

# Genomic Data Analysis on Spark+Hadoop

Ryan Williams

DataEngConf NYC

11/4/2016



# Agenda

- Intro
- Genomics crash course
- Guacamole: somatic mutation calling on Spark
- Other applications / interesting algorithms

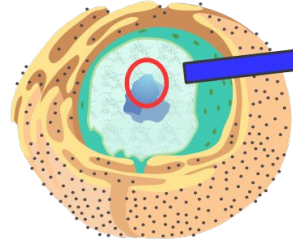
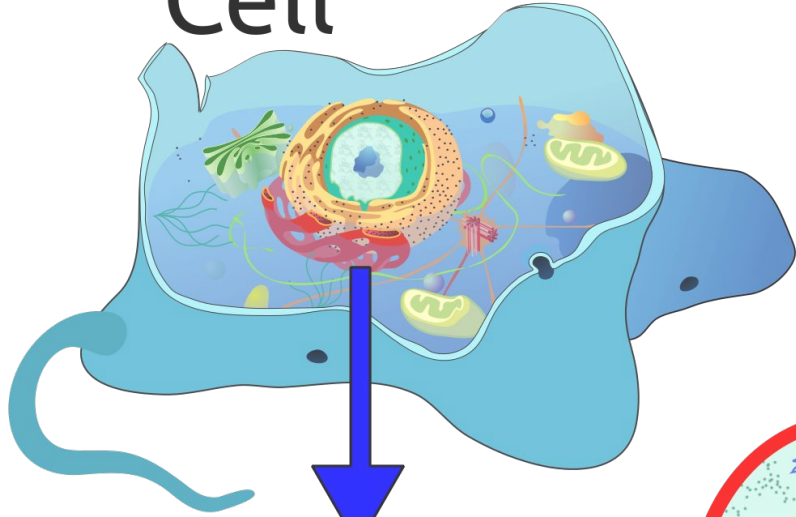
# Hammer Lab

- Computational lab in the department of Genetics and Genomic Sciences at Mount Sinai
- Principal investigator: Jeff Hammerbacher
- Focus on informatics for cancer immunotherapy
- Software developed at [github.com/hammerlab](https://github.com/hammerlab)

