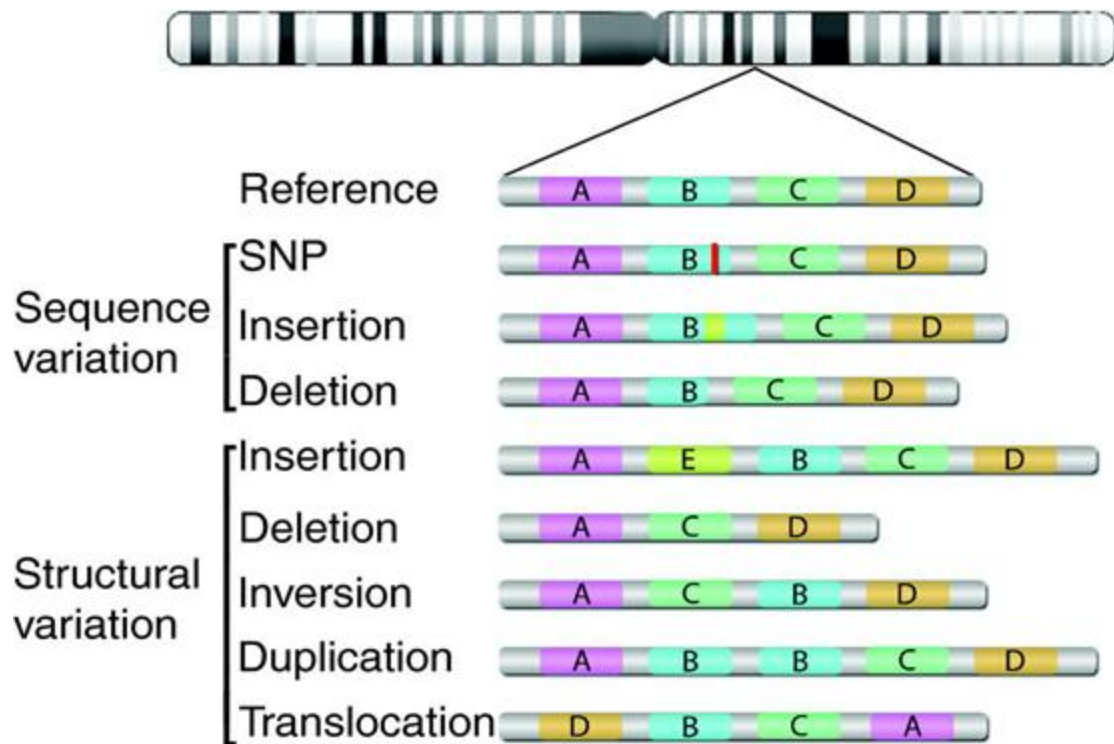


# Variant calling

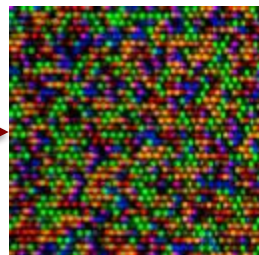
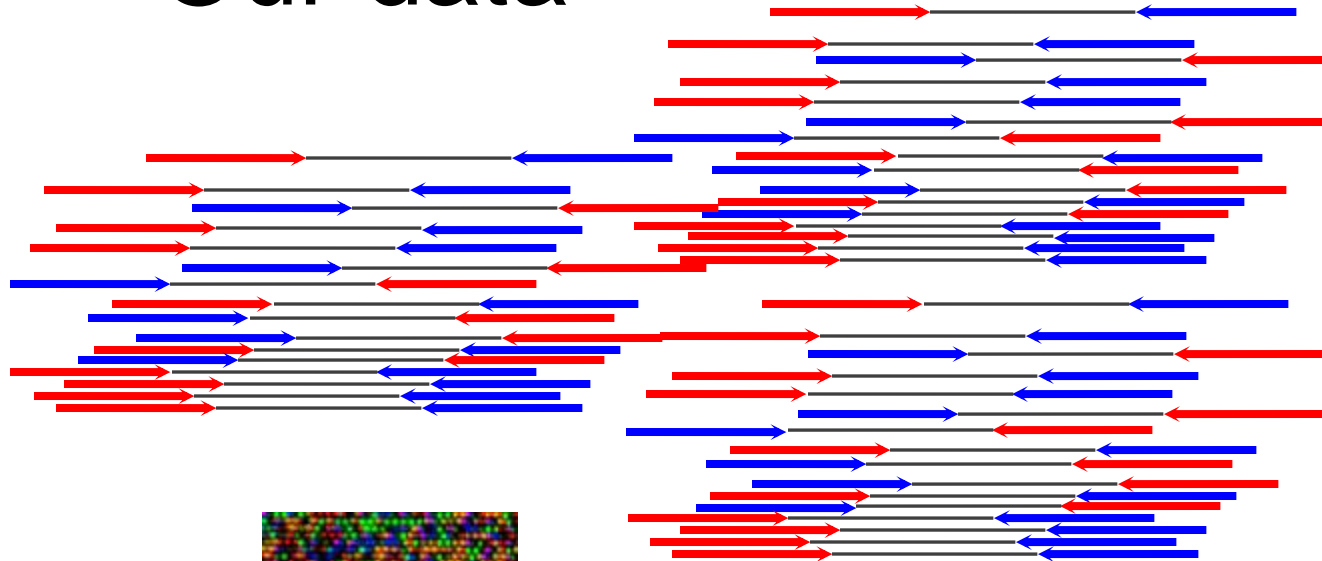
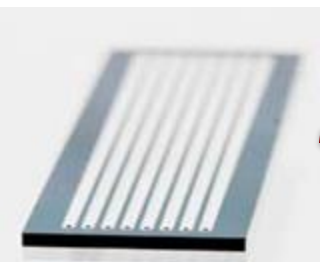
Chris Miller

Some slides adapted from Dave Larson, Aaron Quinlan, Sam Peters, and Alex Paul

# Small Variant Calling

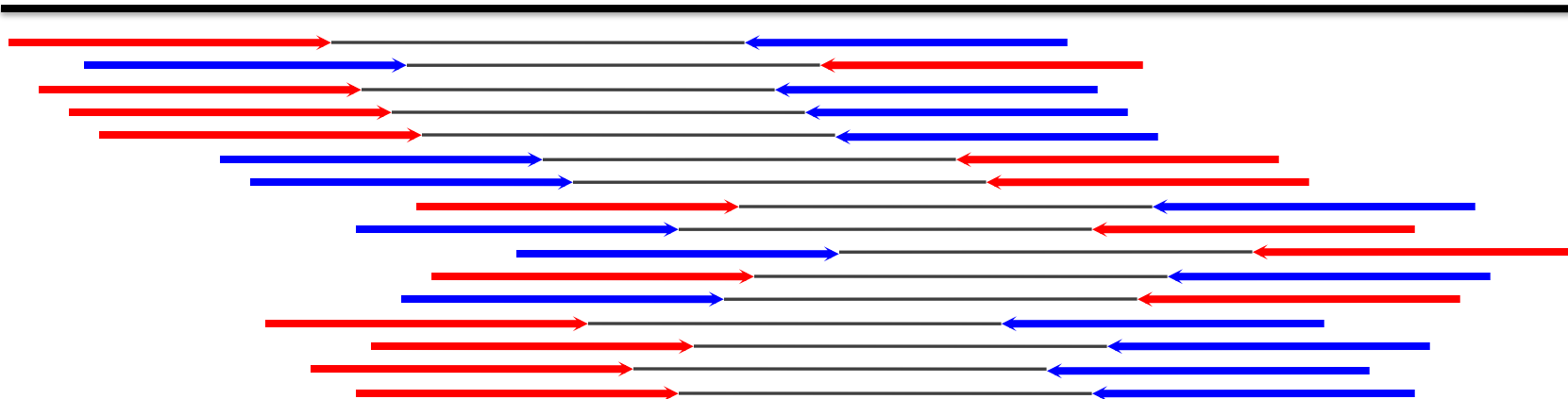


# Our data



# Mapping

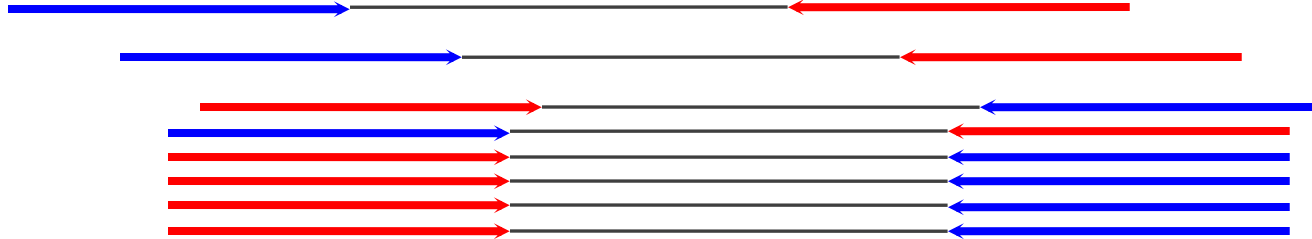
Genome Reference Sequence



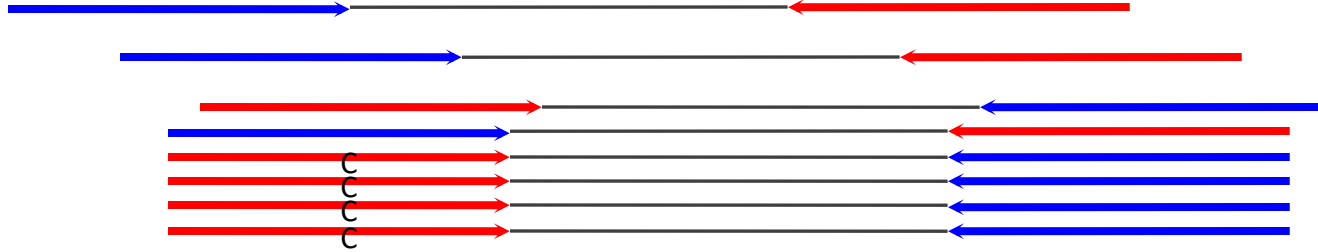
- Single-end reads can be longer, less unique depending on sequence context
- Paired-end reads can span repetitive regions, provide additional information
- Mapping has gotten quite fast, <24 hours for 120 Gbp of sequence
- Split-read alignments are the norm (BWA mem)



