

Genetic variation: what, why, how

Aaron Quinlan

BIMS 6000

04-Sept-2013

quinlanlab.org | arq5x@virginia.edu

Goals

- Understand the origins of genetic mutation
- Be familiar with how mutations become *polymorphic*
- Understand the processes of genetic drift and selection.
- Recognize the different types of genetic variation and their spectrum of functional consequence.

What is genetic variation?

- Differences in DNA content or structure among individuals
- Any two individuals have ~99.5% identical DNA.
- But the human genome is big - each haploid set of 23 chromosomes has 3 billion nucleotides.
- The details matter.

~98-99% identical DNA



~99.5% identical DNA



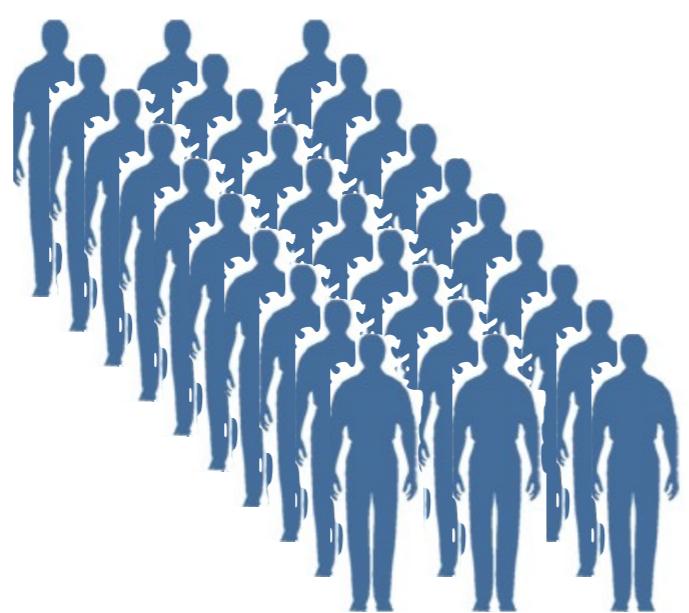
V3073025 [RF] © www.visualphotos.com

CGCAAATTGCCGGATTCTTGTGATGTTAAACGAGATTGCCAGCACC GG TATTCACCATT TT CTT
GTTAACTGCCGTCA GCCTTCTTGACCTCTTCTGTTCATGTGTATTGCTGTCTTAGCCAGACTCCC GTGCCCTTCC
ACCGGGCCTTGAGAGGTACAGGGTCTTGATGCTGTGGCTTCATGCAGGTGTCTGACTCCAGCAACTGCTGCCCTGCCAGG
GTGCAAGCTGAGCACTGGAGTGGAGTTCTGTGGAGAGGCCATGCCTAGAGTGGATGGCCATTGTTCATCTCTGCCCTG
TTGCTGCATGTAACCTAACCAACCAGGCATAGGGAAAGATTGGAGGAAAGATGAGTGAGAGCATCAACTTCTCACAAACCT
AGGCCAGTAAGTAGTGCTTGCTCATCTCCTGGCTGTGATACGTGGCCGCCCTCGCTCCAGCAGCTGGACCCCTACCTGCCGTCT
GCTGCCATCGGAGCCAAAGCCGGCTGTGACTGCTCAGACCAGCCGGCTGGAGGGAGGGCTCAGCAGGTCTGGCTTGGCCCTGGG
AGAGCAGGTGGAAGATCAGGCAGGCCATCGCTGCCACAGAACCCAGTGGATTGGCTAGGTGGATCTTGAGCTCAACAAGCCCTCT
CTGGGTGGTAGGTGCAGAGACGGGAGGGCAGAGCCGCAGGCACAGCCAAGAGGGCTGAAGAAATGGTAGAACGGAGCAGCTGGTGT
GTGTGGGCCACCAGGCCAGGCTCCTGTCTCCCCCAGGTGTGGTGTGCCAGGCATGCCCTCCCCAGCATCAGGTCTCCAGAG
CTGCAGAAGACGACGCCGACTGGATCACACTCTGTGAGTGTCCCCAGTGTGCAGAGGTGAGAGGAGAGTAGACAGTGAGTGG
GTGGCGTCGCCCTAGGGCTTACGGGCCGGTCTCCTGTCTCCTGGAGAGGCTTCATGCCCTCCACACCCCTTTGATCTCCC
TGTGATGTCATCTGGAGCCCTGCTTGCGGTGGCTATAAACGCTCTAGTCTGGCTCCAGGCCTGGCAGAGTCTTCCCAGG
AAGCTACAAGCAGCAAACAGTCTGCATGGTCATCCCCTCACTCCCAGCTCAGAGCCAGGCCAGGGCCCCAAGAAAGGCTCTGG