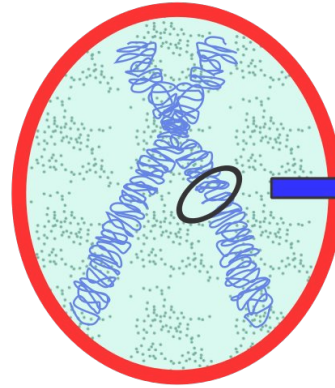


# Genomics / Sequencing Overview

Cell

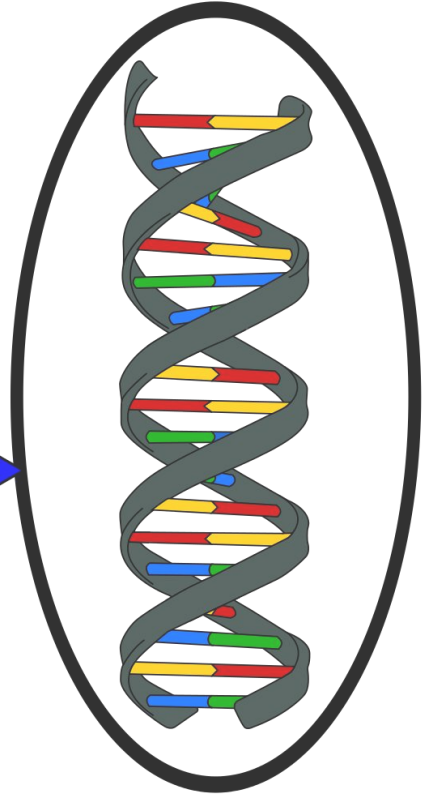


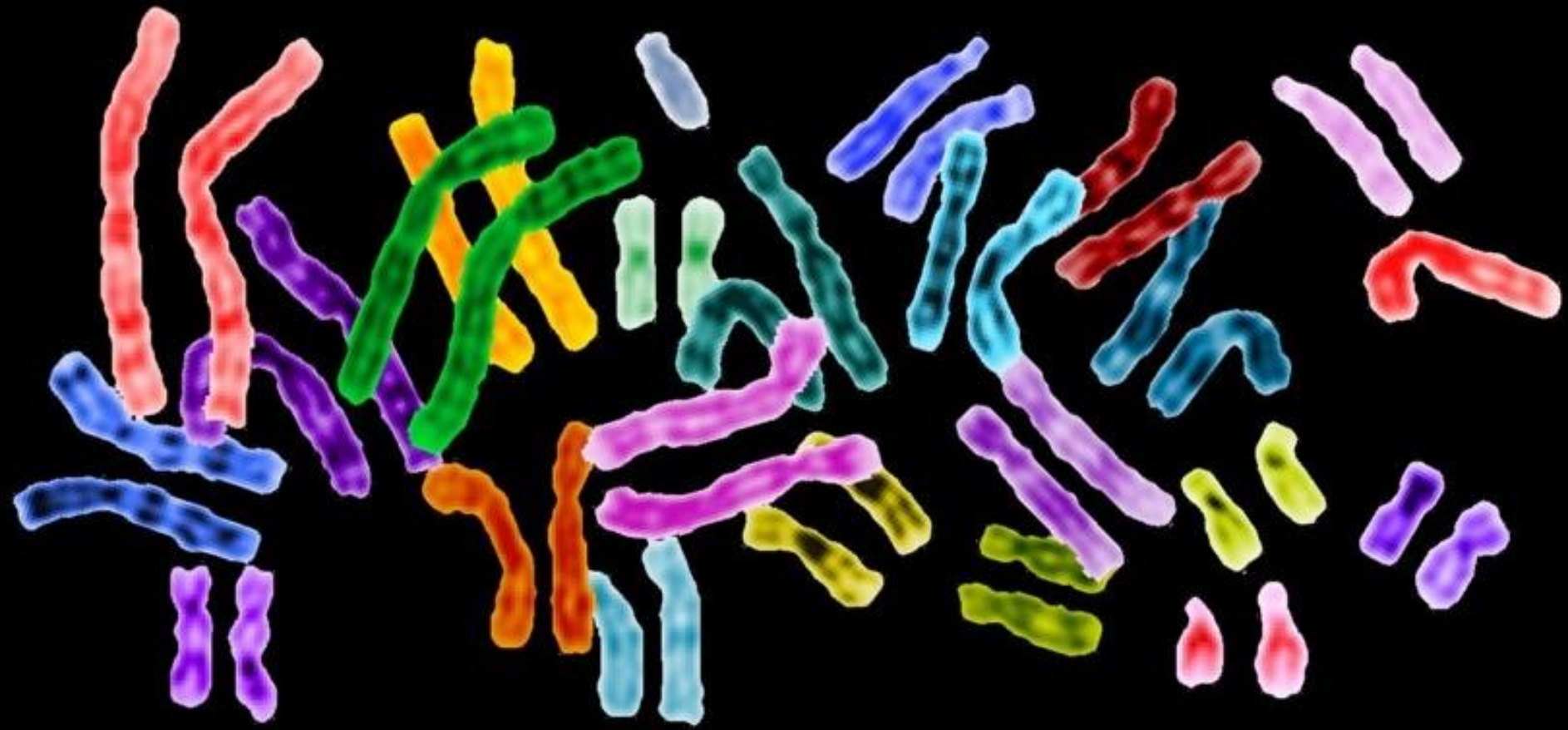
Nucleus



Chromosome

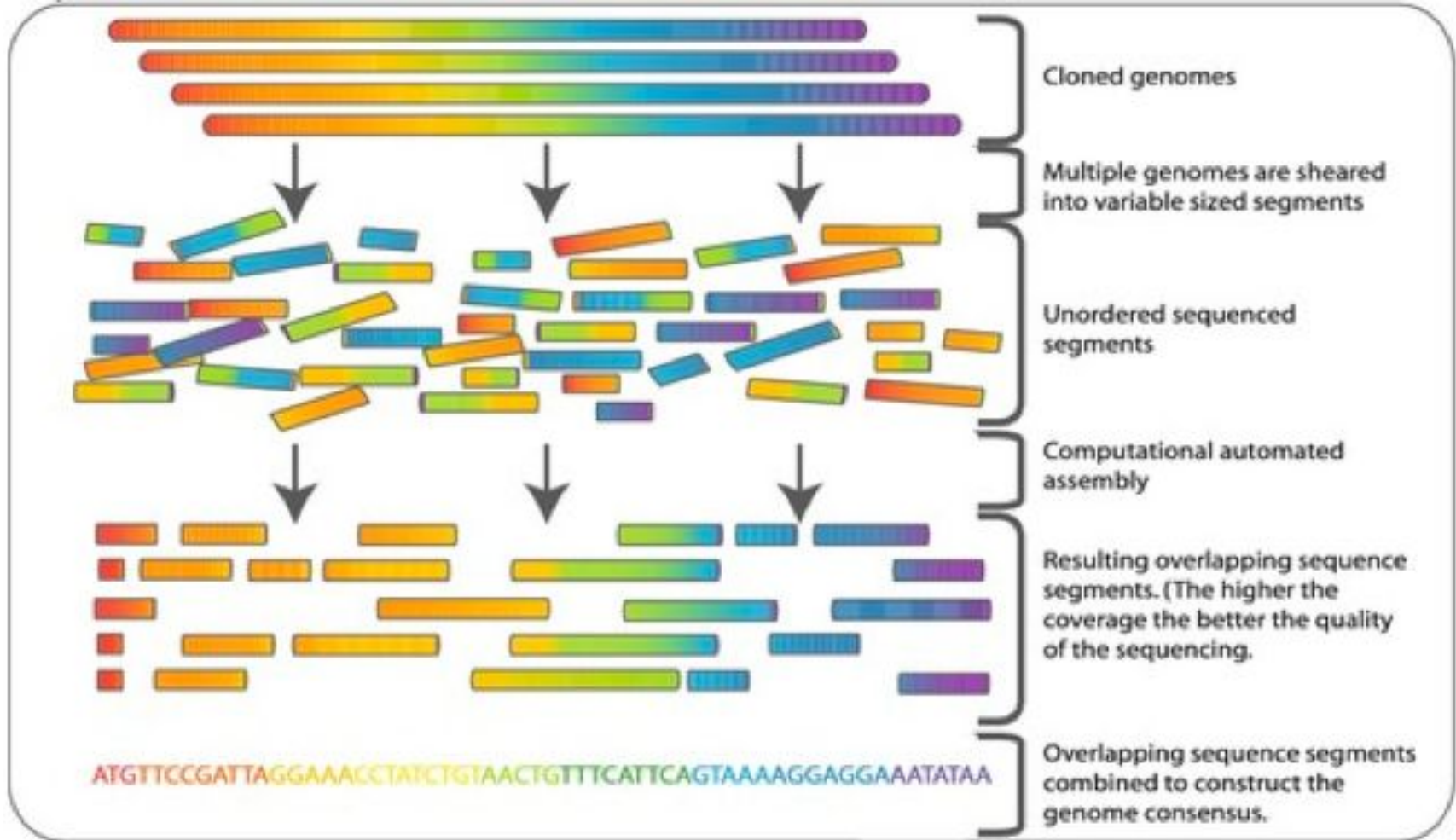
DNA





<http://blogs.plos.org/dnascience/2016/01/21/can-a-quirky-chromosome-create-a-second-human-species/>





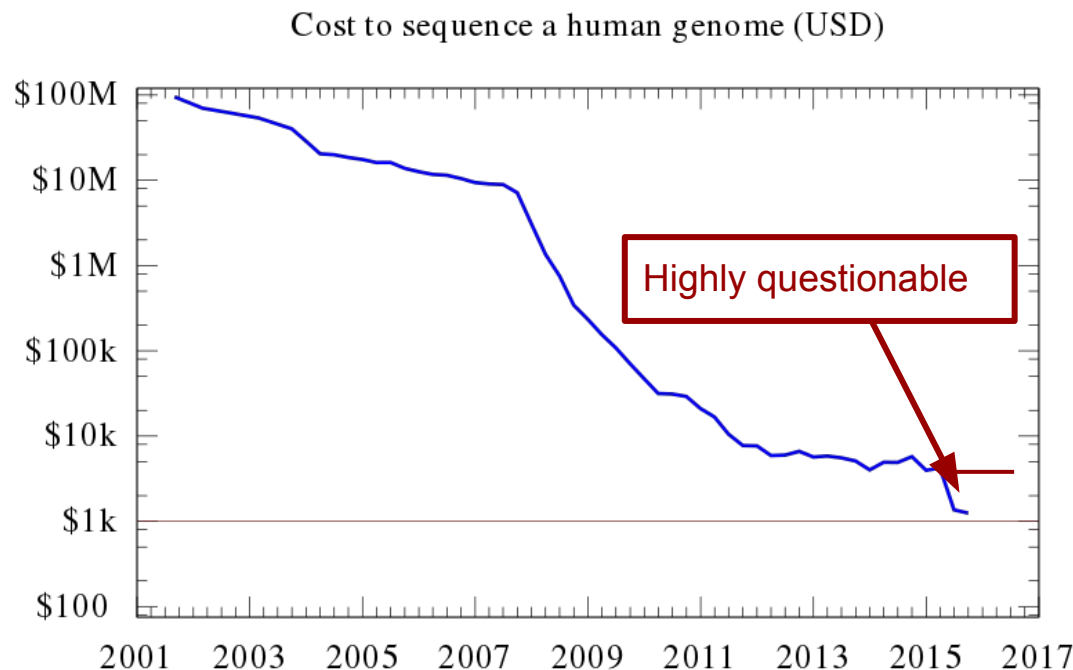
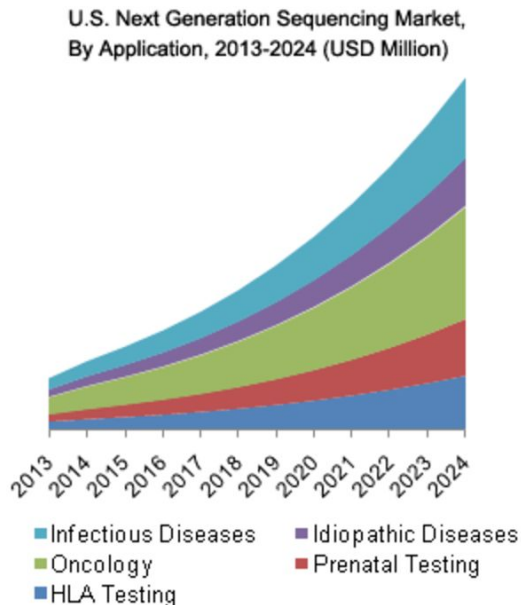
# Storing Human Genomes

- 1 genome  $\approx$  3B base-pairs
- Theory:
  - “2 bits per base-pair” (A, C, G, T)