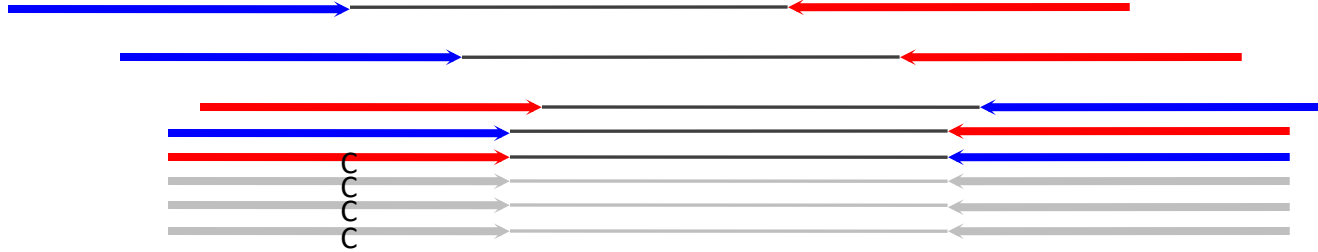
# Duplication



# Duplication



# Duplication

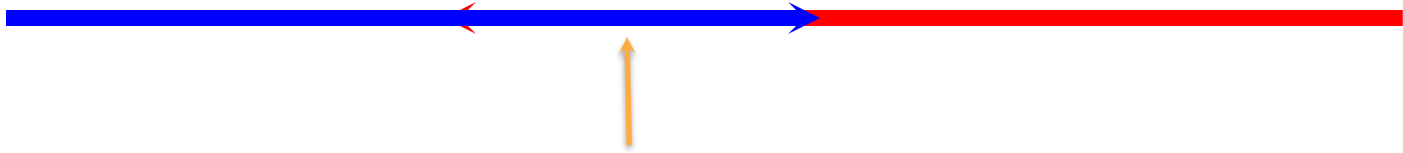


# Overlapping reads

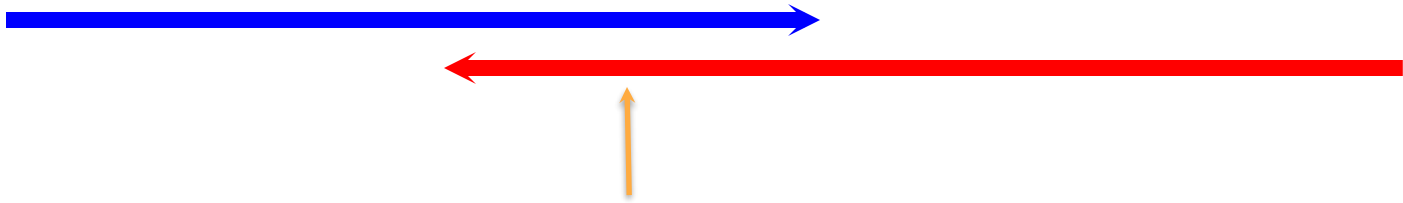




# Overlapping reads



# Overlapping reads



# Overlapping reads



# Every aspect of this process is fraught with error

- Base calling is not perfect: *0.5% error on average*
- Mapping is not perfect: *the reads are short*
- The reference sequence is not perfect

# We have a little help

- Some uncertainty is encapsulated in quality scores
  - the rate at which the data is expected to be wrong
- Each base call (ACTGN) comes with a quality
  - Phred-scaled ( $-10 * \log_{10}$  of quality)
  - A base call with quality of 20 is wrong 1 out of every 100 times.
- Read mapping has quality too
  - These are also Phred-scaled

# Phred quality score calculation

$$Q = -10 \cdot \log_{10}(P_{\text{err}})$$

Error probability ( $P_{\text{err}}$ )	$\log_{10}(P_{\text{err}})$	Phred quality score
1	0	0
0.1	-1	10
0.01	-2	20
0.001	-3	30
0.0001	-4	40

# Goals of a Variant Caller

- Sensitive detect mutations
- Precisely detect mutations
  - Confounded by the error we just talked about
  - FDR must be very low as we're looking across a very large space!

# Goals of a Variant Caller

- Sensitive detect mutations
- Precisely detect mutations
  - Confounded by the error we just talked about
  - FDR must be very low as we're looking across a very large space!
- **An FDR of 0.001 = 3.2 million false positives!**

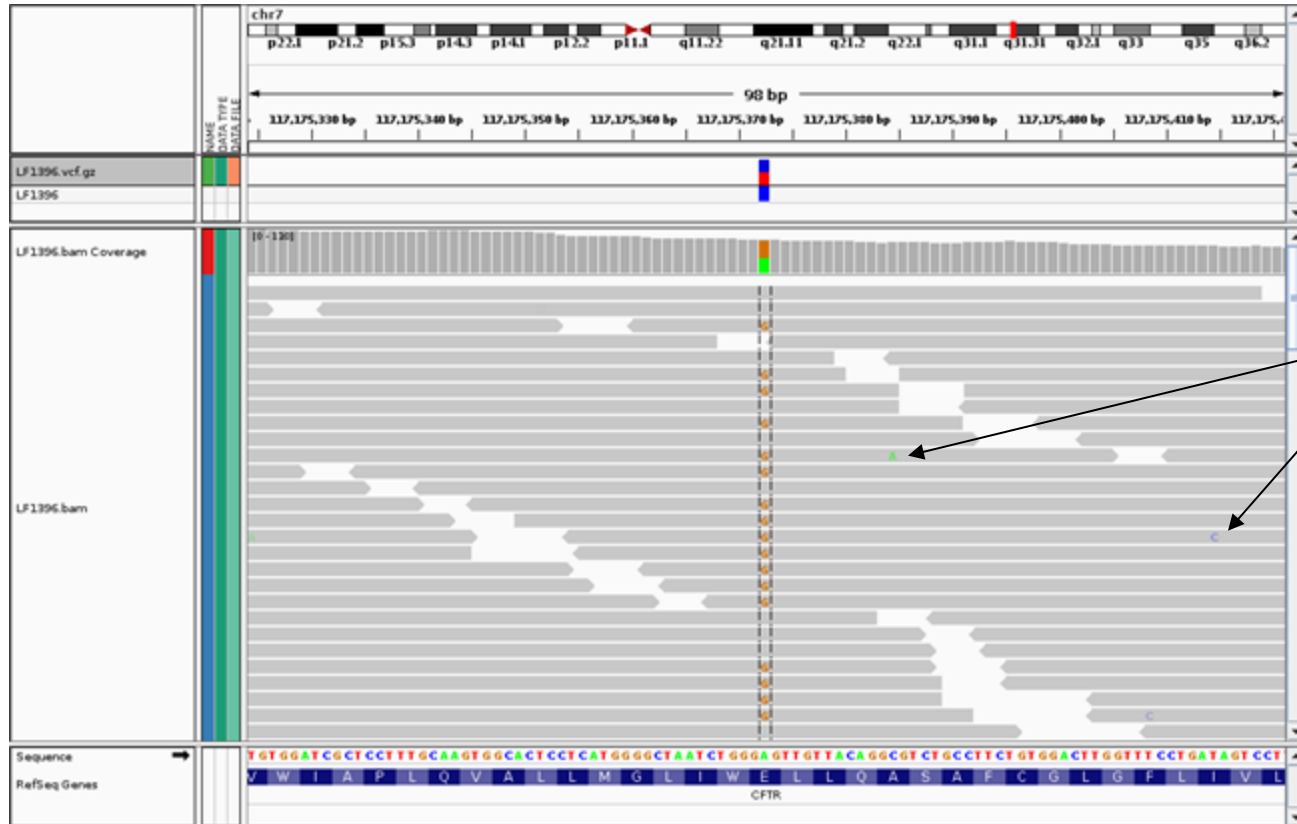


# Homozygous for the "C" allele

Improper  
(too far/too  
close) pairs



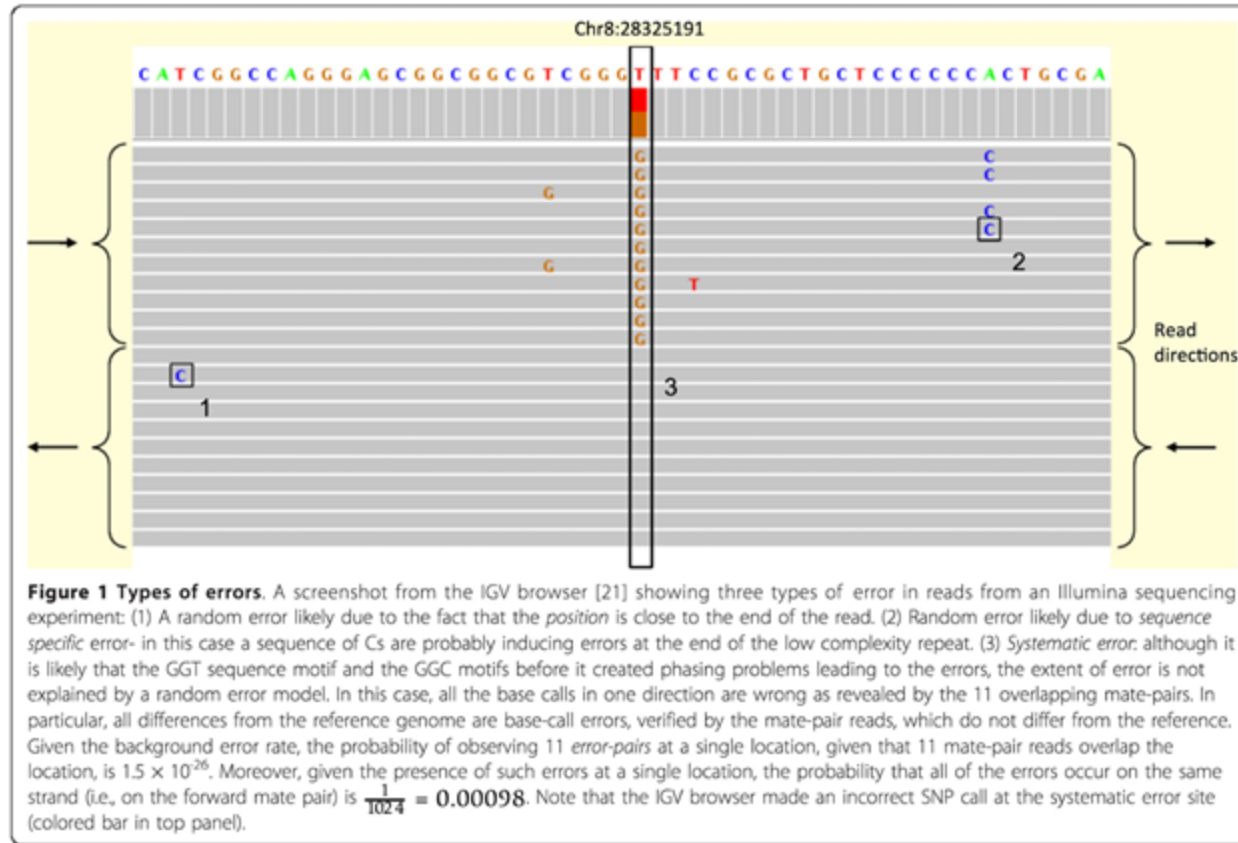
# Sequencing errors fall out as noise (most of the time)



