



Program of Medical Informatics
Faculty of Computers and Information
Mansoura University

Detecting Genetic Variations -I (Genomics)

Sara El-Metwally, Ph.D.
Faculty of Computers and Information,
Mansoura University, Egypt.

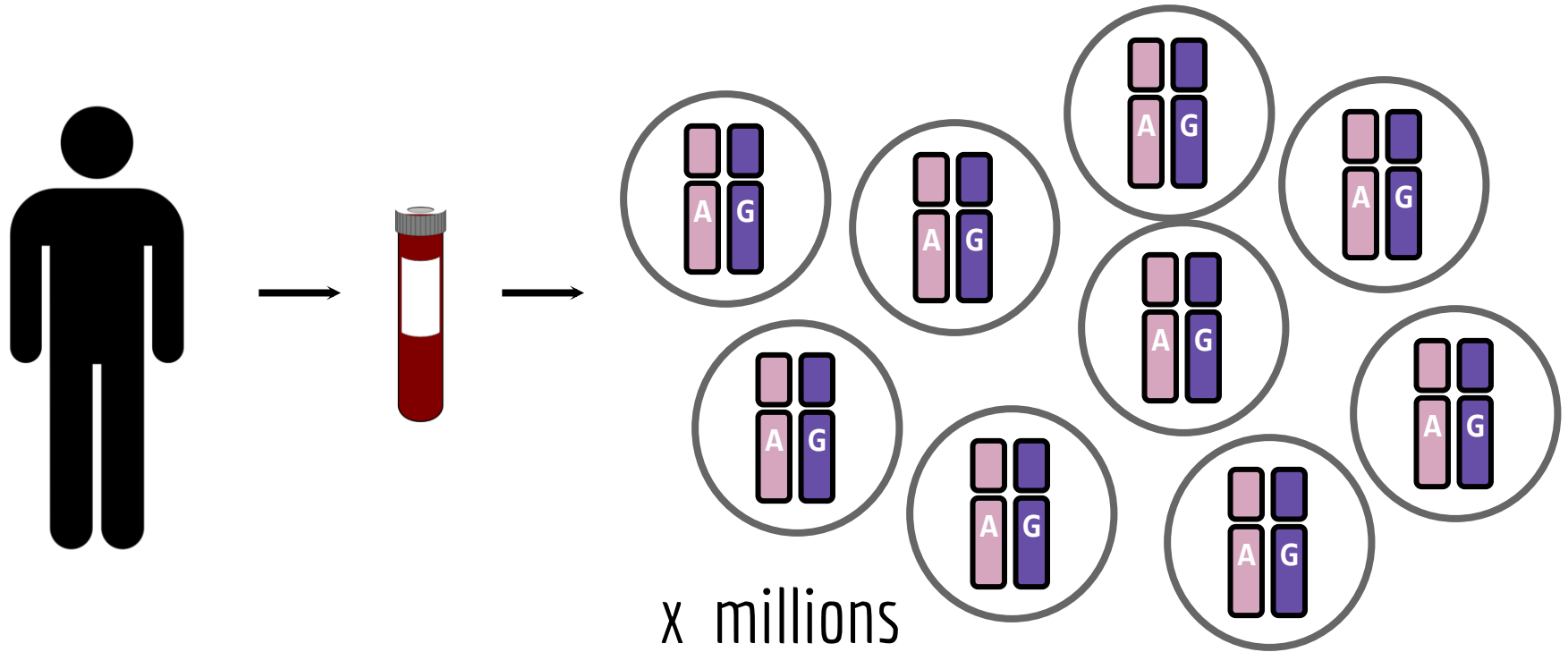
Email: sarah_almetwally4@mans.edu.eg
sara.elmetwally.2007@gmail.com

Office: Faculty of CIS, third floor

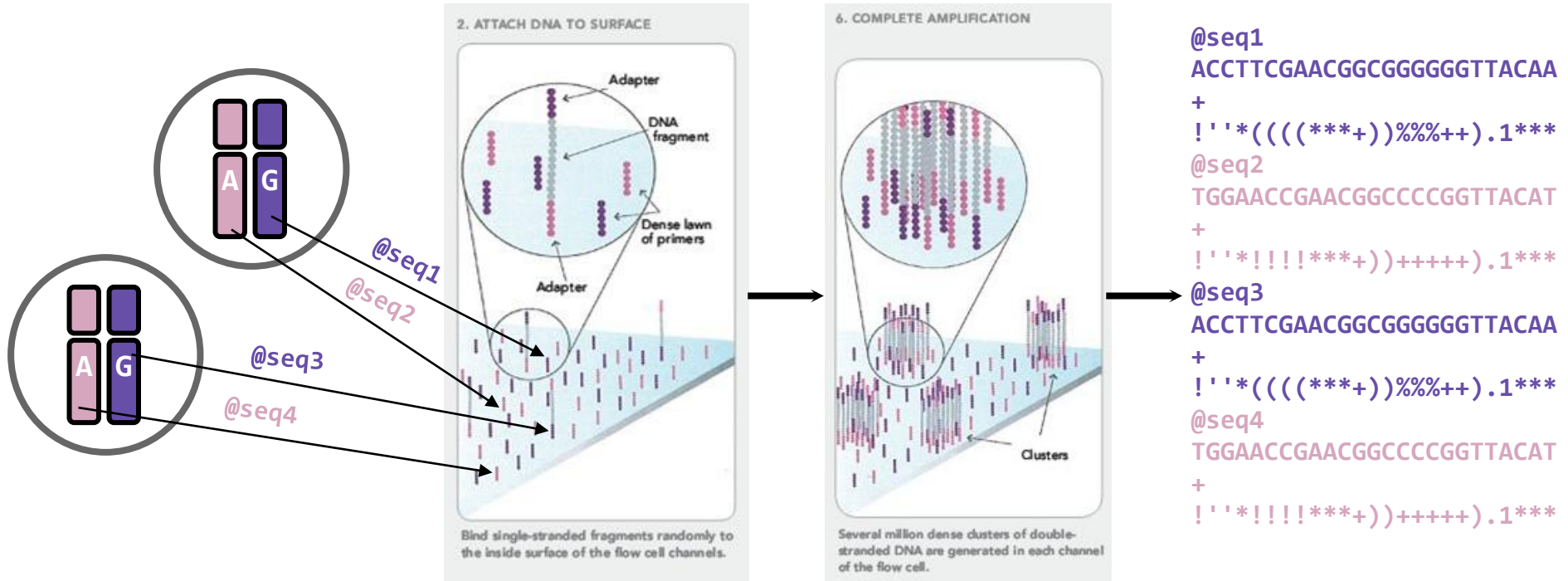
Goal: find all inherited variants in an individual's diploid genome.



