

Whole-Genome Sequencing (WGS)

Aakrosh Ratan

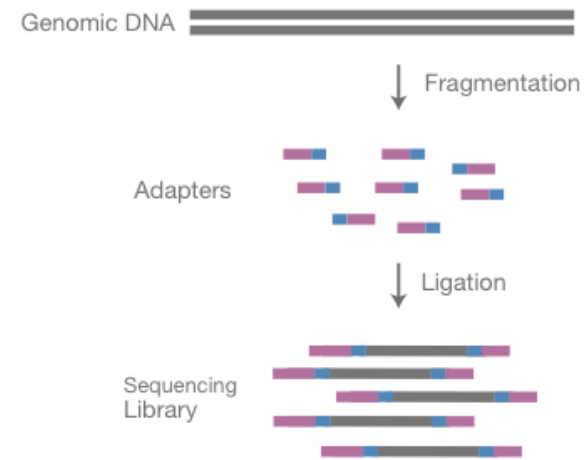
PHS5500: Special Topics in Public Health - Public Health Genomics

14th March, 2016

<http://bims.virginia.edu/faculty/aakrosh-ratan>

ratan@virginia.edu

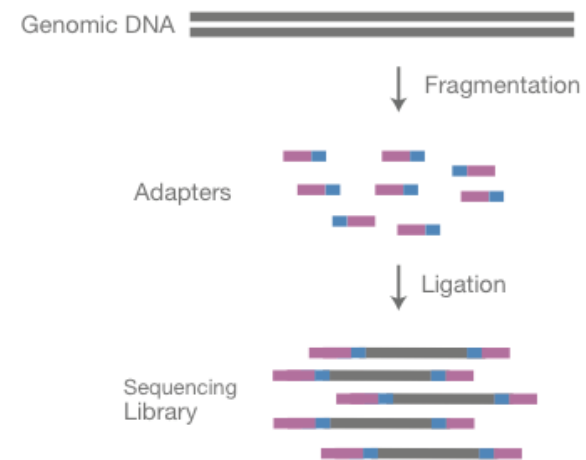
A. Library Preparation



NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

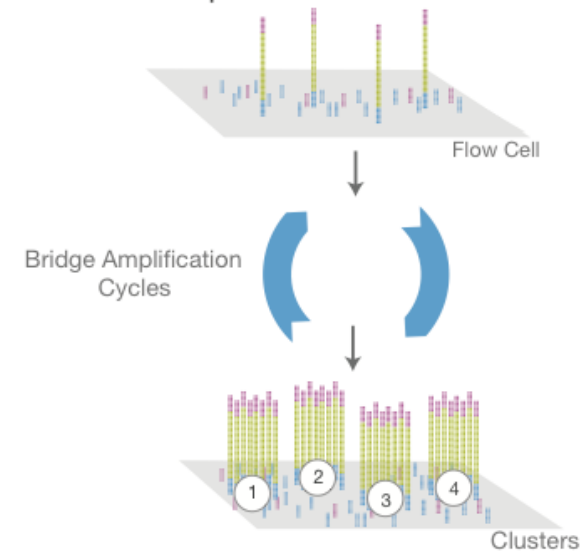
This should be called panel B.

A. Library Preparation



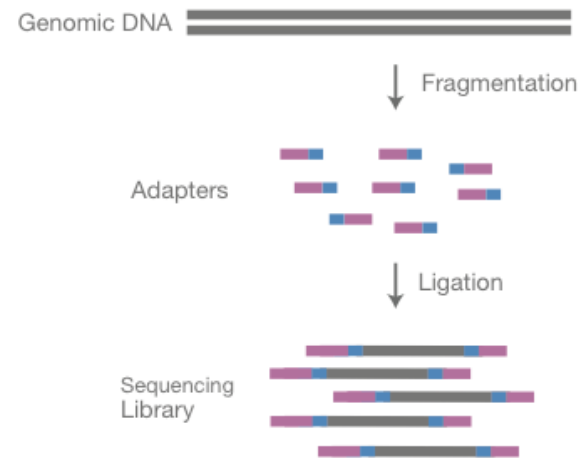
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

A. Cluster Amplification



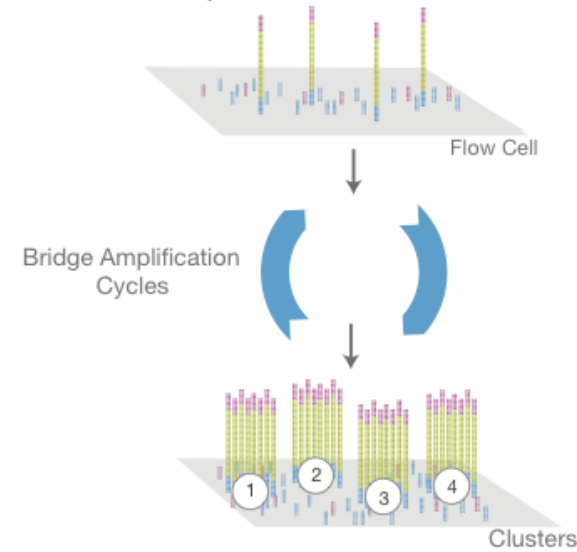
Library is loaded into a flow cell and the fragments hybridize to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

A. Library Preparation



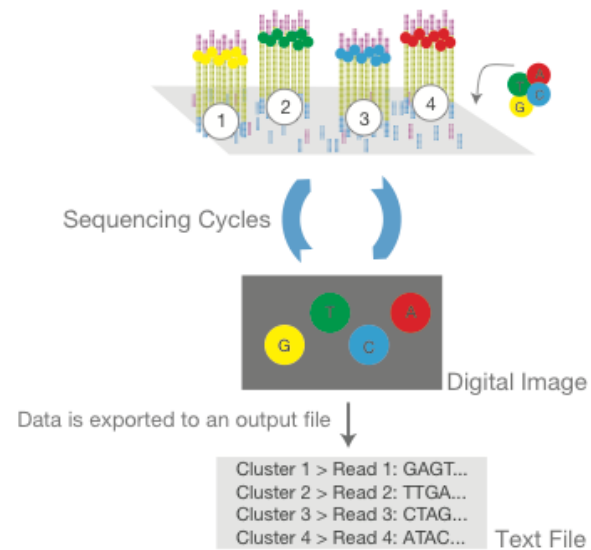
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

A. Cluster Amplification



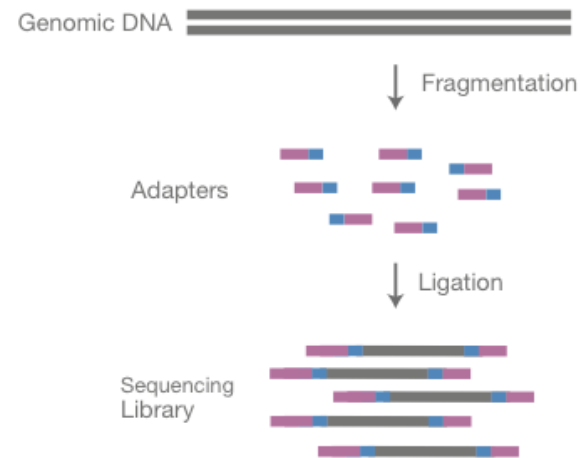
Library is loaded into a flow cell and the fragments hybridize to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

C. Sequencing



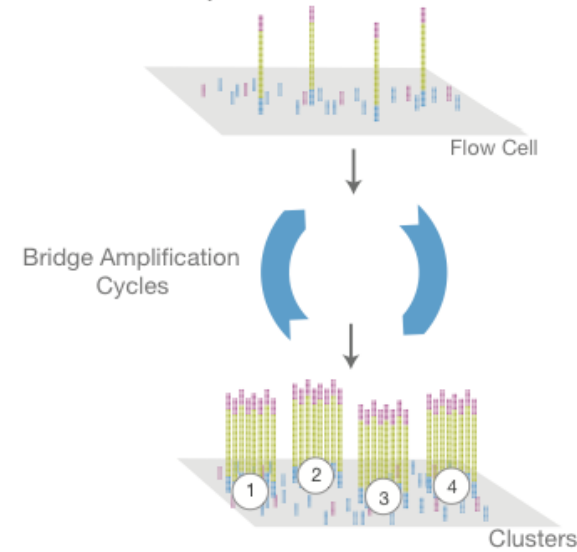
Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

A. Library Preparation



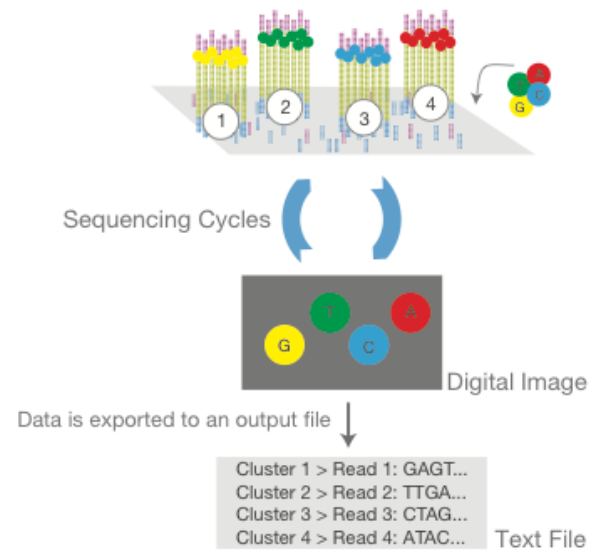
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

A. Cluster Amplification



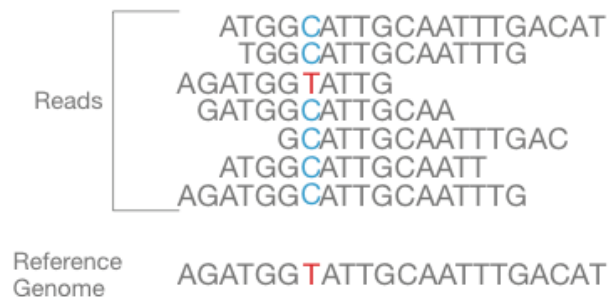
Library is loaded into a flow cell and the fragments hybridize to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