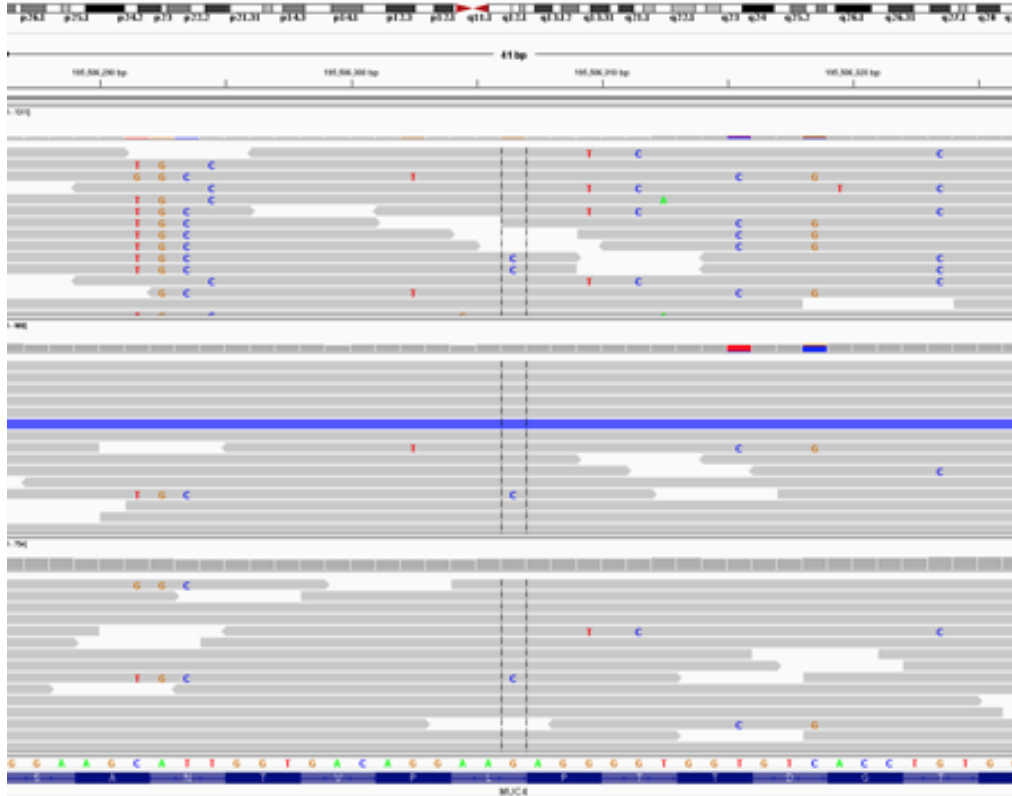
Sequencing errors

It is not always so easy

# Random versus systematic error



# Pileups of many differences from paralogy



RESEARCH ARTICLE | OPEN ACCESS

## FLAGS, frequently mutated genes in public exomes

Casper Shyr, Maja Tarallo-Graovac, Michael Gottlieb, Jessica JY Lee, Clara van Karnebeek and Wyeth W Wasserman

BMC Medical Genomics 2014 7:64 | DOI: 10.1186/s12920-014-0064-y | © Shyr et al.; licensee BioMed Central Ltd. 2014

Received: 16 June 2014 | Accepted: 24 October 2014 | Published: 3 December 2014

Open Peer Review reports



# Calling INDELs is *much* harder than SNPs

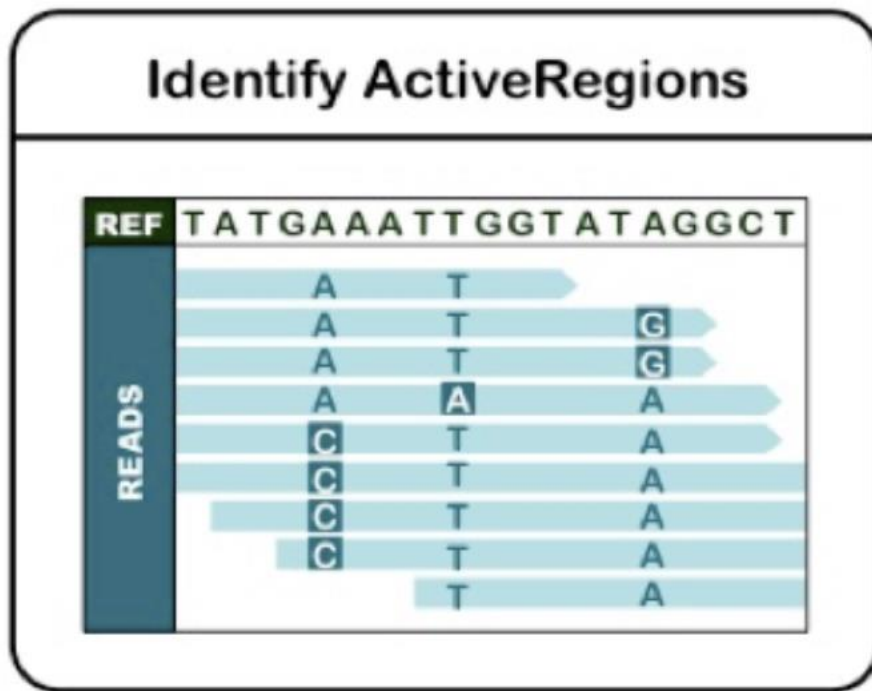


# Exercise

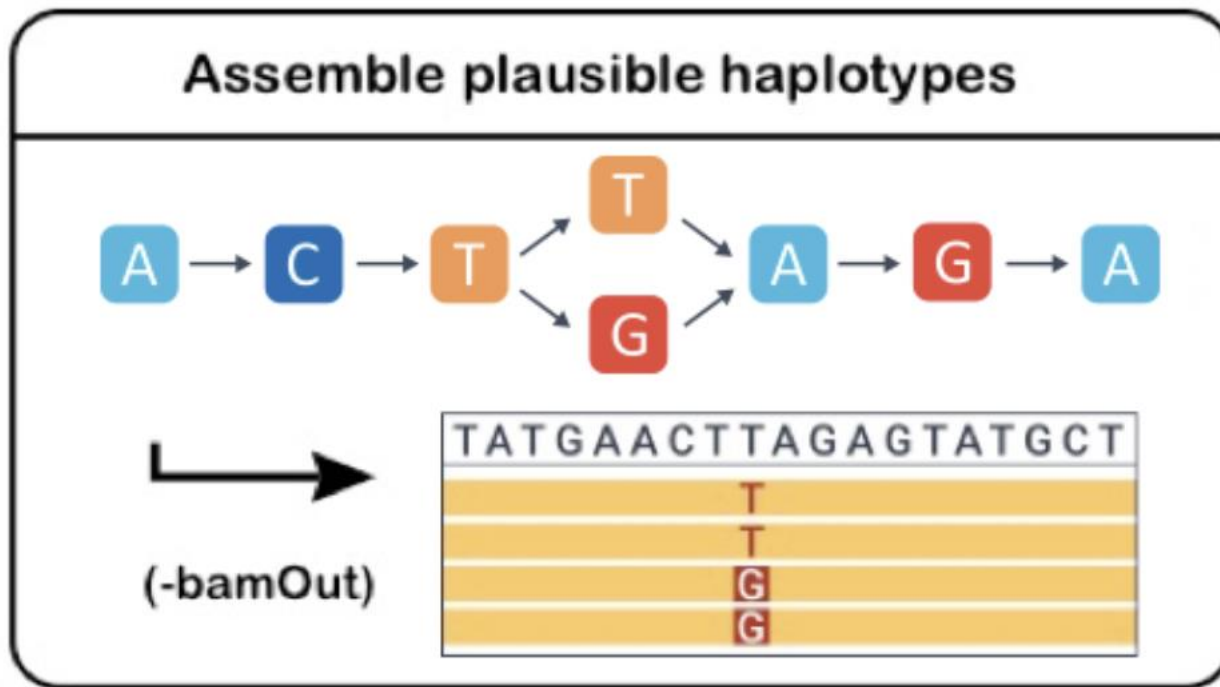
# Germline SNV and Indel Calling



# Call Genotypes Using GATK HaplotypeCaller



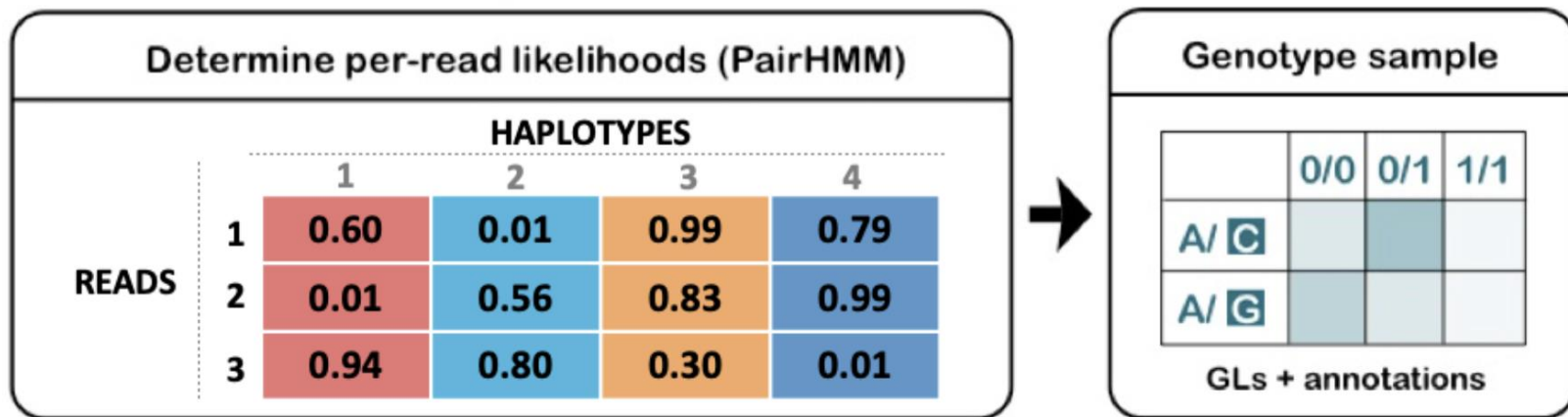
# Call Genotypes Using GATK HaplotypeCaller



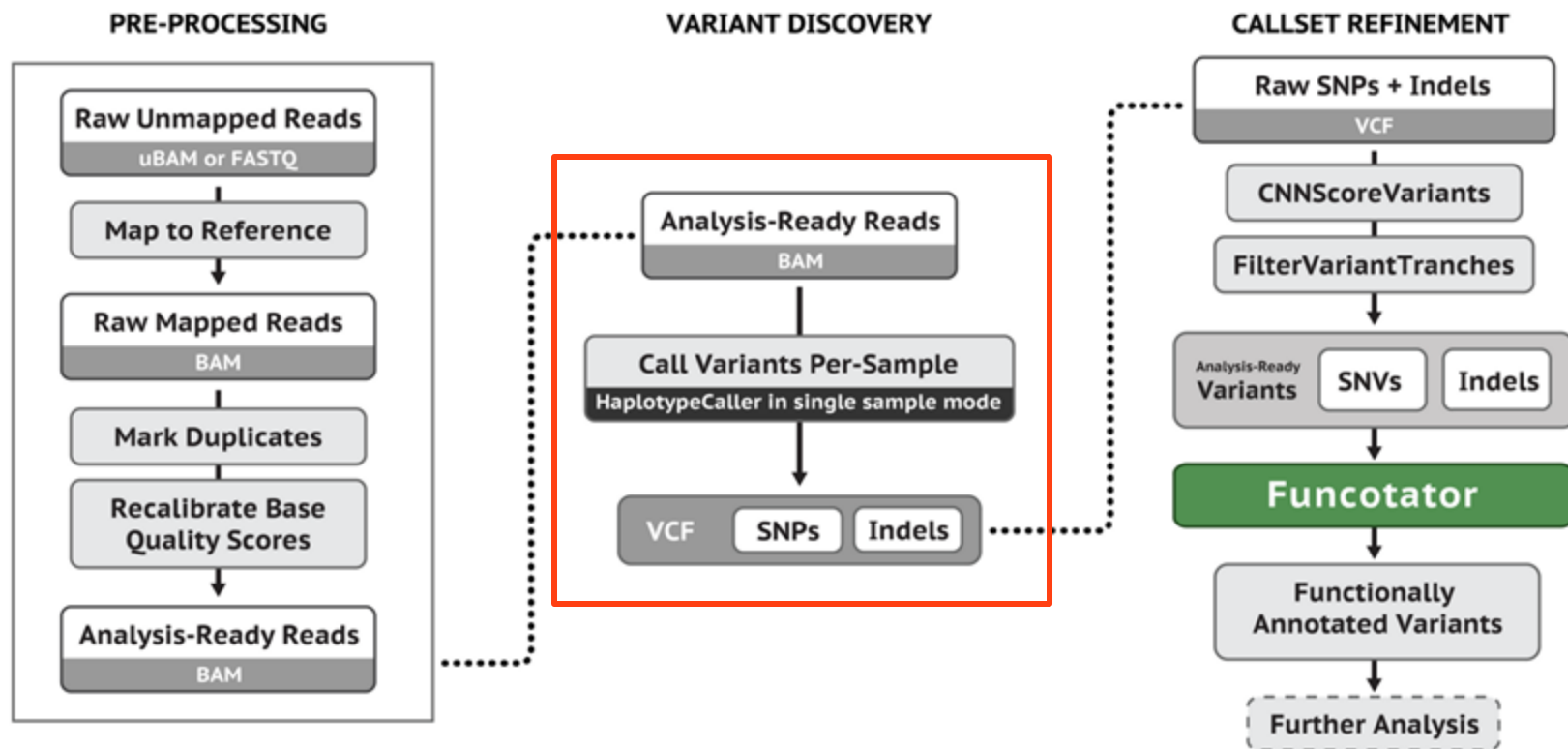
# Indel "realignment"