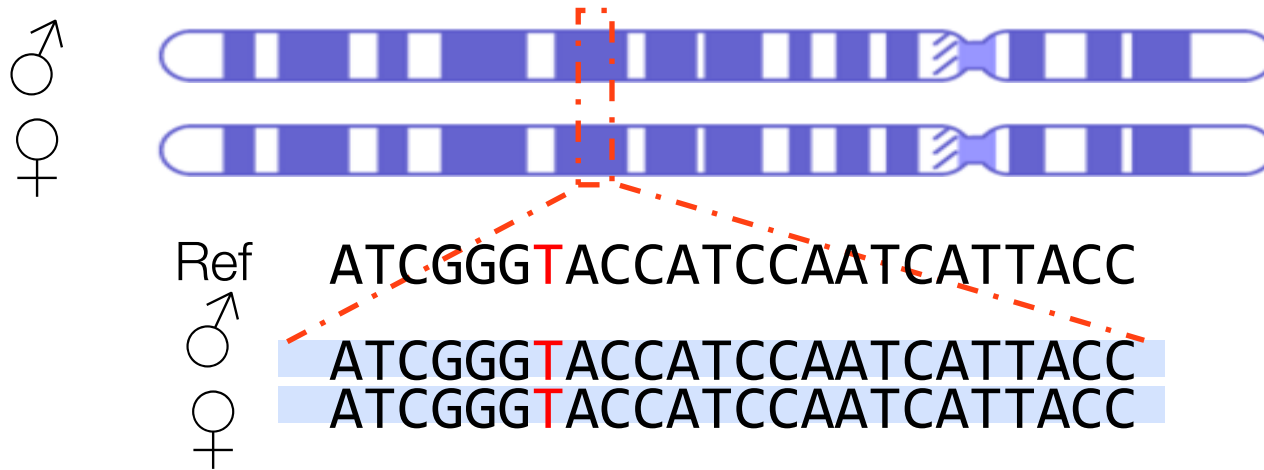
Find inherited genetic variation by sequencing DNA from millions of cells



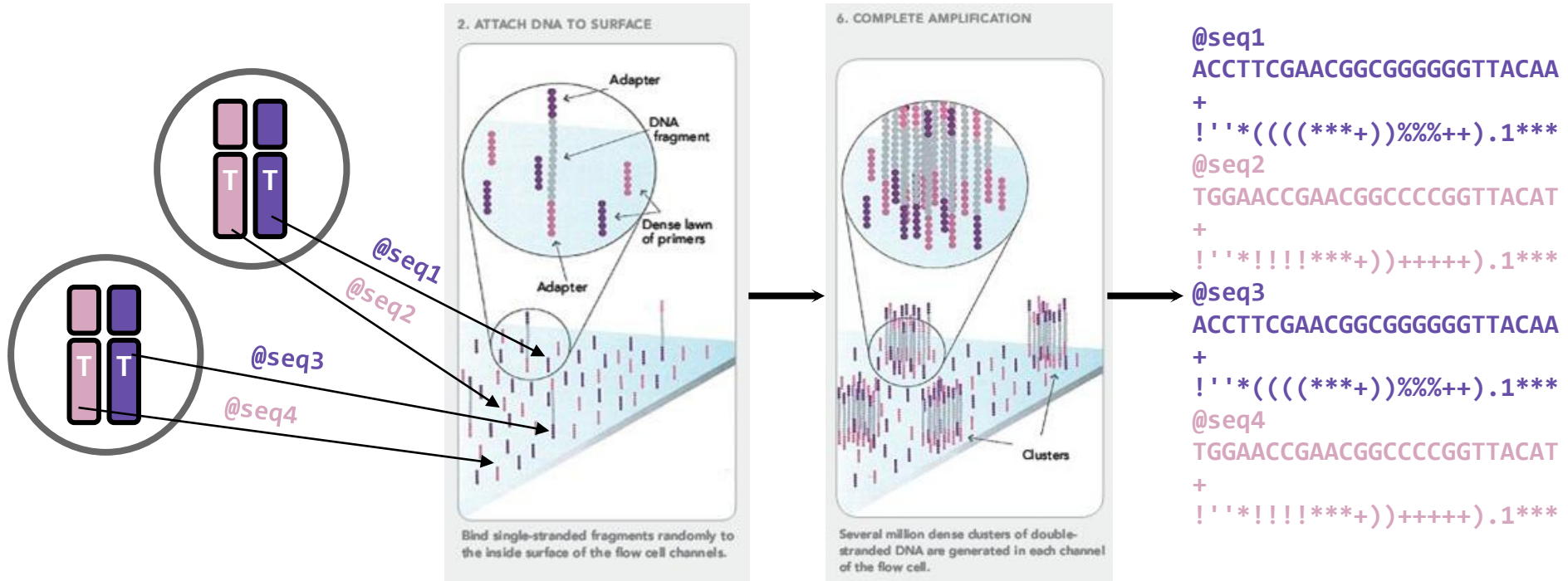
Each DNA cluster is amplified from a single strand from a single haploid chromosome from a single cell.