TGGAGAACCTGTGCATGAAGGCTGTCAACCAGTCCATAGGCAAGCCTGGCTGCCTCAGCTGGTCAGAGGAGAAGGGATGCACTGTTGG
GGAGAAGAGGAAAGTGAGGTTGCCCTGTCTCCTACCTGAGGCTGAGGAAGGAGAAGGGATGCACTGTTGGGAGGCAGCTGTA
ACTCAAAGCCTTAGCCTCTGTTCCCACGAAGGCAGGCCATCAGGCACCAAAGGGATTCTGCCAGCATAGTGCTCCTGGACCAGTGAT
ACACCCGGCACCCCTGCTGGACACGCTGTTGGCCTGGATCTGAGCCTGGTGGAGGTCAAAGCCACCTTGGTTCTGCCATTGCTGC
TGTGTGGAAGTTCACTCCTGCCTTCCCTAGAGCCTCCACCACCCGAGATCACATTCTCACTGCCCTTGTCTGCCAGT
TTCACCAGAAGTAGGCCTTCCGTACAGGCAGCTGCACACTGCCCTGGCGCTGTGCCCTTGTCTGCCCTGGAGACGGT
TTTGTCA TGGCCTGGTCTGCAGGGATCCTGCTACAAAGGTGAAACCCAGGAGAGTGTGGAGTCCAGAGTGTGCCAGGACCCAGGCA
CAGGCATTAGTGCCCGTTGGAGAAAACAGGGGAATCCCGAAGAAATGGTGGGTCTGCCATCCGTGAGATCTCCCAGGTGTGCCGT
TTTCTCTGGAAGCCTTTAAGAACACAGTGGCGCAGGCTGGTGGAGCCGTCCCCCATGGAGCACAGGCAGACAGAACGAGTCCC
CAGCTGTGGCCTCAAGCCAGCCTCCGCTTGAAGCTGGTCTCCACACAGTGCTGGTCCAGCCCTCCAAAGGAAGTAG
TCTGAGCAGCTGTCTGGCTGTCCATGTCAGAGCAACGGCCAAGTCTGGGTCTGGGGGGAGGTGTCA TGGAGCCCTACGA
TTCCCAGTCGTCCTCGTCCTCTGCCCTGTGGCTGCTCGGTGGCGCAGAGGAGGGATGGAGTCTGACACGCCGGCAAAGGCTCCT
CCGGGCCCTCACAGGCCAGGTCTTCCCTAGAGATGCCTGGAGGGAAAAGGCTGAGTGAGGTGGTTGGTGGAAACCCTGGTT
CCCCAGCCCCGGAGACTTAAATACAGGAAGAAAAAGGCAGGACAGAATTACAAGGTGCTGGCCAGGGCGGGCAGGCCCTGCCTC
CTACCCCTGCCCTCATGACCGGAGCCATAGCCCAGGCAGGAGGGCTGAGGACCTCTGGTGGCGGCCAGGGCTCCAGCATGTGCC
TAGGGGAAGCAGGGGCCAGCTGGCAAGAGCAGGGGGTGGCAGAAAGCACCCGGTGGACTCAGGGCTGGAGGGAGGAGGCATCTG
CCCAAGGCCCTCCGACTGCAAGCTCCAGGGCCGCTCACCTGCTCCTGCTCCTGCTGCTGCTGCTTCTCAGCTTCGCTCCTCAT
GCTGCGCAGCTGGCTTGCCTGCCATGCCCTCAGCTGGCGGATGGACTCTAGCAGAGTGGCCAGCCACC GGAGGGTCAACCAC
TGGGAGCTCCCTGGACTGGAGCCGGAGGTGGGGAACAGGGCAAGGAGGAAGGCTGCTCAGGCAGGGCTGGGGAGGCTTACTGTGTC
CAAGAGCCTGCTGGAGGGAAAGTCACCTCCCTCAAACGAGGAGGCCCTGCCTGGAGGGAGGCGACCTTGGAGAGACTGTGTGG
CCTGGGCACTGACTTCTGCAACCACCTGAGCGCGGGCATCCTGTGTCAGATACTCCCTGCTTCTCTAGCCCCCACCCTGCAGAG
CTGGACCCCTGAGCTAGCCATGCTGACAGTCTCAGTTGCACACACGAGCCAGCAGAGGGTTTGTGCCACTTCTGGATGCTAGGG
TTACACTGGAGACACAGCAGTGAAGCTGAAATGAAAAATGTGTTAGTTGCTGTAGTTGTTATTAGACCCTTCCATTGGTTAATT