# Call Genotypes Using GATK HaplotypeCaller

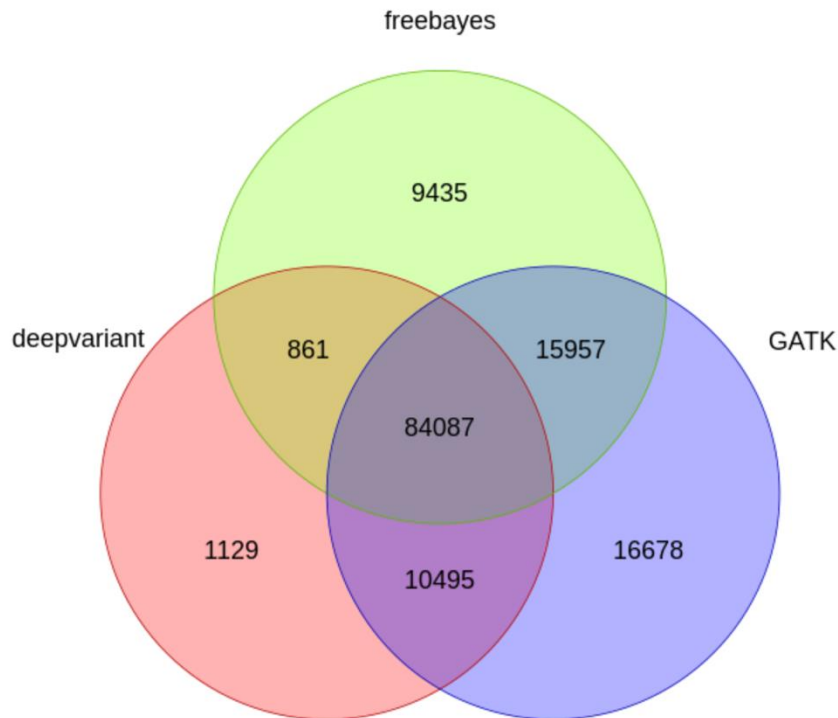


# Main steps for Germline Single-Sample Data



# Other Germline Variant callers

- Freebayes: lighter weight, faster
- DeepVariant: Neural-network based caller from Google ML

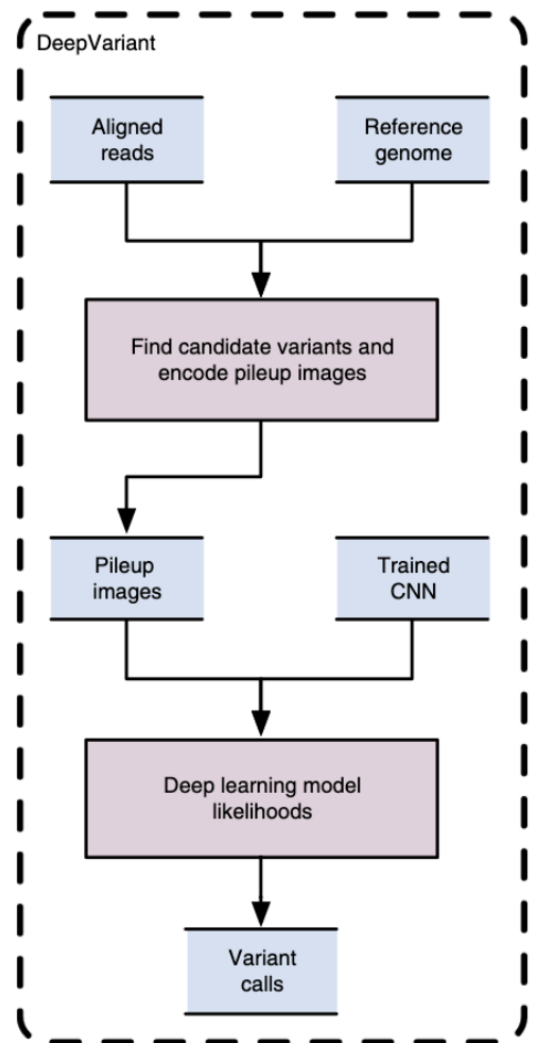


# DeepVariant

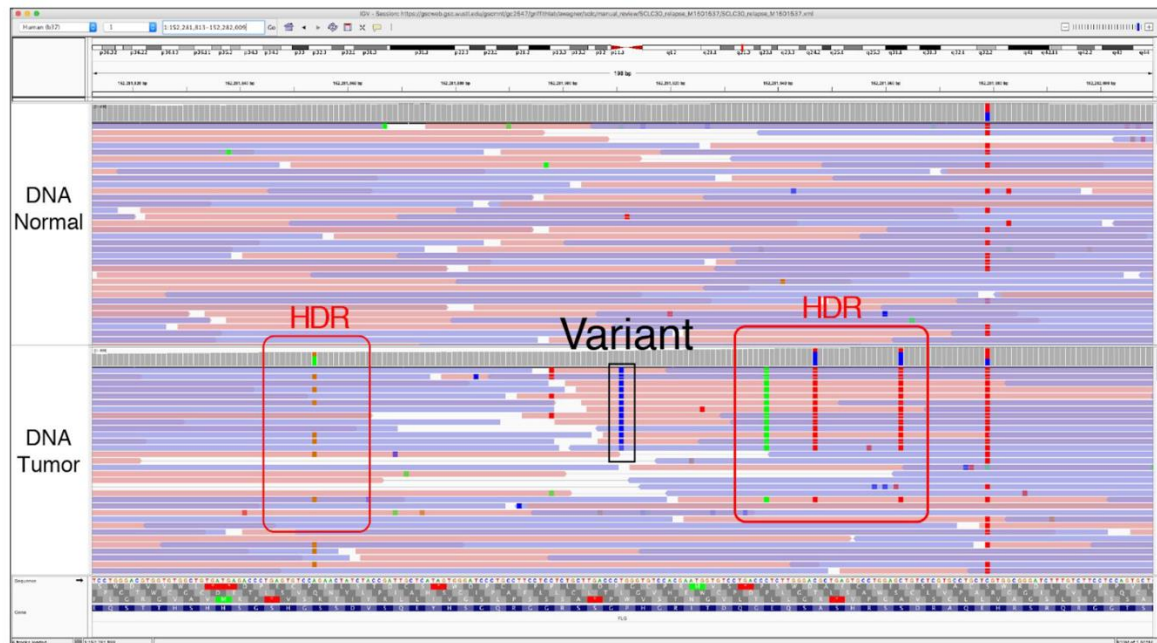


<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6450397/>

<https://www.nature.com/articles/nbt.4235>

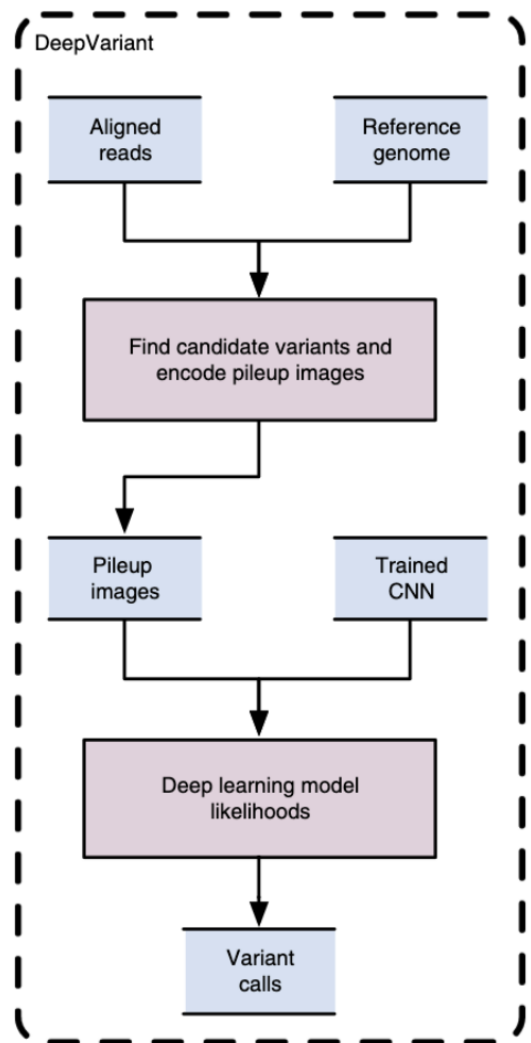


# DeepVariant



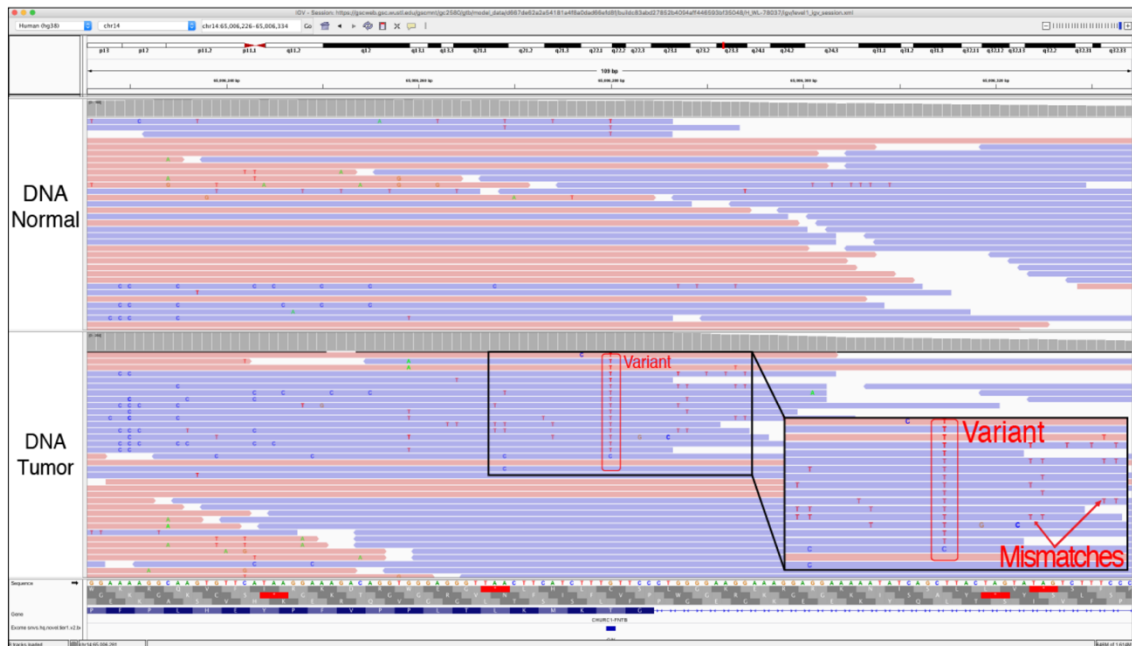
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6450397/>