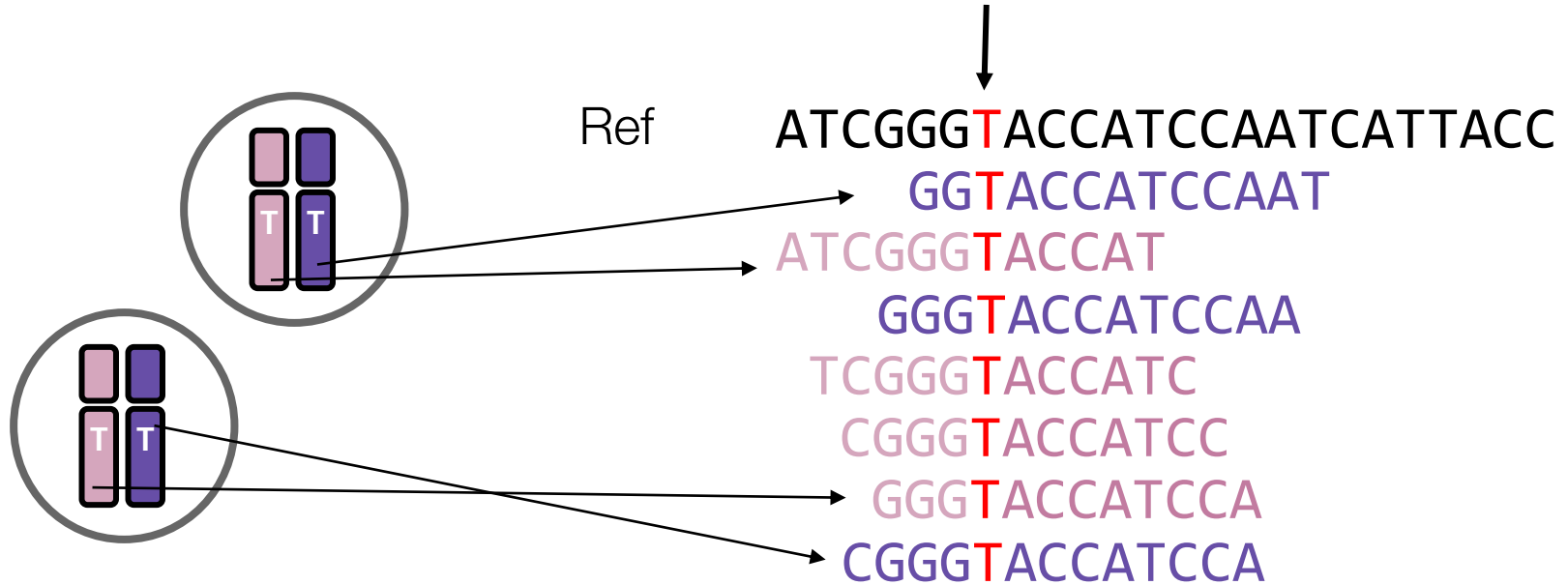
Scenario 1: An individual is homozygous for the "reference" allele.



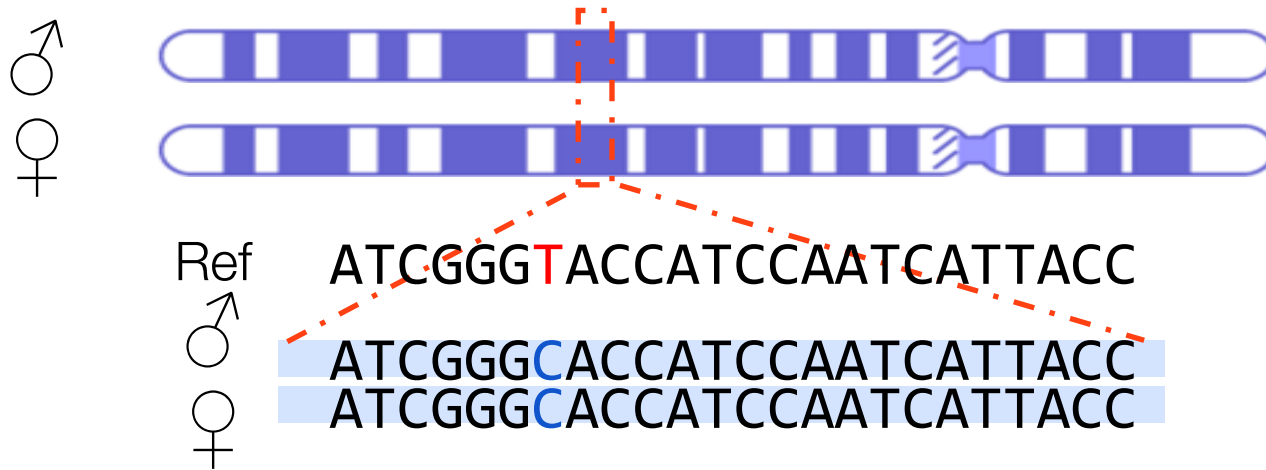
Scenario 1: An individual is homozygous for the "reference" allele.



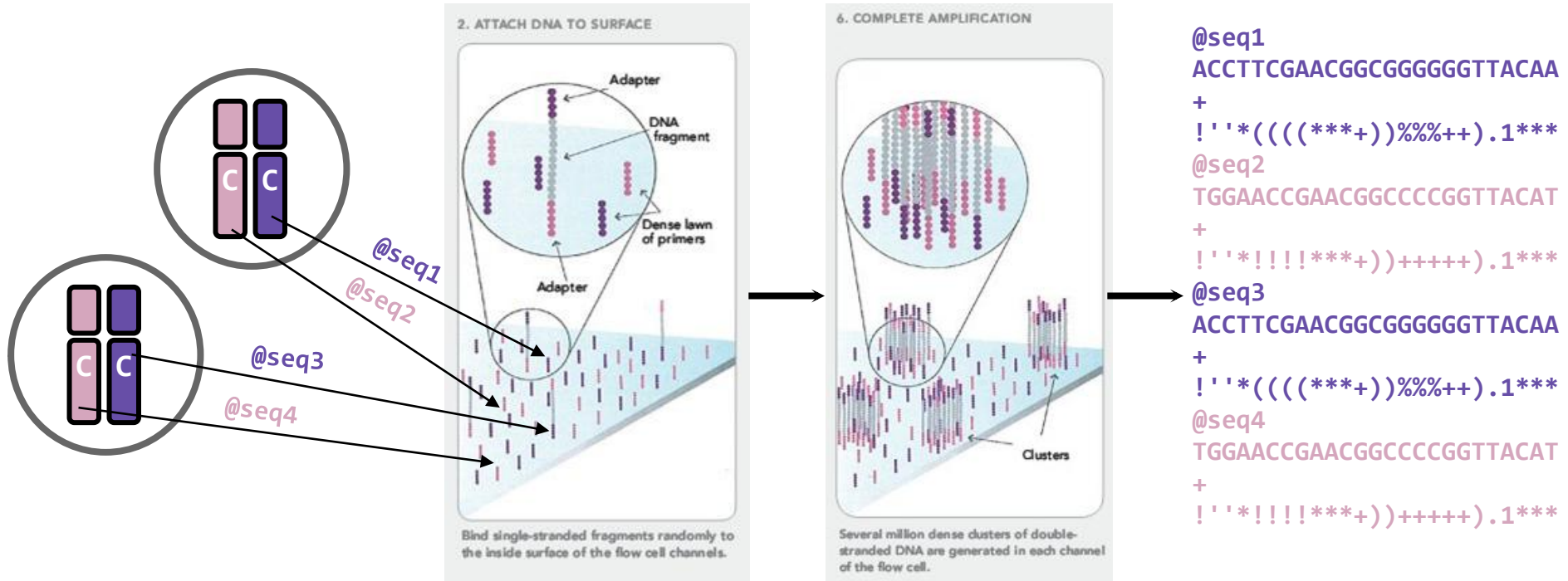
Scenario 1: An individual is homozygous for the "reference" allele.