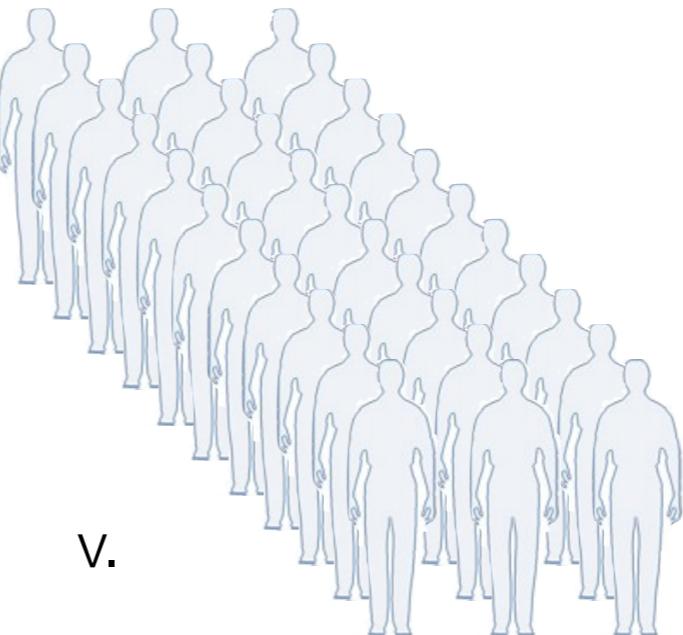
CGCAAATTGCCGGATTCCCTTGCATGTAGTTAACGAGATTGCCAGCACCGGGTATCACCATTTCCTTCA
GTTAACTGCCGTCAGCCTTCTTGACCTCTTCTTCTGTTCATGTGTATTGCTGCTCTTAGCCAGACTCCCCTGCCAGG
ACCGGGCCTTGAGAGGTACAGGGTCTTGATGCTGTTCATCTGCAGGTGCTGACTCCAGCAACTGCTGGCCTGCCAGG
GTGCAAGCTGAGCACTGGAGTTCTGTGGAGAGGGAGCCATGCCTAGAGTGGATGGCCATTGTTCATCTCTGGCCCCCTG
TTGTCTGCATGTAACCTAACCAACCAGGCATAGGGAAAGATTGGAGGAAAGATGAGTGAGAGCATCAACTCTCACAAACCT
AGGCCAGTAAGTAGTGCTTGCTCATCTCCTGGCTGTGATACGTGCCGGCCCTCGCTCCAGCAGCTGGACCCCTACCTGCCGTCT
GCTGCCATCGGAGCCAAAGCCGGCTGTGACTGCTCAGACCAGCCGGCTGGAGGGAGGGCTCAGCAGGTCTGGCTTGGCCCTGGG
AGAGCAGGTGGAAGATCAGGCAGGCCATCGCTGCCACAGAACCCAGTGGATTGGCTAGGTGGATCTTGAGCTCAACAAGCCCTCT
CTGGGTGGTAGGTGCAGAGACGGGAGGGCAGAGCCGCAGGCACAGCCAAGAGGGCTGAAGAAATGGTAGAACGGAGCAGCTGGTGT
GTGTGGGCCACCAGGCCAGGCTCCTGTCTCCCCCAGGTGTGGTGTGCCAGGCATGCCCTCCCCAGCATCAGGTCTCCAGAG
CTGCAGAAGACGACGGCGACTGGATCACACTCTGTGAGTGTCCCAGTGTGCAGAGGTGAGAGGAGAGTAGACAGTGAGTGGGA
GTGGCGTCGCCCTAGGGCTCTACGGGCCGGCTCCTGTCTGGAGAGGGCTTCGATGCCCTCCACACCCCTTGATCTTCCC
TGTGATGTCATCTGGAGCCCTGCTGCTTGCCTGGCTGGCCTATAAGCCTCTAGTCTGGCTCCAAGGCCTGGCAGAGTCTTCCAGGG
AAGCTACAAGCAGCAAACAGTCTGCATGGTCATCCCCCTCACTCCCAGCTCAGAGCCAGGCCAGGGCCCCAAGAAAGGCTCTGG
TGGAGAACCTGTGCATGAAGGCTGTCAACCAGTCCATAGGCAAGCCTGGCTGCCTCCAGCTGGTCGACAGACAGGGCTGGAGAAGG
GGAGAAGAGGAAAGTGAGGTTGCCCTGTCTCCTACCTGAGGCTGAGGAAGGAGAAGGGATGCACTGTTGGGAGGCAGCTGTA
ACTCAAAGCCTAGCCTCTGTTCCCACGAAGGCAGGGCATCAGGCACCAAAAGGGATTCTGCCAGCATAGTGCTCTGGACCAGTGAT
ACACCCGGCACCCCTGCTGGACACGCTGTTGGCTGGATCTGAGCCCTGGTGGAGGTCAAAGCCACCTTGGTCTGCCATTGCTGC
TGTGTGGAAGTTCACTCCTGCCCTTCCCTAGAGCCTCCACCCCCGAGATCACATTCTCACTGCCCTTGTCTGCCAGT
TTCACCAGAAGTAGGCCTTCCCTGACAGGCAGCTGCACCACTGCCTGGCGCTGTGCCCTTGTCTGCCAGGGCTGGAGACGGTG
TTTGTCATGGCCTGGCTGCAGGGATCCTGCTACAAAGGTGAAACCCAGGAGAGTGTGGAGTCCAGAGTGTGCCAGGACCCAGG
CAGGCATTAGTGCCGTTGGAGAAAACAGGGGAATCCGAAGAAATGGTGGCTGGCCATCCGTGAGATCTCCAGGTGTGCCGT
TTTCTCTGGAAGCCTTTAACAGAACACAGTGGCGCAGGCTGGTGGAGCCCTCCACAGTGCTGGTCCAGCCCTCCAGGAAGTAG
CAGCTGTGTGGCCTCAAGCCAGCCTCCGCTCTGAAGCTGGTCTCCACACAGTGCTGGTCCAGCCCTCCAGGAAGTAG
TCTGAGCAGCTTGTCTGGCTGTCCATGTCAGAGCAACGGCCAAGTCTGGTCTGGGGGGAGGTGTCACTGGAGCCCCCTACGA
TTCCCAAGTCGTCCTCGTCTGCCCTGCTGGCTGCTGCCGTGGCGAGAGGAGGGATGGAGTCTGACACGCCGGCAAAGGCTCCT
CCGGGCCCTCACCAGCCCCAGGTCTTCCAGAGATGCCCTGGAGGGAAAAGGCTGAGTGAGGGTGGTGGAGGAAACCCCTGGTTC
CCCCAGCCCCGGAGACTTAAATACAGGAAGAAAAAGGCAGGACAGAATTACAAGGTGCTGGCCAGGGCGGGCAGCGGCCCTGCC
CTACCCCTGCGCCTCATGACCGGAGCCATAGCCCAGGCAGGAGGGCTGAGGACCTCTGGTGGCGGCCAGGGCTTCCAGCATGTGCC
TAGGGGAAGCAGGGGCCAGTGGCAAGAGCAGGGGGTGGCAGAAAGCACCCGGTGGACTCAGGGCTGGAGGGAGGAGGCGATCTG
CCCAAGGCCCTCCGACTGCAAGCTCCAGGGCCGCTCACCTGCTCCTGCTGCTGCTTCTGCTGCTTCCAGCTTGCTCCTCAT
GCTGCGCAGCTTGGCCTTGCGATGCCCGAGCTTGGCGGATGGACTCTAGCAGAGTGGCCAGCCACCGGAGGGTCAACCACCTCC
TGGGAGCTCCCTGGACTGGAGCCGGAGGTGGGAACAGGGCAAGGAGGAAAGGCTGCTCAGGCAGGGCTGGGAAGCTTACTGTGTC
CAAGAGCCTGCTGGAGGGAAAGTCACCTCCCTCAAACAGAGGAGCCCTGCCTGGGGAGGGCAGCTTGGAGACTGTGTGGGG
CCTGGGCACTGACTTCTGCAACCACCTGAGCGGGCATCCTGTGTGCAGATACTCCCTGCTTCTAGCCCCCACCCTGCAGAG
CTGGACCCCTGAGCTAGCCATGCTCTGACAGTCTCAGTTGCACACACGAGCCAGCAGAGGGTTTGCCACTCTGGATGCTAGGG
TTACACTGGGAGACACAGCAGTGAAGAAAAATGTGTTGCTGTAGTTGTTATTAGACCCCTTCCATTGGTTAAATT