C. Sequencing



Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

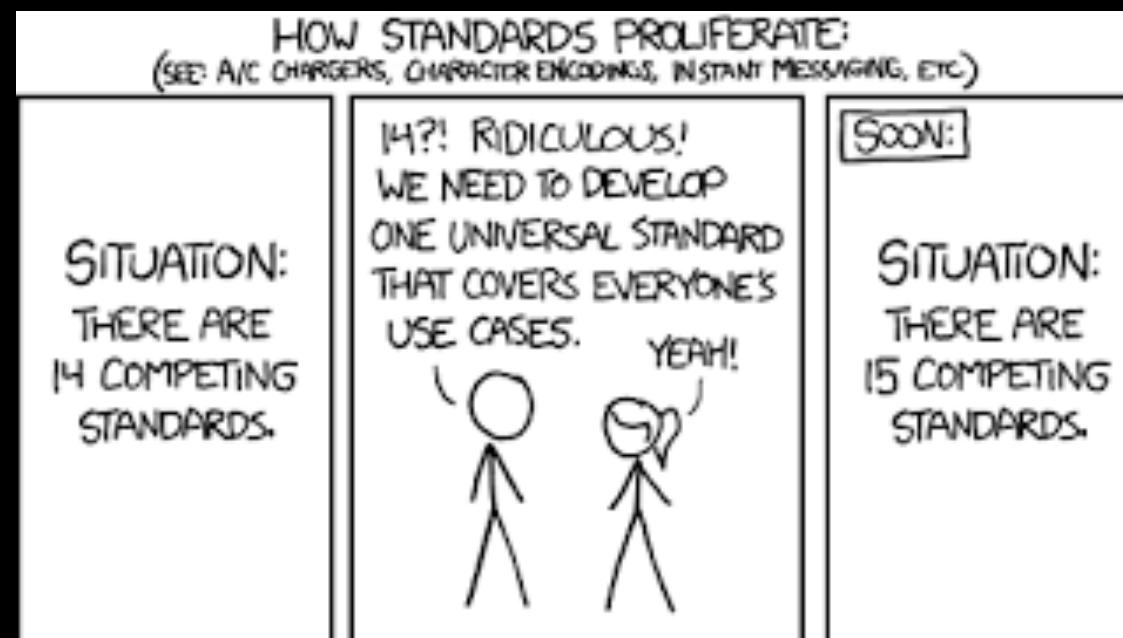
D. Alignment & Data Analysis



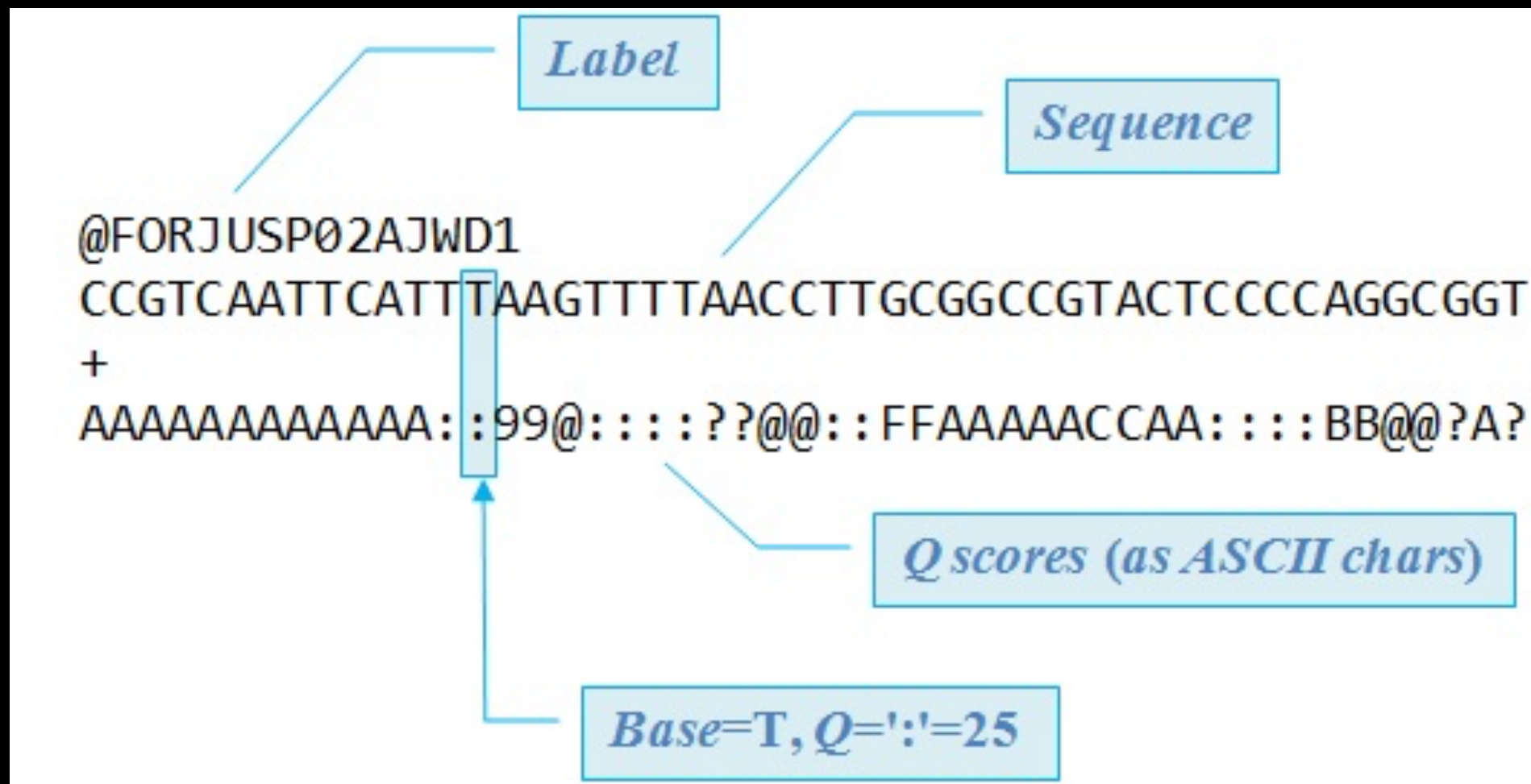
Reads are aligned to a reference sequence with bioinformatics software. After alignment, differences between the reference genome and the newly sequenced reads can be identified.



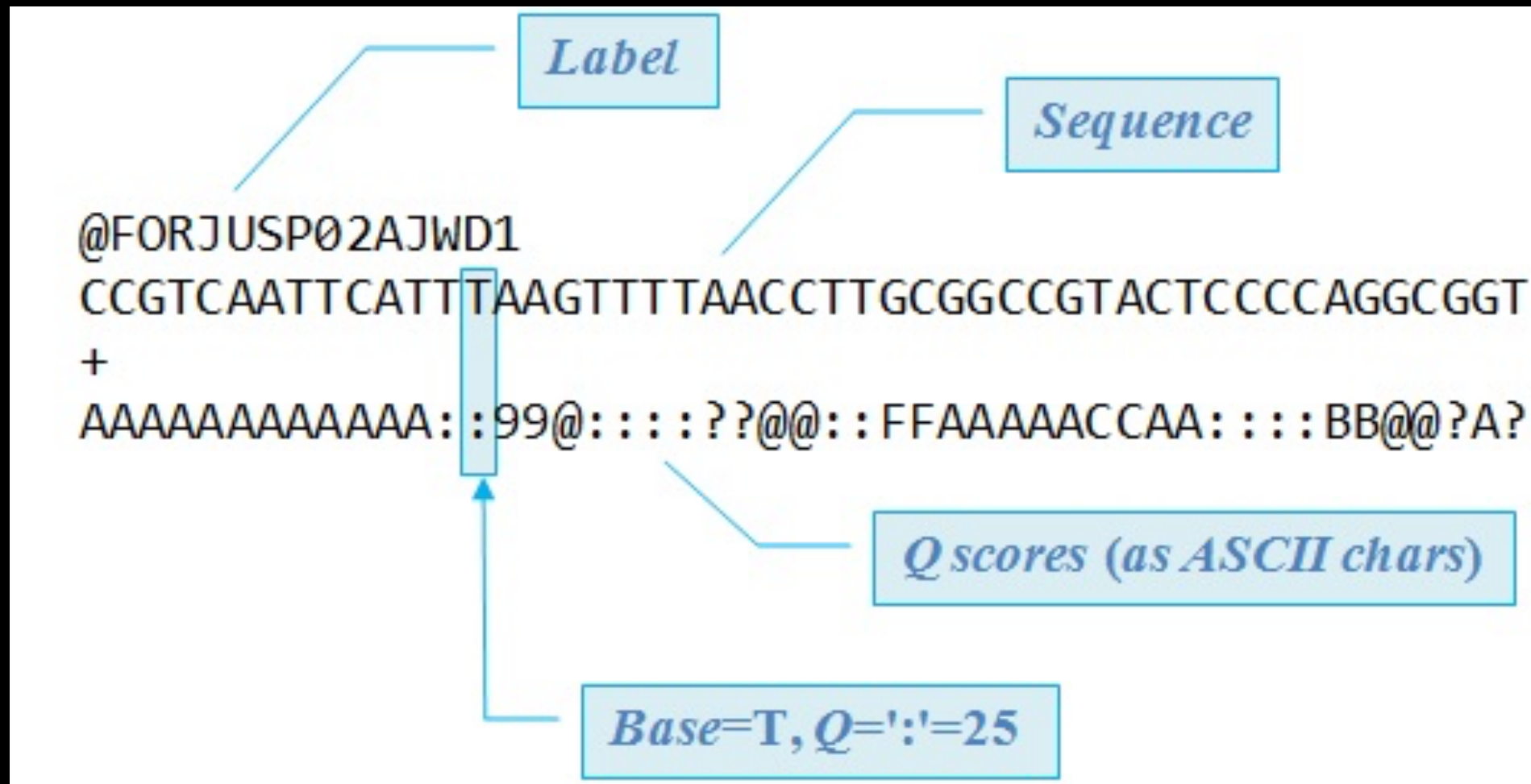
Paired-End Sequencing



Source: <https://xkcd.com/927/>



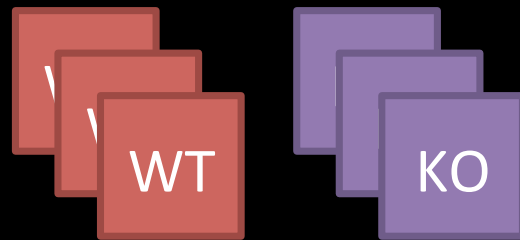
Source: http://drive5.com/usearch/manual/fastq_files.html



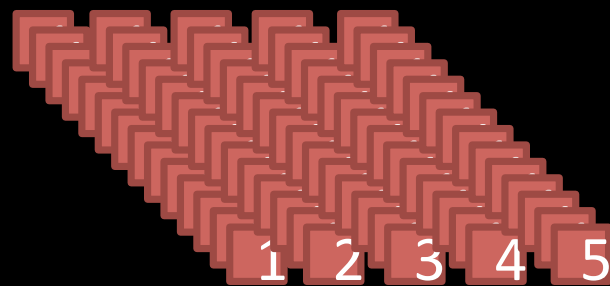
Source: http://drive5.com/usearch/manual/fastq_files.html

$$A = 1 - E = 1 - 10^{-\left(\frac{Q}{10}\right)}$$

phred Q25 ~ 0.9968



Small experiment; we can examine things in detail



Large experiment; we can't truly look at all details

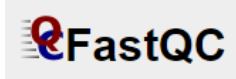
Quality Control

QC should tell you what to look at more closely. It should NOT be used as an automated filter.

(Baby) Steps

- Decide what is “normal”
- Calculate the same metric in your datasets
- Check from deviance from normal
- Trigger an alarm (visual alarm, email...) to notify user to look at the data more closely
- Summarize, Visualize and Flag

Raw sequence:



PRINSEQ

FastQ screen


Mapped sequence:



QualiMap

Bismark (specialised)

Application specific:

- **RNA-Seq**
SeqMonk RNA-Seq QC
RNASeQC
- **Small RNA**
SeqMonk QC
- **ChIP**
ChIPQC package
- **Bi-Sulphite Sequencing**
Bismark
- **Hi-C**
The logo for HiCUP, featuring a red and white checkered pattern next to the text 'HiCUP'.

Data Visualisation:

SeqMonk



Intergrated Genome Viewer (IGV)

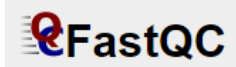


Compile many QC analyses
into a single report:



Few Existing Tools for QC

Raw sequence:



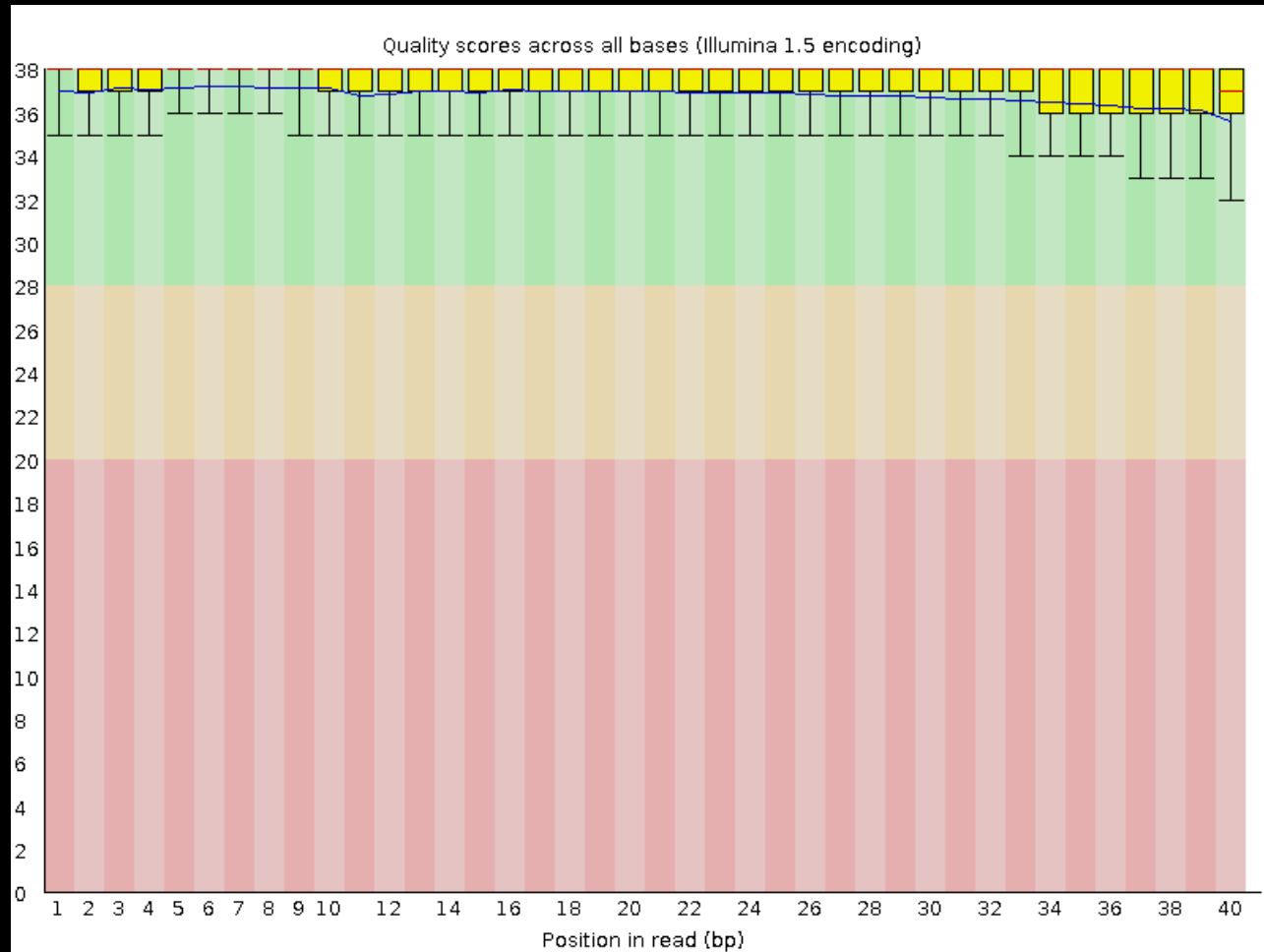
PRINSEQ

FastQ screen

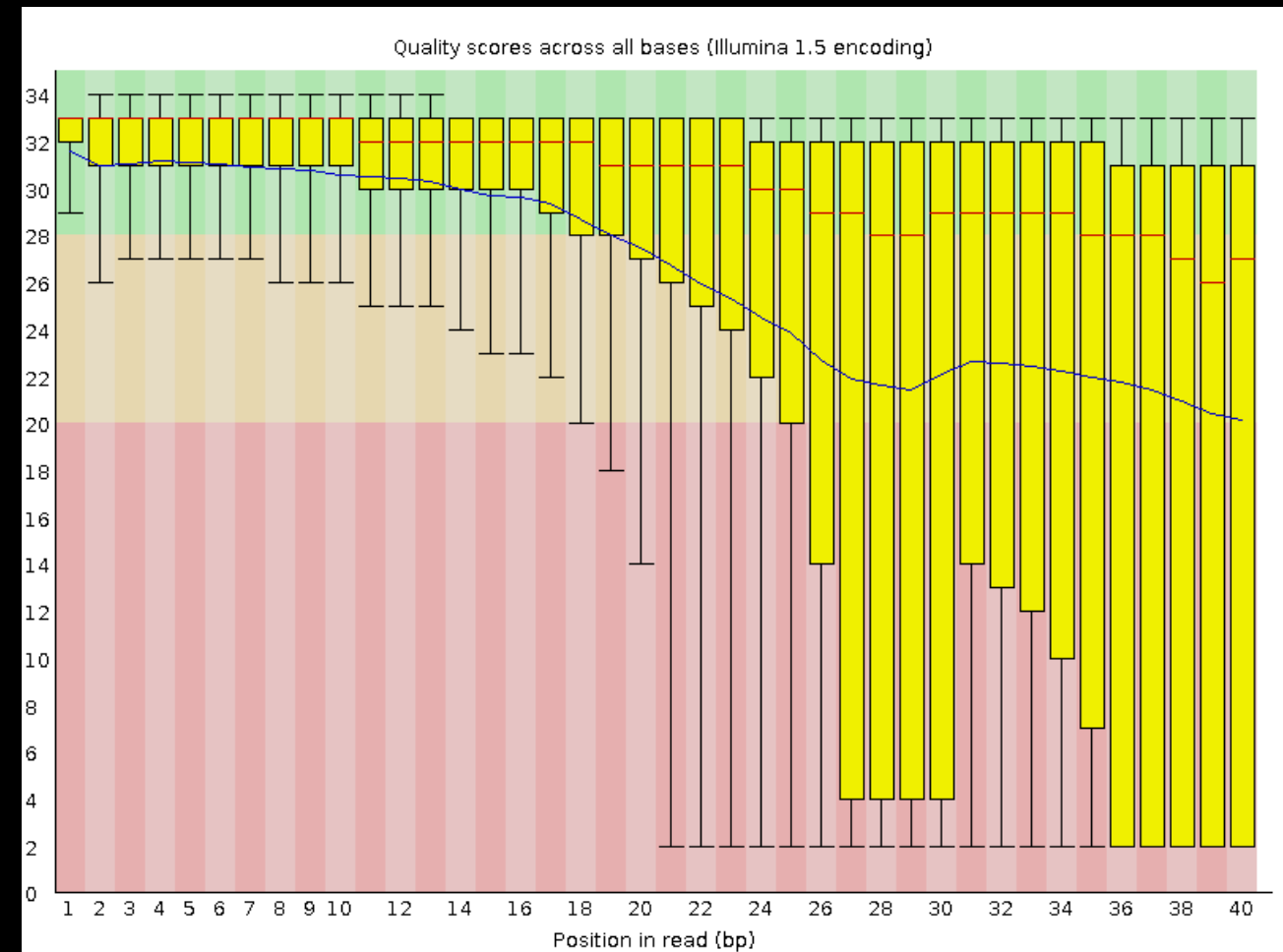
Few Existing Tools for QC

Per base sequence quality

Good



Bad

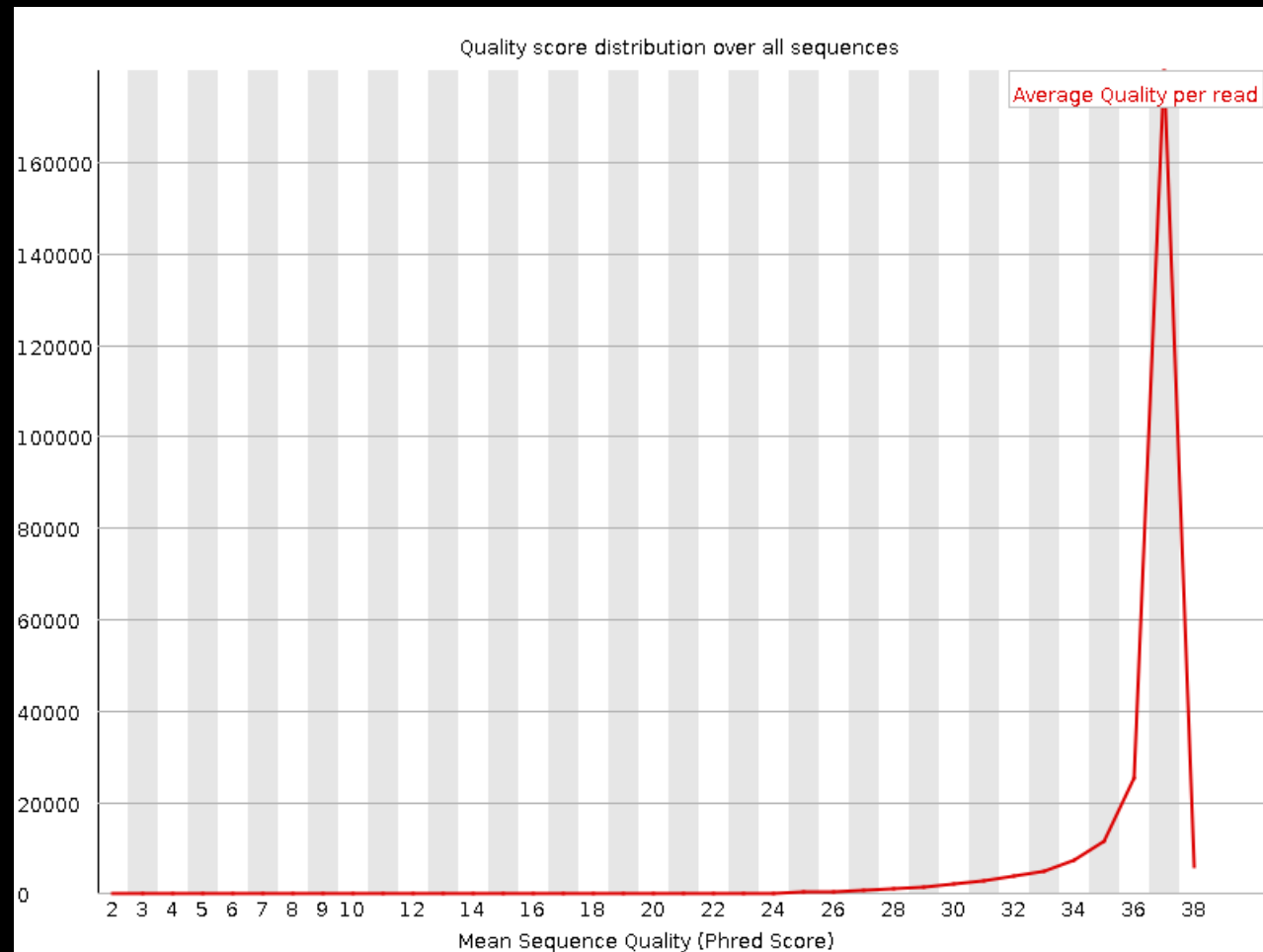


FASTQC

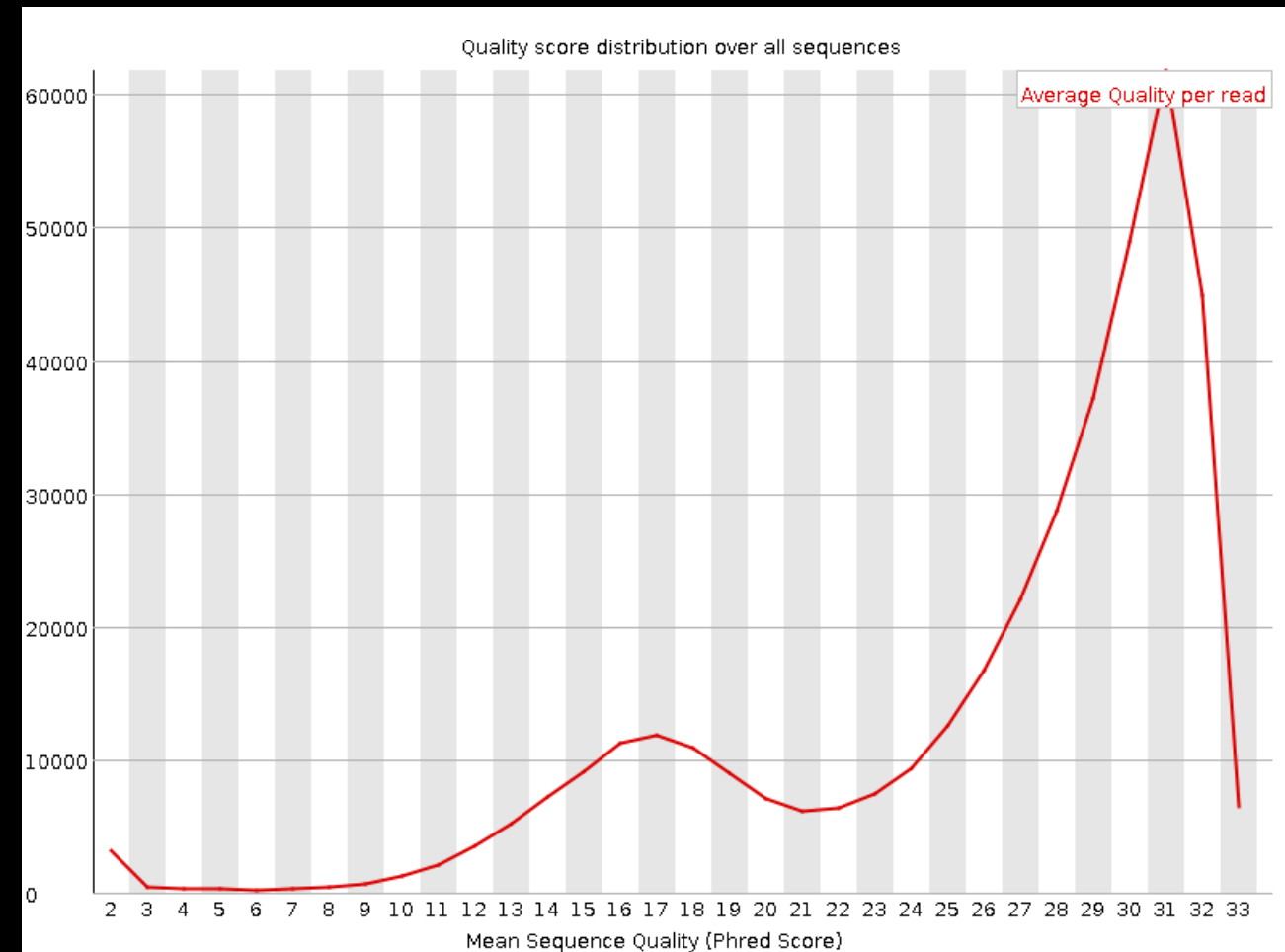
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Per sequence quality scores