  - $\Rightarrow$  1 genome  $\approx$  750MB
  - $<1\%$  unique, person to person
  - 7BN genomes  $\approx$  50PB
- Reality:
  - 1BN 100bp “reads”
    - $\Rightarrow$  100BN sequenced base-pairs
    - Cover the genome at average depth 30 (“30x coverage”)
  - 2-bit base, 1-byte quality score  $\Rightarrow$  **100GB / genome**
  - **100-100k genomes  $\Rightarrow$  10TB-10PB**



# Sequencing Human Genomes

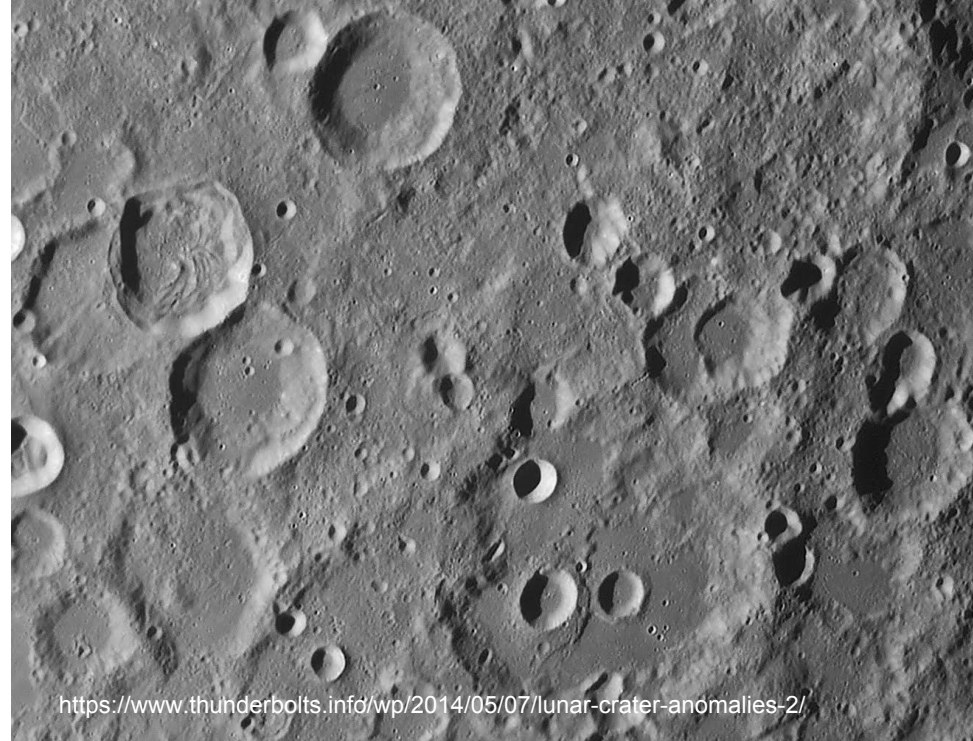
- Human Genome Project, 1990-2003
- 1000 Genomes Project, 2008-2012
- 100k Genomes Project, 2012-???