Why do we care?

- Understand the relationship between genotype and phenotype.



Cases
(have disease)

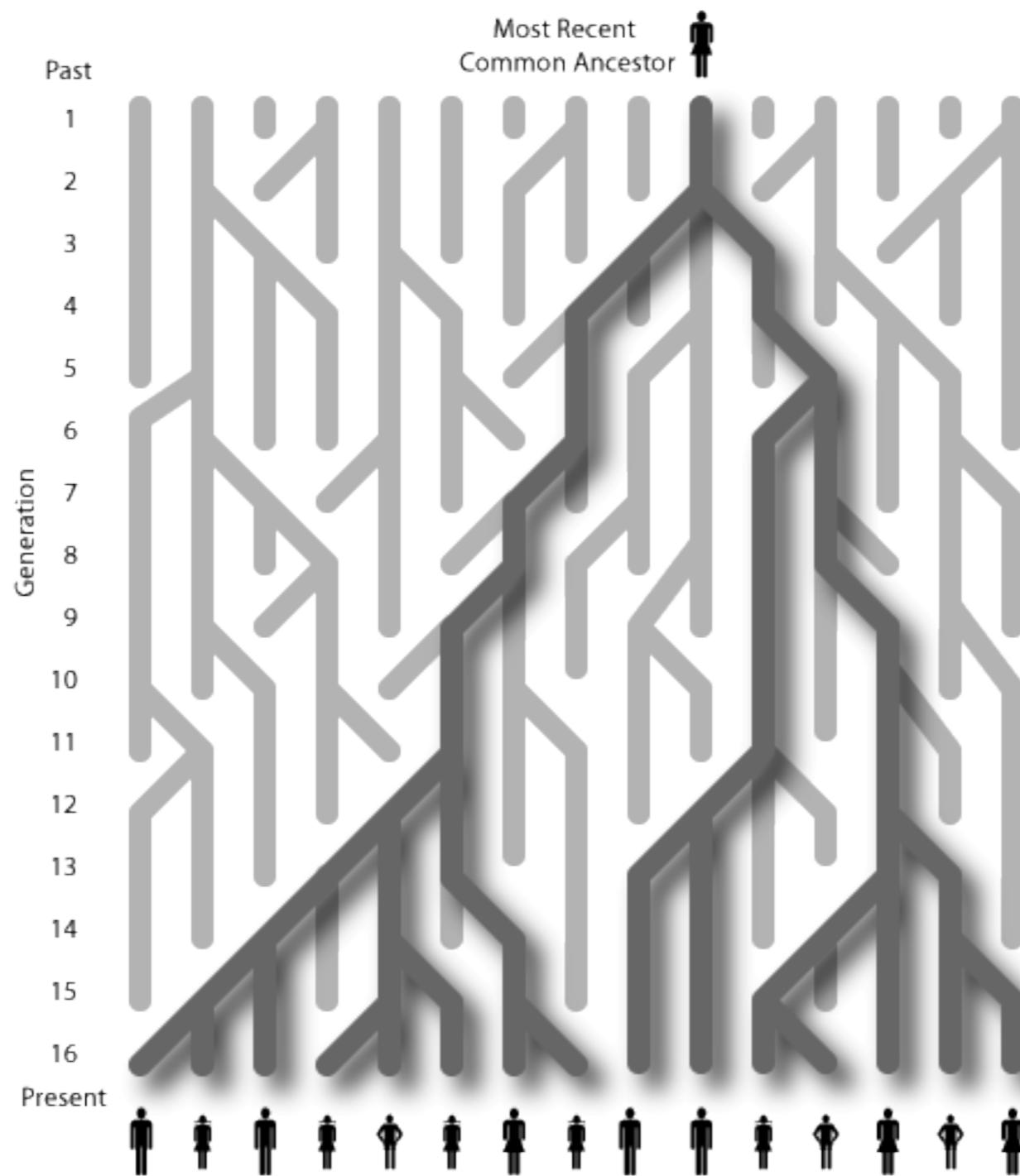


Controls
(no disease)

Complex diseases
(multiple genes contribute to risk)

Why do we care?

- Bread crumbs of evolution



Why do we care?

- How, when, where does our genome evolve?



Types of genetic variation

ctc**c**gag
ctc**t**gag

Single-nucleotide
polymorphisms
(SNPs)

ctc--ag
ctc**tg**ag

Insertion-deletion
polymorphisms
(INDELs)

ctcag
ctc  ag

Structural
variants
(SVs)

“spelling mistakes”

*“extra or missing
letters”*

*“extra, missing
or reordered
chapters”*

Properties of genetic variation

Single-nucleotide (SNPs)	ctcc c gag	ctc -- ag	ctc ag
	ctct t gag	ctc tg ag	ctc  ag
Size	1bp	1-100bp	100bp-1Mb+
Frequency	3 million / genome	300K / genome	3,000 / genome
Detection Difficulty	Easy	Medium	Hard

How different are we?