Good



Bad

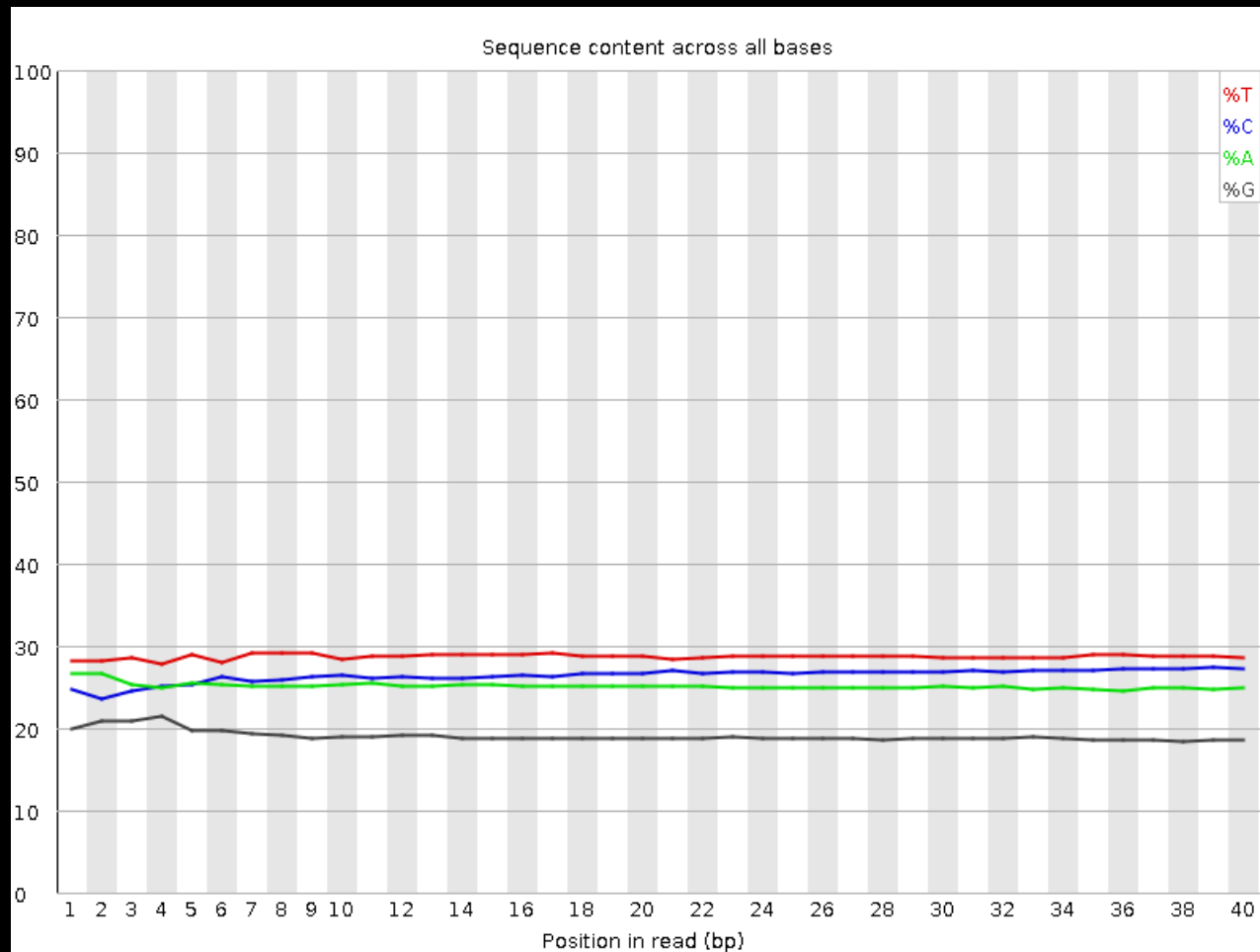


FASTQC

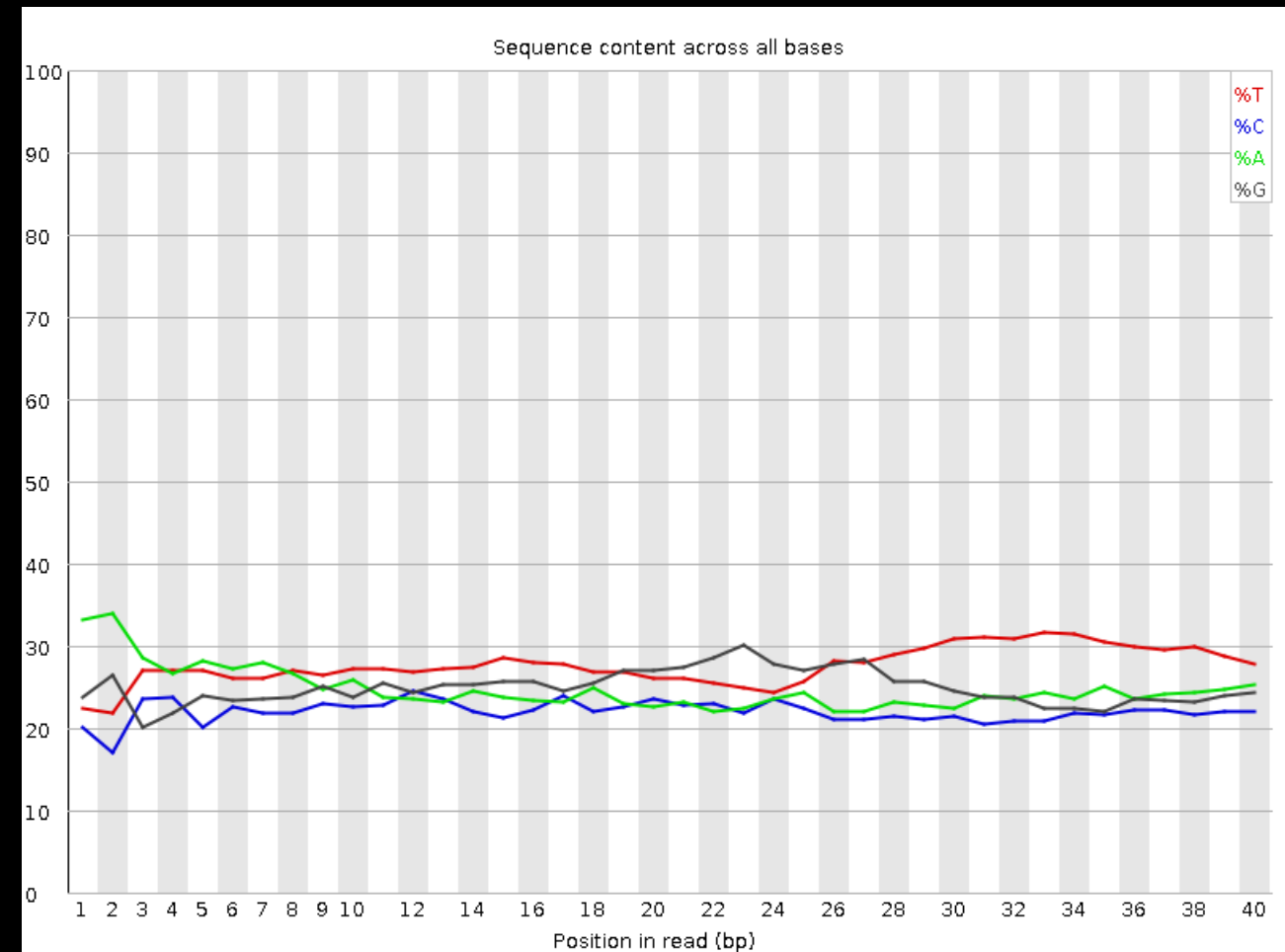
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Per base sequence quality

Good



Bad

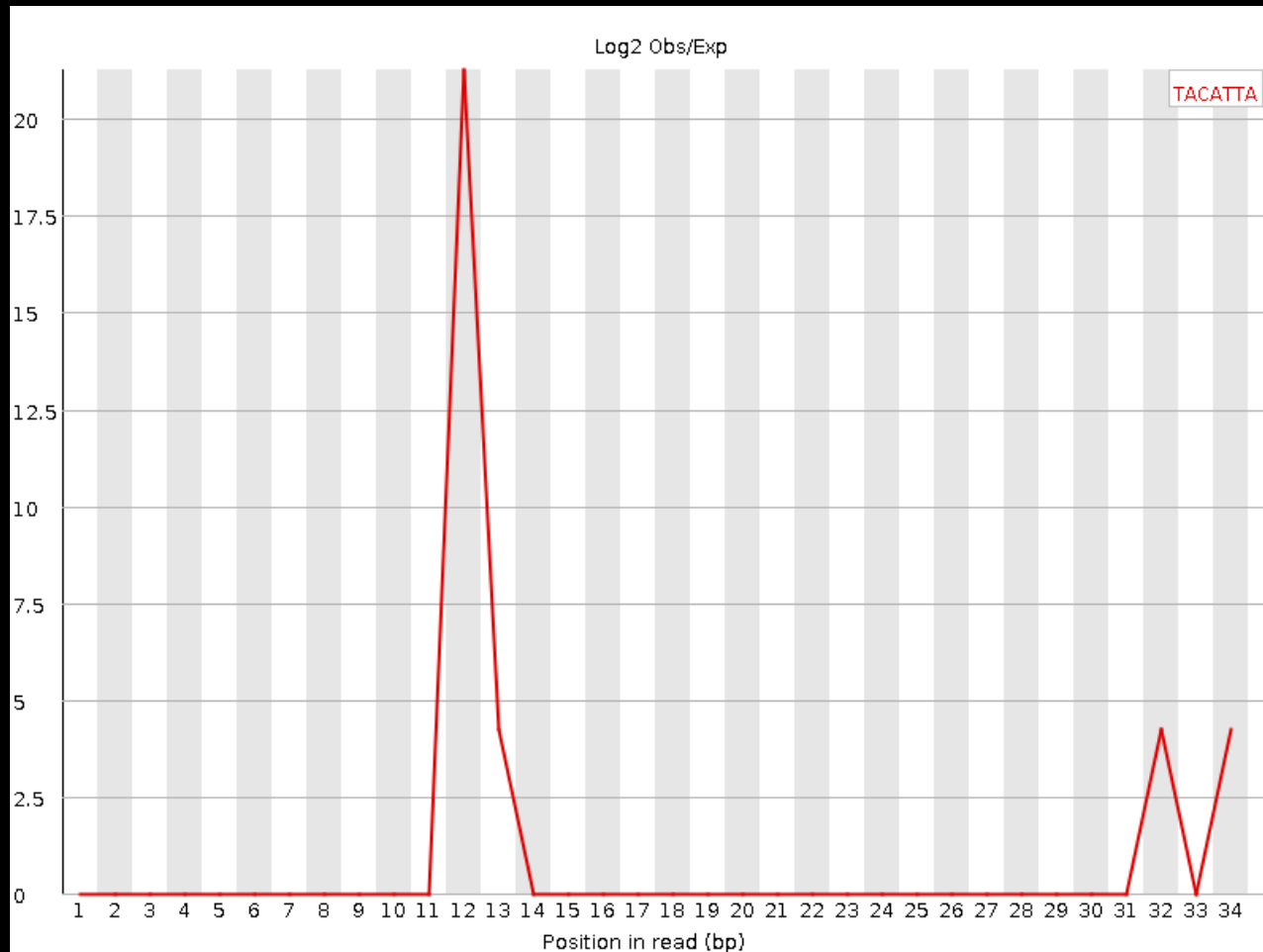


FASTQC

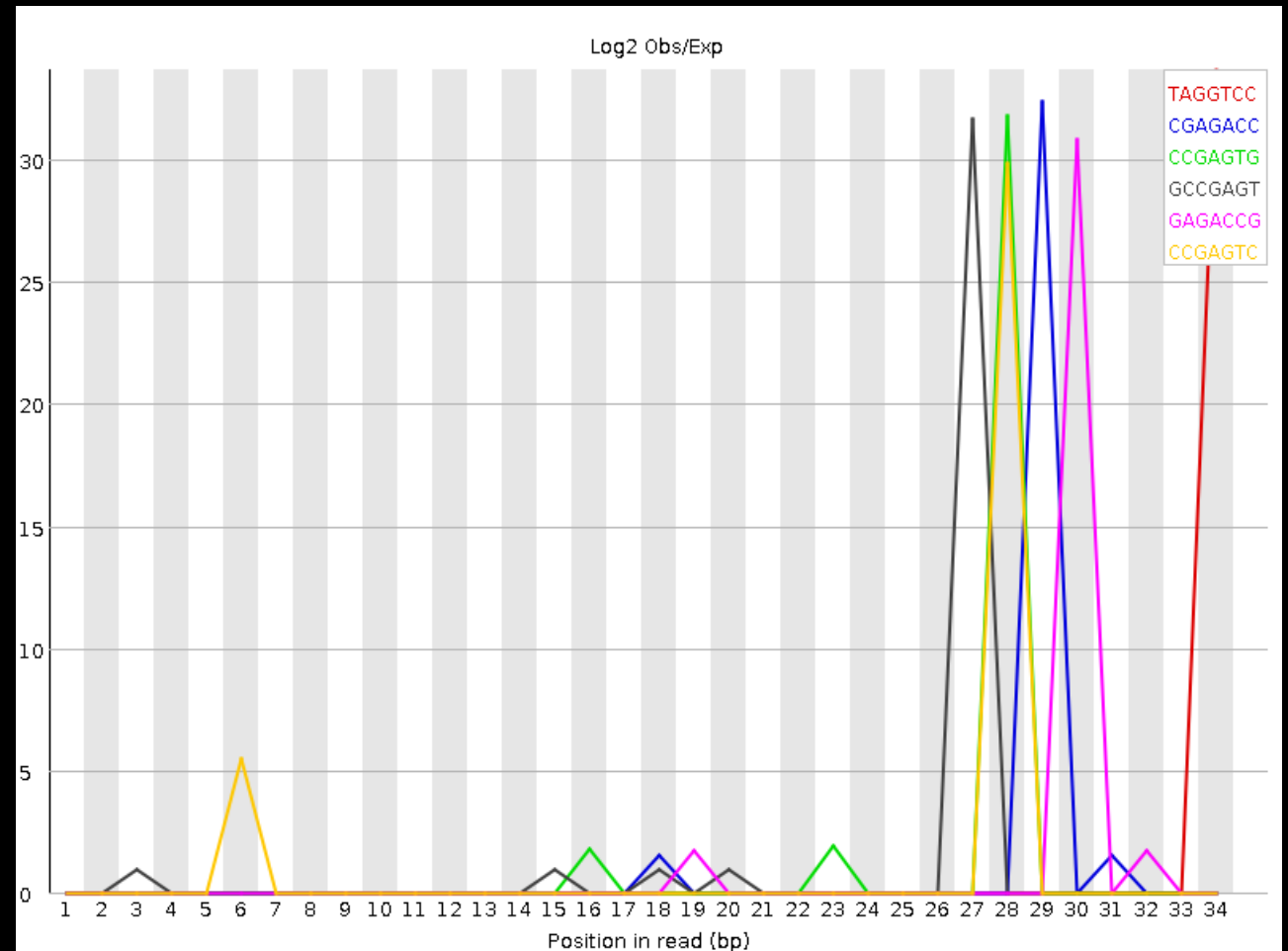
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

K-mer Content

Good (Not that good!!!)



Bad



FASTQC

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Several other metrics can be checked (based on application)

- Per tile sequence quality (if some specific region on the sequencer was enriched for bad sequences)
- Per base N content (if some specific base position was enriched for ambiguous base-calls)
- Sequence length distribution (Adapter contamination, could also signal degraded DNA)

Several other metrics can be checked (based on application)

- Per tile sequence quality (if some specific region on the sequencer was enriched for bad sequences)
- Per base N content (if some specific base position was enriched for ambiguous base-calls)
- Sequence length distribution (Adapter contamination, could also signal degraded DNA)

And of course, there is additional QC after every step in the pipeline



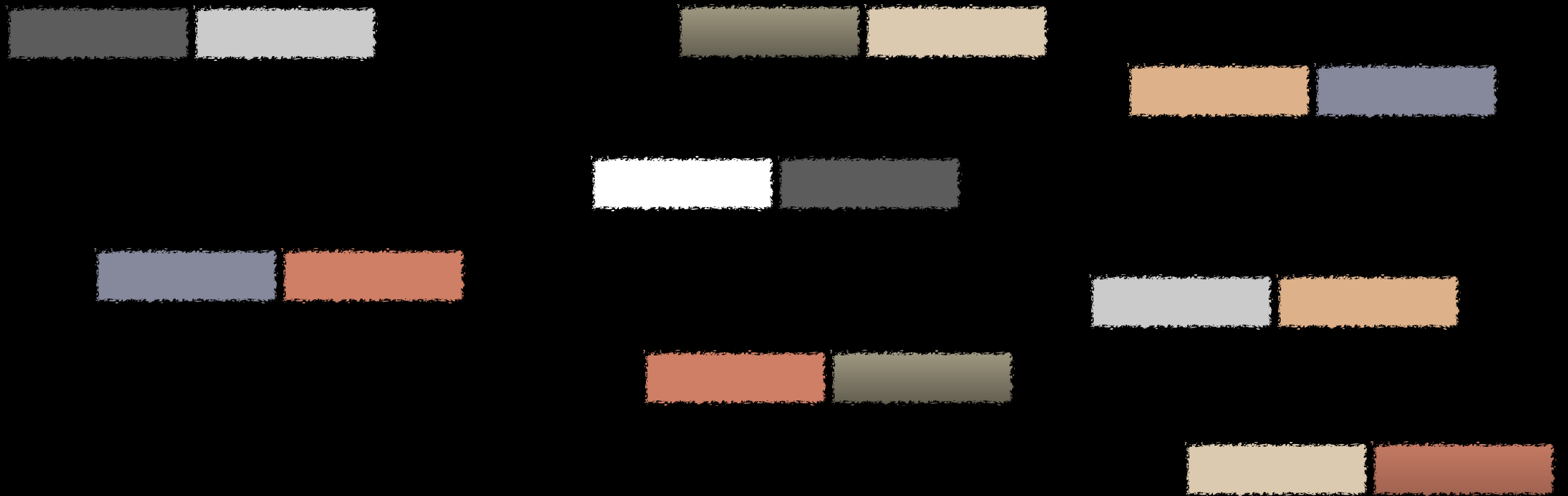
So let us sequence Mr. H.



