

SNP and INDEL discovery

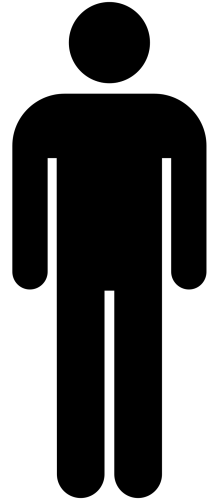
CSHL Advanced Sequencing Technologies 2022

11/17/2022

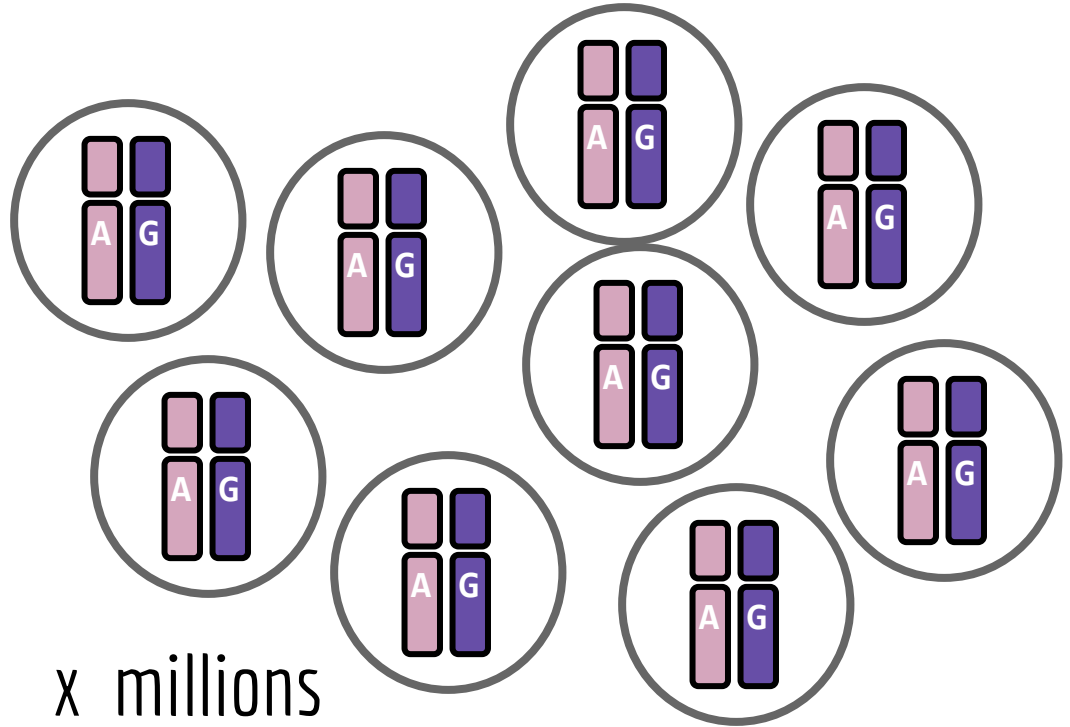
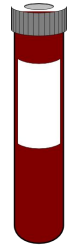
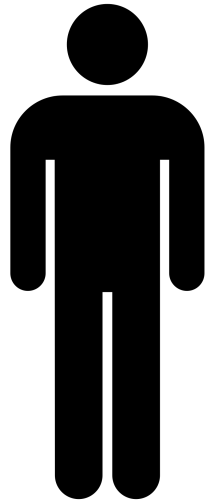
Joshua Mincer, Jason Kunisaki



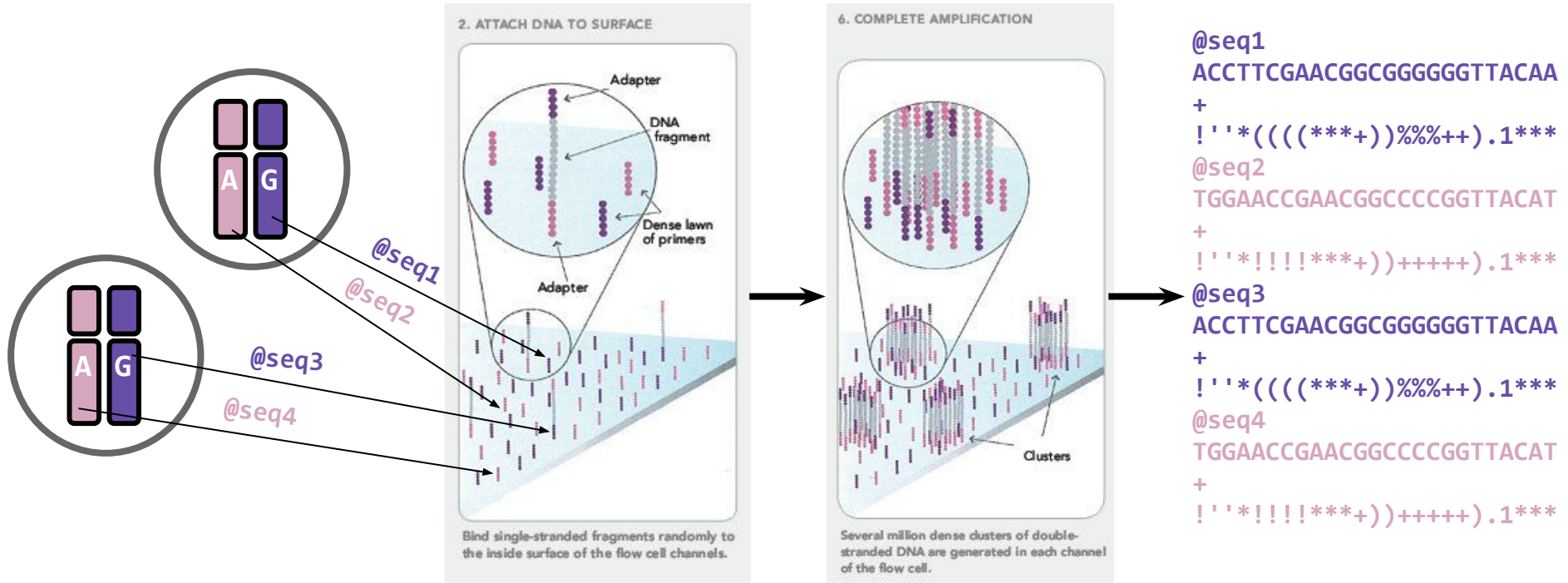
Goal: find all inherited variants in an individual's diploid genome.



Find inherited genetic variation by sequencing DNA
from millions of cells

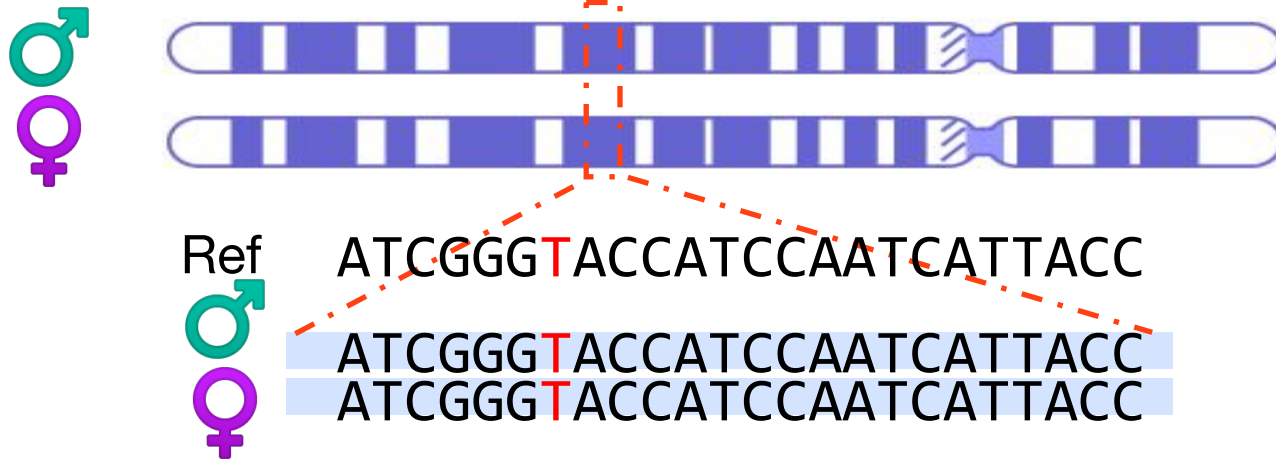


Each DNA cluster is amplified from a single strand from a single haploid chromosome from a single cell.

