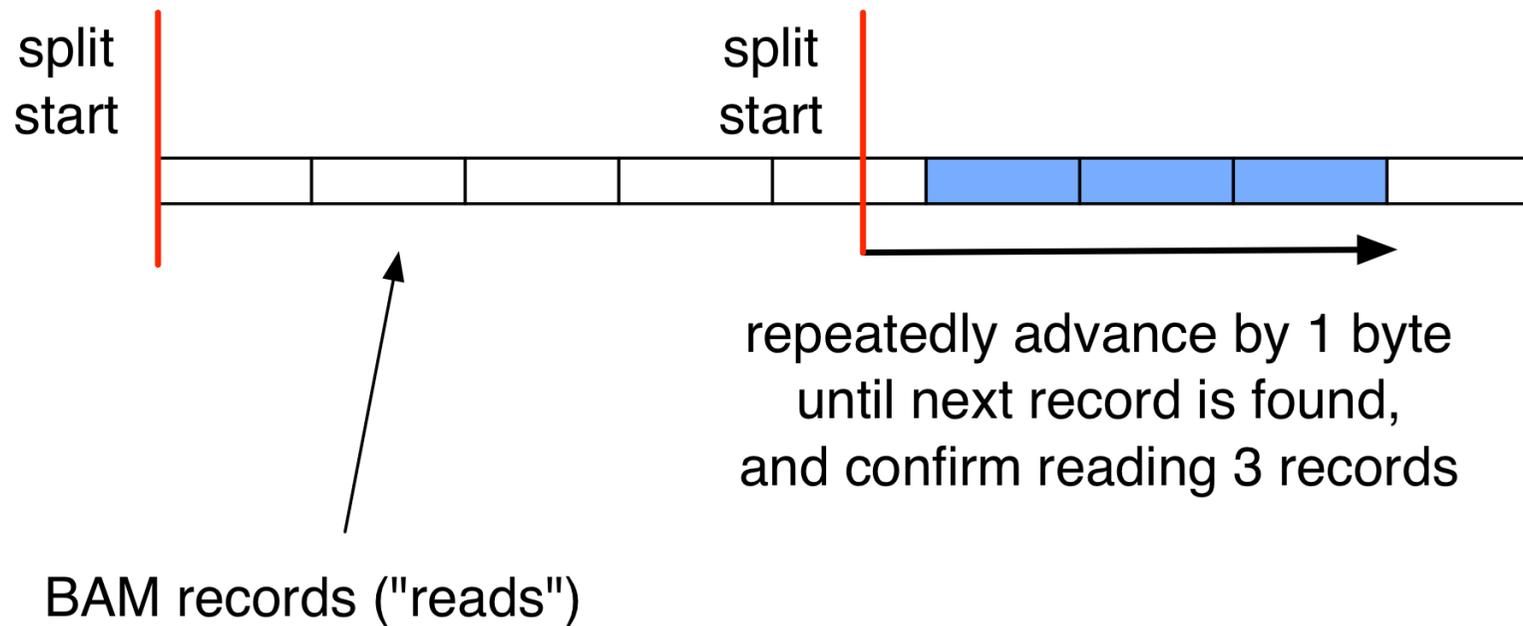
Bioinformatics File Formats

- Hand crafted
- Poorly specified
- Text based
- Unsplittable (in the Hadoop sense)