Single-nucleotide
(SNPs)

ctc**c**gag
ctc**t**gag

Insertion-deletions
(INDELs)

ctc--ag
ctc**tg**ag

Structural
variants (**SVs**)

ctcag
ctc  ag

How different are we?

The genomes of any two humans are ~99.5% identical

Single-nucleotide
(SNPs)

ctc**c**gag
ctc**t**gag

Insertion-deletions
(INDELs)

ctc--ag
ctc**tg**ag

Structural
variants (**SVs**)

ctcag
ctc  ag

How different are we?

The genomes of any two humans are ~99.5% identical

But, the genome is big.

Single-nucleotide
(SNPs)

ctcc**c**gag
ctc**t**gag

Insertion-deletions
(INDELs)

ctc--ag
ctc**tg**ag

Structural
variants (**SVs**)

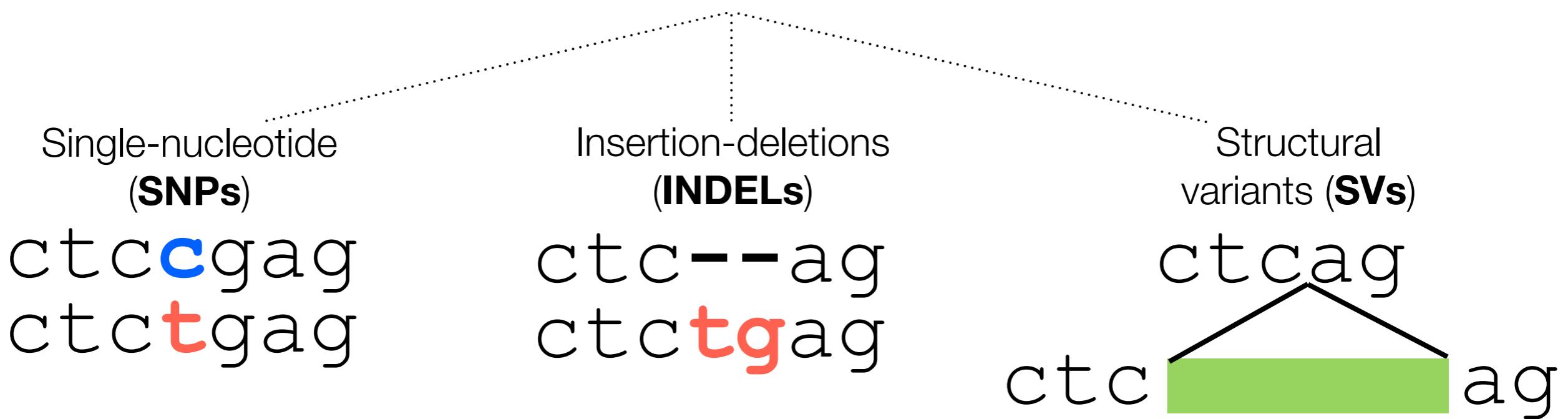
ctcag
ctc  ag

How different are we?

The genomes of any two humans are ~99.5% identical

But, the genome is big.

~21,000,000 different base pairs.



It all starts with mutation

acctccgagta

a toy population of 10 identical chromosomes

Mutation creates genetic diversity

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctccgagta

acctc**T**gagta

mutation:
private to this chromosome

From mutation to polymorphism

acctccgagta

acctccgagta

acctccgagta

acctc**T**gagta

acctccgagta

acctc**T**gagta

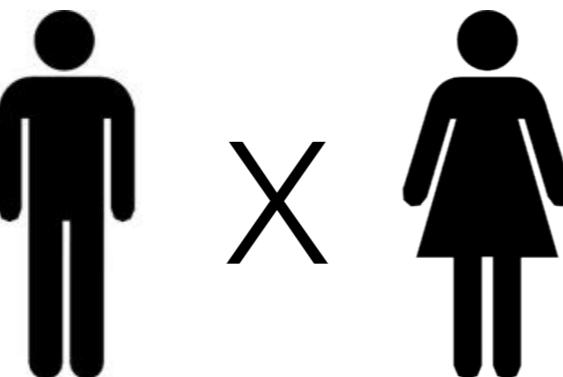
acctccgagta

acctc**T**gagta

acctccgagta

acctc**T**gagta

de novo mutation: the birth of new variation



♂ ctcc**c**gag ♀ ctcc**c**gag
♂ ctcc**c**gag ♀ ctcc**c**gag

Example: Mom and dad are homozygous for the same alleles.



**New mutation occurs in
father's or mother's germ cell**

♂ ctcc**c**gag → ♂ ctcc**c**gag



♂ ctct**t**gag
♀ ctcc**c**gag

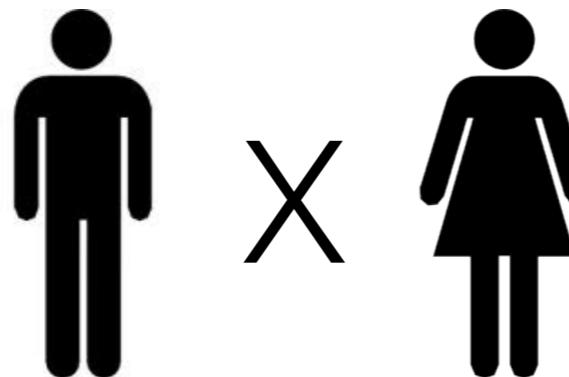
Note: This is a derivative chromosome of the one the father inherited from his parents. The mutation occurred in his gamete (sperm) and was passed on to the child.