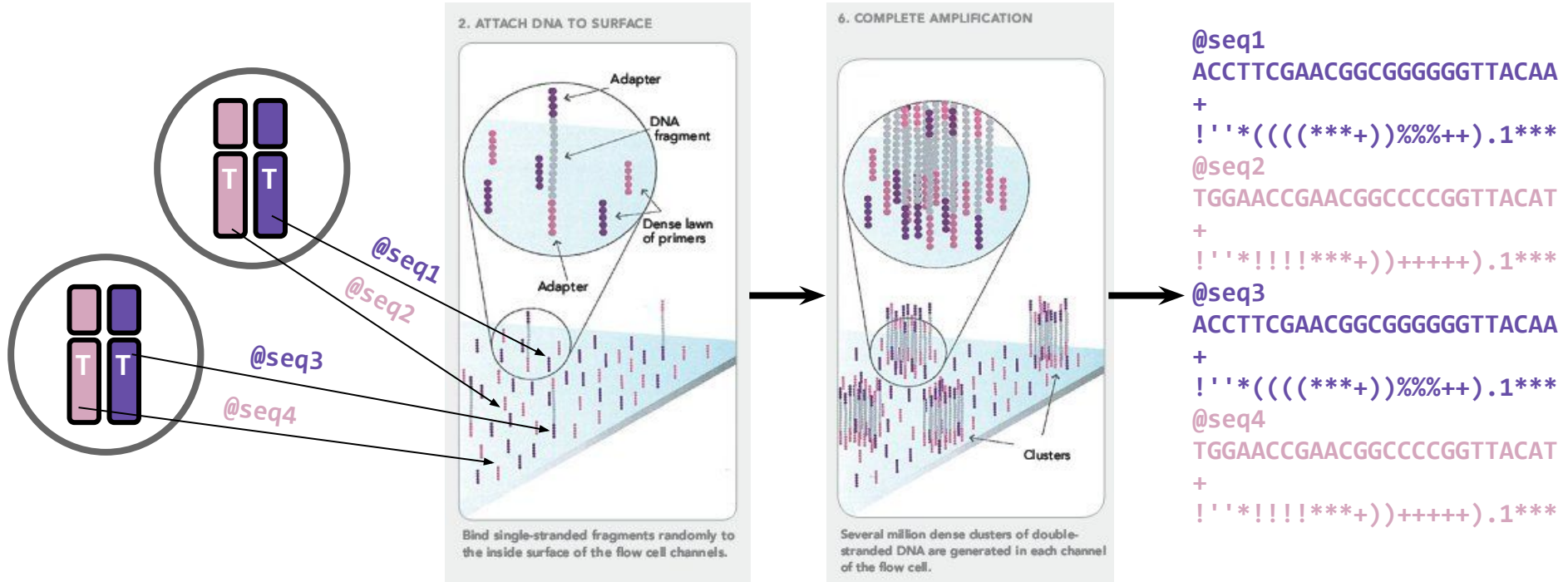


```
@seq1
ACCTTCGAACGGCGGGGGTTACAA
+
!' '*((( (**+))%%++) .1***
@seq2
TGGAACCGAACGGCCCCGGTTACAT
+
!' '*!!!! (**+))++++).1***
@seq3
ACCTTCGAACGGCGGGGGTTACAA
+
!' '*((( (**+))%%++) .1***
@seq4
TGGAACCGAACGGCCCCGGTTACAT
+
!' '*!!!! (**+))++++).1***
```

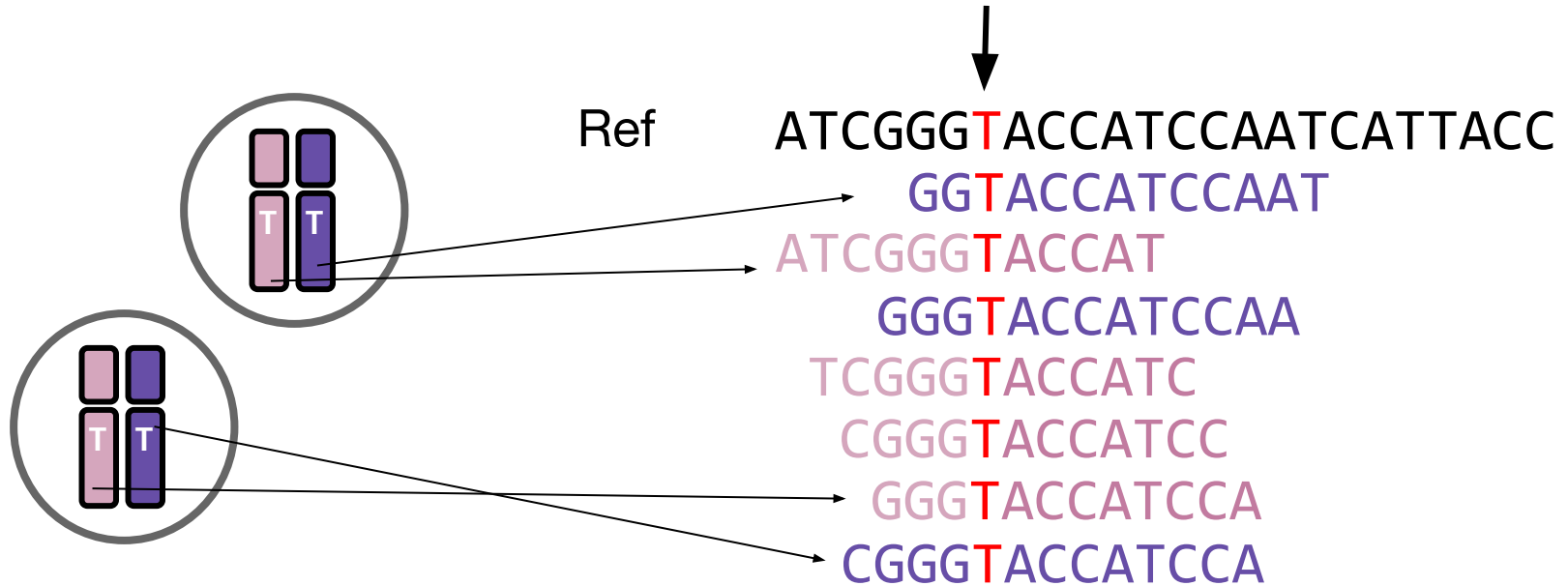
Scenario 1: An individual is homozygous for the "reference" allele.



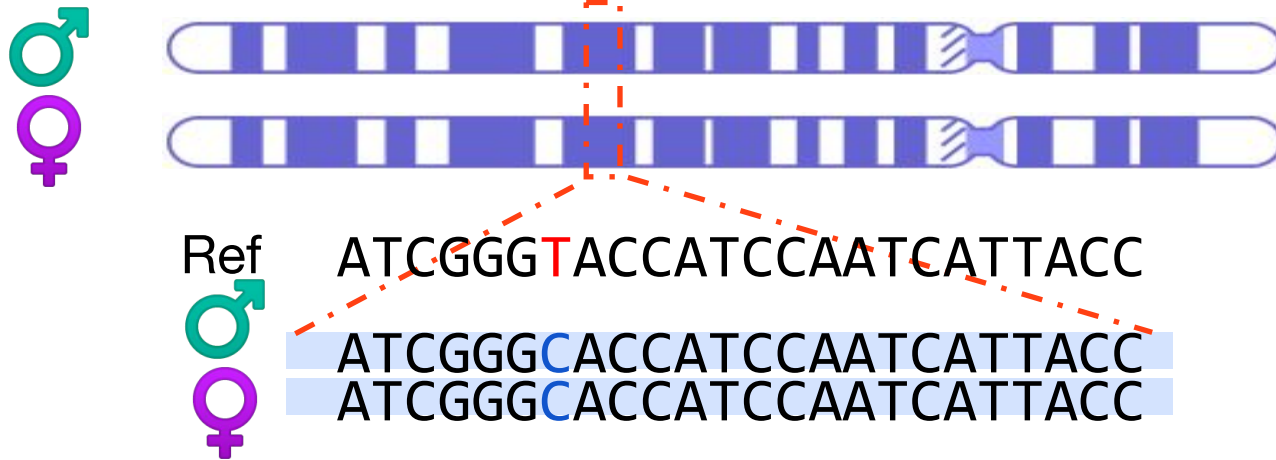
Scenario 1: An individual is homozygous for the "reference" allele.