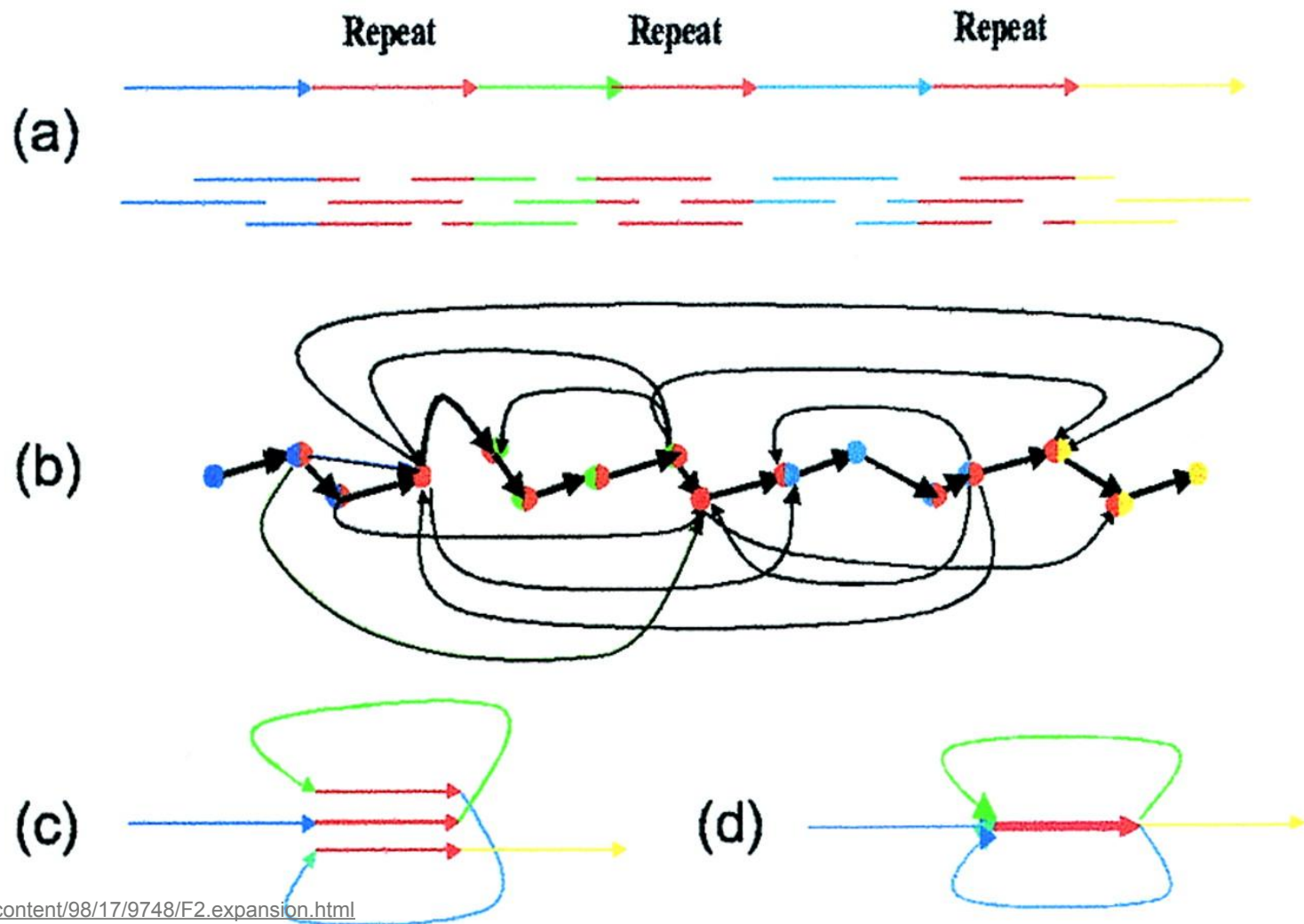


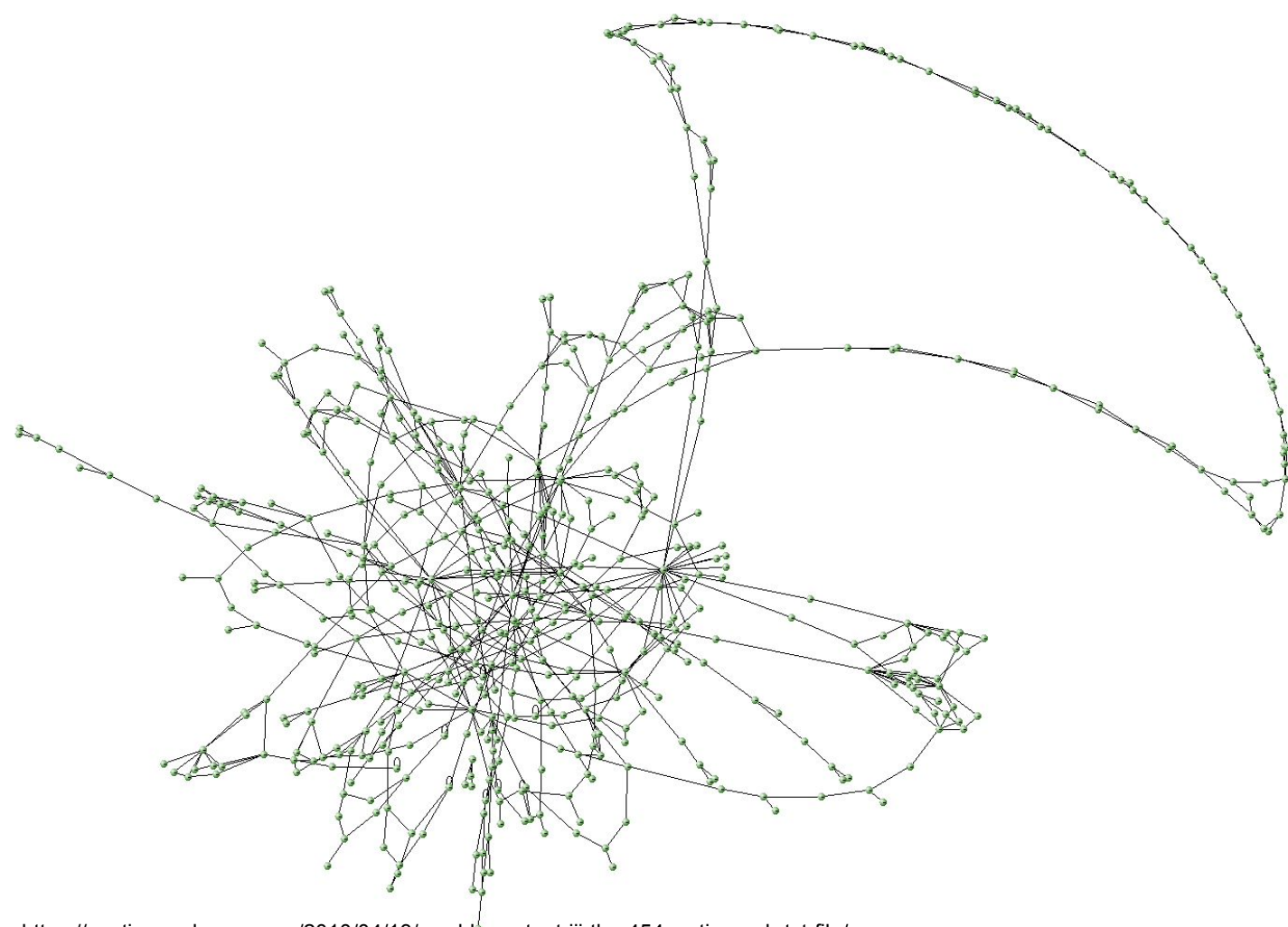
# Genome Alignment / Assembly

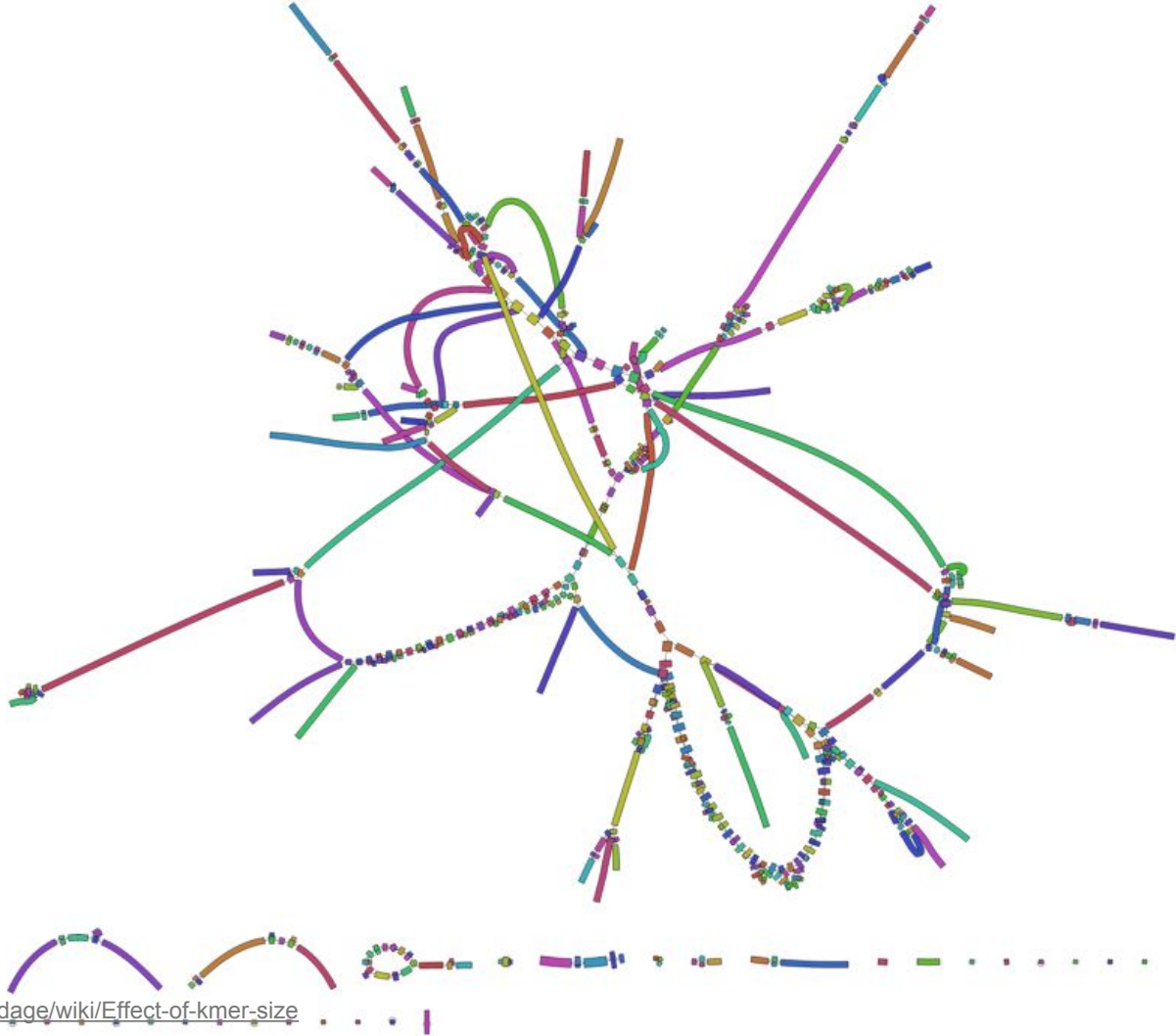
# Genome structure makes things difficult

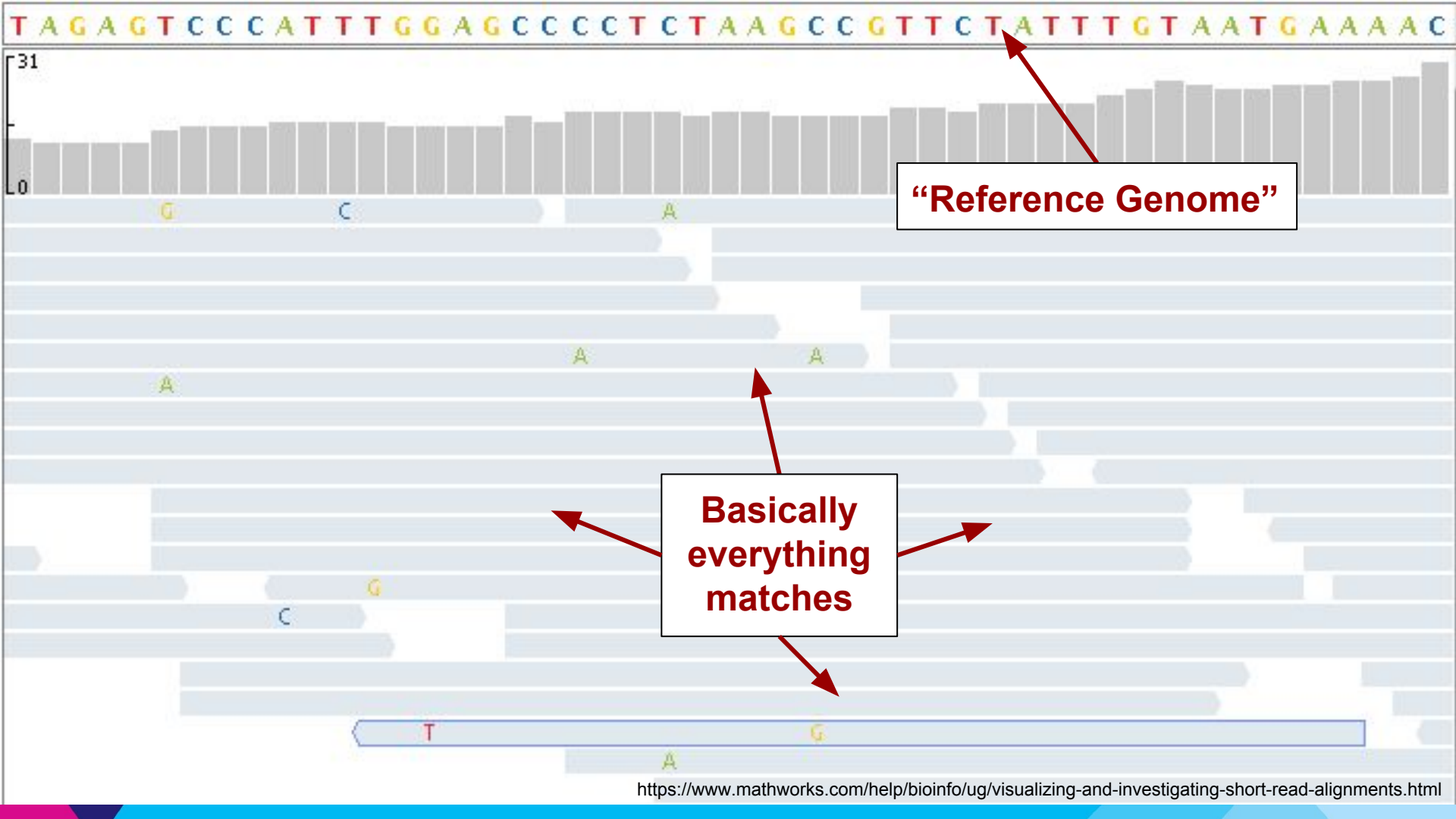
- Excessive repetitiveness
- 20% retrotransposons
  - L1: 7000bp, 100k copies
- Pseudogenes
- Impossible to resolve with “short reads”