Using reads that are 2 bps long.



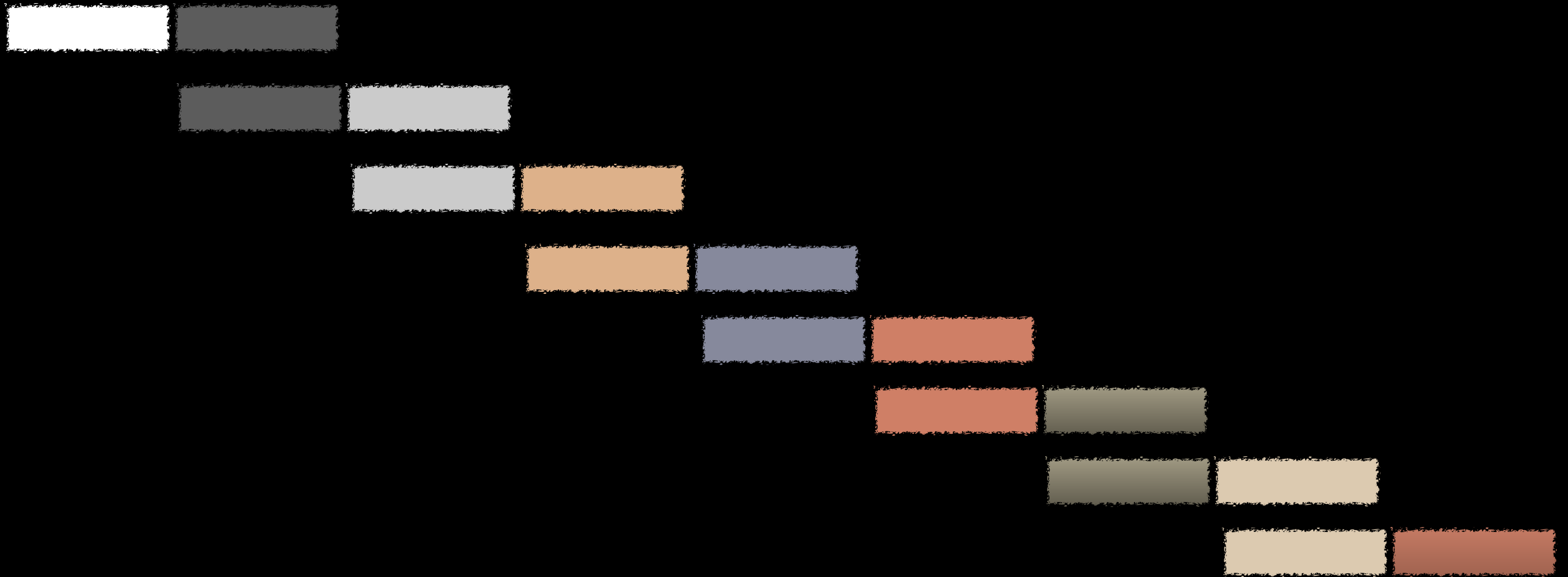
We actually do not know his genome
sequence!!!



De Novo Assembly



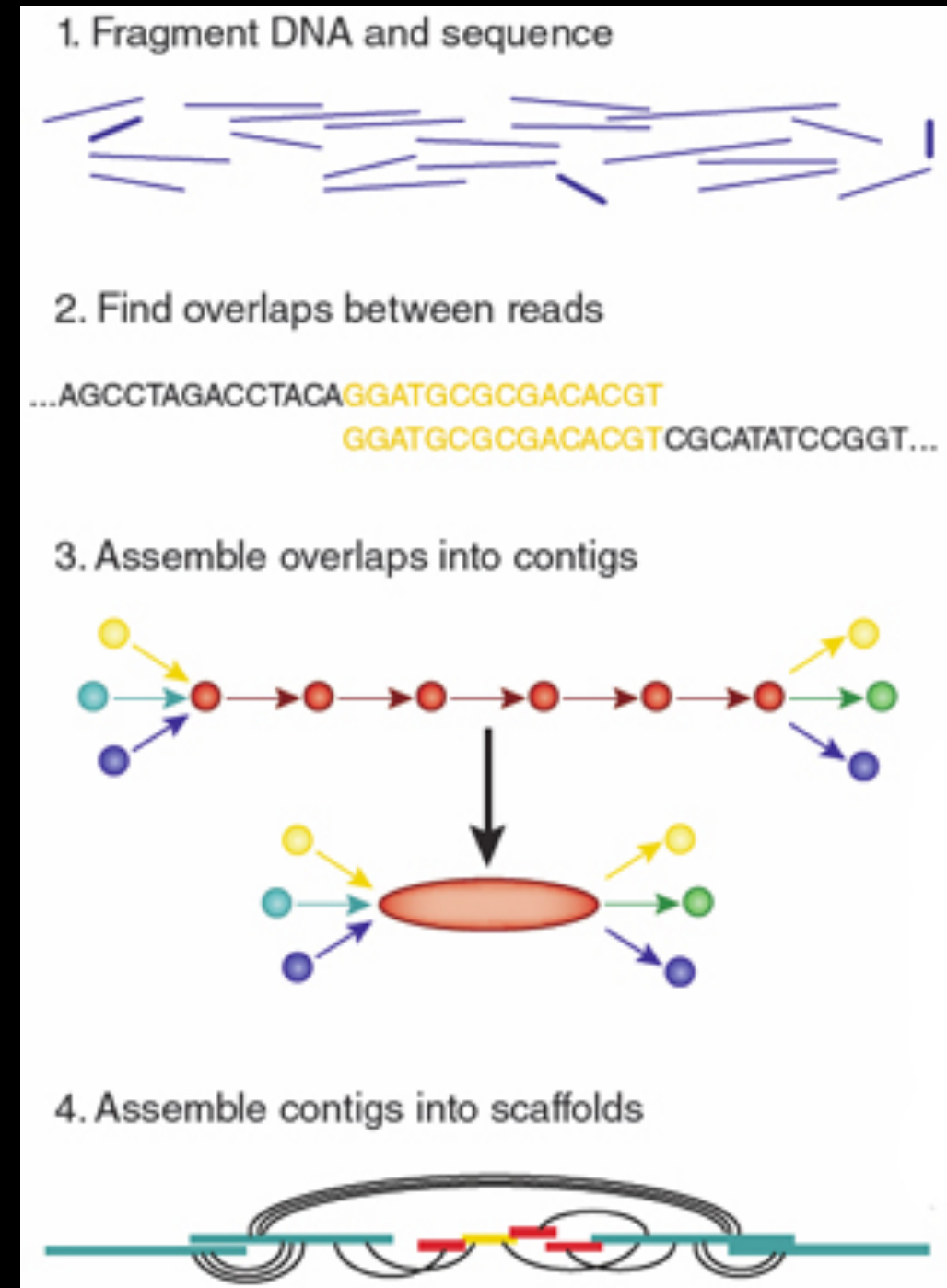
Overlapping Fragments



De Novo Assembly

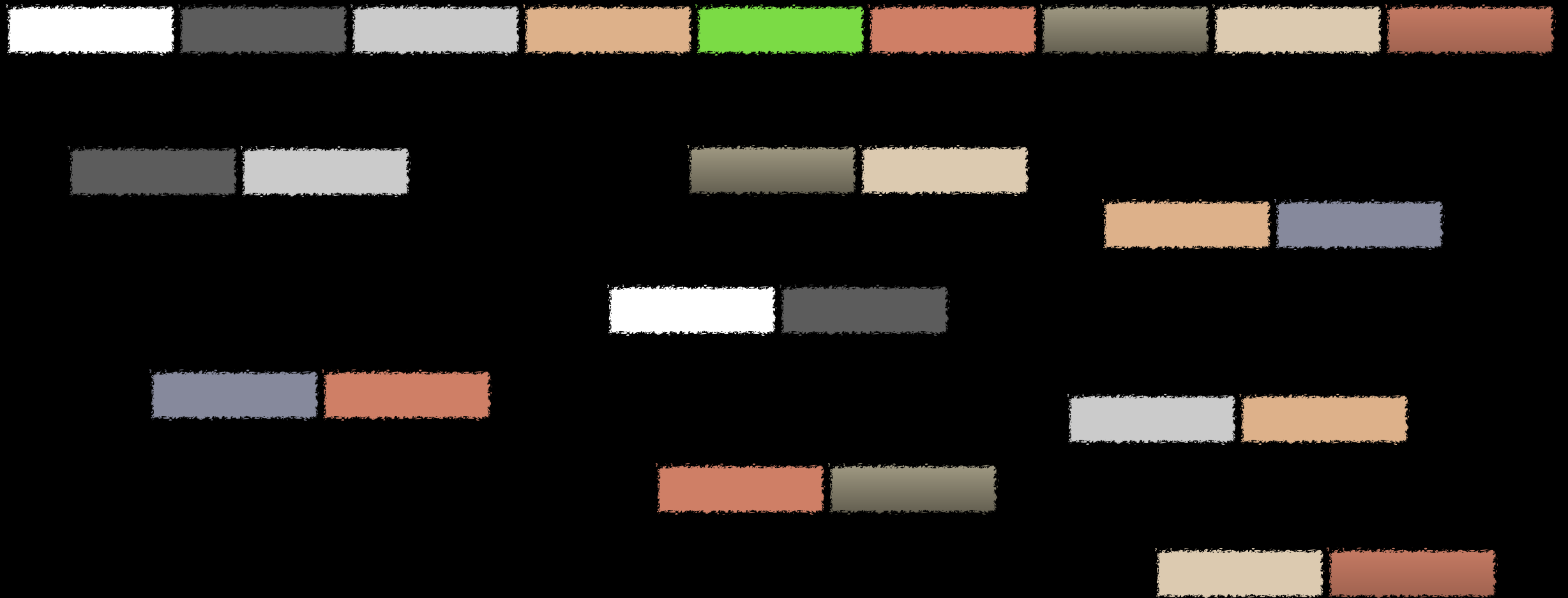
De Novo Assembly

- Repeats complicate assemblies.
- Typically require large amounts of memory for mammalian sized genomes
- Several approaches:
 - Overlap graphs
 - De Bruijn graphs
- Some de novo assemblers for short-reads
 - Velvet, ABySS, Forge, SOAPdenovo...

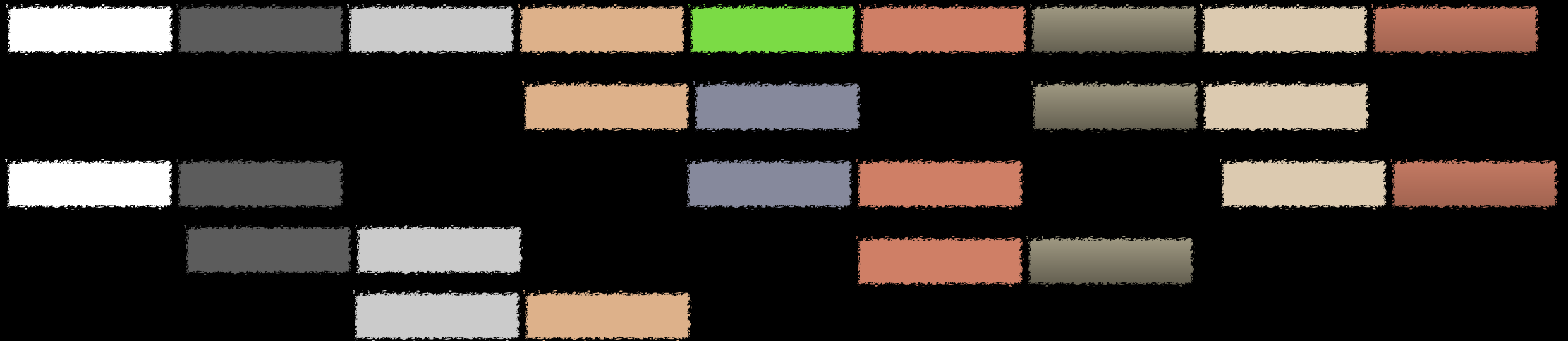




But wait !!! We have a reference genome
from Mr. T.