BAM files



Example: splitting BAM files

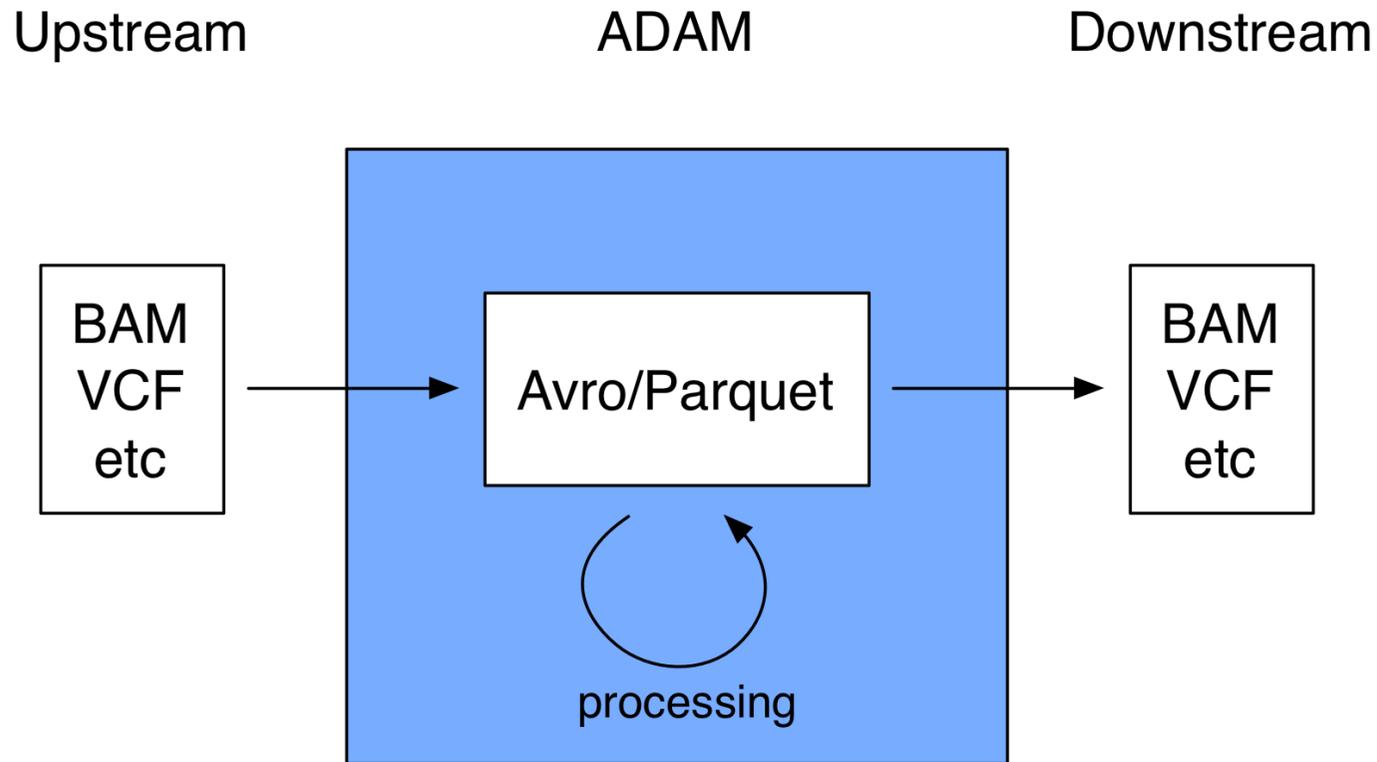


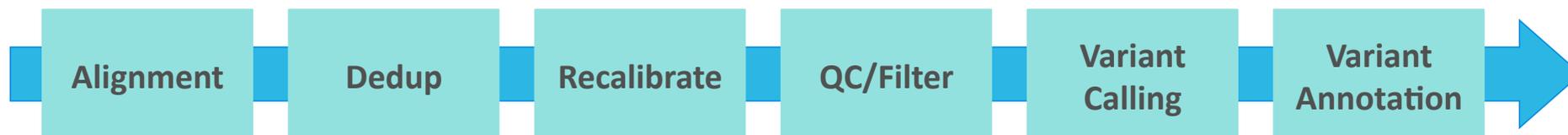
Example: splitting BAM files (BGZF compression)