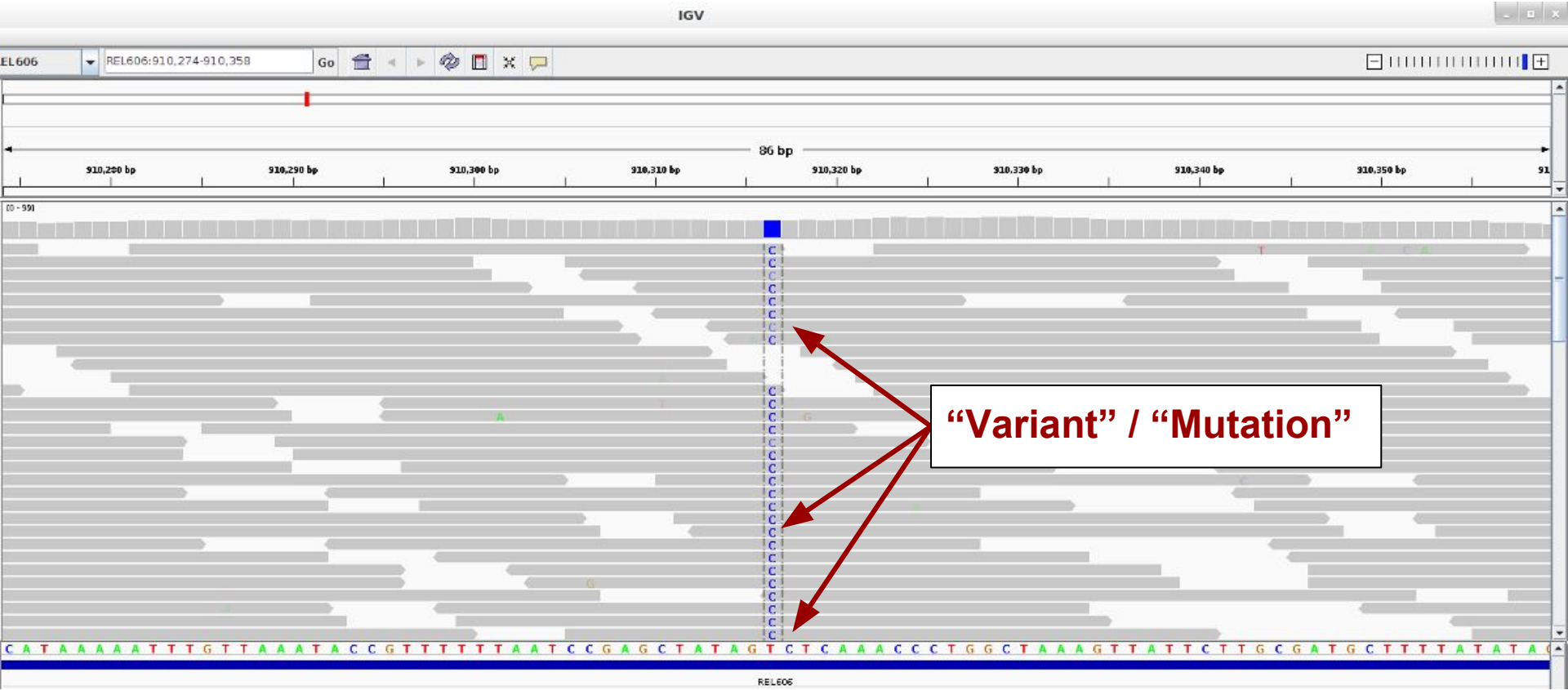




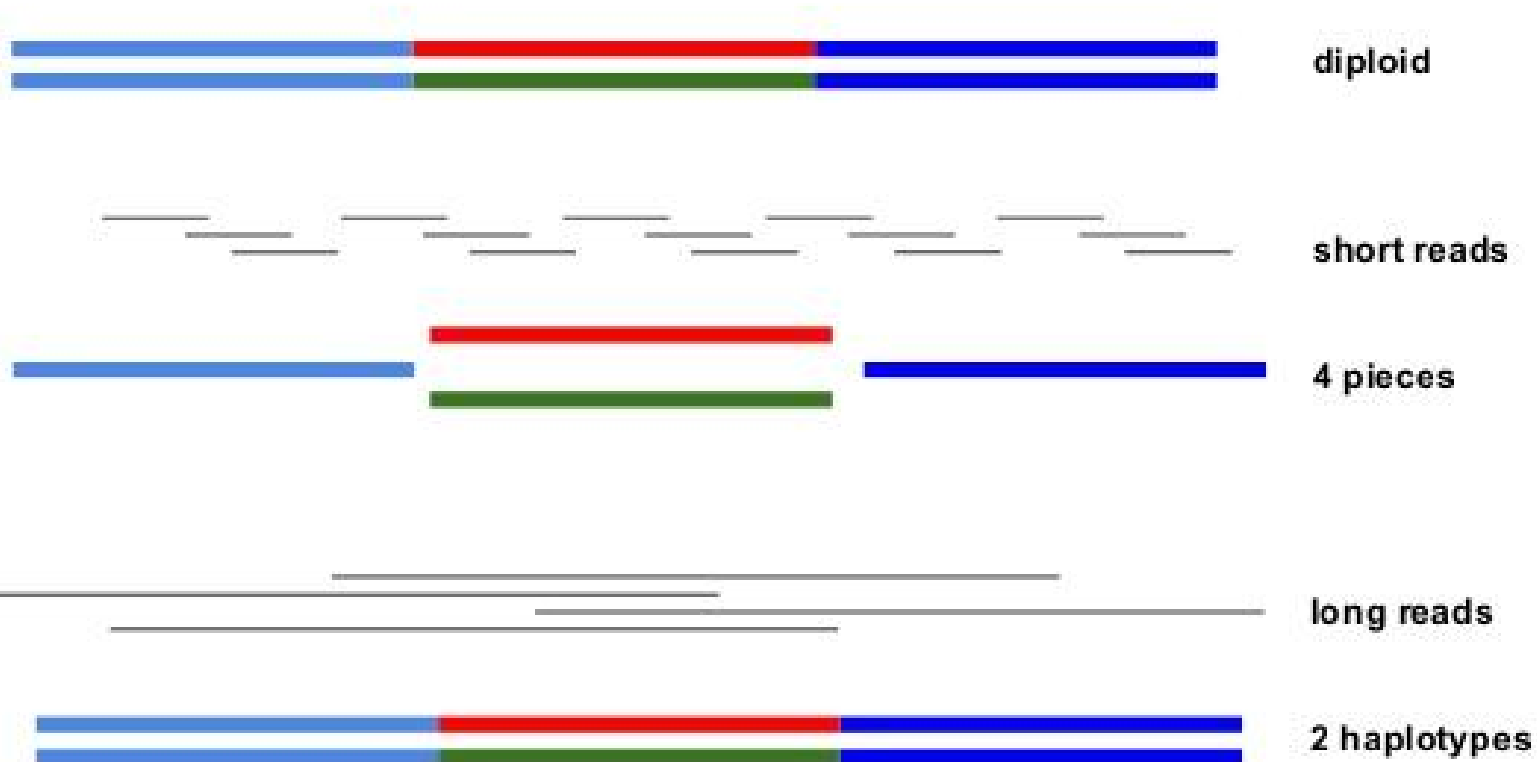


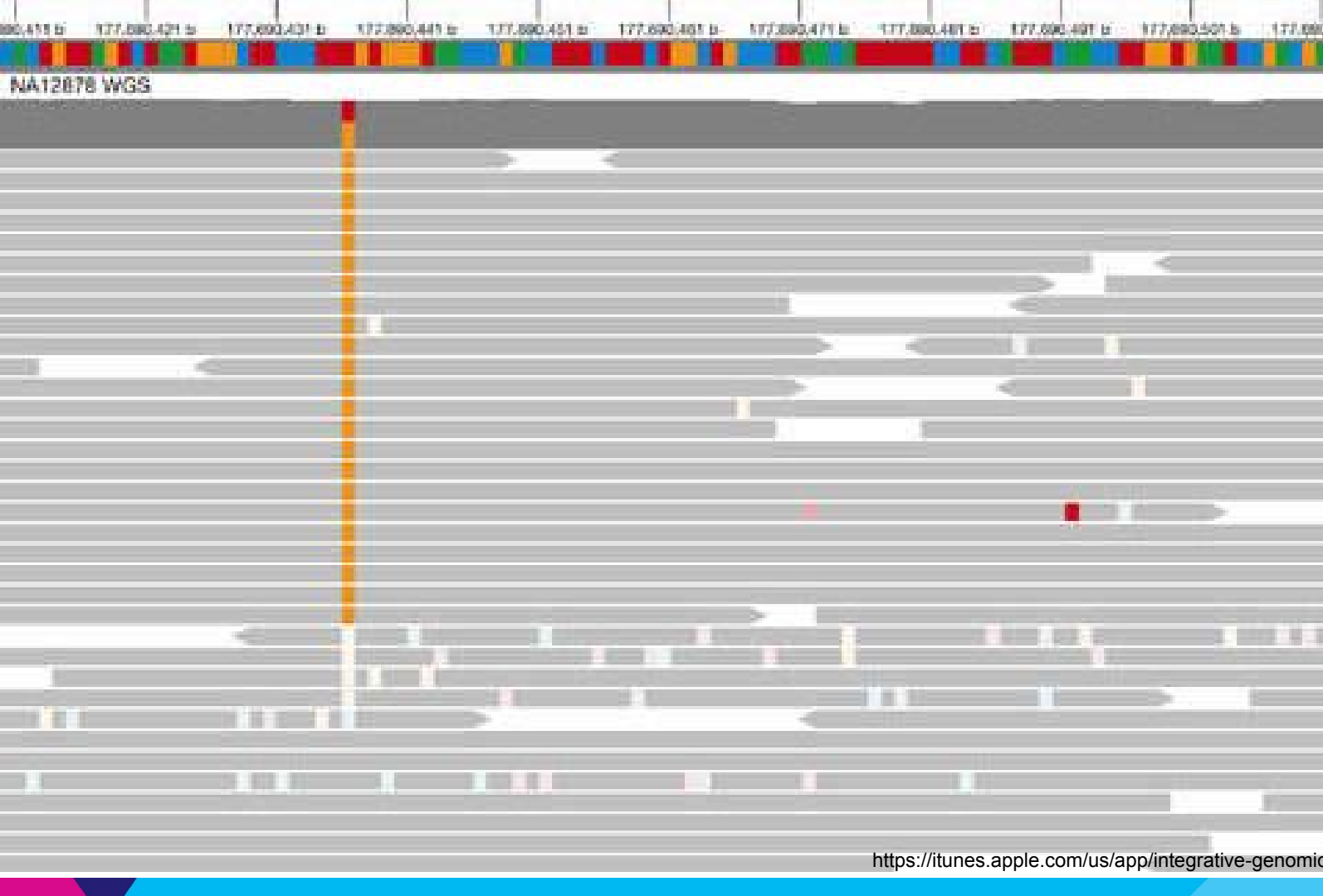


# Variant Calling



# Heterozygosity



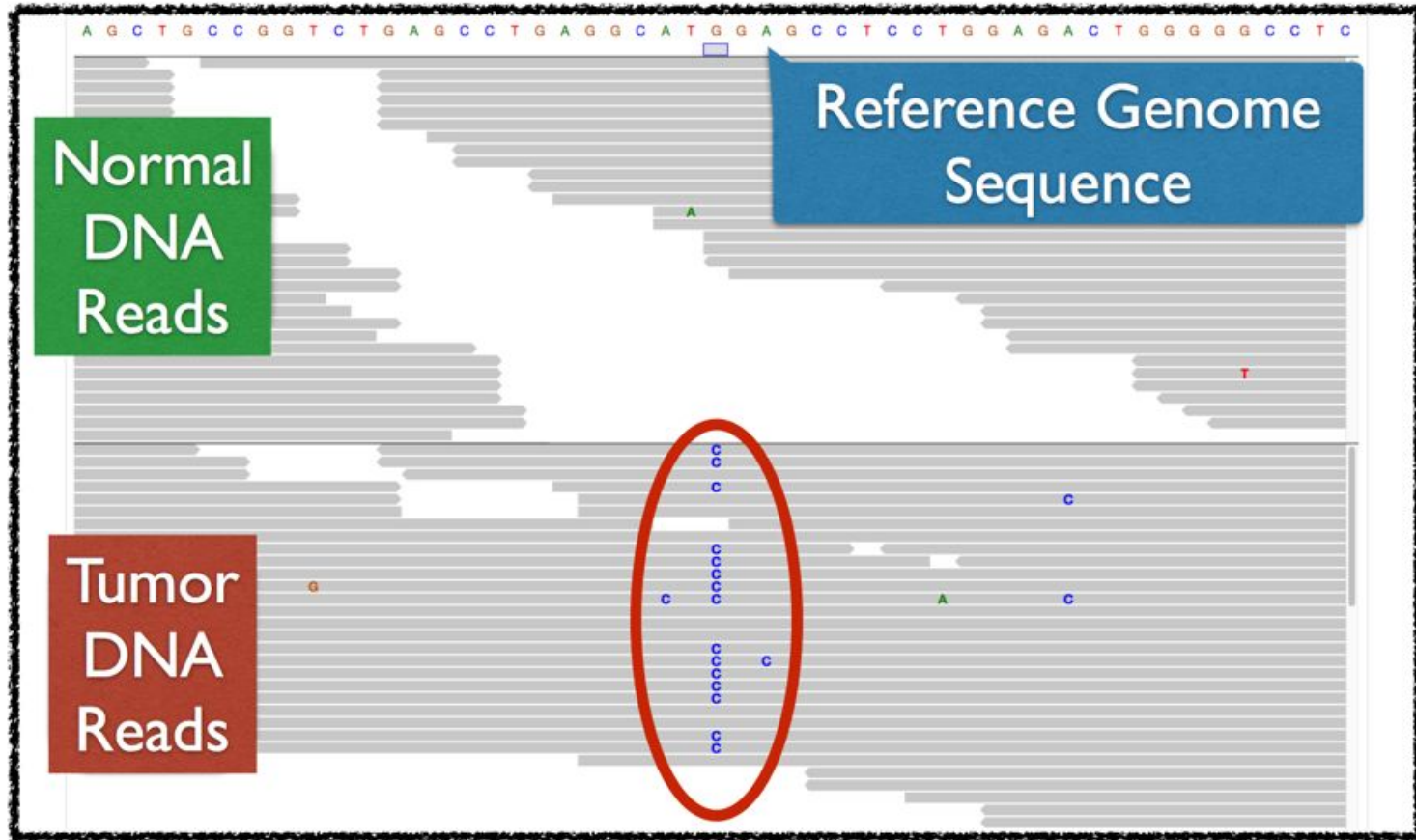


# Somatic Mutation Calling

# Somatic Mutation Calling

- Cells come from two populations
  - e.g. “normal” and tumor cells
- Find mutations specific to cancer cells

# Somatic Mutation Calling

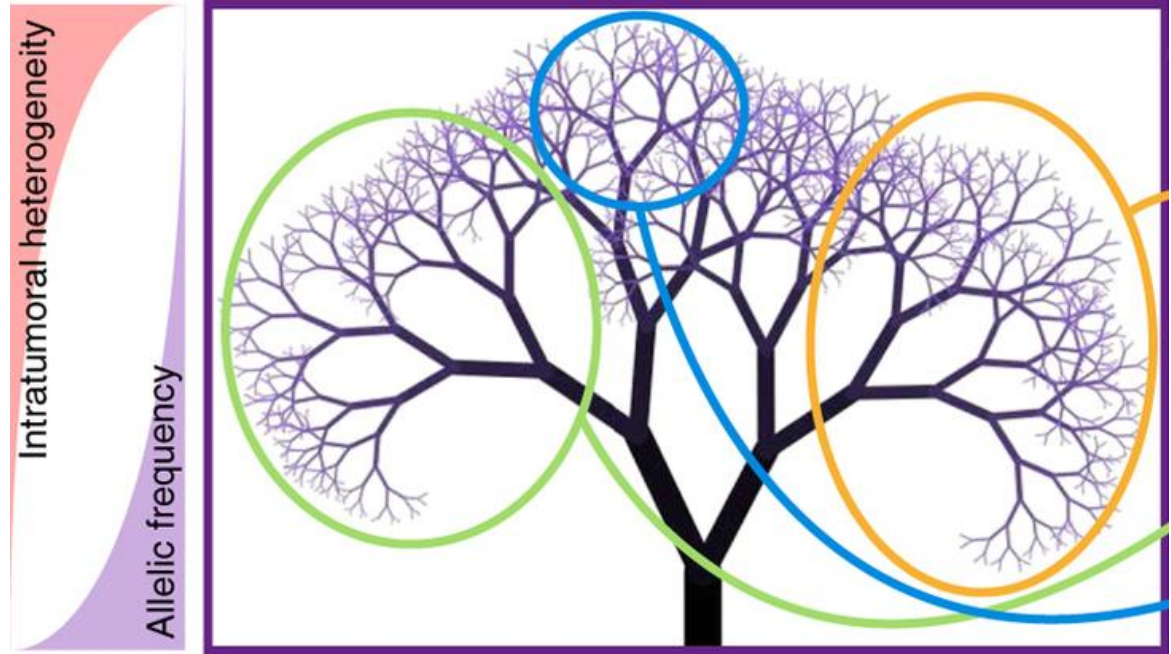






# Somatic Mutation Calling

- Underdetermined!
  - Sub-clonality
  - Tumor sample purity



# Agreement on somatic variant calls across tools is surprisingly poor

## Exome sequencing

### SNVs

	EBCall	Mutect	Seurat	Shimmer	Somatic Sniper	Strelka	Varscan 2	Virmid
EBCall	1.00	0.51	0.60	0.36	0.50	0.43	0.57	0.48
Mutect	0.20	1.00	0.47	0.18	0.25	0.26	0.30	0.33
Seurat	0.08	0.16	1.00	0.10	0.26	0.09	0.19	0.12
Shimmer	0.32	0.41	0.66	1.00	0.39	0.39	0.46	0.40
Somatic Sniper	0.08	0.10	0.31	0.07	1.00	0.06	0.20	0.09
Strelka	0.52	0.80	0.79	0.53	0.52	1.00	0.69	0.68
Varscan	0.21	0.28	0.52	0.20	0.50	0.21	1.00	0.24
Virmid	0.37	0.64	0.72	0.35	0.47	0.43	0.51	1.00

Krøigård, A.B. et al., 2016. Evaluation of Nine Somatic Variant Callers for Detection of Somatic Mutations in Exome and Targeted Deep Sequencing Data. Plos One, 11(3), p.e0151664. Available at: <http://dx.plos.org/10.1371/journal.pone.0151664>.