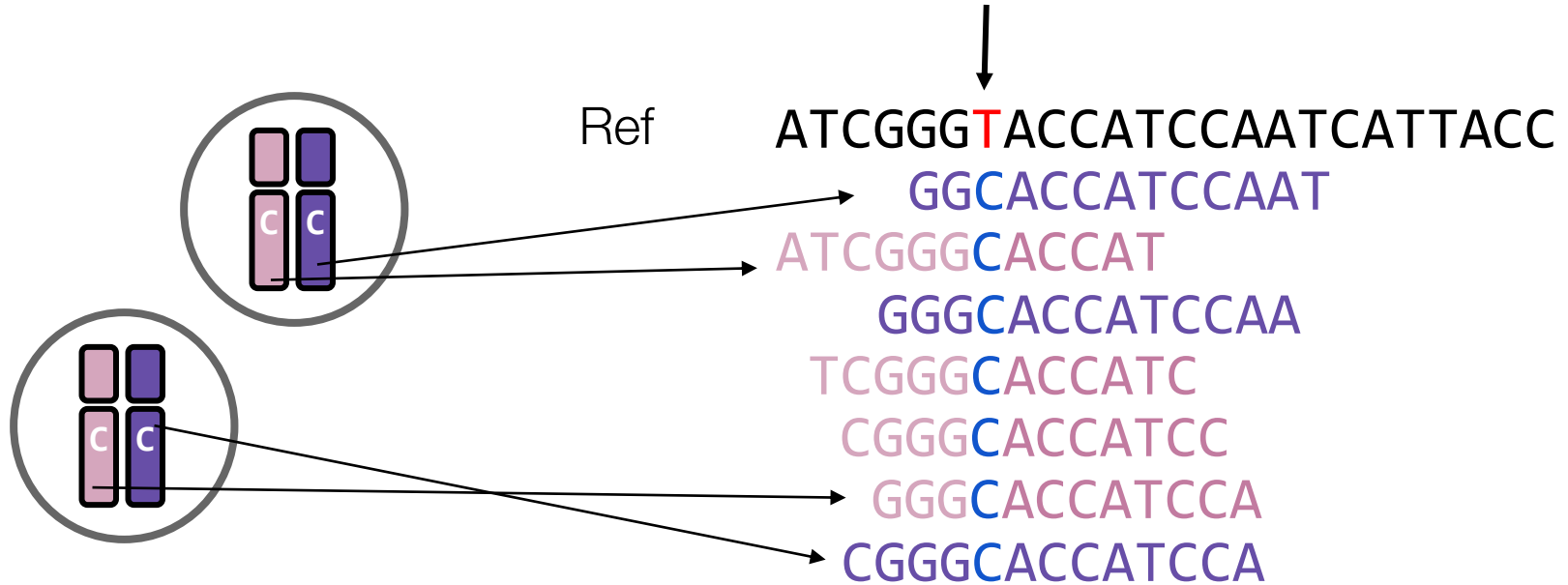
Scenario 2: An individual is homozygous for an "alternate" allele.



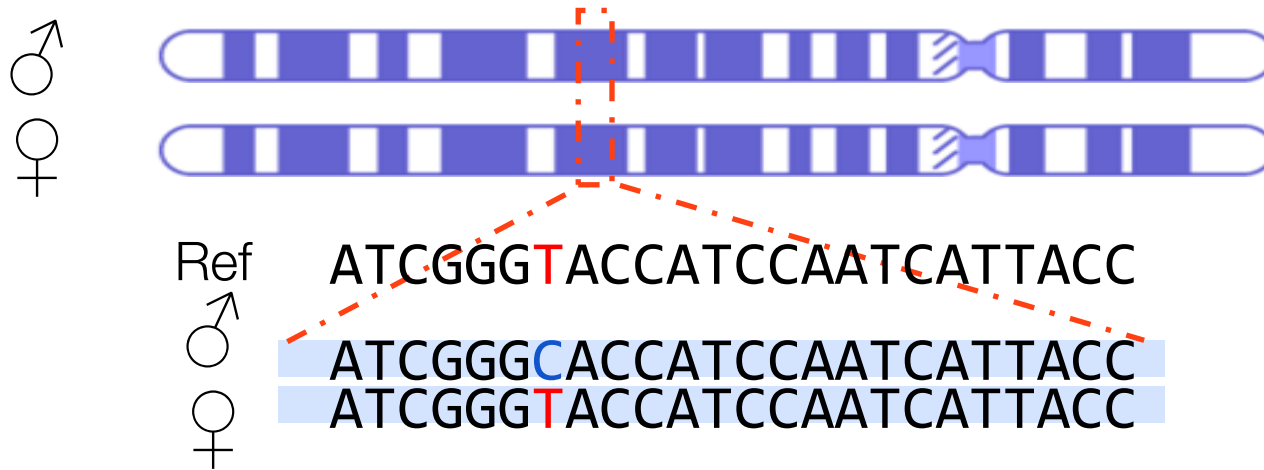
Scenario 2: An individual is homozygous for an "alternate" allele.