<https://www.nature.com/articles/nbt.4235>



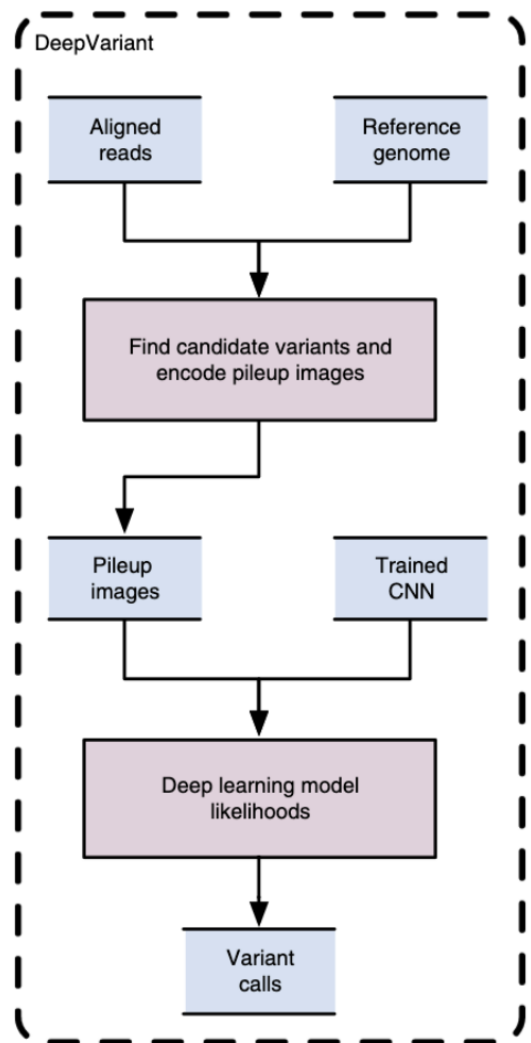


# DeepVariant



<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6450397/>

<https://www.nature.com/articles/nbt.4235>



The screenshot displays the Integrative Genomics Viewer (IGV) interface. At the top, the URL bar shows the session link. The main view consists of several horizontal tracks:

- Reference:** Shows the genomic coordinates from 149,754,000 to 149,754,800 bp.
- DNA Normal:** Displays sequencing reads from the normal sample, mostly aligned to the reference.
- DNA Tumor:** Displays sequencing reads from the tumor sample. A specific variant is highlighted with a red box, indicating a G-to-A transition at position 1:149,754,652.

An information panel on the right provides detailed metadata for the selected variant:

```

Read name = K00193.86.HPLSYBBXX.1.2109.21065
Sample = H_NP-234000-M1501537
Library = H_NP-234000-M1501537-g1-lb1
Read group = 2996181311
Read length = 151bp

Mapping - Primary @ BAMRD
Reference span = 1:149,754,791 (+) - 151bp
Cigar = 151M
Clipping = None
Mate is mapped = yes
Insert start = 1:149,754,896 (-)
Insert size = 197
Flag is pair
Pair orientation = FRXC

PG = MarkDuplicates
AQ = 2996181311
NM = 1
AS = 148
XS = 146
Hidden tags: MD

Location = 1:149,754,652
Base = A @ QV 41
    
```



<https://www.nature.com/articles/nbt.4235>

# DeepVariant

Here are 3 examples that we would consider canonical easy-to-classify loci, and that DeepVariant calls confidently and correctly:



The variant above is a "2", which means both chromosomes match the variant allele, so this locus represents a homozygous alternate locus.



DeepVariant correctly classifies the variant above as a "1", which means that one of the two alleles matches the variant allele, i.e. it is heterozygous.



# Somatic Mutation Calling

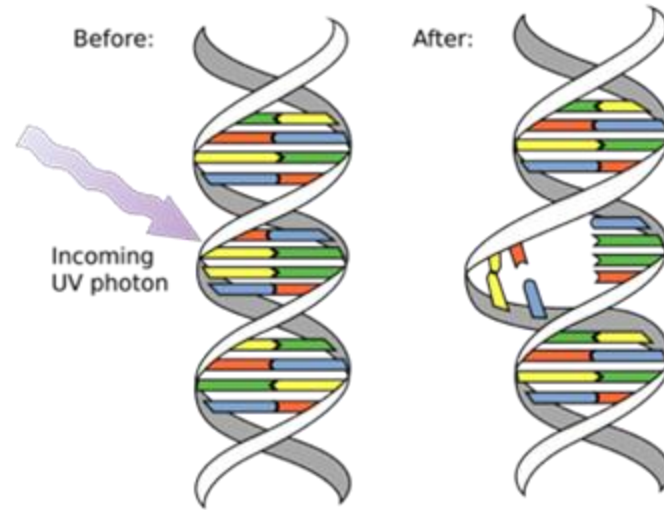
# Cancer is a disease of the genome

- Cancer is caused by **somatic** mutations

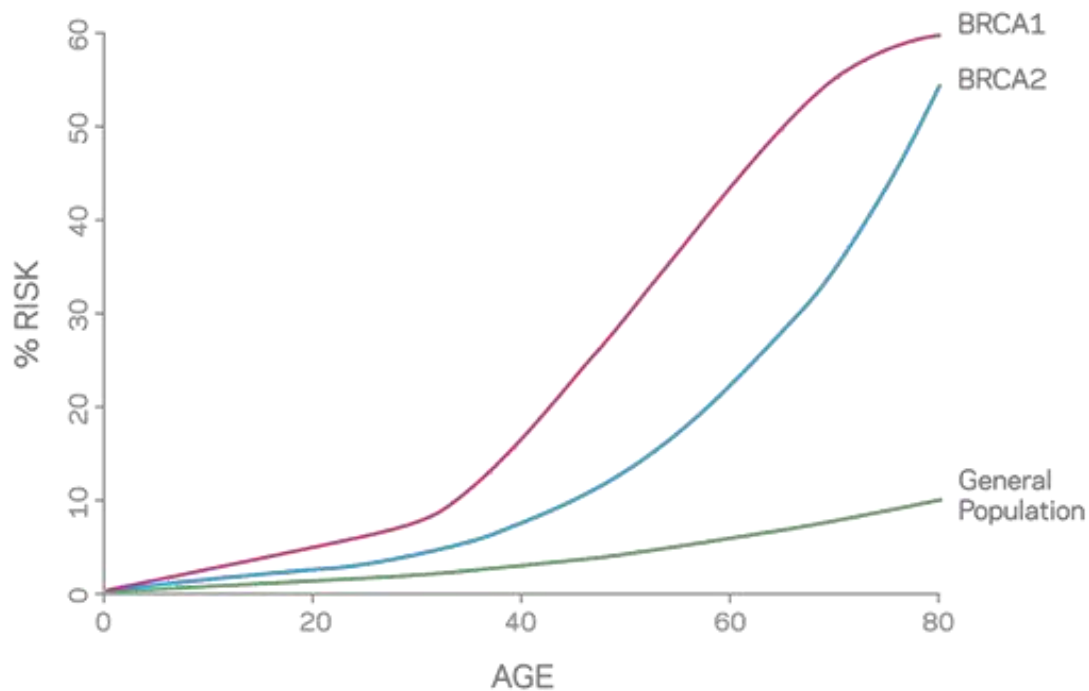


# Cancer is a disease of the genome

- Cancer is caused by **somatic** mutations
- These mutations are introduced into the genome of a cell (errors in DNA copying, UV light, chemicals)
- Most cancers require around 3 driver mutations