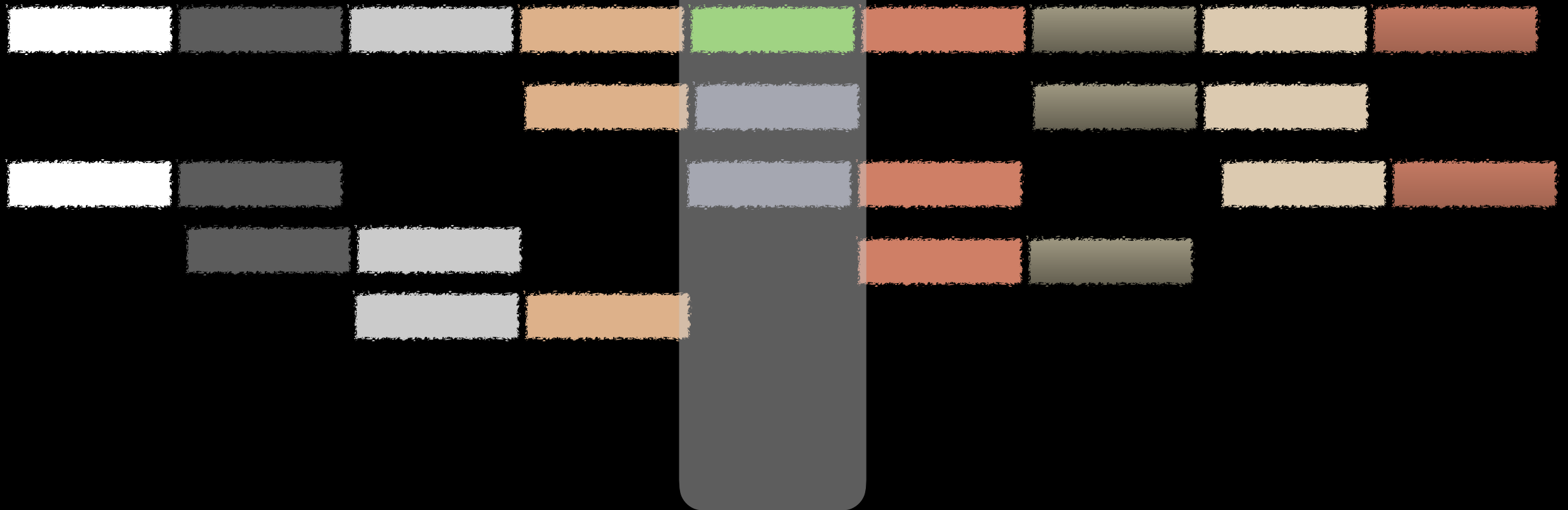
So now we need to find the best placement for each sequence from Mr. H.



So now we need to find the best placement for each sequence from Mr. H.



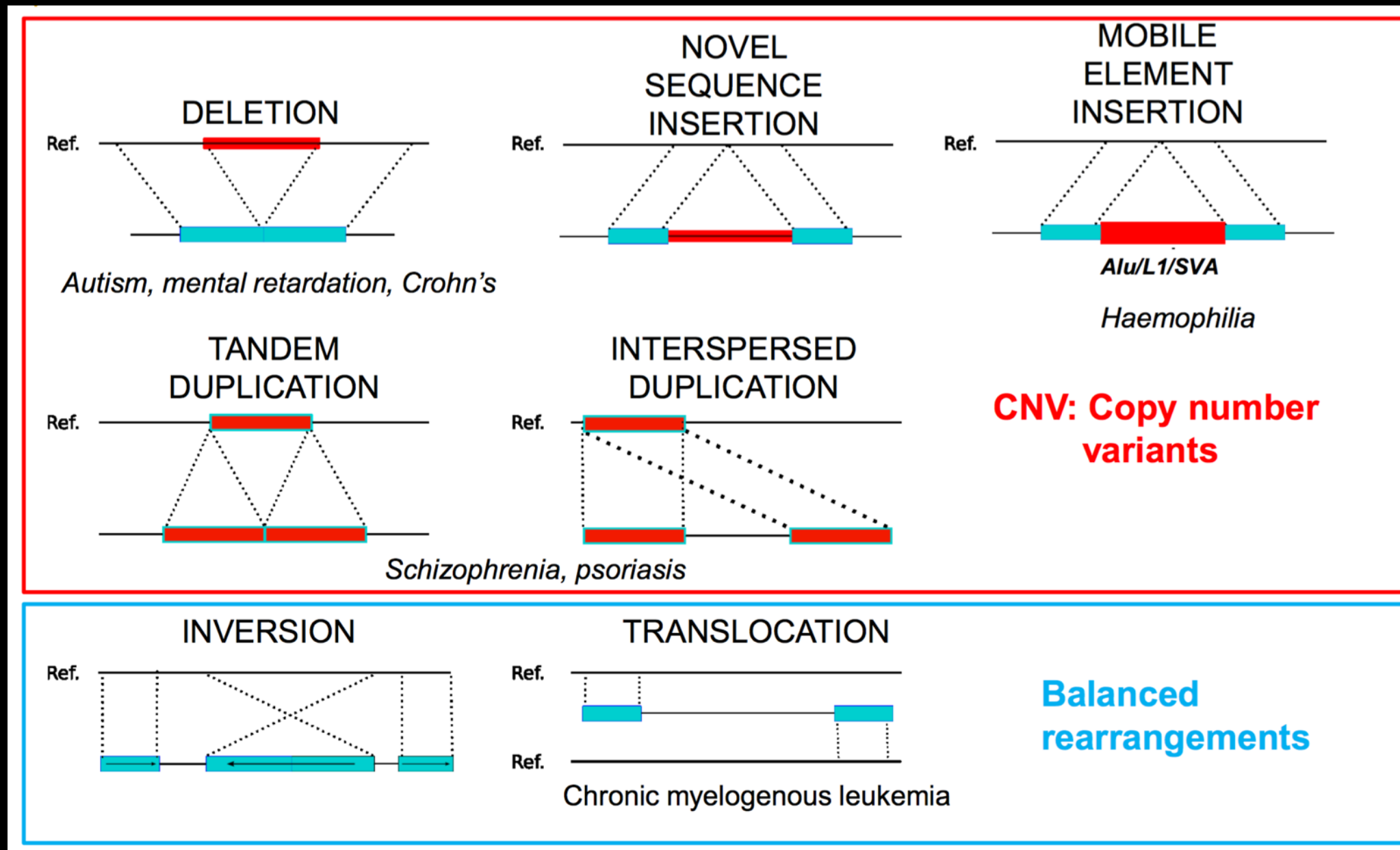
SNP



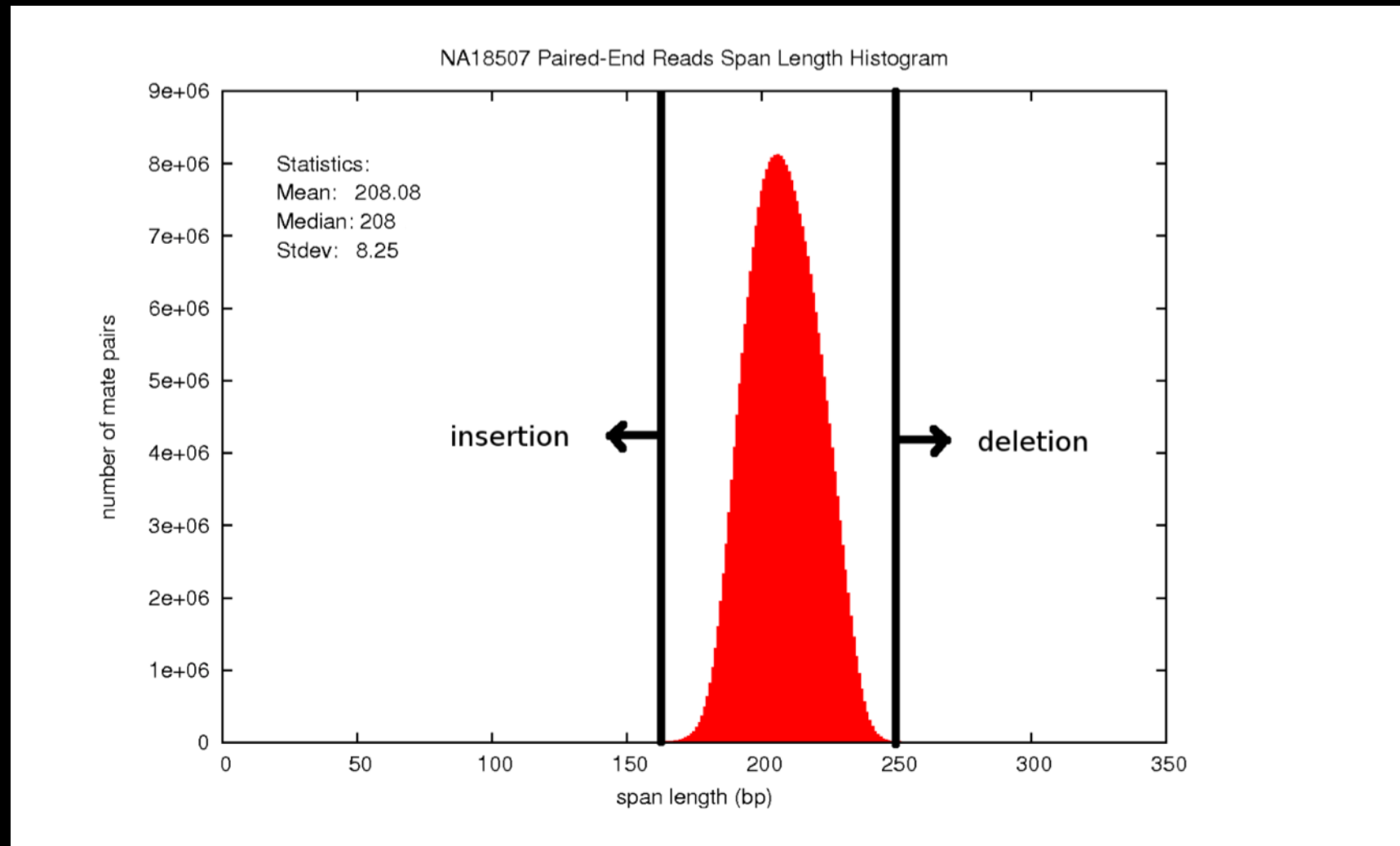
So now we need to find the best placement for each sequence from Mr. H.

Aligning to a reference

- Typically faster and requires less resources
- SNPs and other variations are more easily placed and identified
- Large fraction of sequence that does not align is either really divergent or not present in the reference
- Several approaches
 - Seed and extend
 - BWT ...
- Some alignment tools:
 - BWA, Bowtie, LASTZ, LAST, BLAST ...