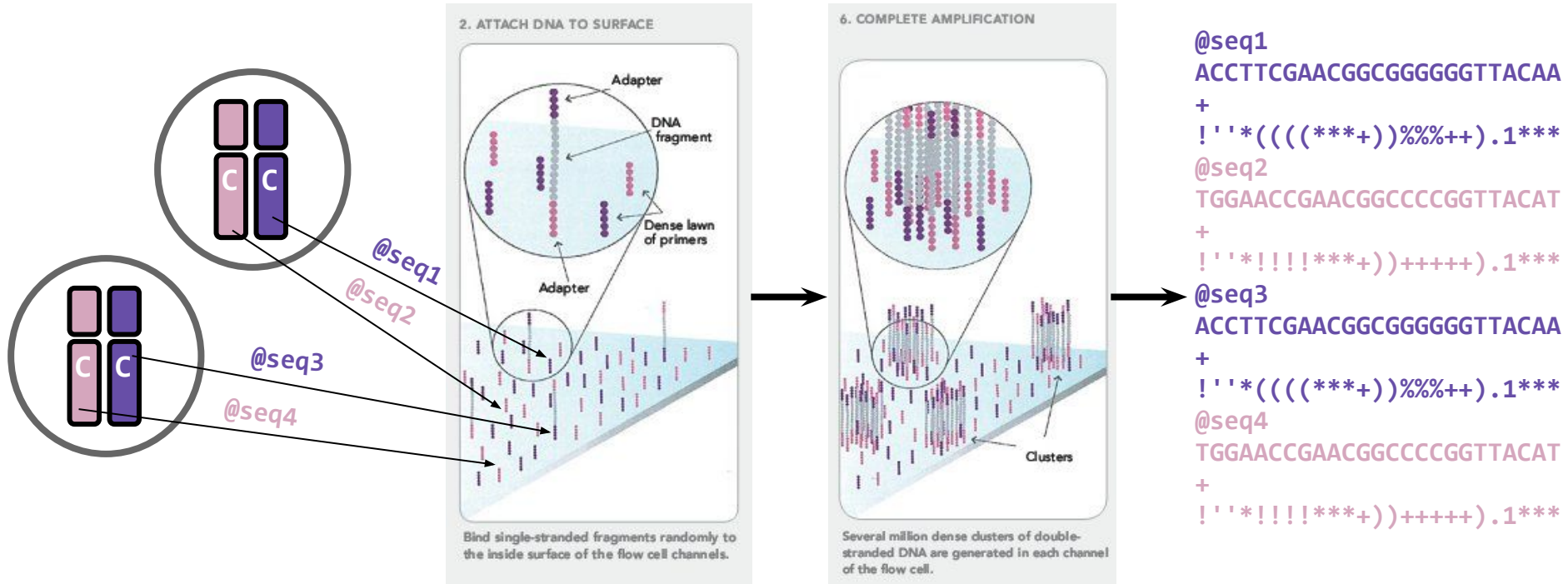
Scenario 1: An individual is homozygous for the "reference" allele.



Scenario 2: An individual is homozygous for an "alternate" allele.

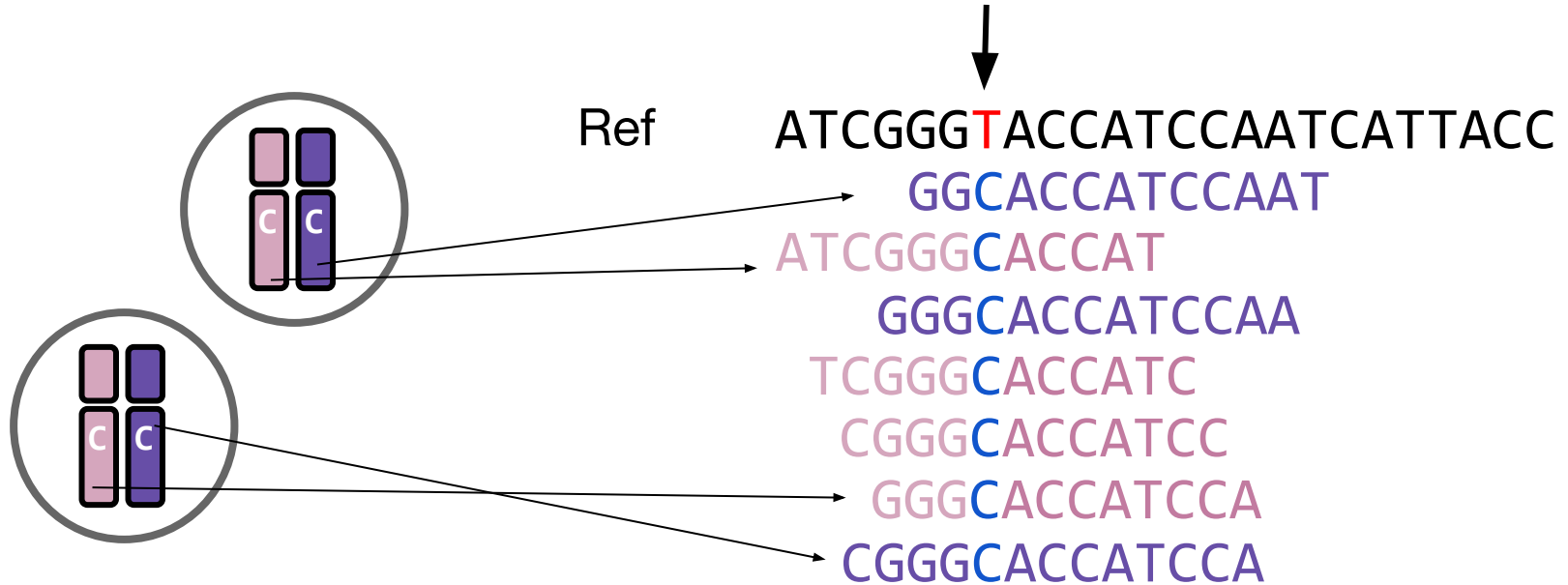


Scenario 2: An individual is homozygous for an "alternate" allele.

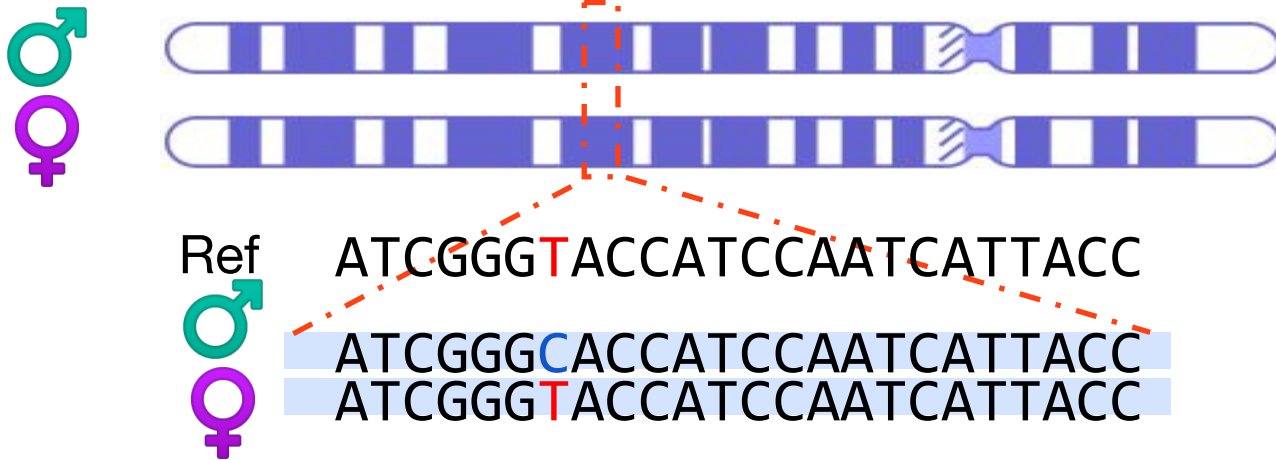


```
@seq1  
ACCTTCGAACGGCGGGGGTTACAA  
+  
!' '*((( (**+))%%++) .1***  
@seq2  
TGGAACCGAACGGCCCCGGTTACAT  
+  
!' '*!!!! (**+))++++).1***  
@seq3  
ACCTTCGAACGGCGGGGGTTACAA  
+  
!' '*((( (**+))%%++) .1***  
@seq4  
TGGAACCGAACGGCCCCGGTTACAT  
+  
!' '*!!!! (**+))++++).1***
```