# Germline Predisposition

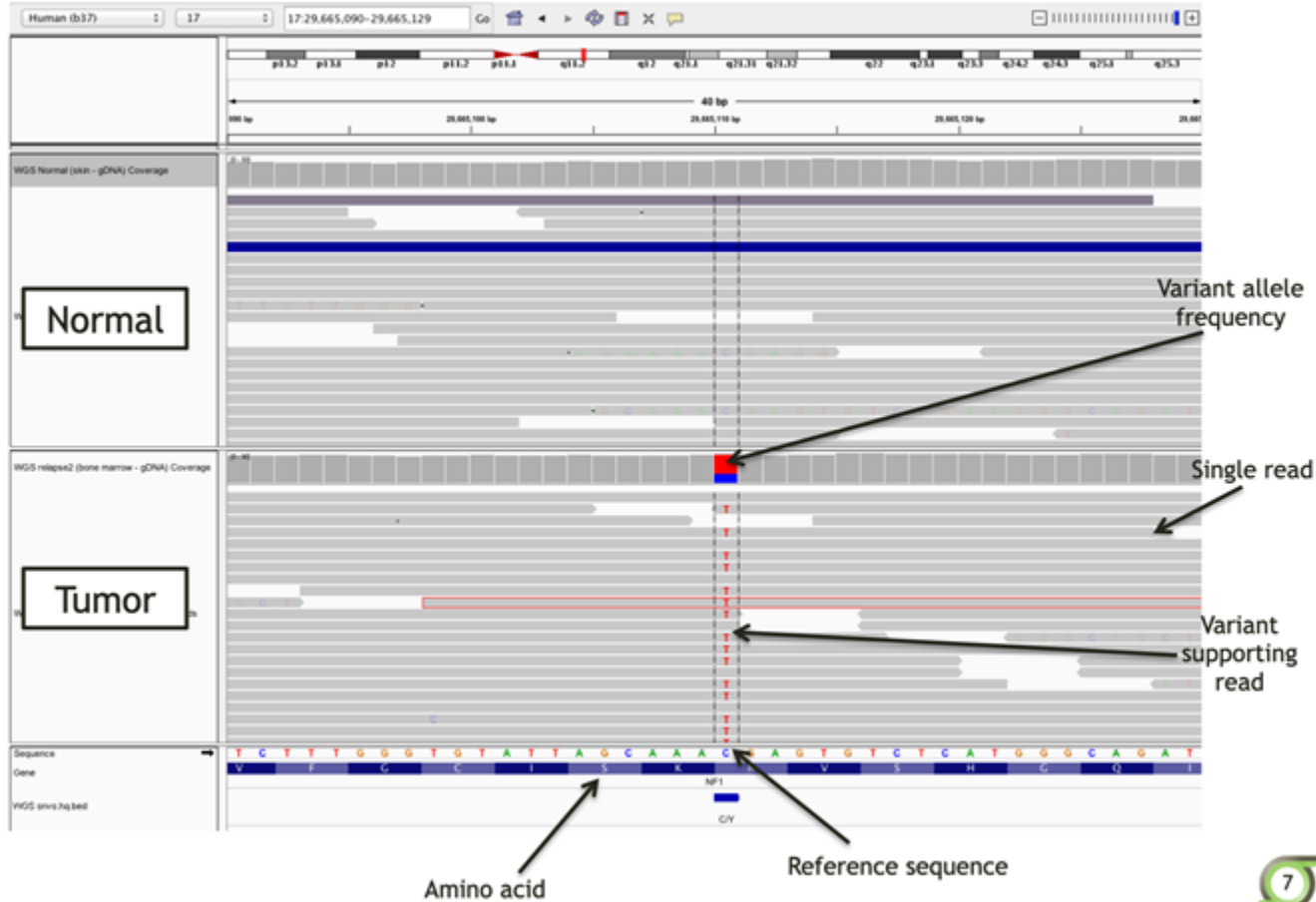


# Cancer Sequencing

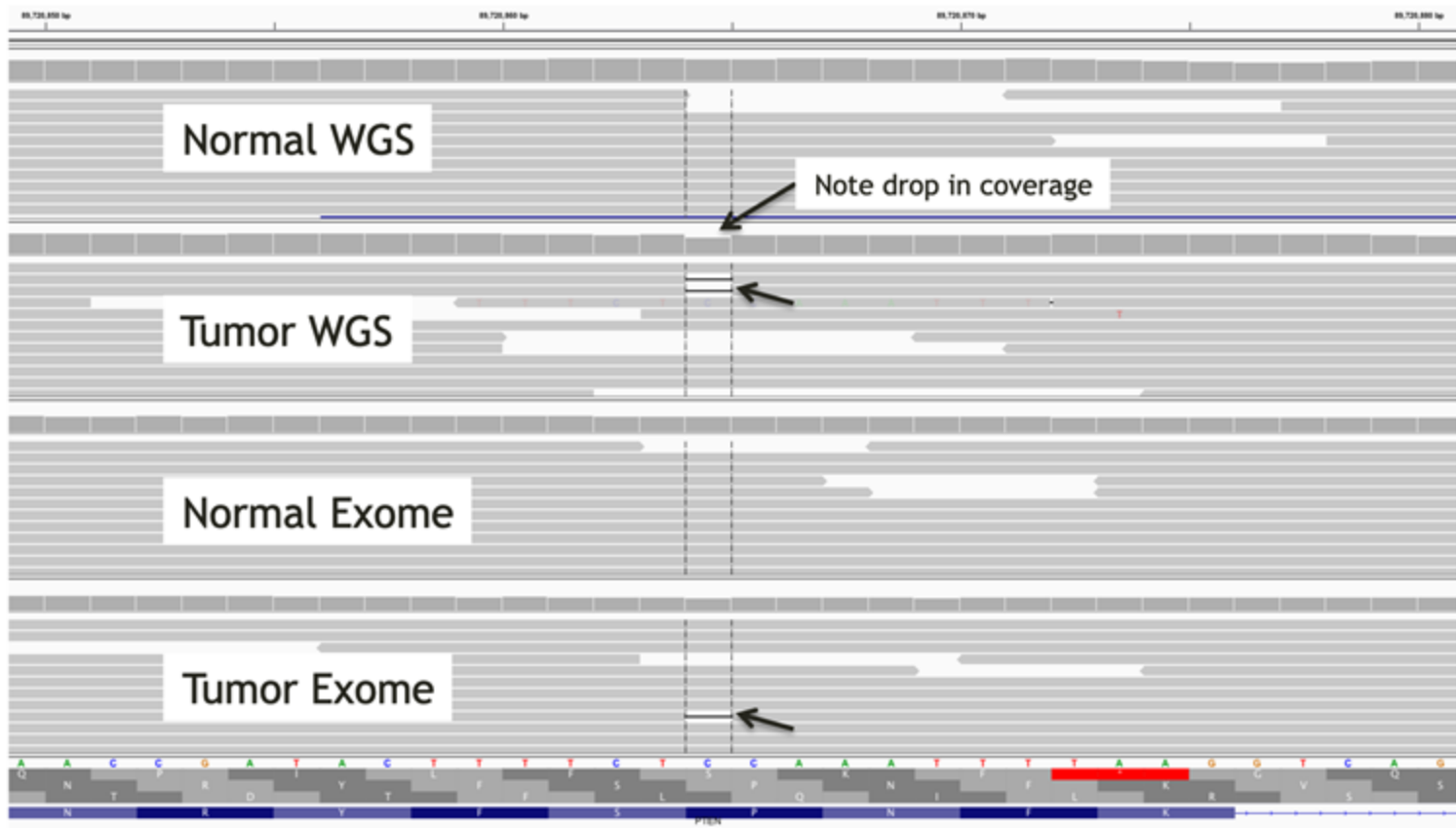
- In cancer, we have to (at least) double sequencing costs
- Uses both a tumor sample and a matched normal
- We compare them to find somatic mutations



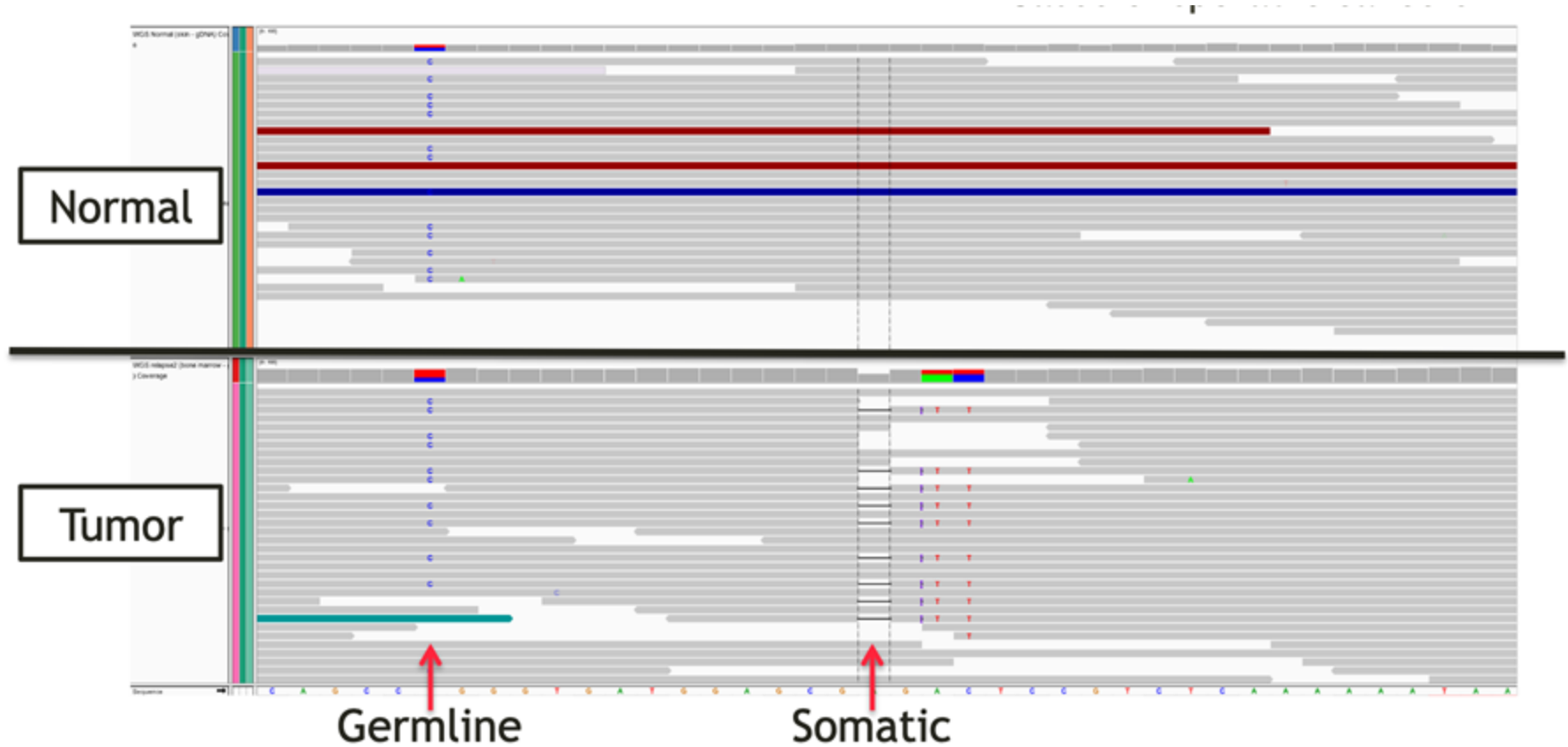
# What do somatic variants look like?



# Indels



# Germline vs Somatic



VAF = Variant reads / Total reads



Normal

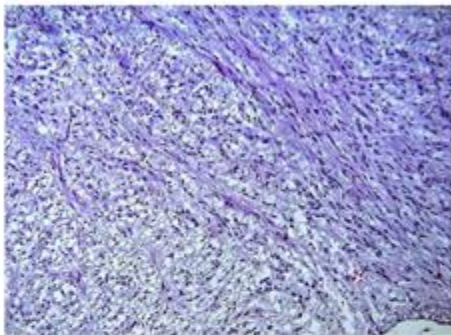
VAF = 0/20 = 0%

Tumor

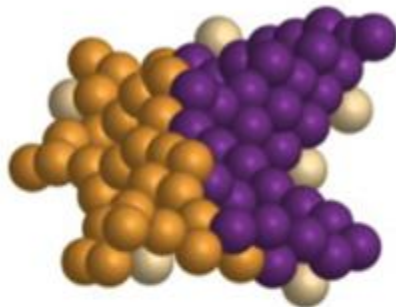
VAF = 14/20 = 70%

Check on Mutect

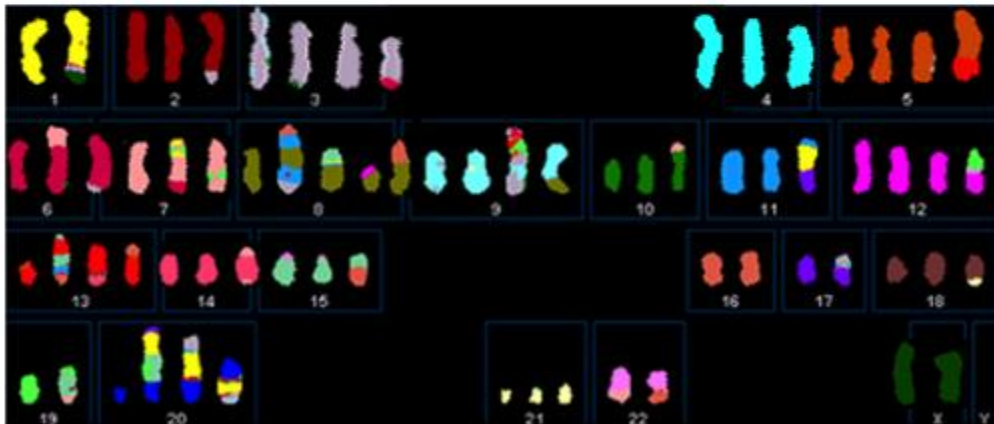
# Tumors are often impure, heterogeneous, and aneuploid



Tumors are often impure  
(contain normal cells)