Kid is heterozygous owing to *de novo mutation*.
(C/T)

How existing (germline) variation is inherited

♂
♀

denotes from which
parent the chrom. was
inherited



ctcc**c**gag
ctct**t**gag

ctcc**c**gag
ctct**t**gag



or

♂ ctcc**c**gag
♀ ctcc**c**gag

♂ ctcc**c**gag
♀ ctct**t**gag

♂ ctcc**c**gag
♂ ctct**t**gag

♂ ctct**t**gag
♀ ctct**t**gag

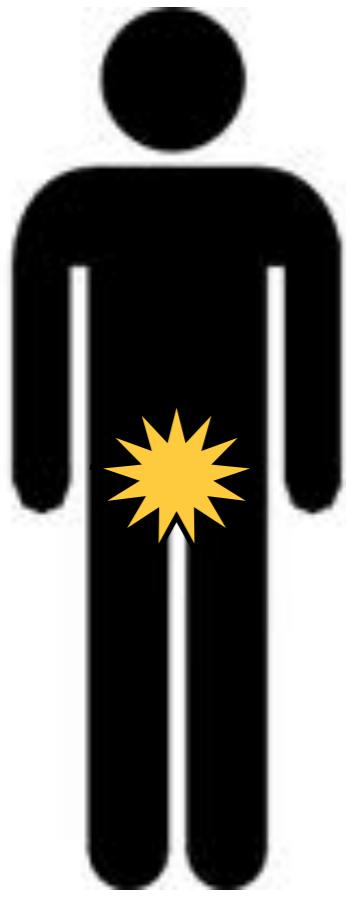
Kid is homozygous
(C/C)

Kid is heterozygous
(C/T)

Kid is homozygous
(T/T)

Example: Mom and dad are heterozygous; that is, the zygote from which they developed was comprised of a sperm and egg with two different alleles

somatic mutations



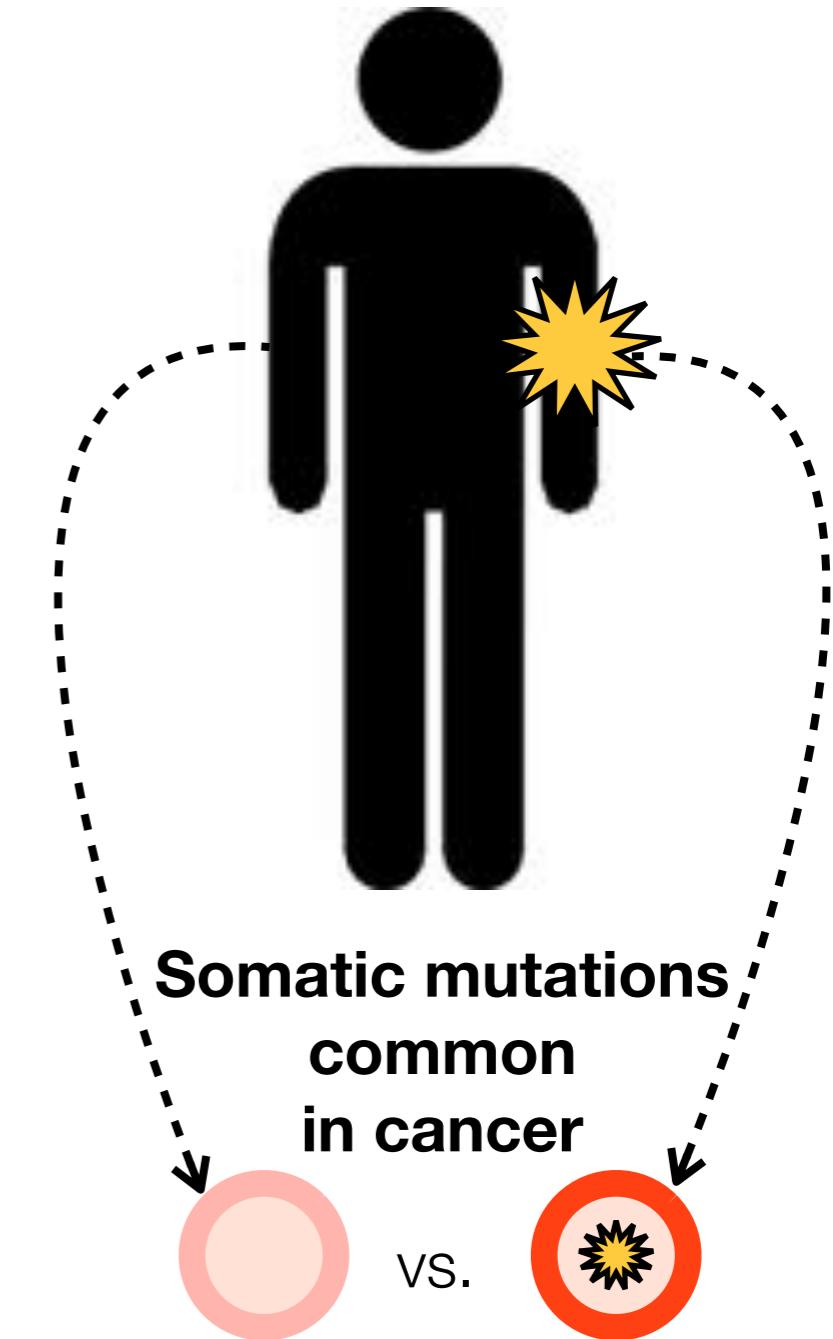
Germline mutation

- occur in sperm or egg.
- are heritable



Somatic mutation

- non-germline tissues.
- are not heritable



compare DNA from cancer cells to healthy cells from same individual

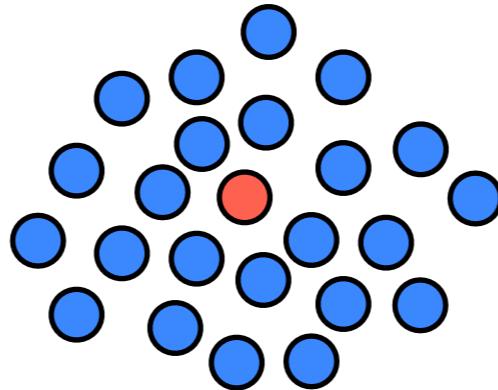
Selection and genetic drift

All other chromosomes

acctcc**c**gagta

acctct**t**gagta

Chromosome with new allele



What if the mutation is:

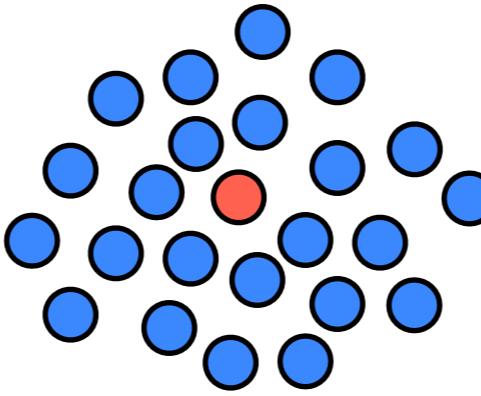
✗

Selection and genetic drift

All other chromosomes
acctcc**c**gagta

acctct**t**gagta

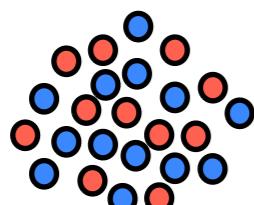
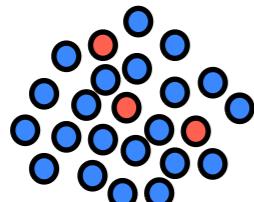
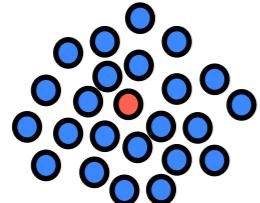
Chromosome with new allele



What if the mutation is:

beneficial

Time
↓



✗

Positive
selection

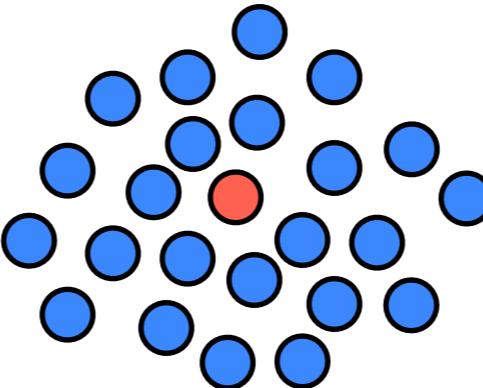
Selection and genetic drift

All other chromosomes

acctcc**c**gagta

acctct**t**gagta

Chromosome with new allele

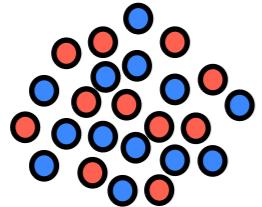
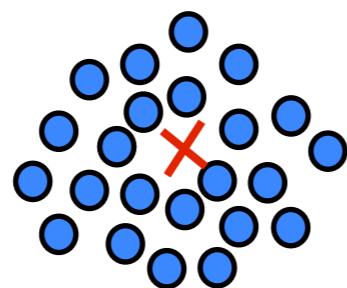
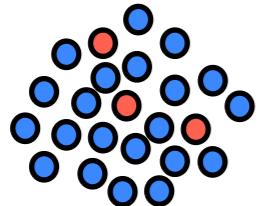
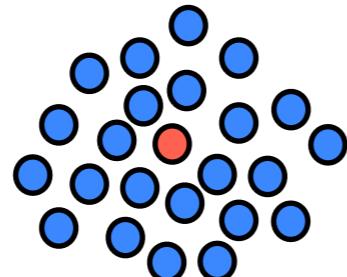
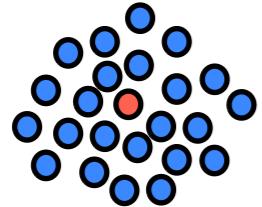


What if the mutation is:

beneficial

deleterious

Time

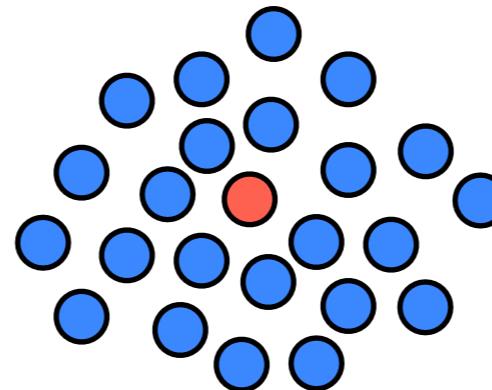


**Positive
selection**

**Negative
selection**

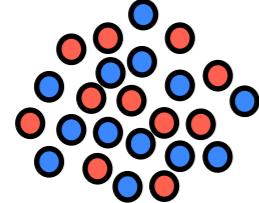
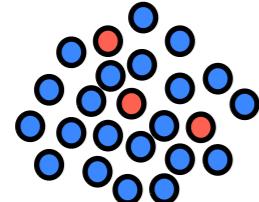
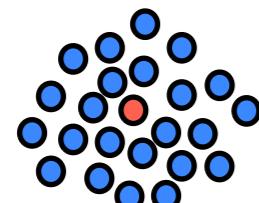
Selection and genetic drift