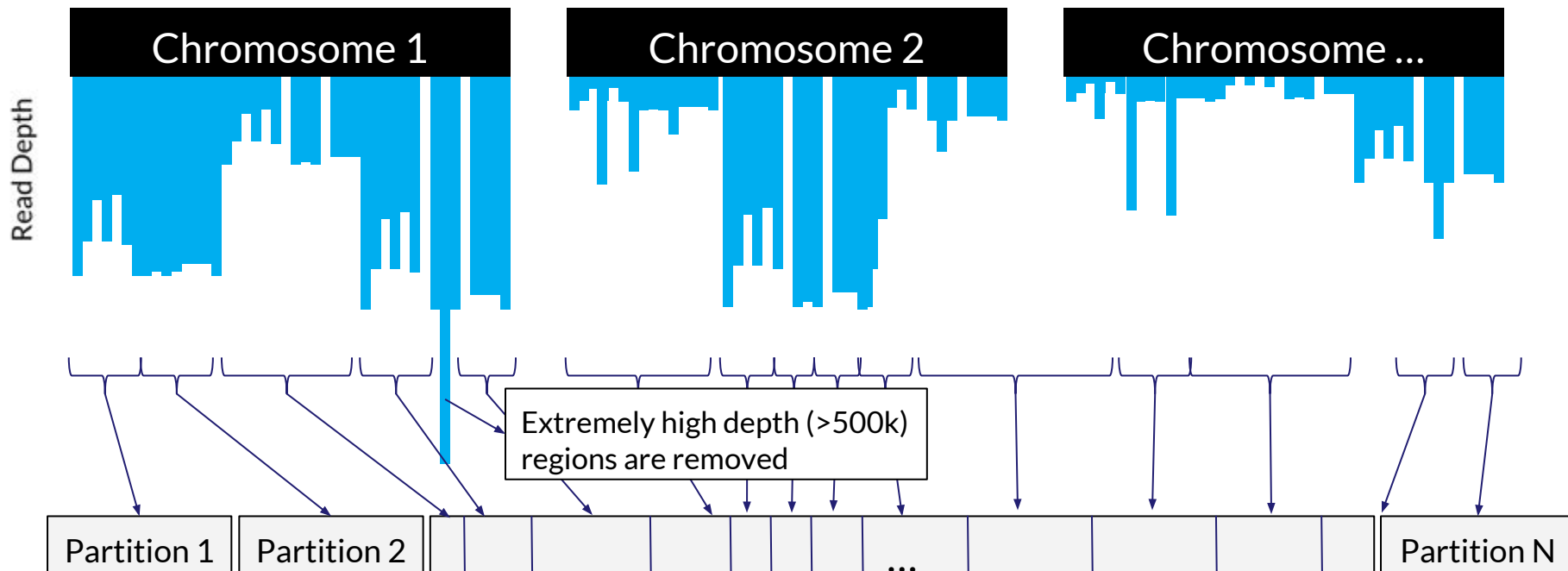
# Guacamole: Somatic Mutation Calling with Apache Spark

# A typical Guacamole analysis

1. Partition the genome
2. Partition reads according to (1)
3. Build pileups at each site
4. Apply user-supplied function at each pileup
5. Write output

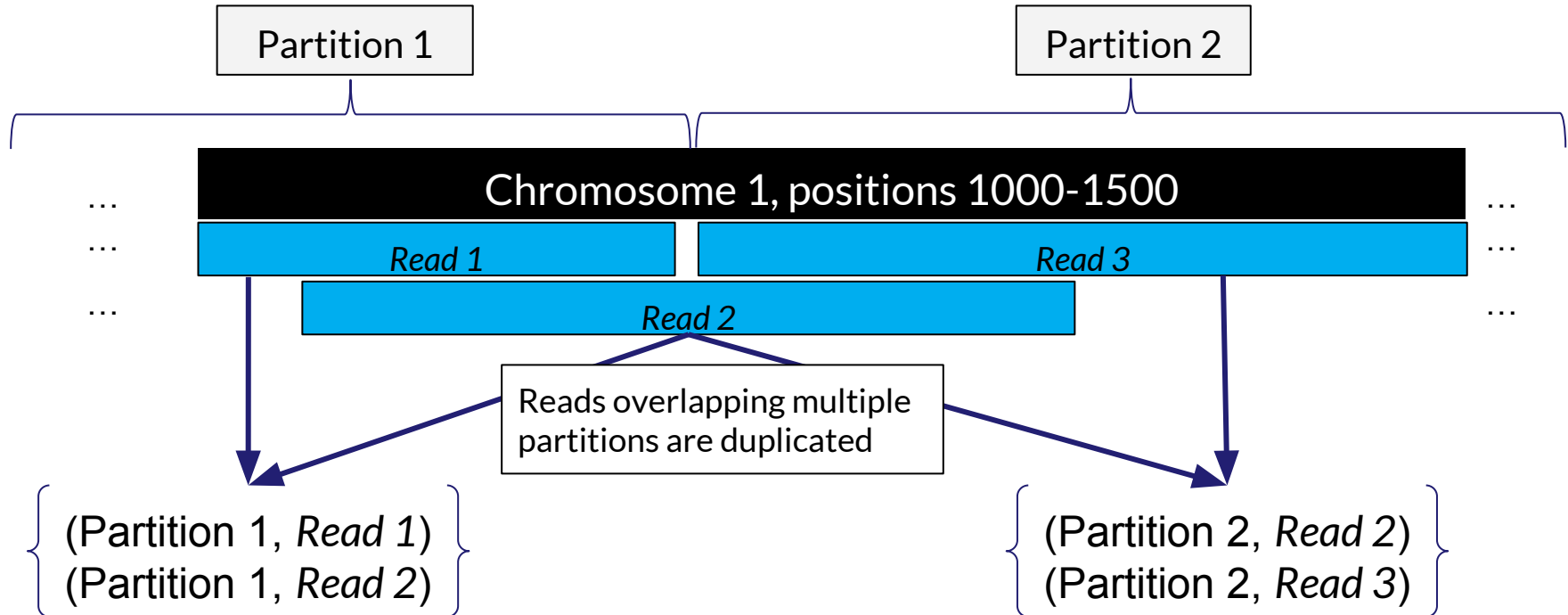
# Step 1: Partition the genome

- Partition the genome into intervals, balancing the number of reads overlapping each partition
- Each interval will correspond to one Spark partition