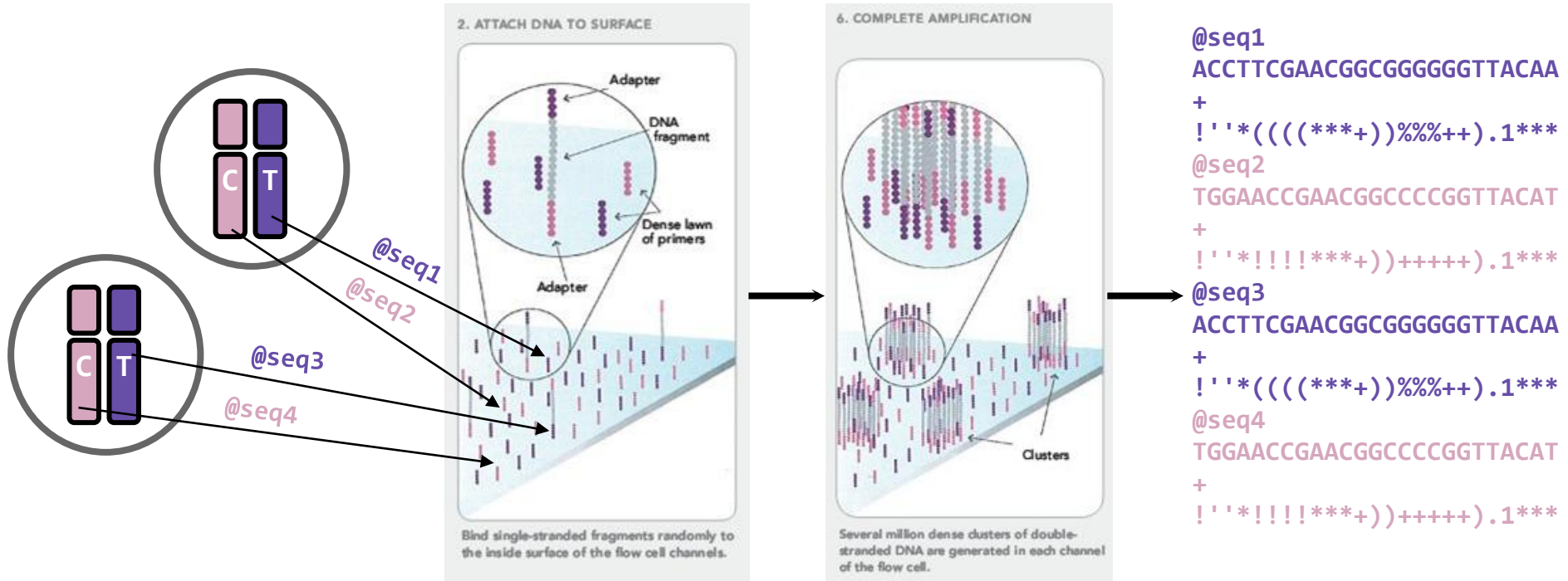
Scenario 2: An individual is homozygous for an "alternate" allele.



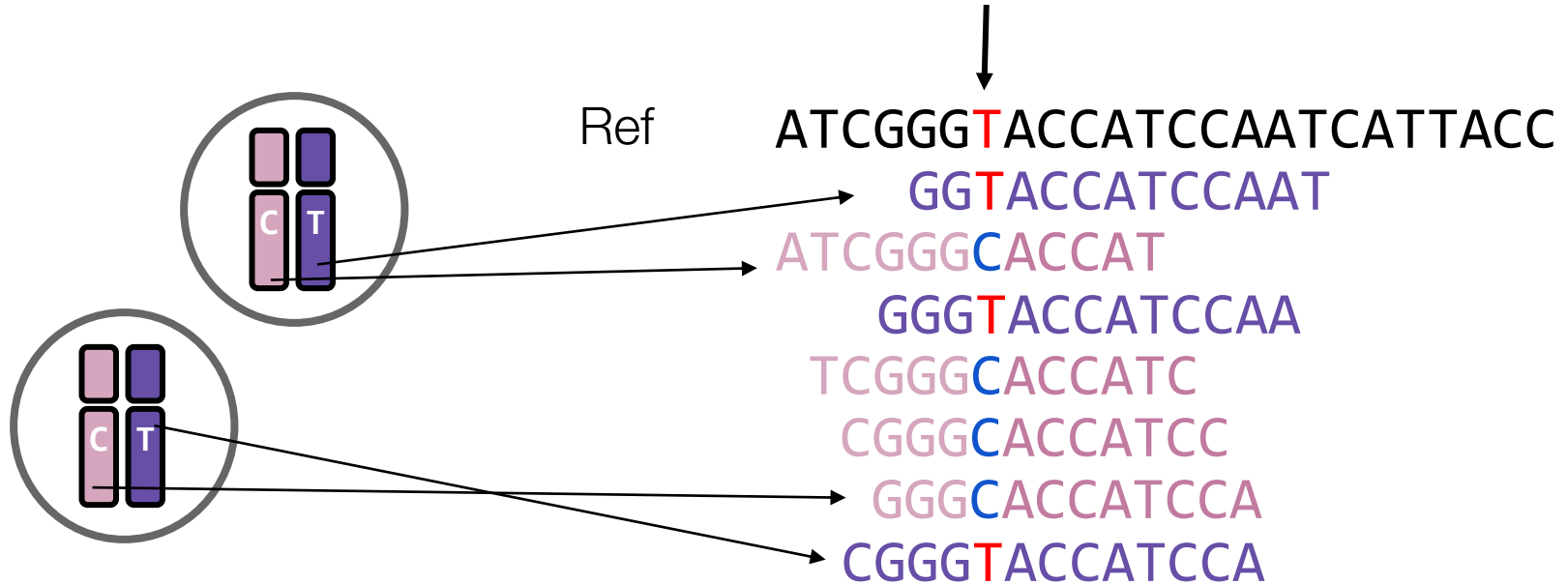
Scenario 3: An individual is heterozygous for an "alternate" allele.



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



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



Why might finding heterozygous variants be harder?

The binomial distribution: adventures in coin flipping

Sequencing read's position ?