header includes block length (as gzip extension),
so easy to index and split

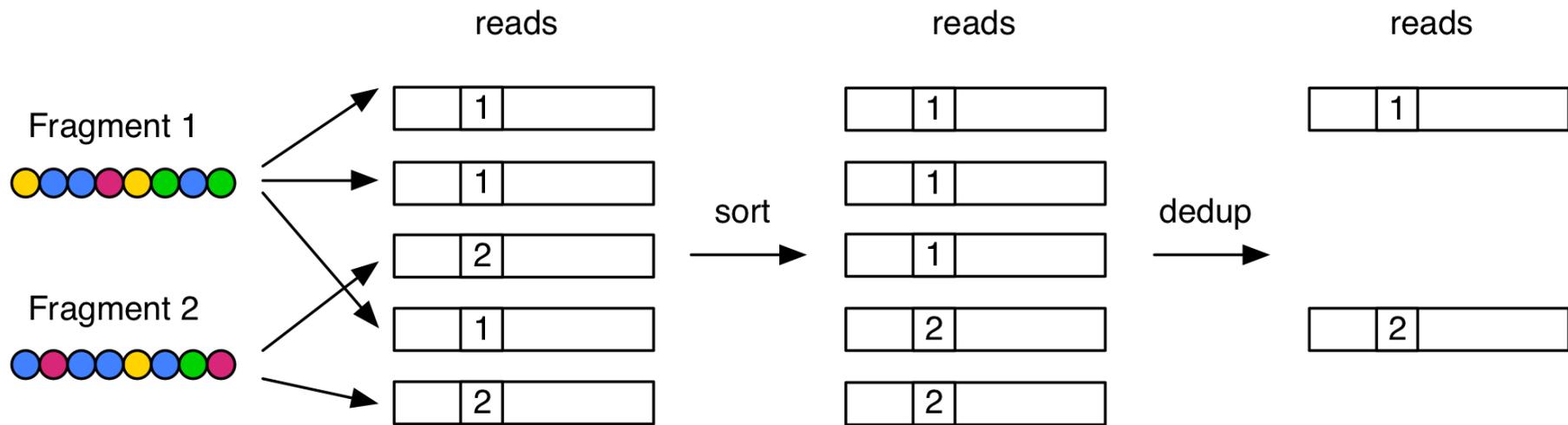
Why not use Hadoop formats?





Dedup

Mark Duplicates



Method

```
/**  
 * Main work method. Reads the BAM file once and collects sorted information about  
 * the 5' ends of both ends of each read (or just one end in the case of pairs).  
 * Then makes a pass through those determining duplicates before re-reading the  
 * input file and writing it out with duplication flags set correctly.  
 */
```

```
protected int doWork() {  
    // build some data structures  
    buildSortedReadEndLists(useBarcodes);  
    generateDuplicateIndexes(useBarcodes);  
  
    final SAMFileWriter out =  
        new SAMFileWriterFactory().makeSAMOrBAMWriter(outputHeader, true, OUTPUT);  
    final CloseableIterator<SAMRecord> iterator = headerAndIterator.iterator;  
    while (iterator.hasNext()) {  
        final SAMRecord rec = iterator.next();  
        if (!rec.isSecondaryOrSupplementary()) {  
            if (recordInFileIndex == nextDuplicateIndex) {  
                rec.setDuplicateReadFlag(true);  
                // Now try and figure out the next duplicate index  
                if (this.duplicateIndexes.hasNext()) {  
                    nextDuplicateIndex = this.duplicateIndexes.next();  
                } else {  
                    // Only happens once we've marked all the duplicates  
                    nextDuplicateIndex = -1;  
                }  
            } else {  
                rec.setDuplicateReadFlag(false);  
            }  
        }  
        recordInFileIndex++;  
        if (!this.REMOVE_DUPLICATES || !rec.getDuplicateReadFlag()) {  
            out.addAlignment(rec);  
        }  
    }  
}
```

Code

cloudera

```
@Option(shortName = "MAX_FILE_HANDLES",
        doc = "Maximum number of file handles to keep open when spilling " +
              "read ends to disk. Set this number a little lower than the " +
              "per-process maximum number of file that may be open. This " +
              "number can be found by executing the 'ulimit -n' command on " +
              "a Unix system.")
public int MAX_FILE_HANDLES_FOR_READ_ENDS_MAP = 8000;
```



Spark Implementation

```
JavaPairRDD<String, Iterable<GATKRead>> keyedReads = ...;
```

```
JavaPairRDD<String, PairedEnds> keyPairs =  
    keyedReads.flatMapToPair(keyedRead -> { ... });
```

```
JavaPairRDD<String, Iterable<PairedEnds>> keyedPairs =  
    keyPairs.groupByKey(numReducers);
```

```
JavaRDD<GATKRead> markedDups = markPairedEnds(keyedPairs,  
    scoringStrategy, finder, header);
```