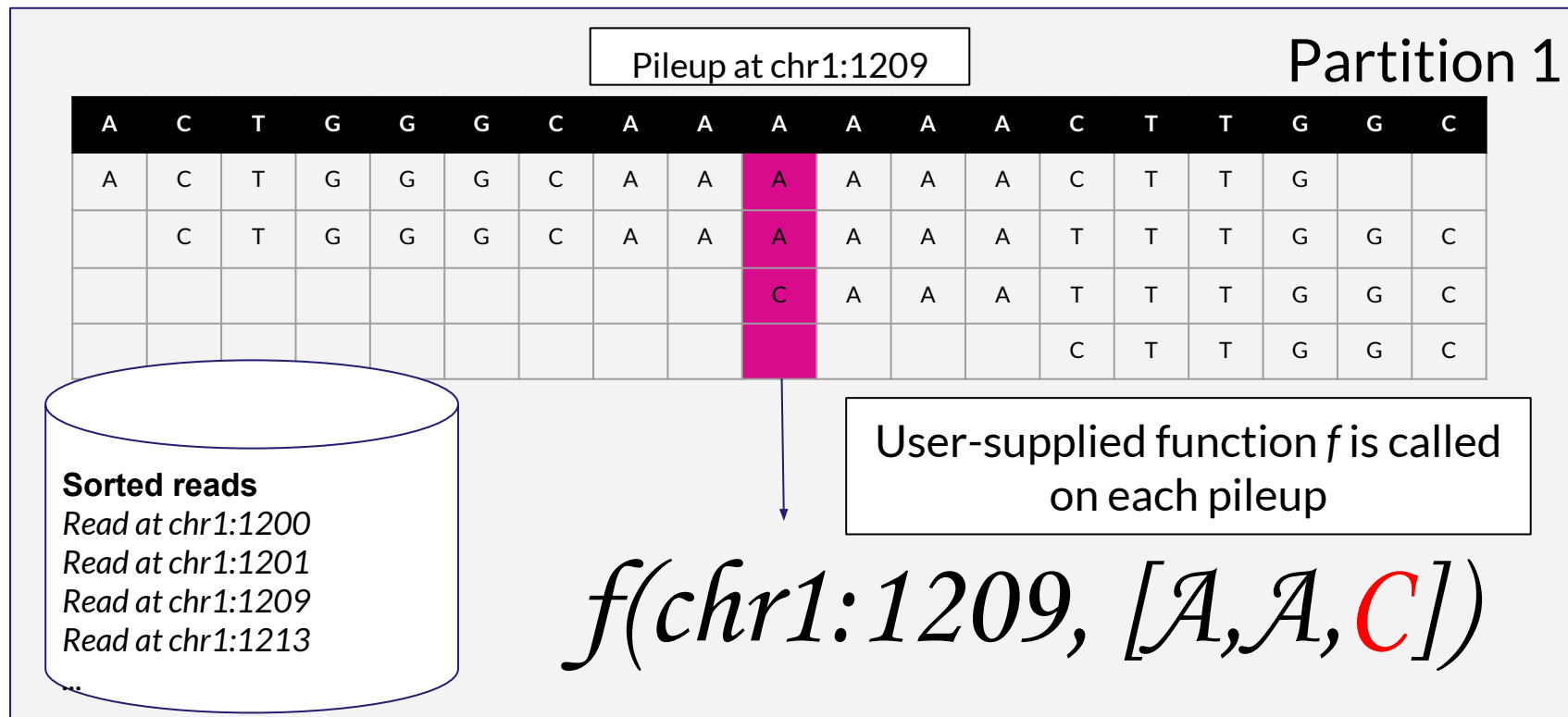
## Step 2: Partition reads

- All-to-all shuffle of reads based on the genomic partition in Step 1
- A copy of each read goes to each partition it overlaps





# Step 3: each partition streams through reads to generate pileups



# Benchmarking

# Testing cluster

Hardware	
Nodes	100
Cores	2400
Memory	12.5 TB
Storage	3.6 PB

Software	
Spark	1.6.1
Hadoop	2.6.0-cdh5.5.1
OS	CentOS 7.2.1511



# Guacamole speed

Guacamole	
<b>2 WGS samples (DREAM Synth4)</b>	22 minutes
<b>3 Whole Genome Samples (AOCS-034)</b>	31 minutes
<b>10 Whole Exome Samples (PT189)</b>	52 minutes

Mutect	
<b>2 WGS samples (DREAM Synth4) chromosome 1</b>	158 minutes
We compare to chromosome 1 single-node runtime because Mutect is typically parallelized by chromosome, and chromosome 1 takes the longest	

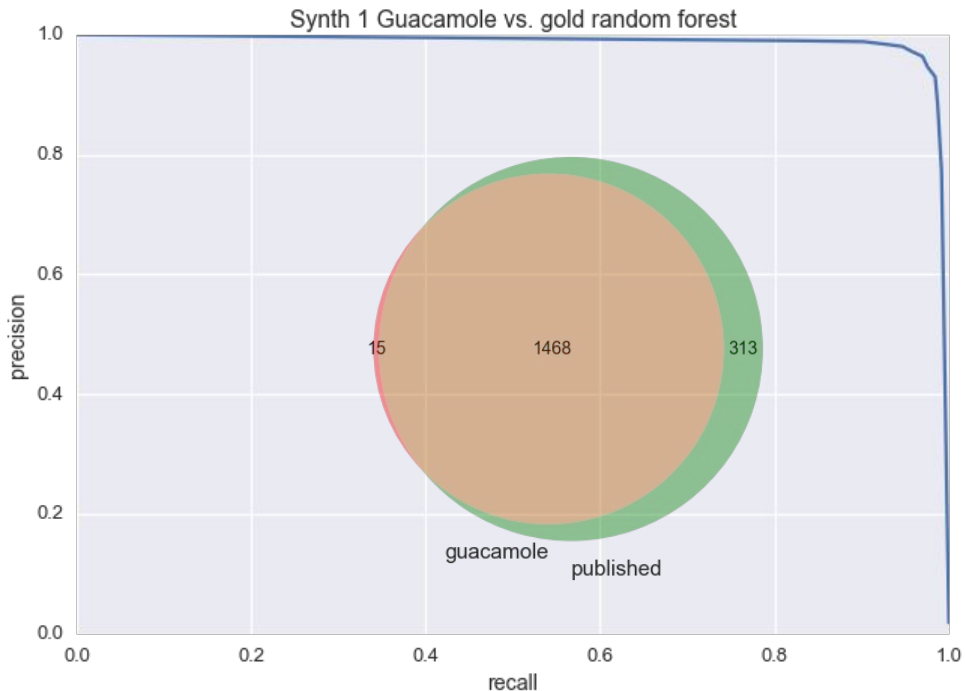
# WIP: Model-based optimization

## Features

- Raw likelihood
- Difference of ref and alt likelihood
- Variant allele fractions
- Allele depths
- Strand bias

## Methodology

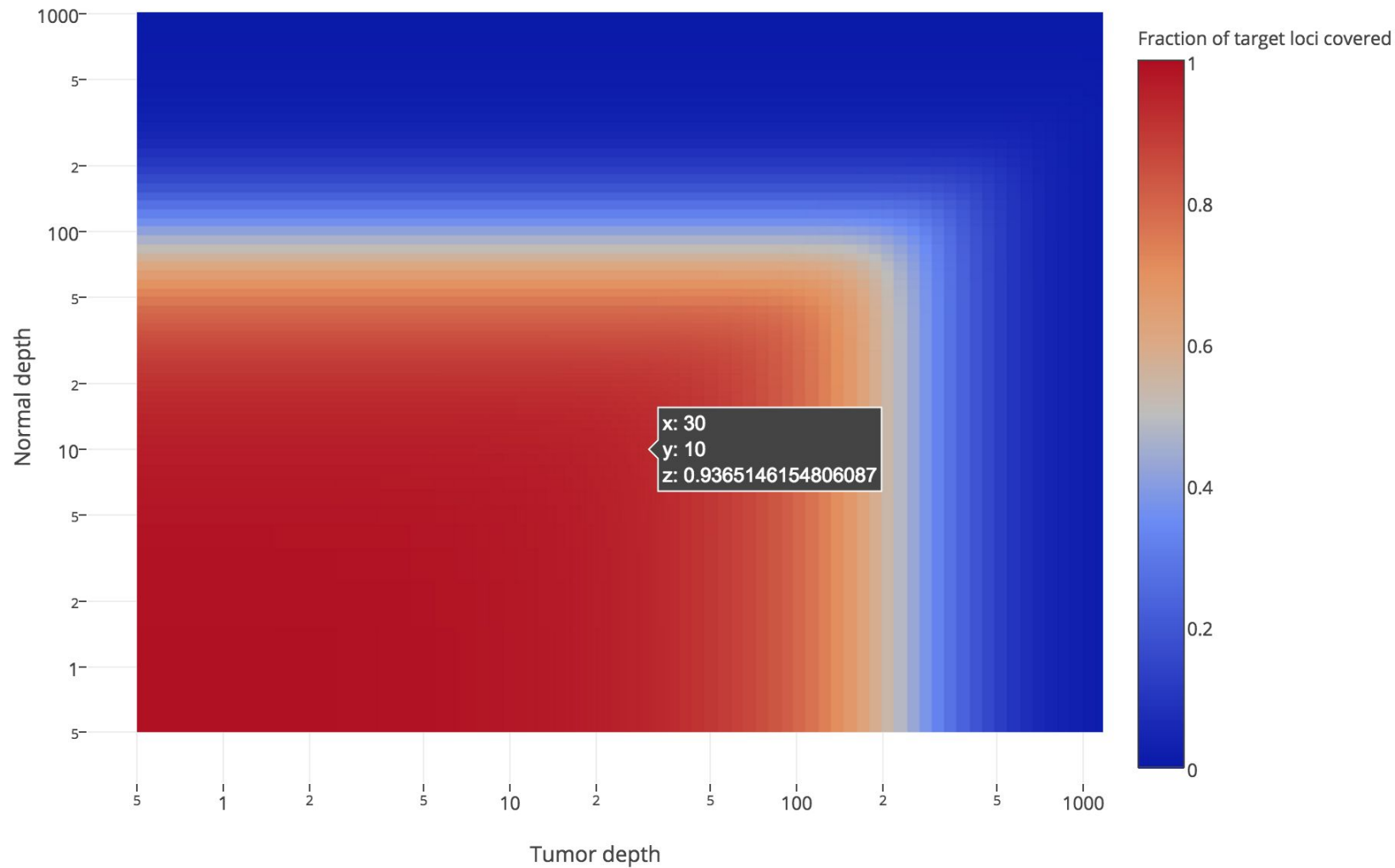
- Random forest
- 1:1 train/test split



# Other applications: QC Analysis

# Exome Sequencing

- “Exome”: just the genes. 1% of genome.
- Question: how much of the exome was covered with at least X reads in one sample and Y reads in the other.





# Distributed 2D-Prefix-Sum

- Demo / viz
- Spark implementation: hammerlab/magic-rdds

**Thanks!**