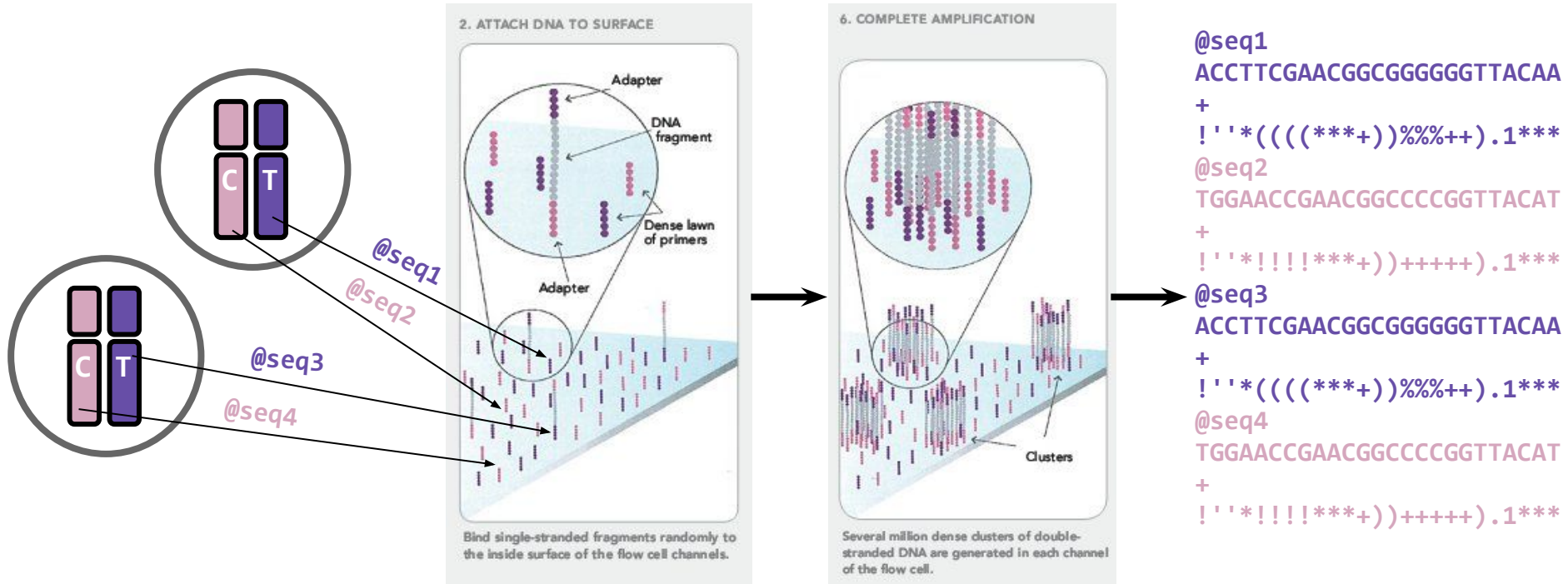
Scenario 2: An individual is homozygous for an "alternate" allele.



Scenario 3: An individual is heterozygous for an "alternate" allele.

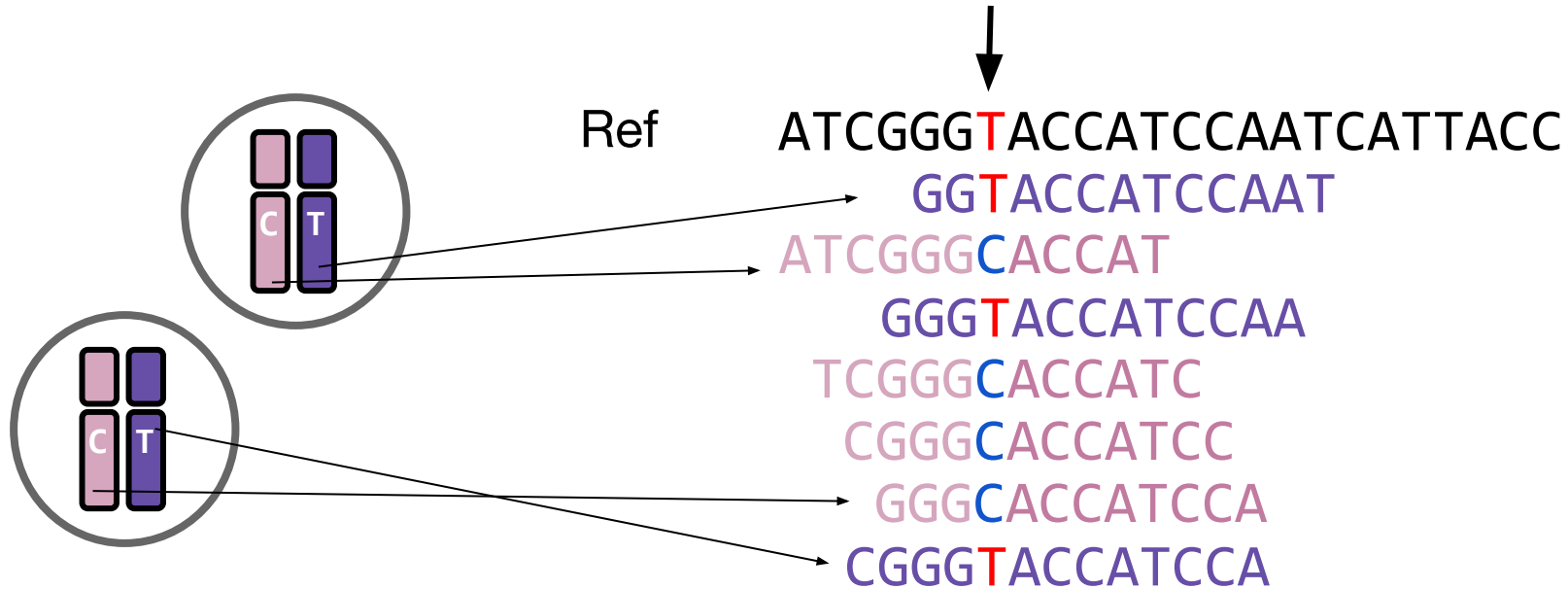


Scenario 3: An individual is heterozygous for an "alternate" allele.



```
@seq1  
ACCTTCGAACGGCGGGGGTTACAA  
+  
!' '*((( (**+))%%++) .1***  
@seq2  
TGGAACCGAACGGCCCCGGTTACAT  
+  
!' '*!!!! (**+))++++).1***  
@seq3  
ACCTTCGAACGGCGGGGGTTACAA  
+  
!' '*((( (**+))%%++) .1***  
@seq4  
TGGAACCGAACGGCCCCGGTTACAT  
+  
!' '*!!!! (**+))++++).1***
```

Scenario 3: An individual is heterozygous for an "alternate" allele.



Why might finding heterozygous variants be harder?

The binomial distribution: adventures in coin flipping



$$P(\text{heads}) = 0.5$$



$$P(\text{tails}) = 0.5$$

Thinking about allele sampling with the binomial distribution

The **binomial distribution** with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent yes (e.g., "heads" or "reference allele") or no (e.g., "tails", or "alternate allele") experiments, each of which yields success with probability p .