Lessons Learned

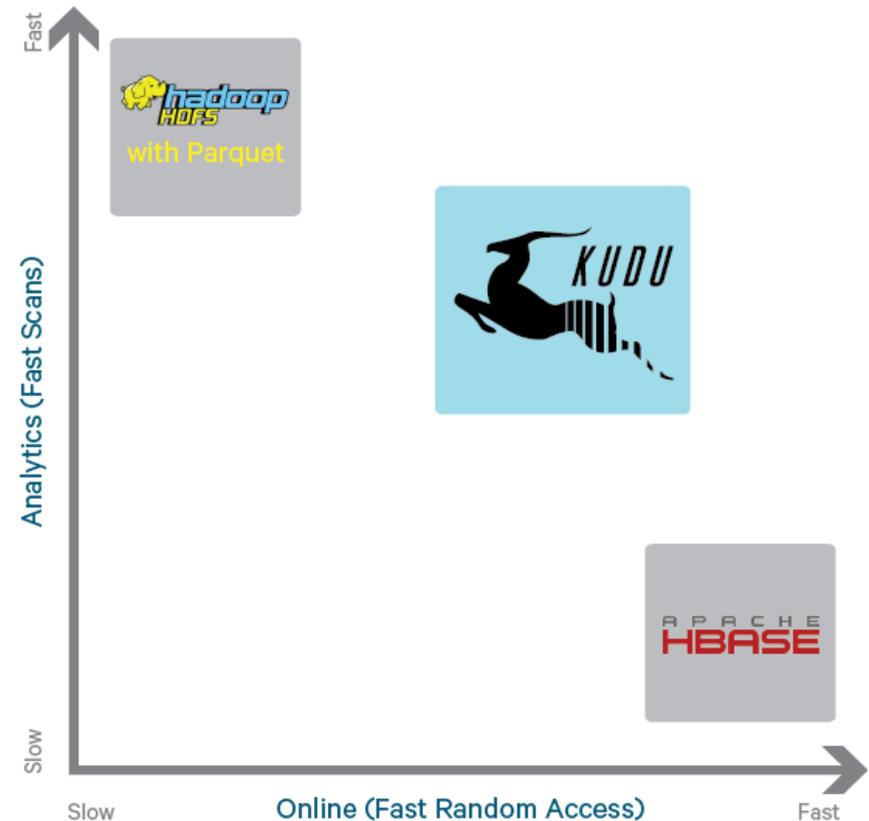
1. Figure out how to read and write existing formats efficiently
2. Spark is a great API, but developers need to understand consequences of e.g. the shuffle, serialization cost
3. Work with domain experts, on existing projects, if possible

Future developments

Kudu for Variant Stores

- Kudu fills gap between HDFS and HBase
- Fast scans and updateable
- Add new annotations to genomics data (variants) without rewriting whole dataset
- Key = genome position
- Range partitioning

cloudera



Hail

- Scalable variant analytics in Spark
- Command-line tools like PLINK
- Parquet-based storage by default, other storage possibilities like Kudu

Links

- ADAM
 - <https://github.com/bigdatagenomics/adam>
- GATK4
 - <https://github.com/broadinstitute/gatk>

Acknowledgements

UCBerkeley

Matt Massie
Frank Nothaft
Michael Heuer

Tamr

Timothy Danford

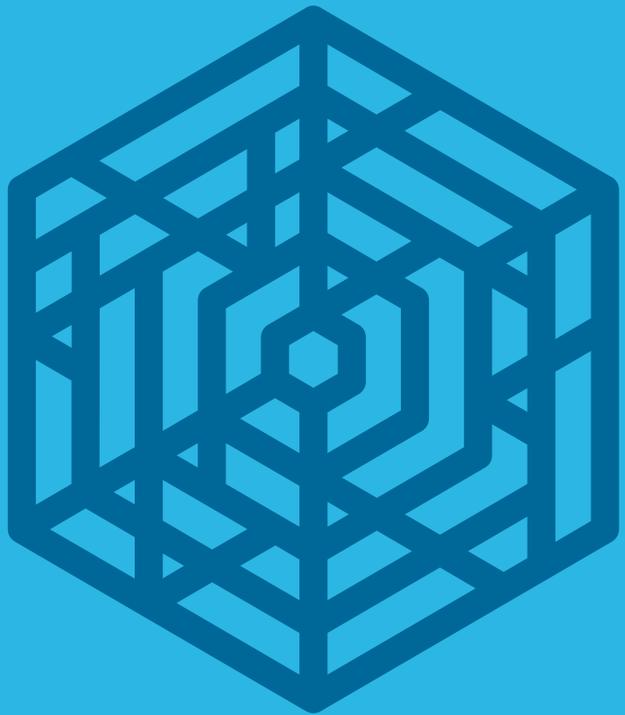
MSSM

Jeff Hammerbacher
Ryan Williams

Cloudera

Uri Laserson

Sandy Ryza
Sean Owen



cloudera

Thank you

@tom_e_white

tom@cloudera.com