Classes of other variants

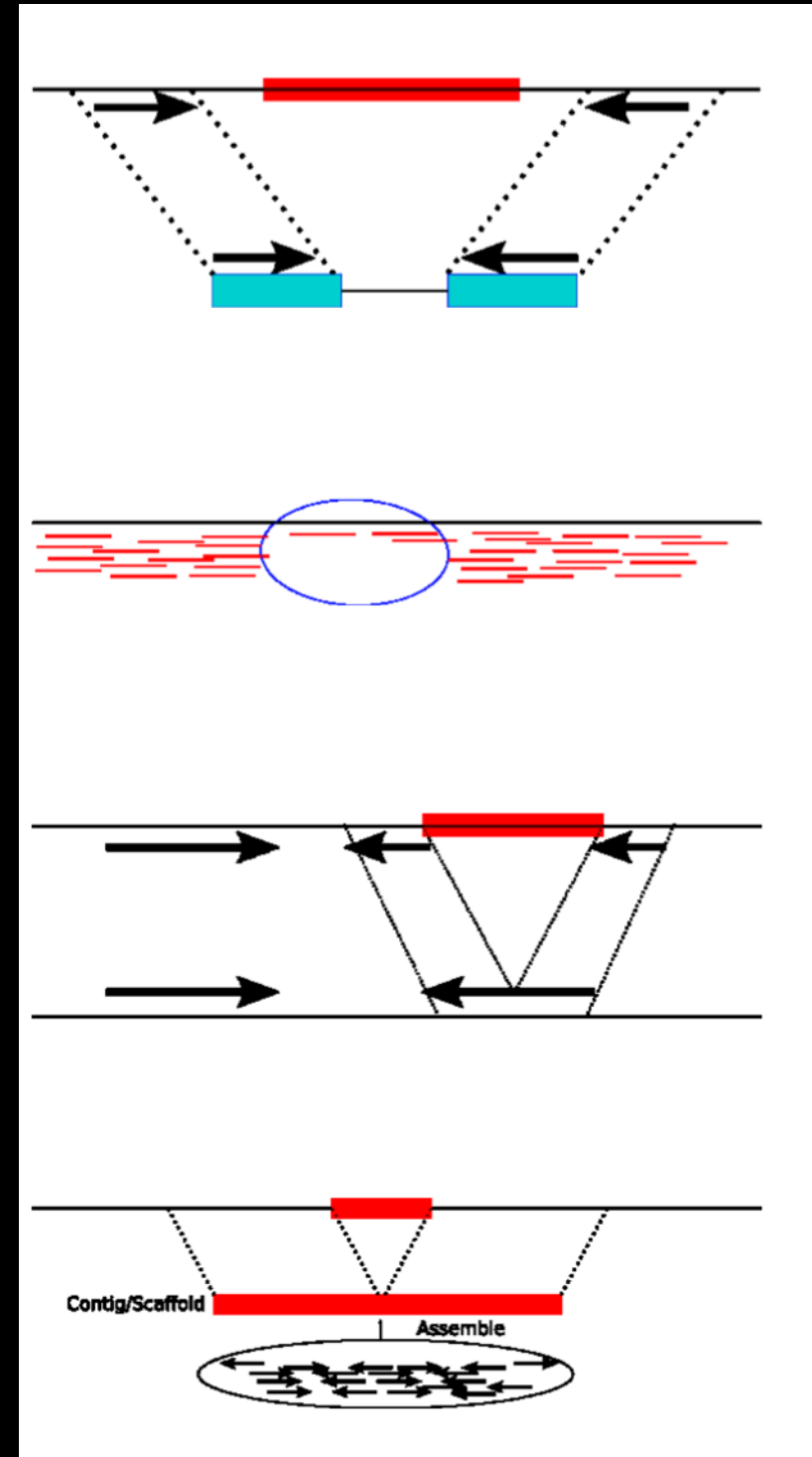


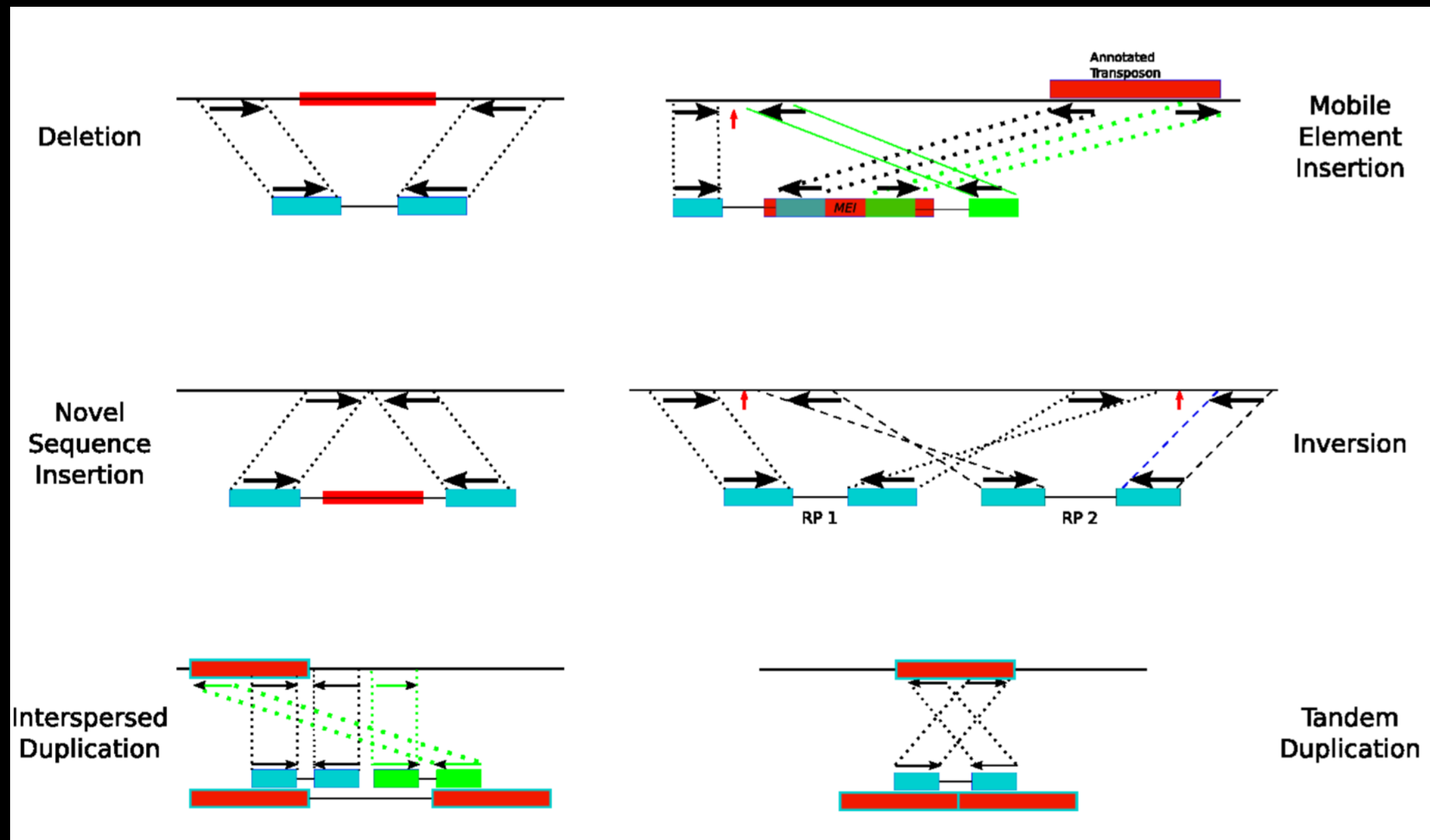
Concordant = read pairs that map in expected orientation & size
Discordant = read pairs that map different than what is expected

Insert Length Distribution

Sequence signatures of structural variation

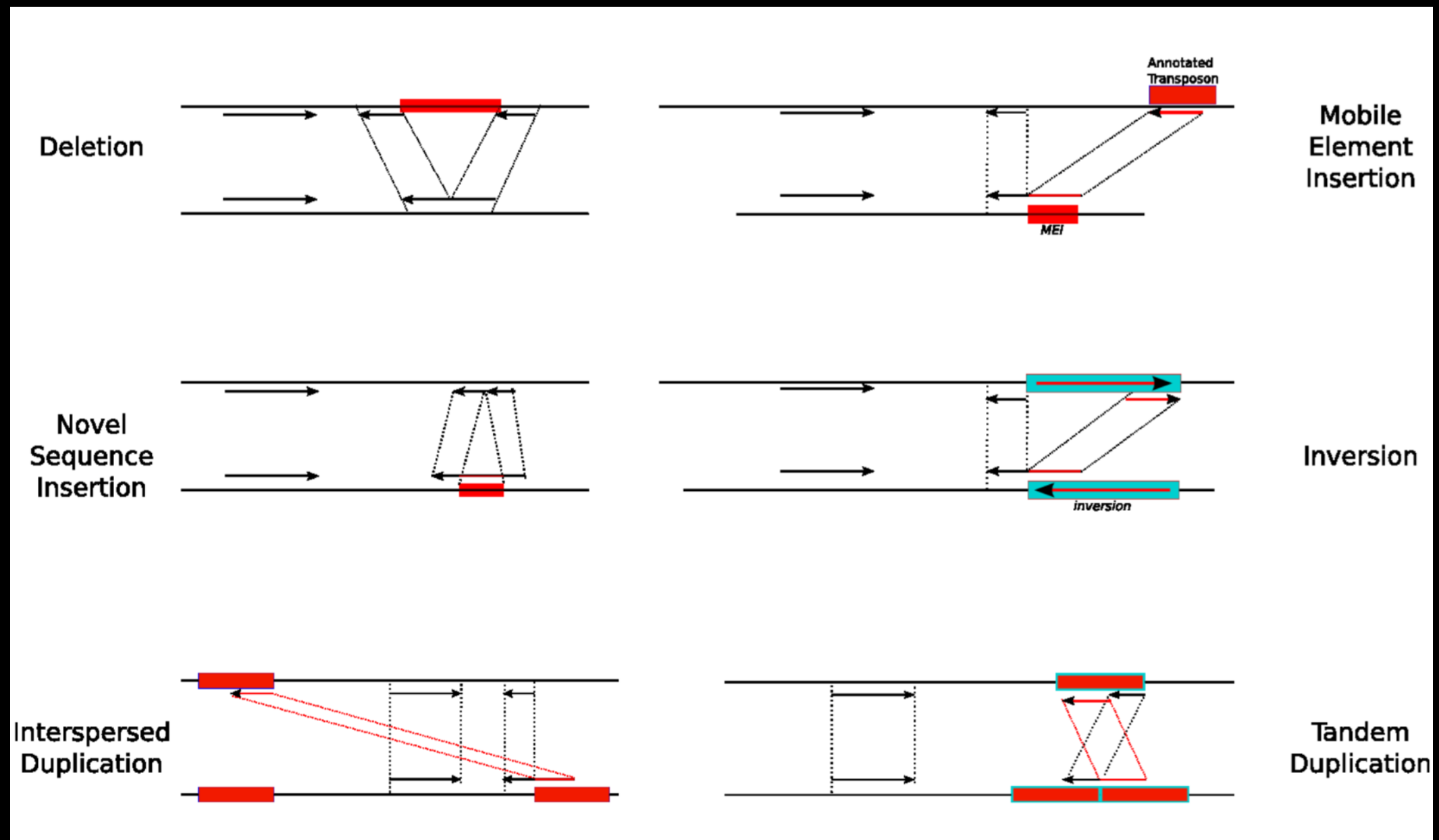
- Read pair analysis
 - Deletions, small novel insertions, inversions, transposons
 - Size and breakpoint resolution dependent to insert size
- Read depth analysis
 - Deletions and duplications only
 - Relatively poor breakpoint resolution
- Split read analysis
 - Small novel insertions/deletions, and mobile element insertions
 - 1bp breakpoint resolution
- Local and de novo assembly
 - SV in unique segments
 - 1bp breakpoint resolution





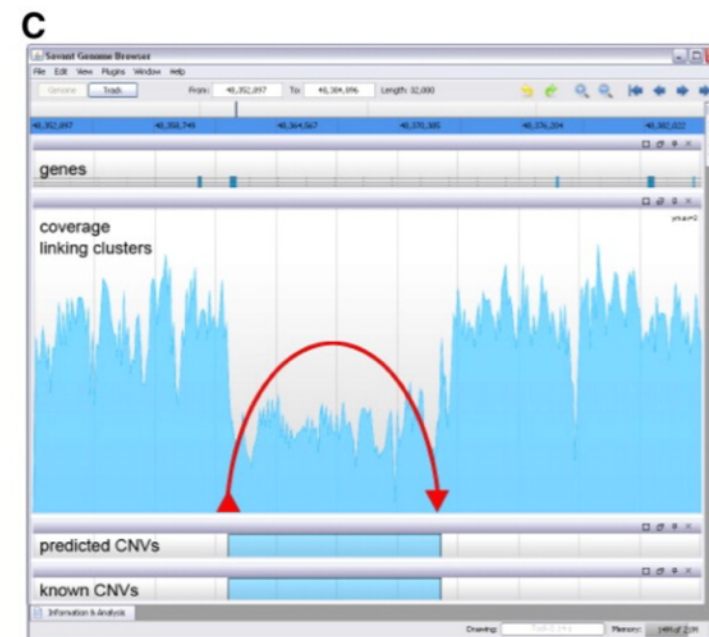
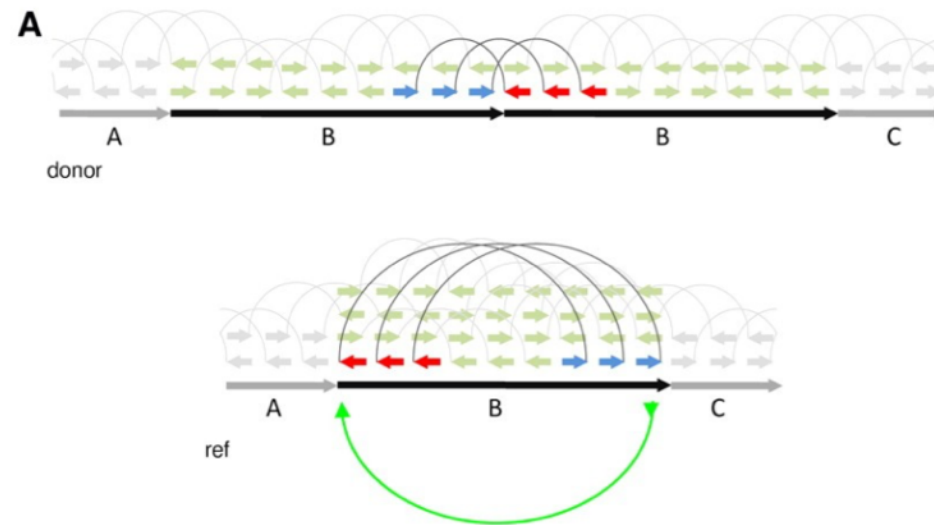
Patterns of SVs from Paired-End reads