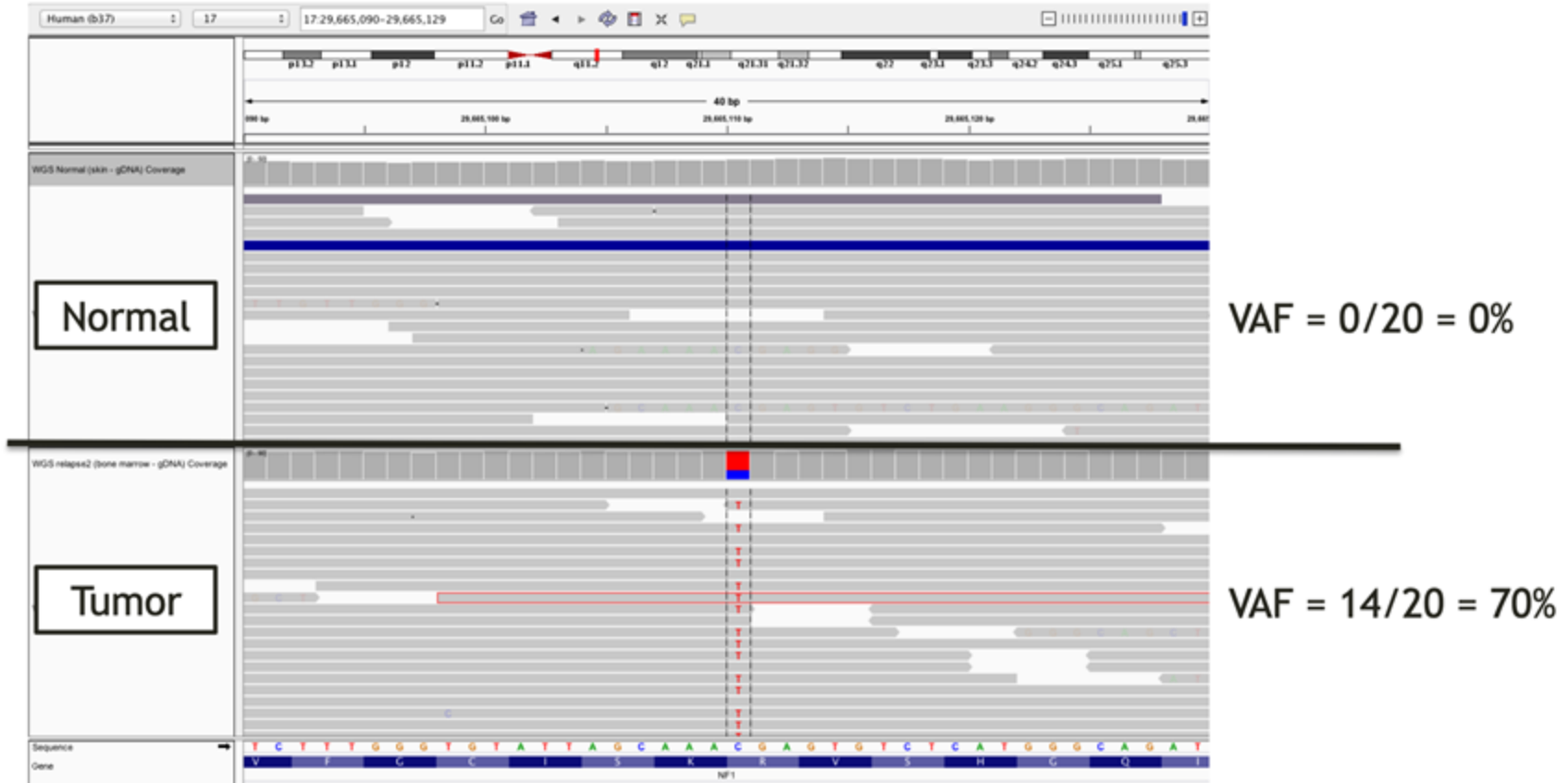
Tumors are often genetically  
diverse collections of cells



Tumors may be aneuploid

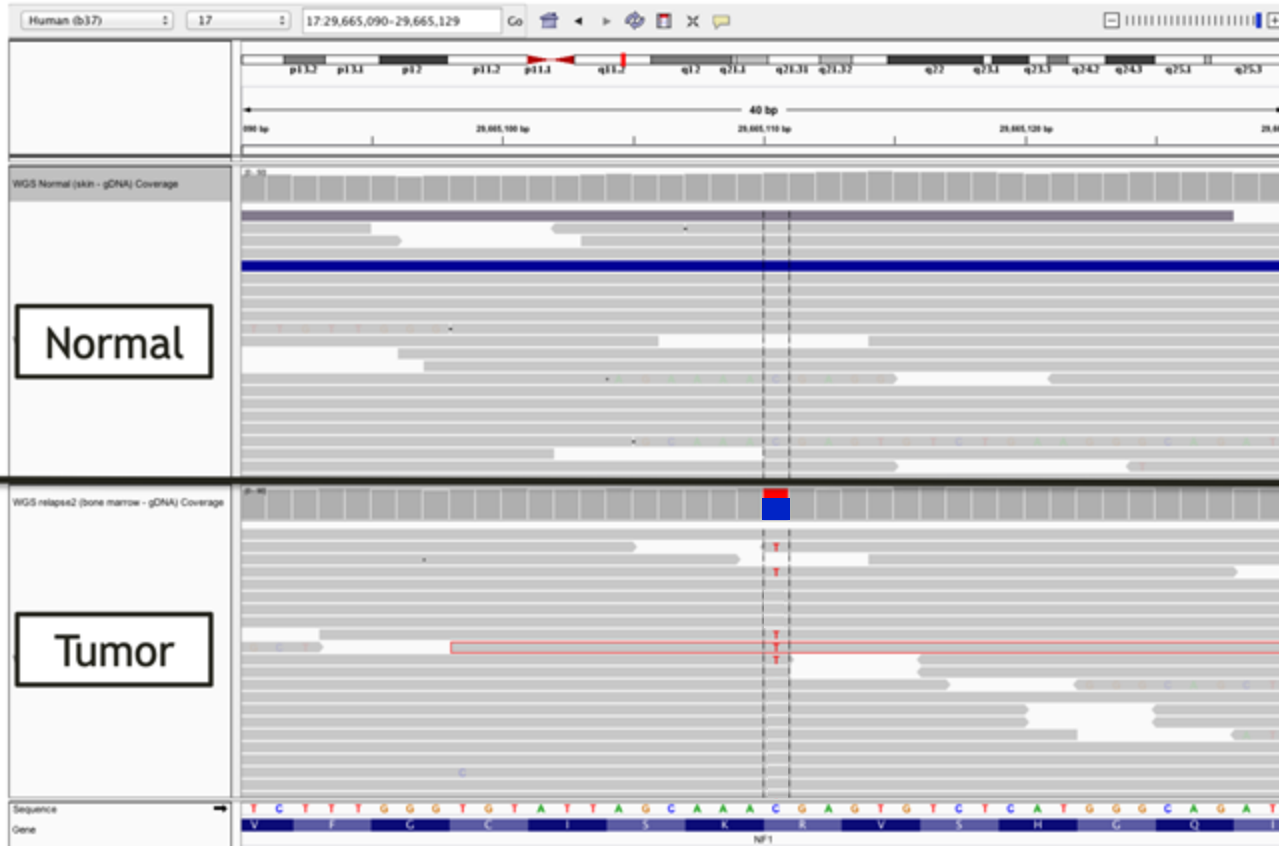
# How does purity influence VAF?