The probability of getting exactly k successes in n trials is given by the probability mass function:

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

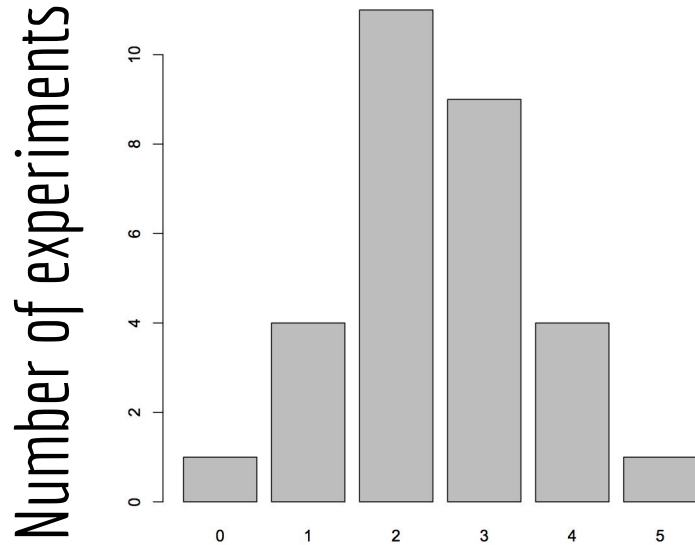
What is the probability of seeing $k=1$ tails in $n=3$ flips of a fair coin with the probability of a tail (p) = 0.5?

3 choose $1 = 3$; $0.5^1 = 0.5$; $(1-0.5)^{(3-1)} = 0.25$. So... $3*0.5*0.25 = \mathbf{0.375}$

In R, the function would be: `dbinom(1, size=3, prob=0.5)`

What is the distribution of tails (alternate alleles) do we expect to see after 5 tosses (sequence reads)?

What is the distribution of tails (alternate alleles) do we expect to see after 5 tosses (sequence reads)?



Number of "tails"

R code:

```
barplot(table(rbinom(30, 5, 0.5)))
```

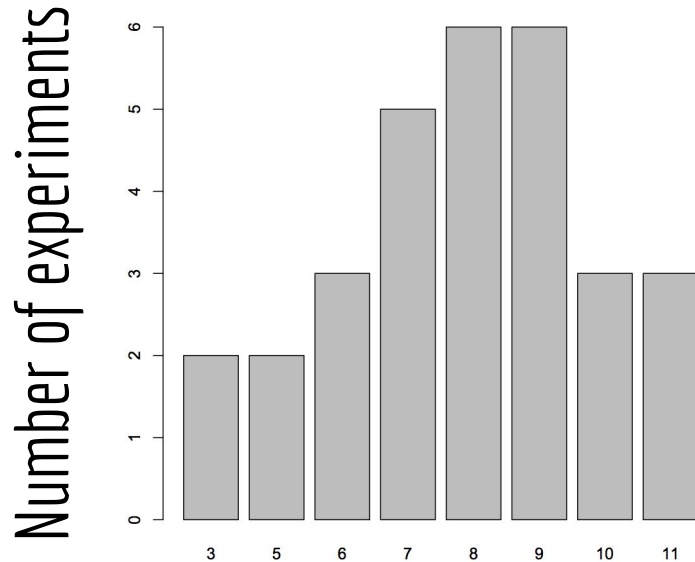
30 experiments (students tossing coins)

5 tosses each

Probability of Tails

What is the distribution of tails (alternate alleles) do we expect to see after 15 tosses (sequence reads)?

What is the distribution of tails (alternate alleles) do we expect to see after 15 tosses (sequence reads)?



Number of "tails"

R code:

```
barplot(table(rbinom(30, 15, 0.5)))
```

30 experiments (students tossing coins)

15 tosses each

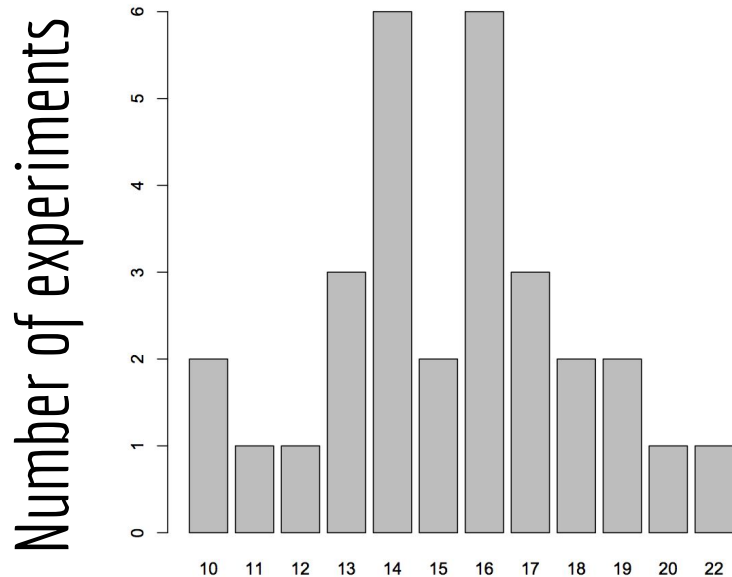
Probability of Tails

What is the distribution of tails (alternate alleles) do we expect to see after 30 tosses (sequence reads)?

Record your result in the following spreadsheet:

<https://docs.google.com/spreadsheets/d/1i8sA1KMeYc9UhwTnKg0tLFjCy8x5LIsBITcXrz5La94/edit?usp=sharing>

What is the distribution of tails (alternate alleles) do we expect to see after 30 tosses (sequence reads)?



Number of "tails"

R code:

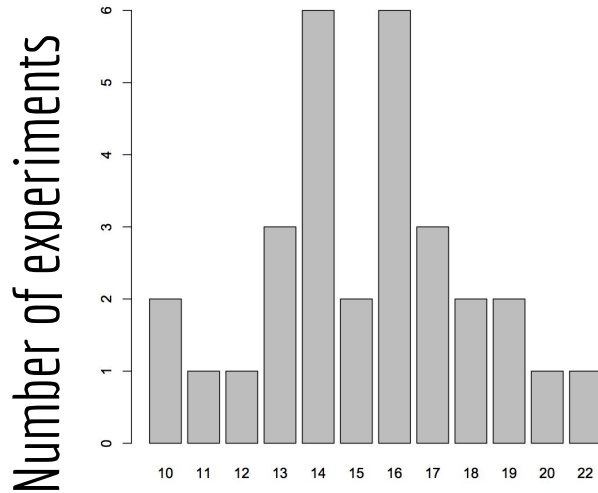
```
barplot(table(rbinom(30, 30, 0.5)))
```

30 experiments (students tossing coins)

30 tosses each

Probability of Tails

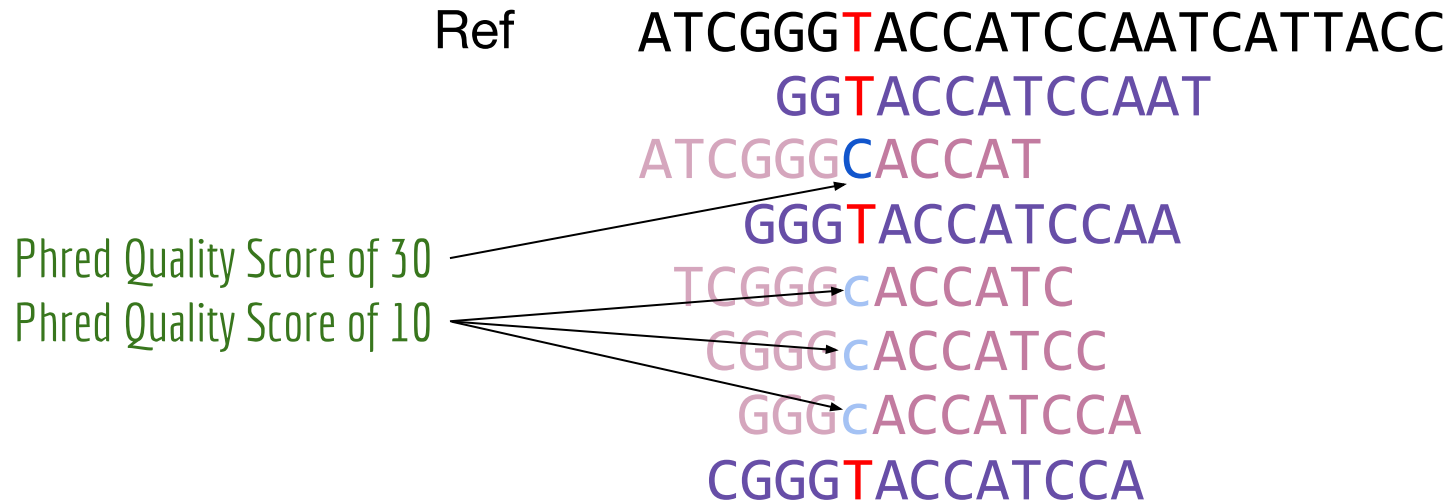
So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome



Number of "alternate alleles"

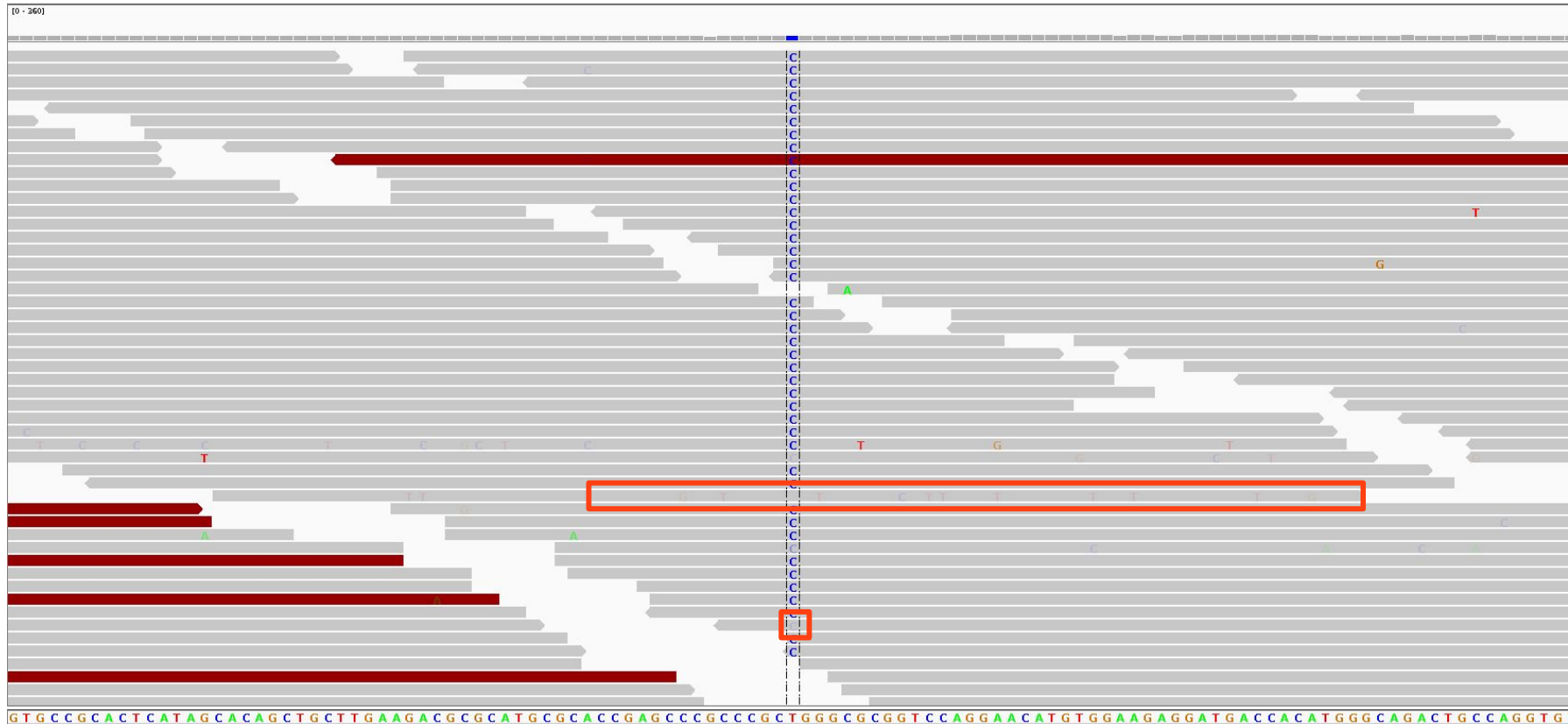
This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to find the majority of heterozygous alleles

Depth tackles the allele sampling issue and lower quality scores



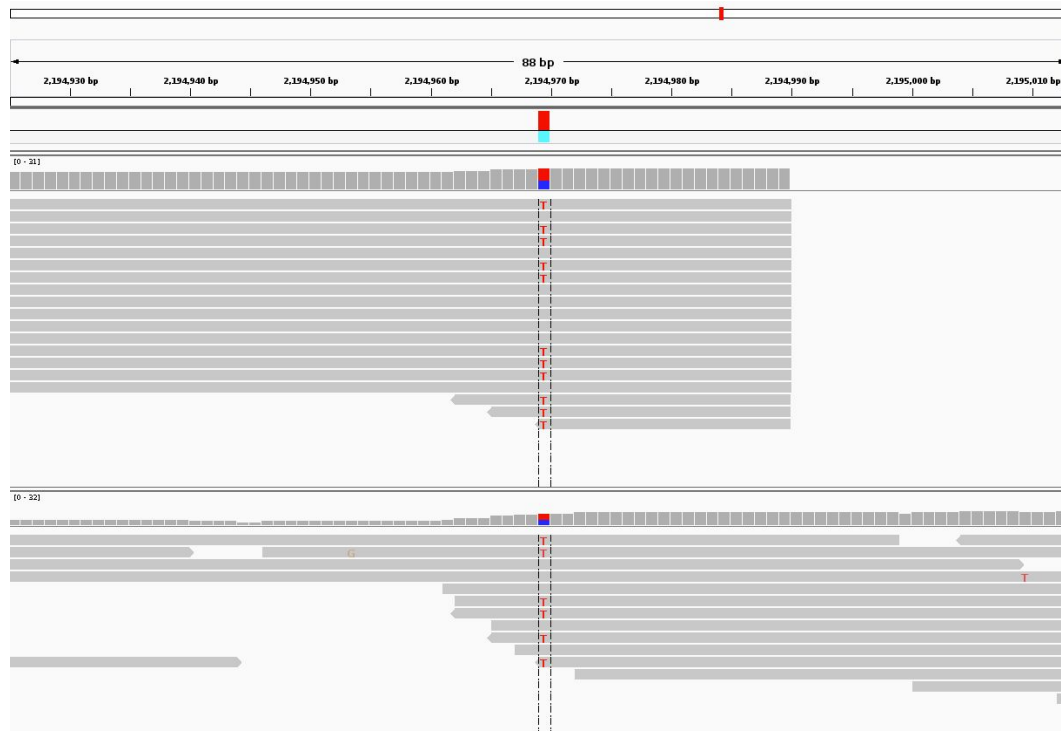
Some real examples of SNPs in IGV: validating variants via manual review