BreakDancer, GenomeSTRiP, VariationHunter, HYDRA



Patterns of SV from split-reads

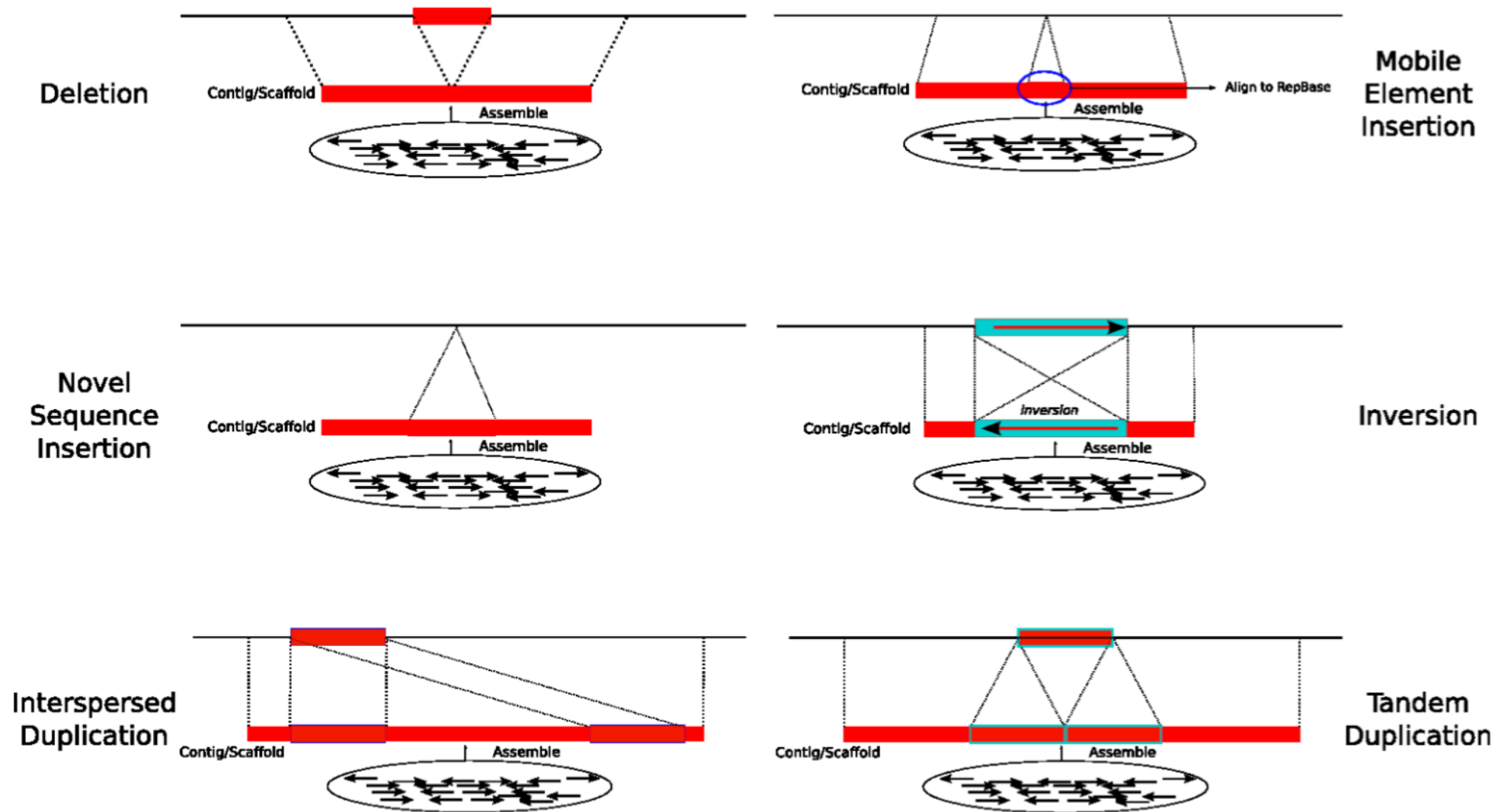
PINDEL, SPLITREAD, indelMINER



Medvedev et al., Genome Res, 2010

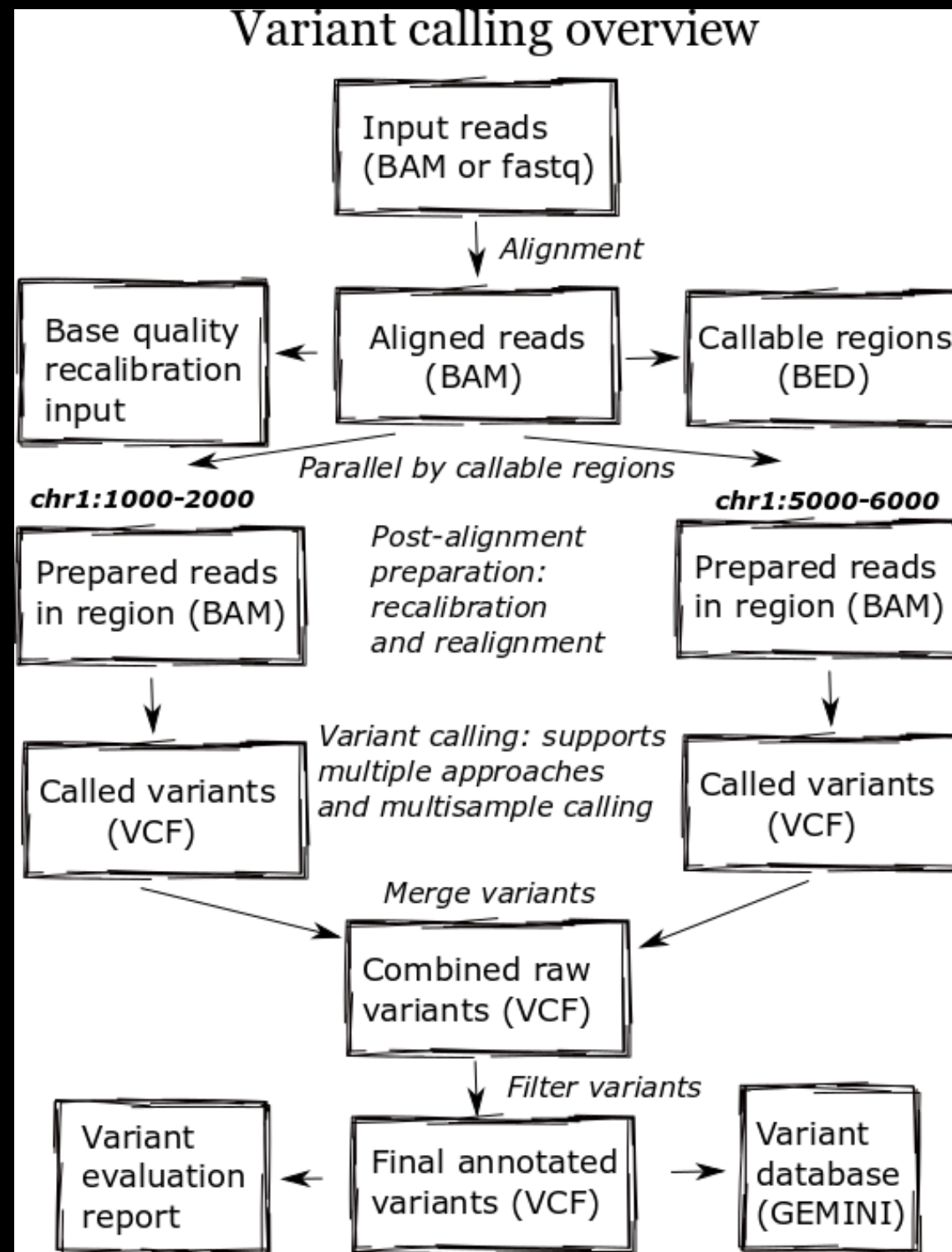
Using multiple signals

CNVer, LUMPY



Using Genome Assembly

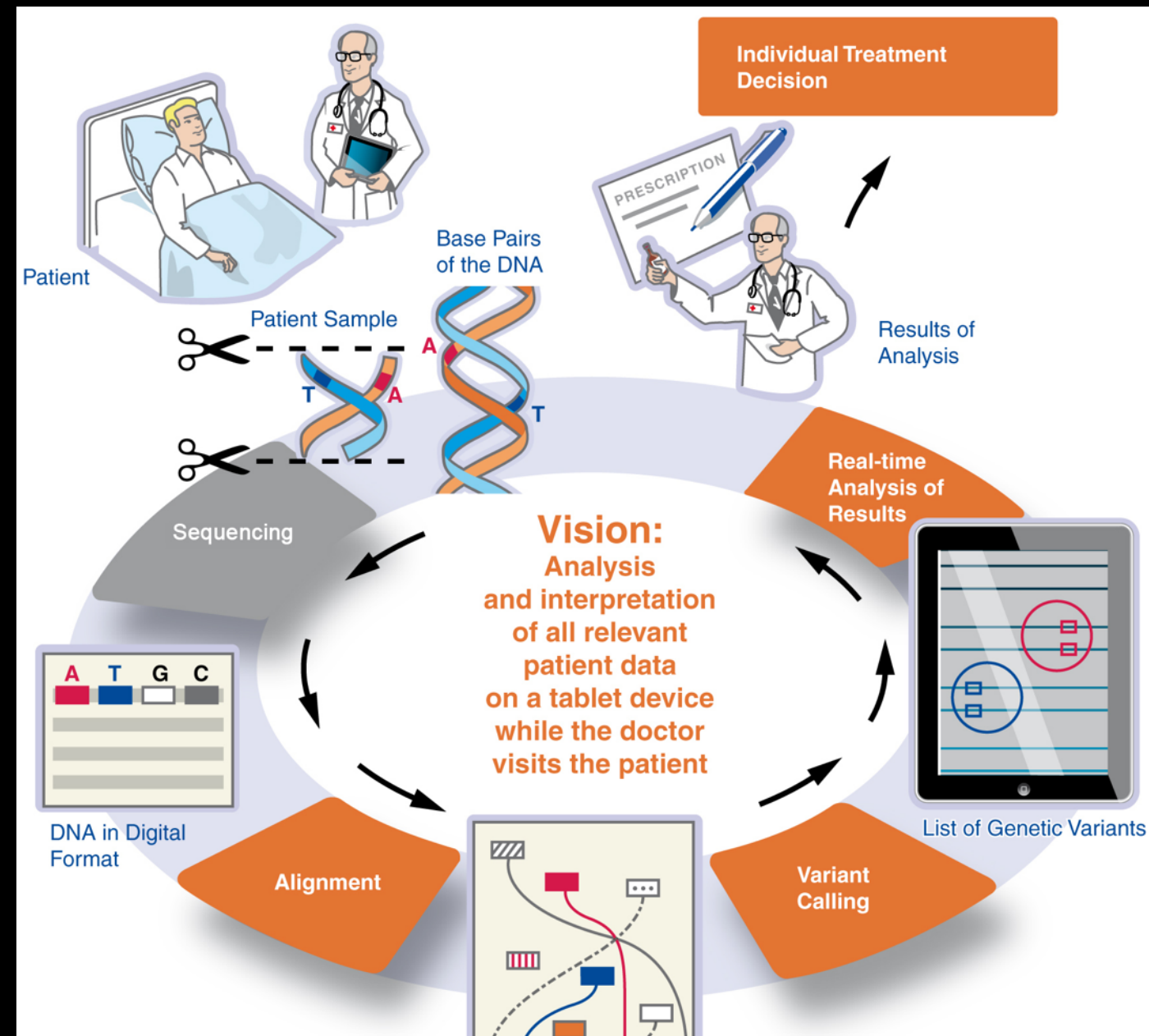
NovelSeq

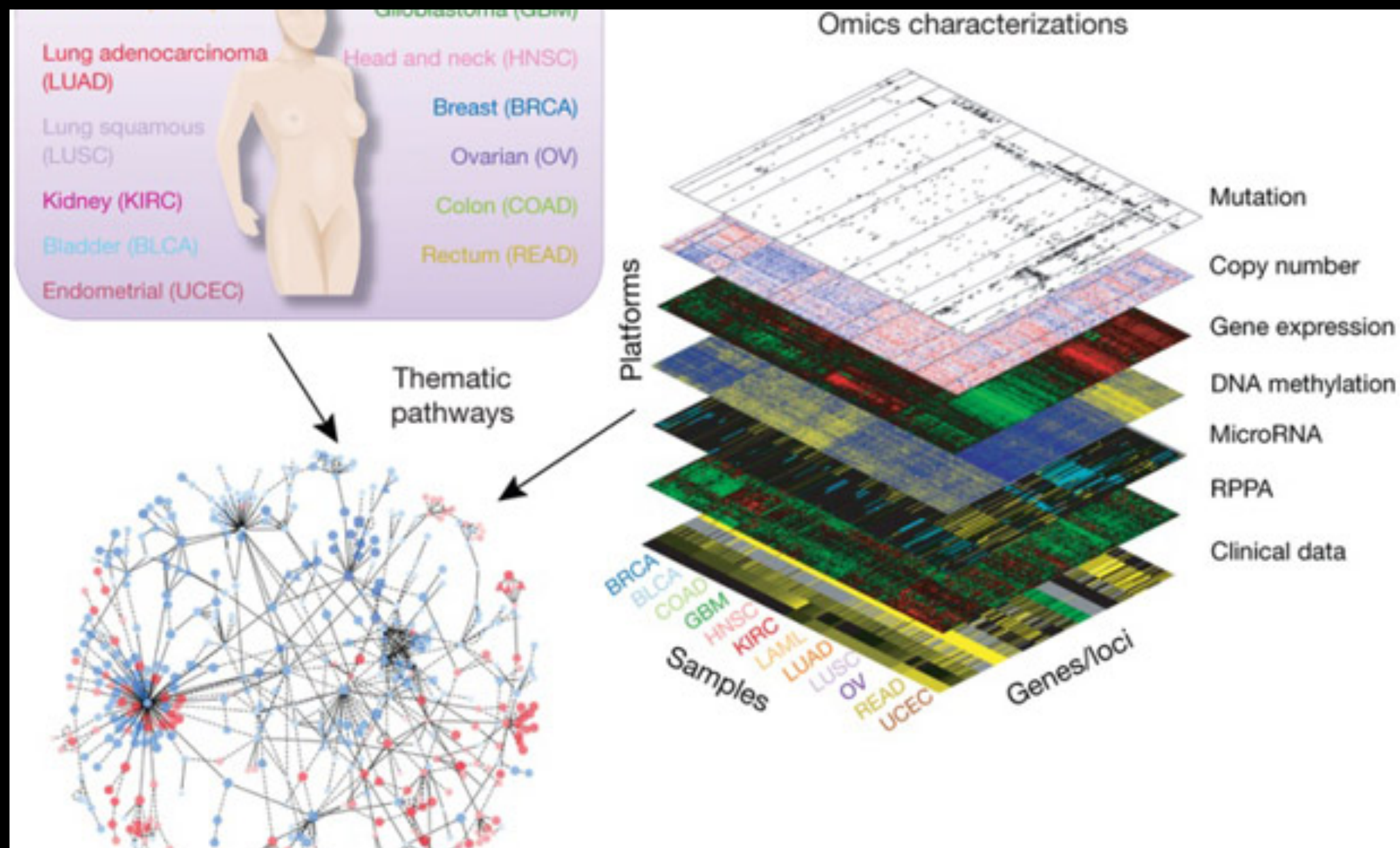


Typical Analysis Pipeline in human genome sequencing

Again, why sequence genomes?

- From genome sequences to know genome variation between individuals and study
 - Disease
 - Drug response
 - Biomes
 - Energy
 - Agriculture ...





Combining Datatype to enable precision medicine, improvements in agriculture, energy ...



HAPPY PI DAY