T



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



C



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

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:

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

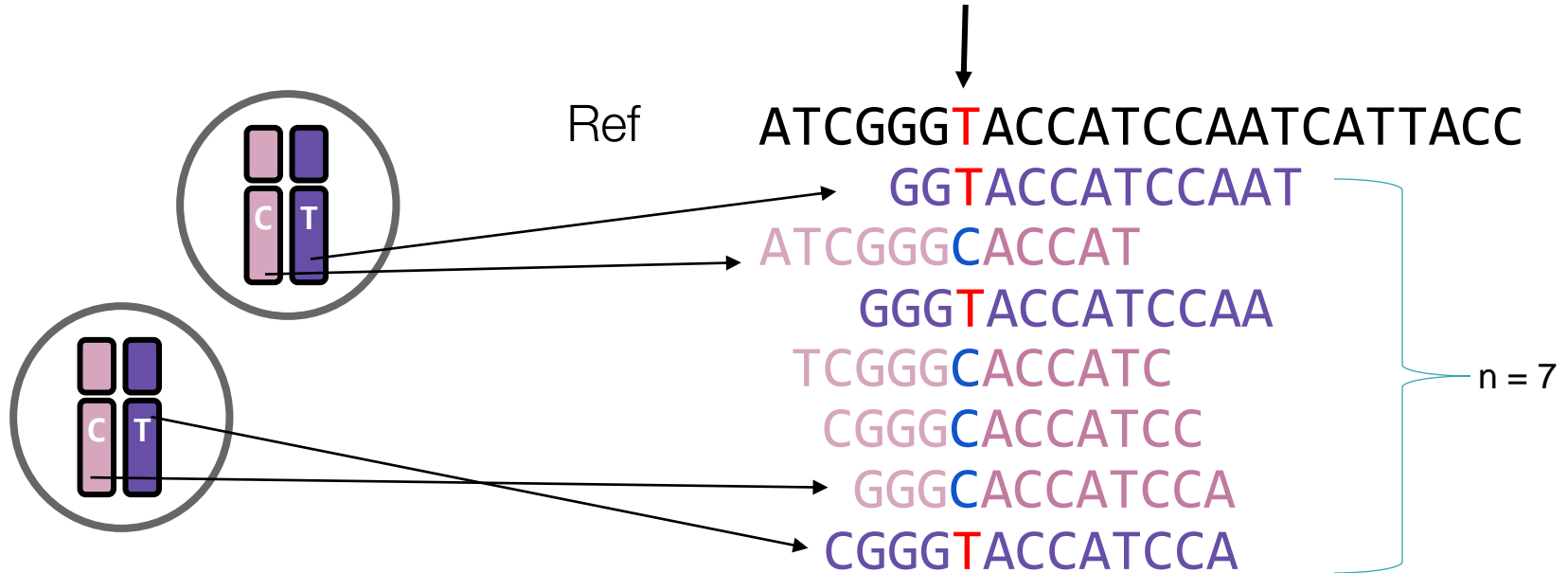
Notes:

n : number of reads that cover a particular heterozygous position

If **X** is a binomial random variable for modelling the number of seeing **T/C** at this particular position with respect to the whole number of reads that cover this position. Then:

$\Pr(X=k)$: represents the probability of seeing exactly **k 'Cs** or **k 'Ts** depending on the defining random variable **X** .

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



$n = 7$
 $\Pr [X = 2],$ where X represents C's

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:

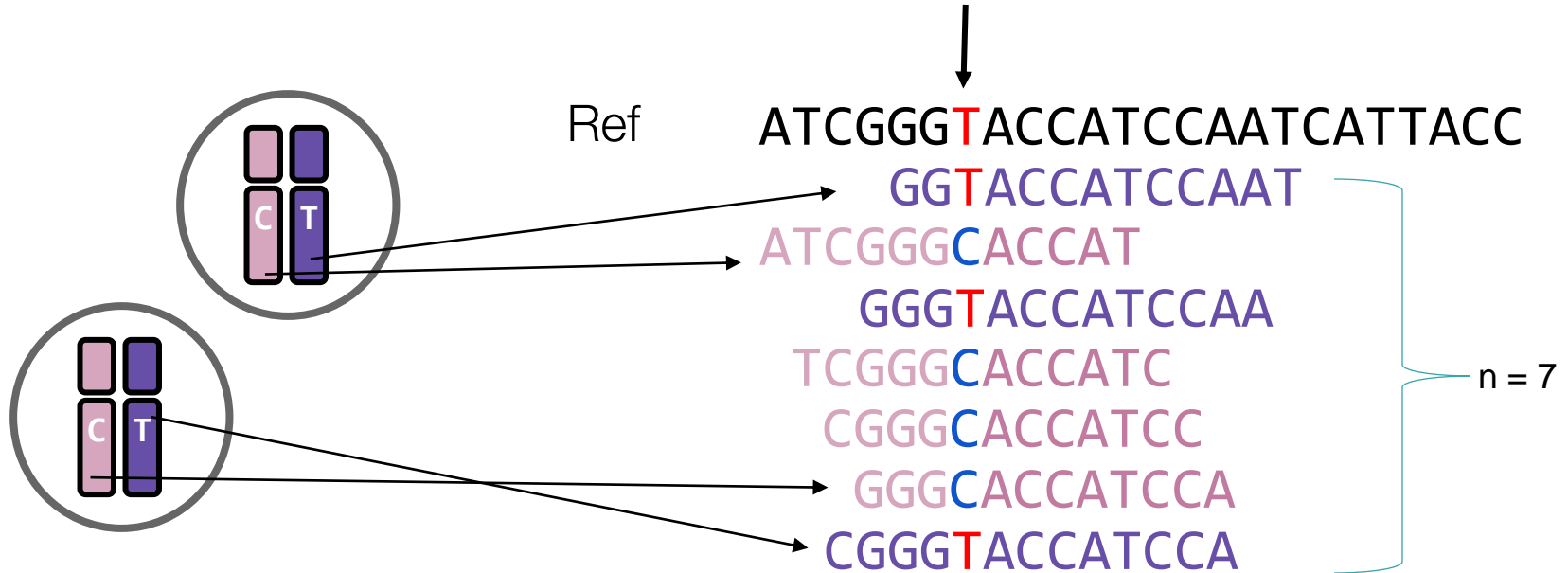
$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

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?