Homozygous for the "C" allele



Heterozygous for the alternate allele

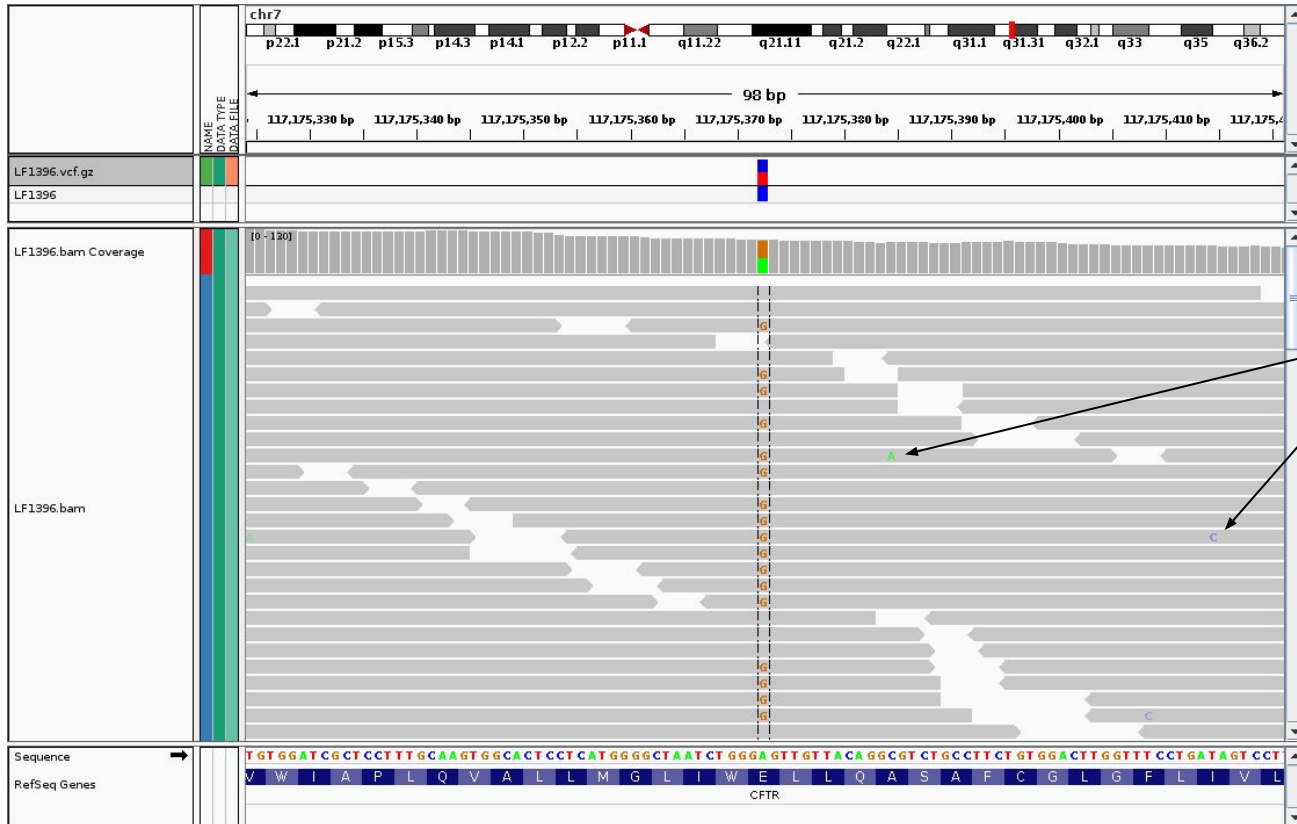
Individual 1



Individual 2

Which genotype prediction would you have more confidence in?

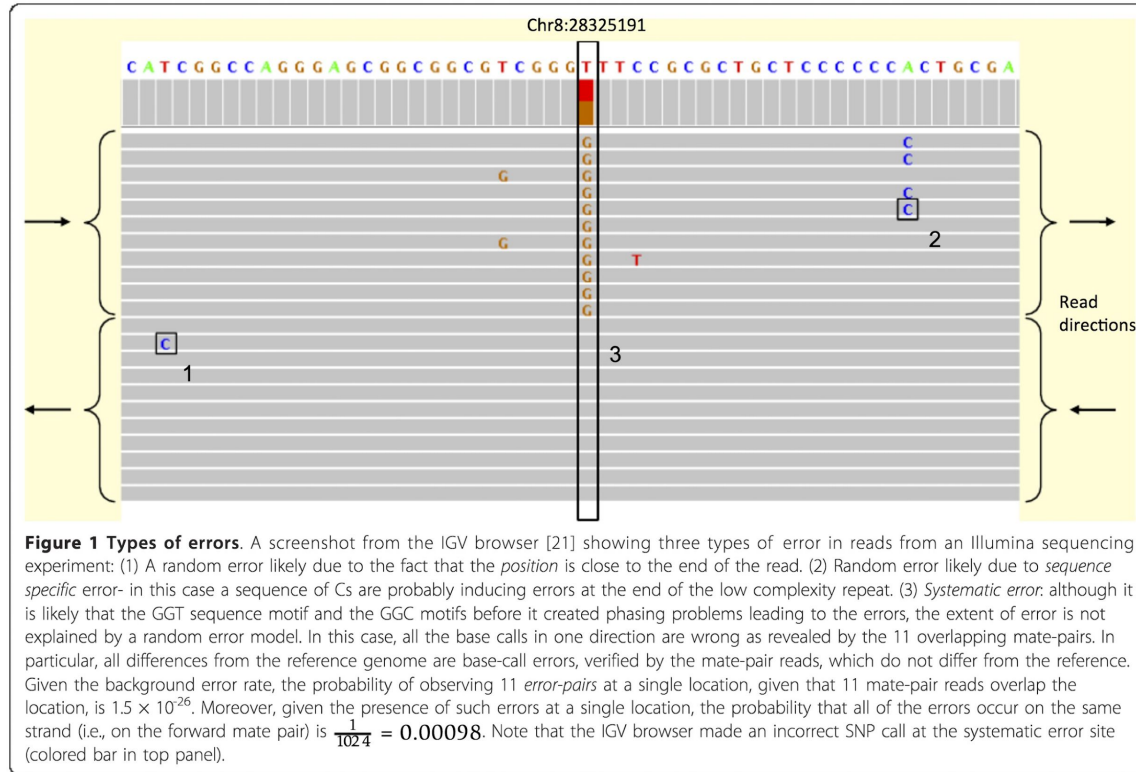
Sequencing errors fall out as noise (most of the time)