VAF = Variant reads / Total reads



# How does purity influence VAF?

VAF = Variant reads / Total reads



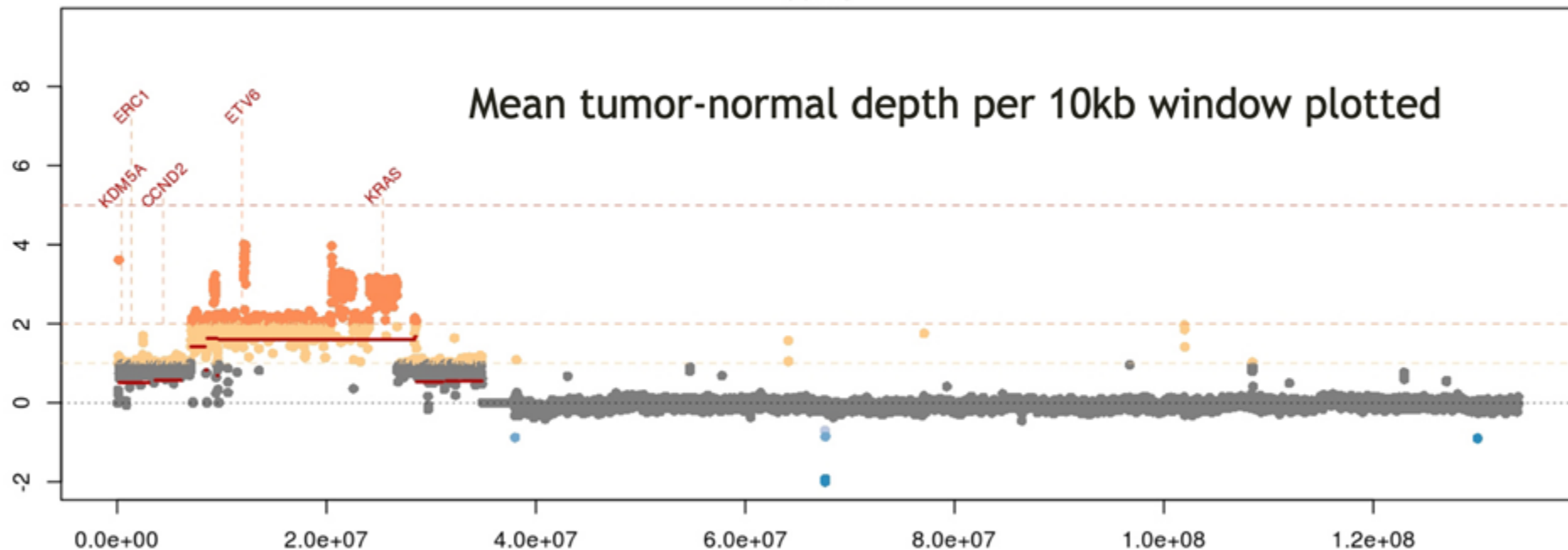
$$\text{VAF} = 0/20 = 0\%$$

$$\text{VAF} = 5/20 = 25\%$$



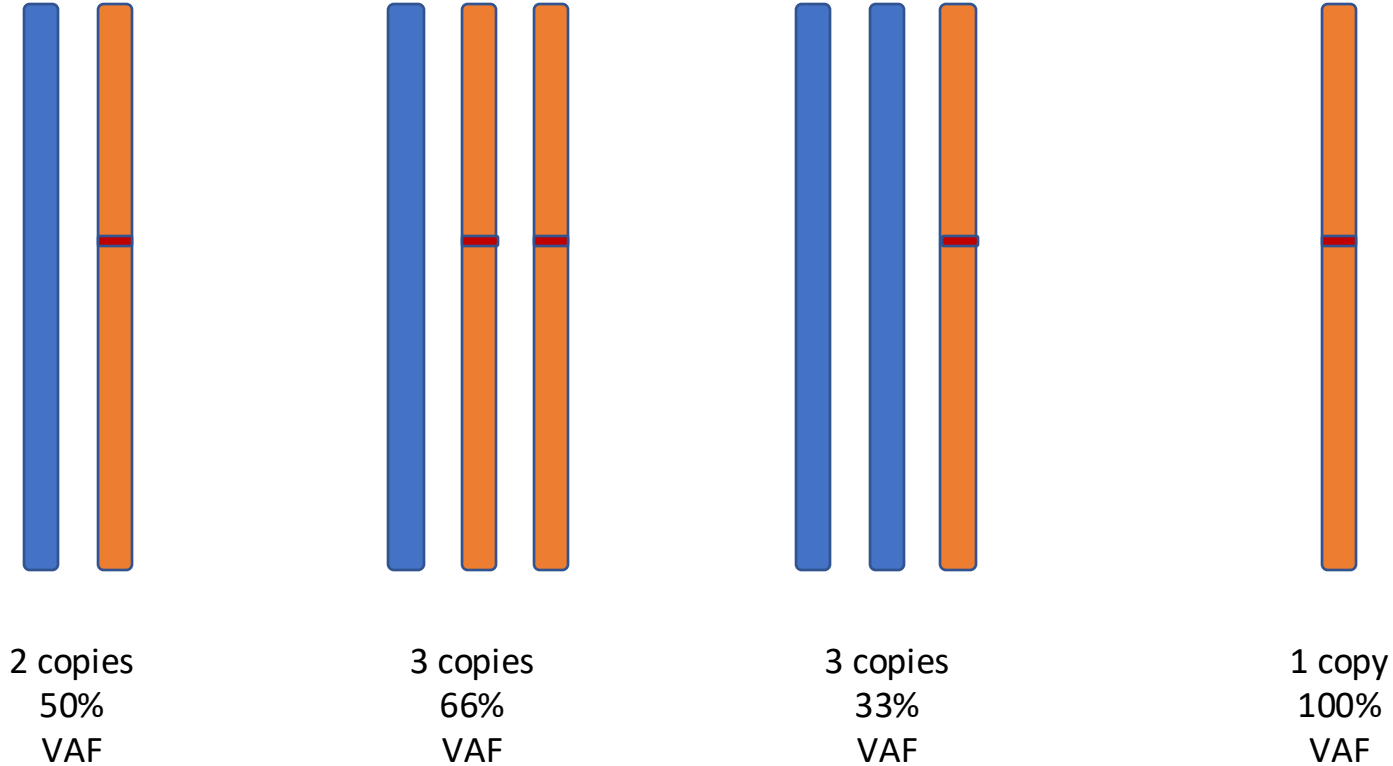


### Gains

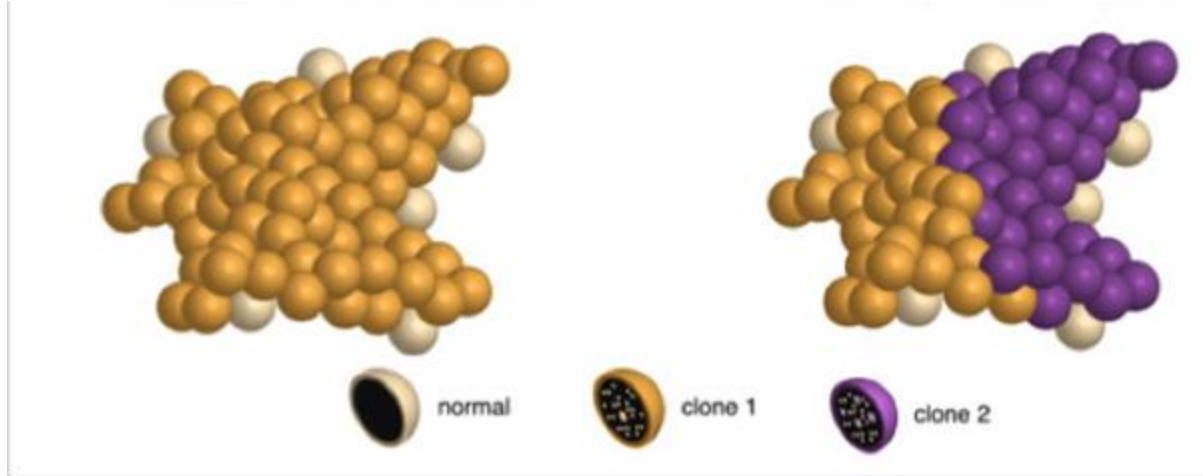


KRAS amplification in a metastatic breast cancer

# How does copy number influence VAF?

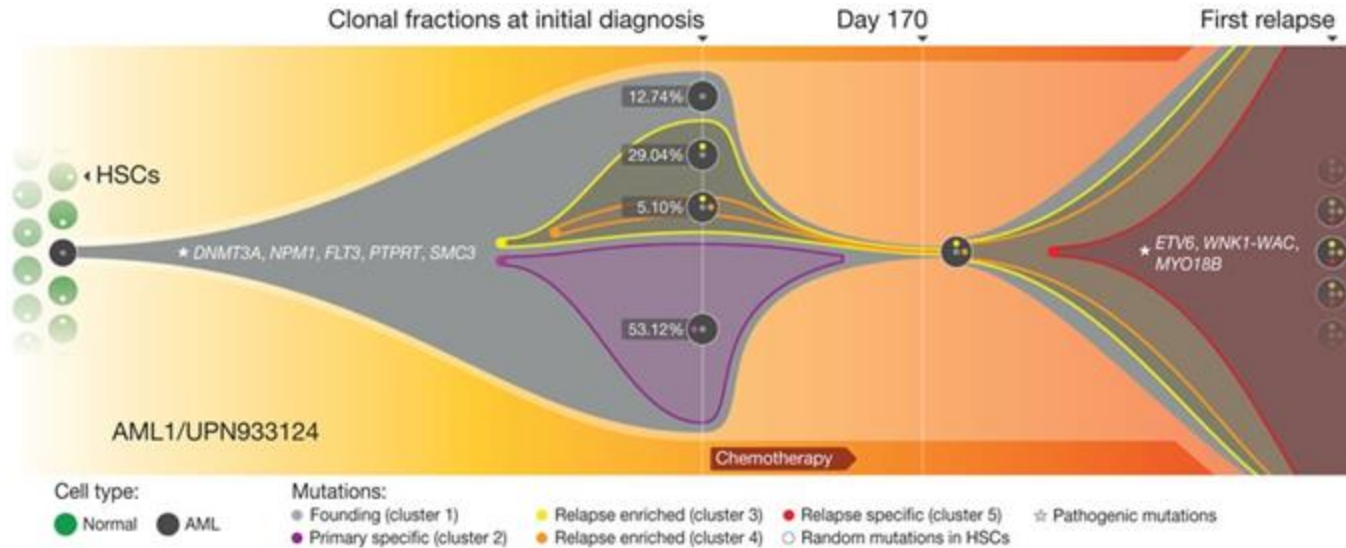


# How does clonality influence VAF?

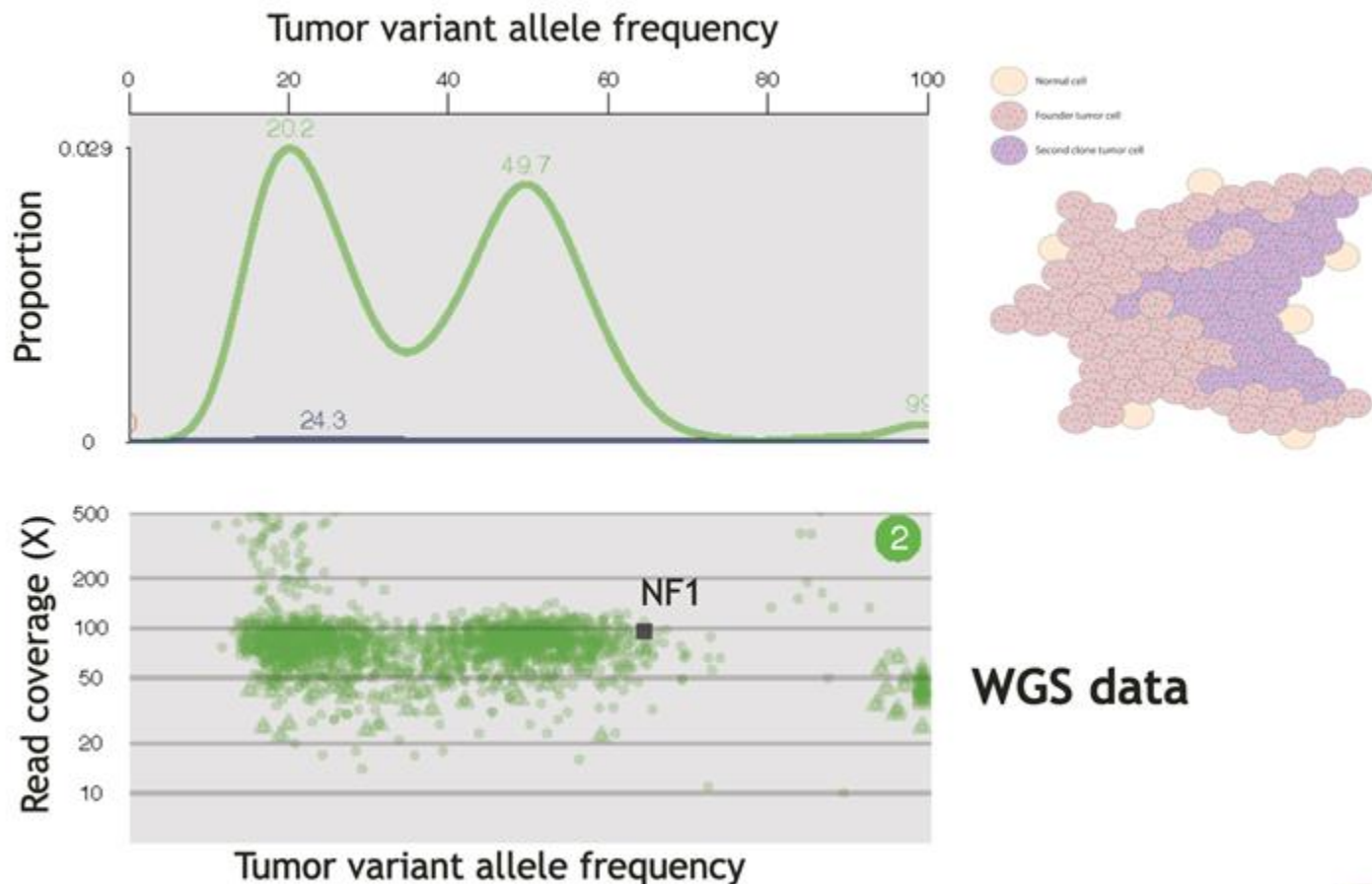


- Subclones contain genetically diverse populations of cells
- Evolution occurs at the molecular and cellular levels
- The growth rates for subclones are often different

# Clonal evolution in relapsed AML



# Dominant clone vs. sub-clonal (and driver vs. passenger)



# Somatic variant calling is harder

- There are more factors to consider, a wider range of possibilities, and often, more sketchy samples

# Somatic Variant detection callers

- Mutect
- Strelka
- Varscan
- Pindel
- Lancet
- Deep Somatic
- VarDict
- Seurat
- Shimmer
- more...

Lots of choices!

# Use of multiple variant callers can improve sensitivity and accuracy

Performance of caller Intersections