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

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

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



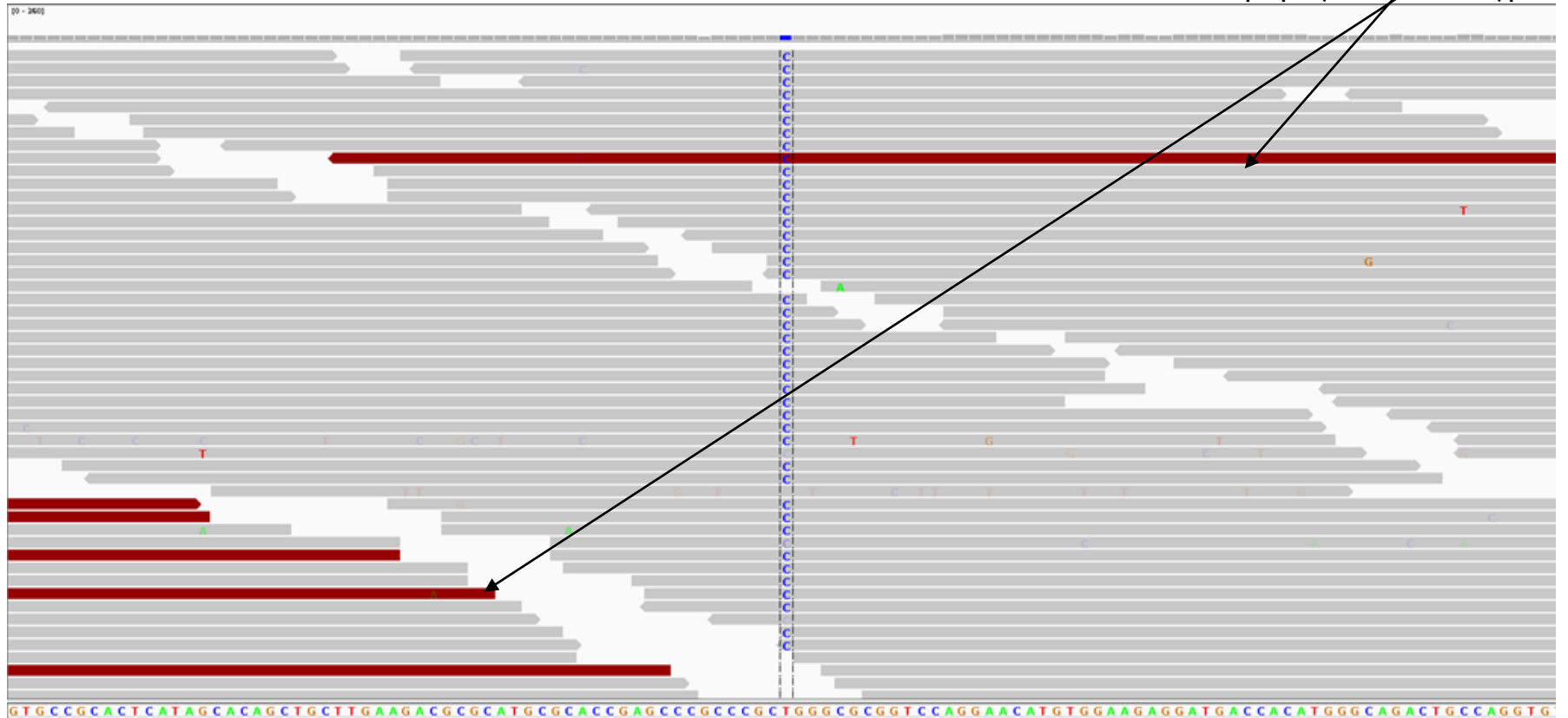
$n = 7$

$\Pr [X = 2],$ where X represents C's, $p =$ probability of success 0.57

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

Some real examples of SNPs in IGV

Homozygous for the "C" allele



Heterozygous for the alternate allele

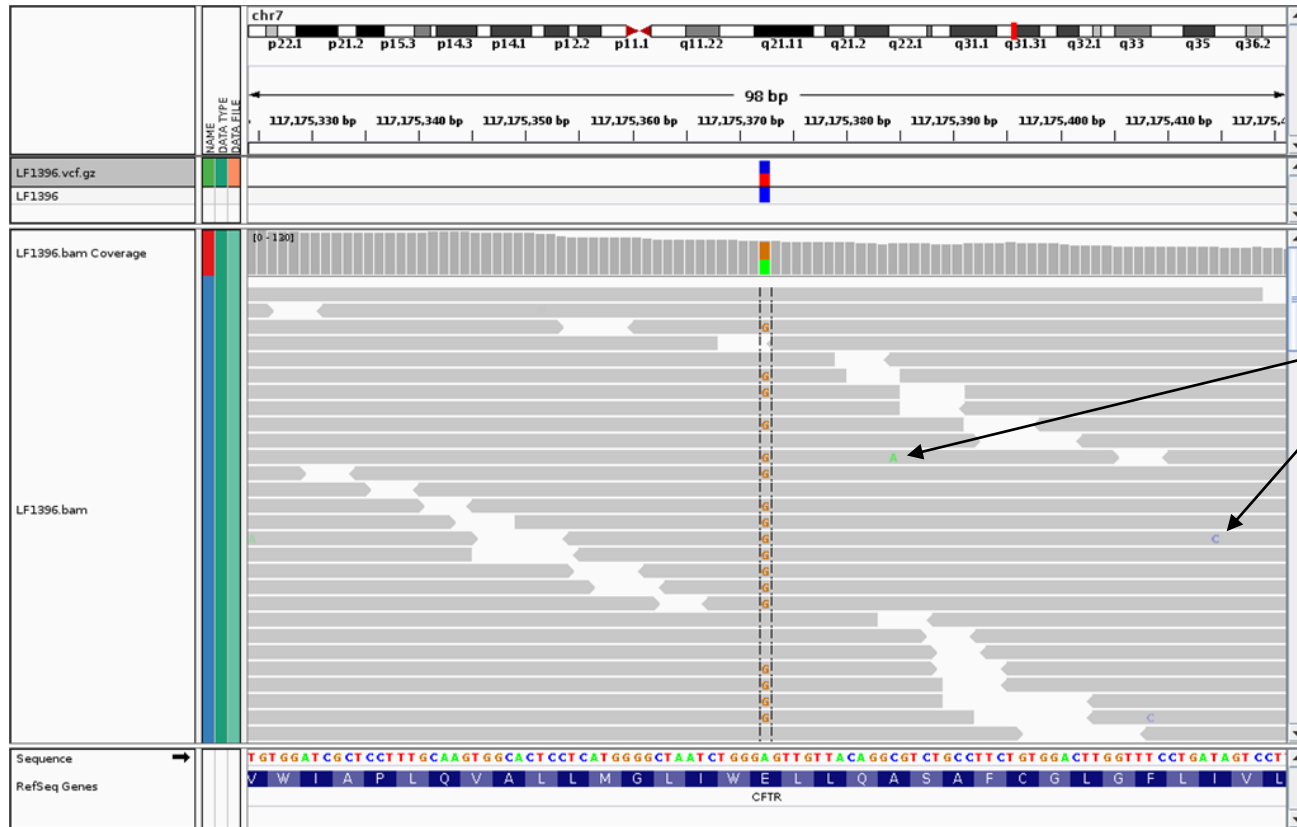
Individual 1

Individual 2



Which genotype prediction would you have more confidence in?