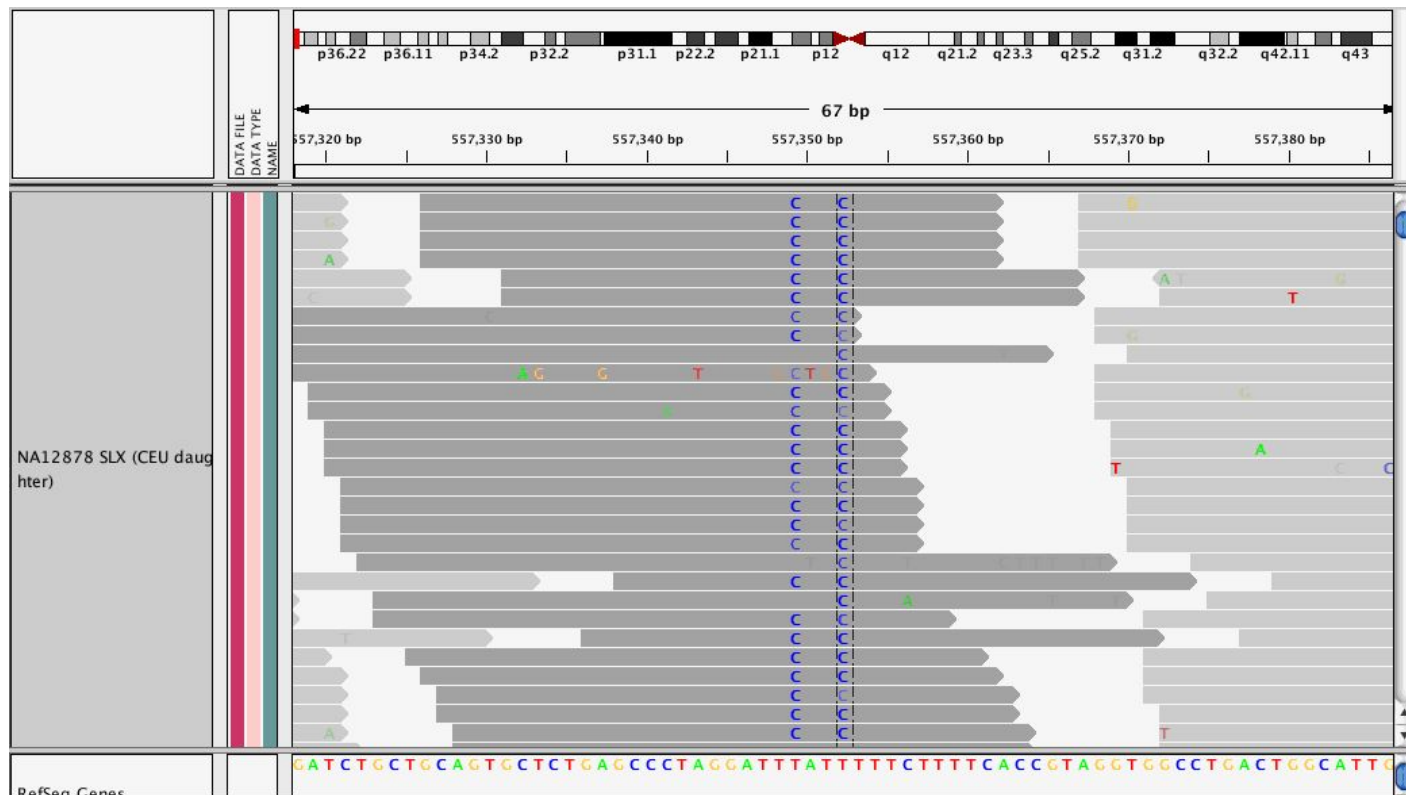
Sequencing errors

It is not always so easy

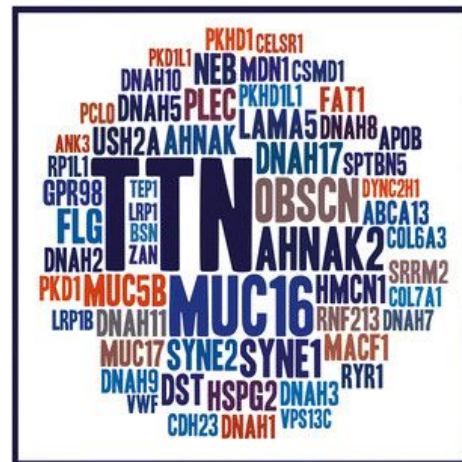
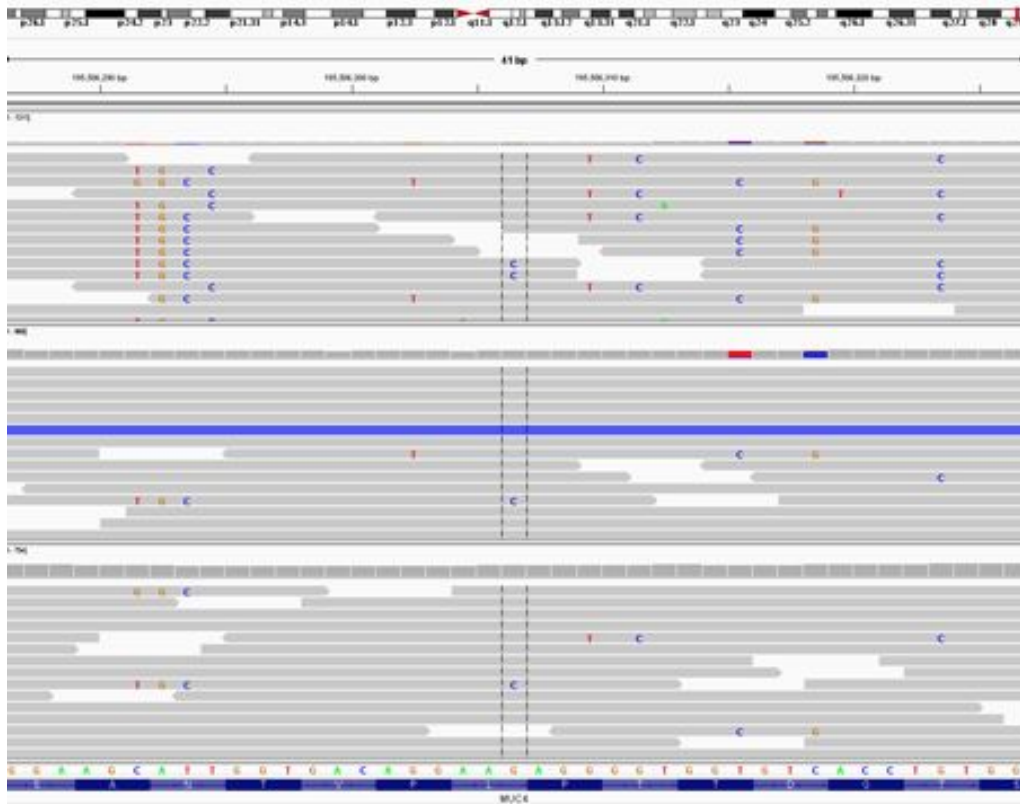
Random versus systematic error



Strand bias from PCR



Pileups of many differences from paralogy



RESEARCH ARTICLE | OPEN ACCESS

FLAGS, frequently mutated genes in public exomes

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

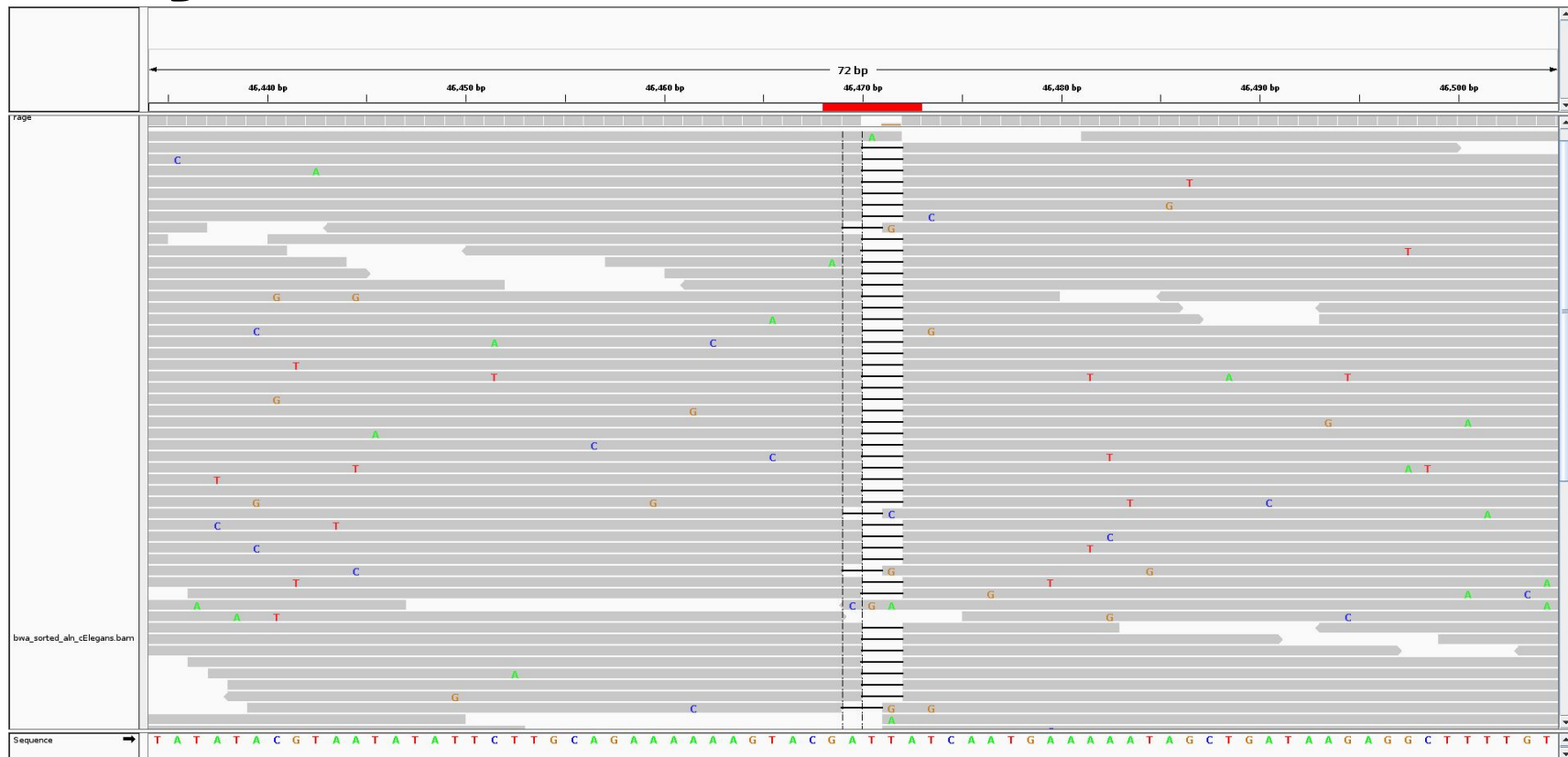
BMC Medical Genomics 2014 7:64 | DOI: 10.1186/s12920-014-0064-y | © Shyr et al.; licensee BioMed Central Ltd. 2014

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

[Open Peer Review reports](#)



Calling INDELs is much harder than SNPs



Some excellent resources to learn about manual review

Griffith Lab guides to manual review in IGV:

- <https://rnabio.org/module-02-alignment/0002/04/01/IGV/>
- [Standard operating procedure for somatic variant refinement of sequencing data with paired tumor and normal samples](#)