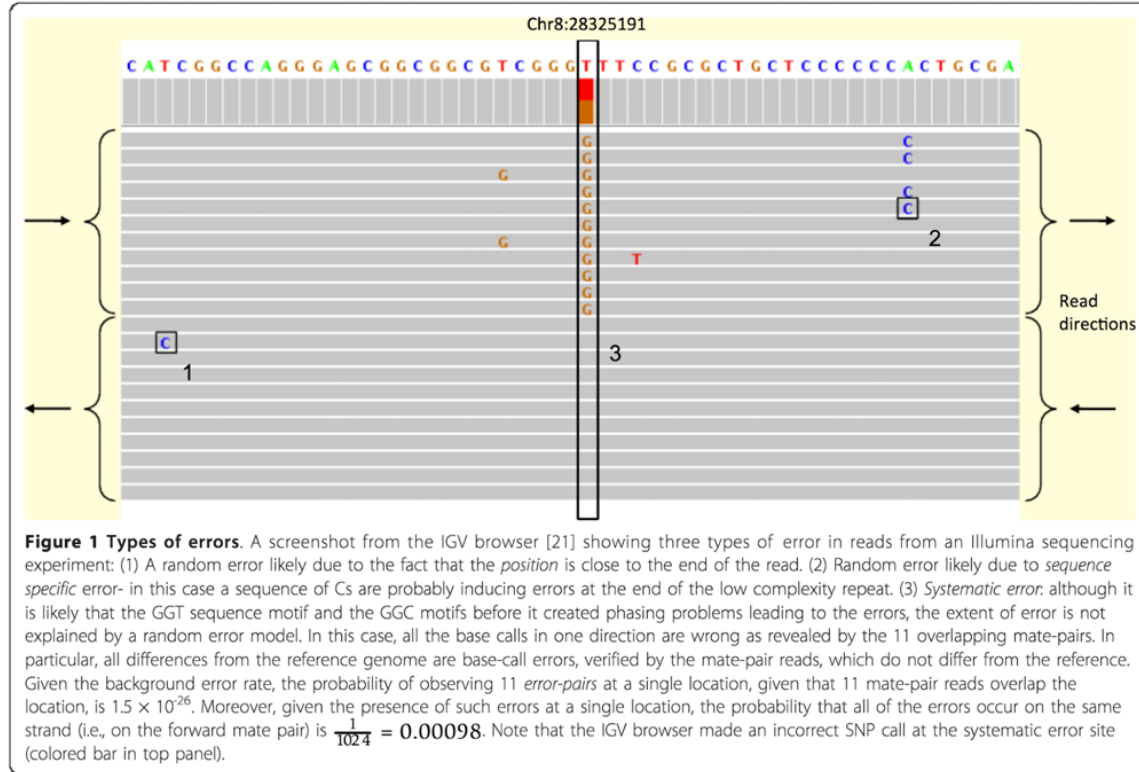
Sequencing errors fall out as noise (most of the time)



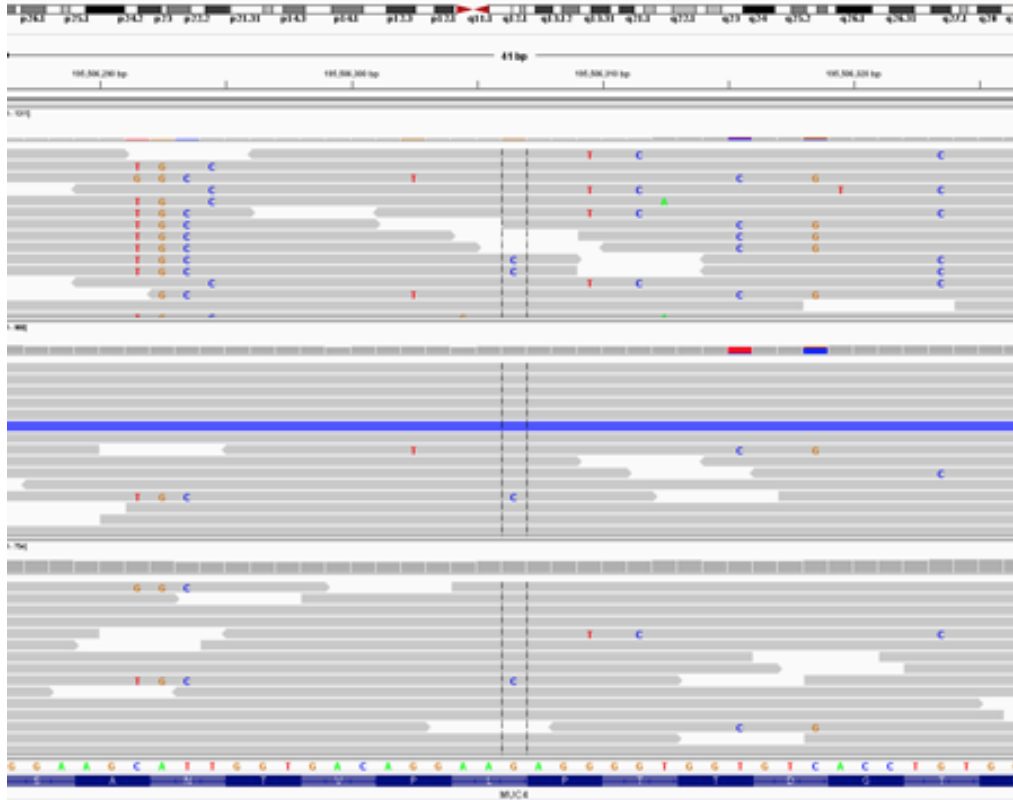
Error rate in this window= total number of errors / total numbers of bases

$$= 2 / (50 * 98) = 0.00041$$

Random versus systematic error



Pileups of many differences from paralogy



Calling INDELs is _much_ harder than SNPs