All other chromosomes
acctcc**c**gagta
acctct**t**gagta
Chromosome with new allele



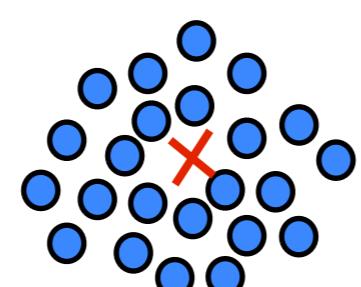
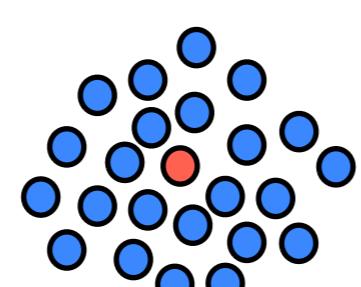
What if the mutation is:

beneficial



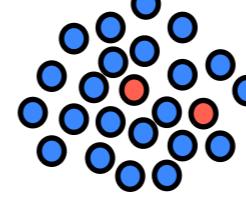
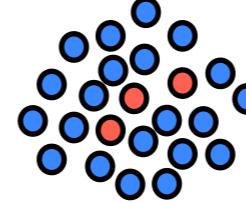
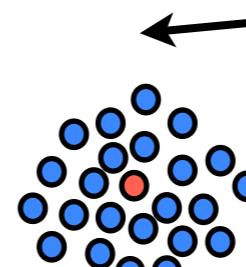
Positive selection

deleterious

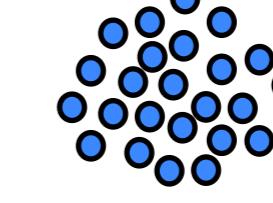
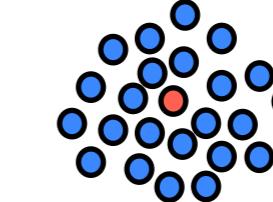
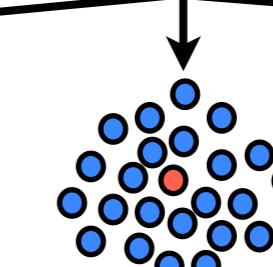


Negative selection

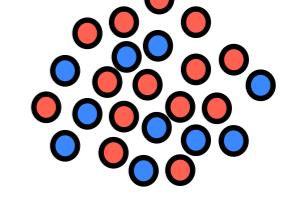
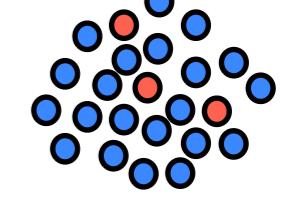
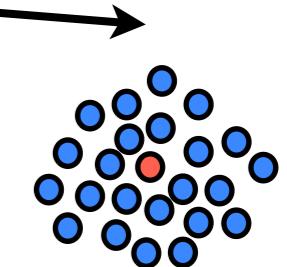
neutral



freq. =



freq. ↓



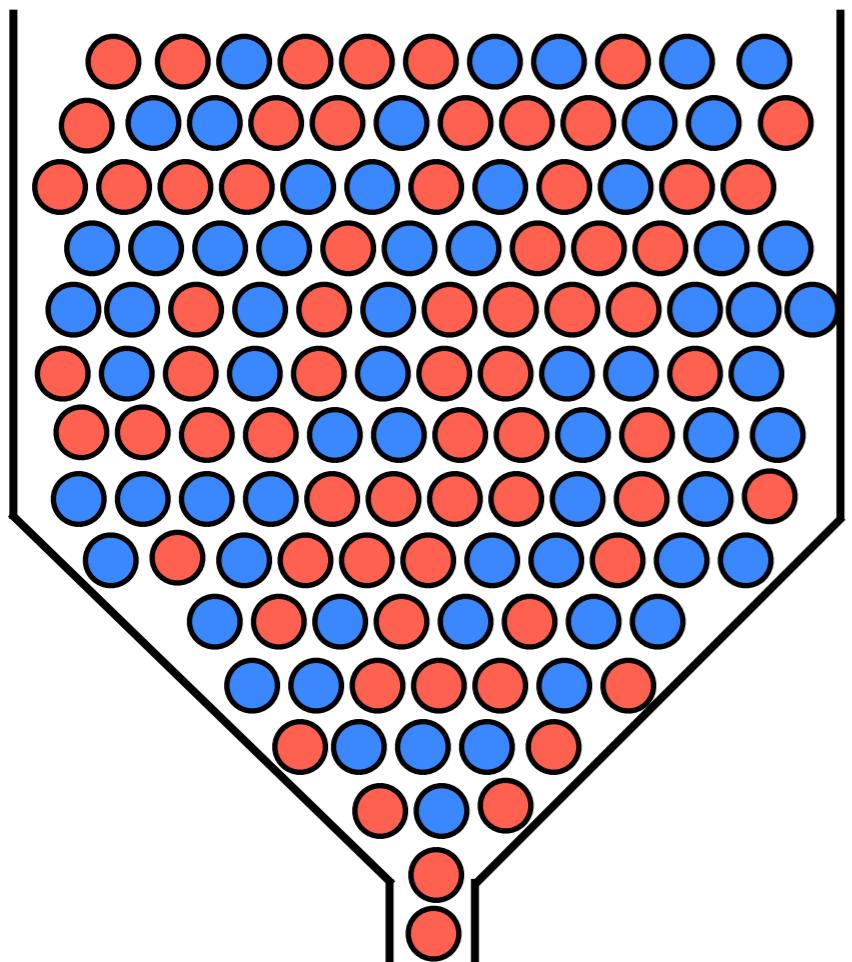
freq. ↑

Genetic “drift” - random process

Time
↓

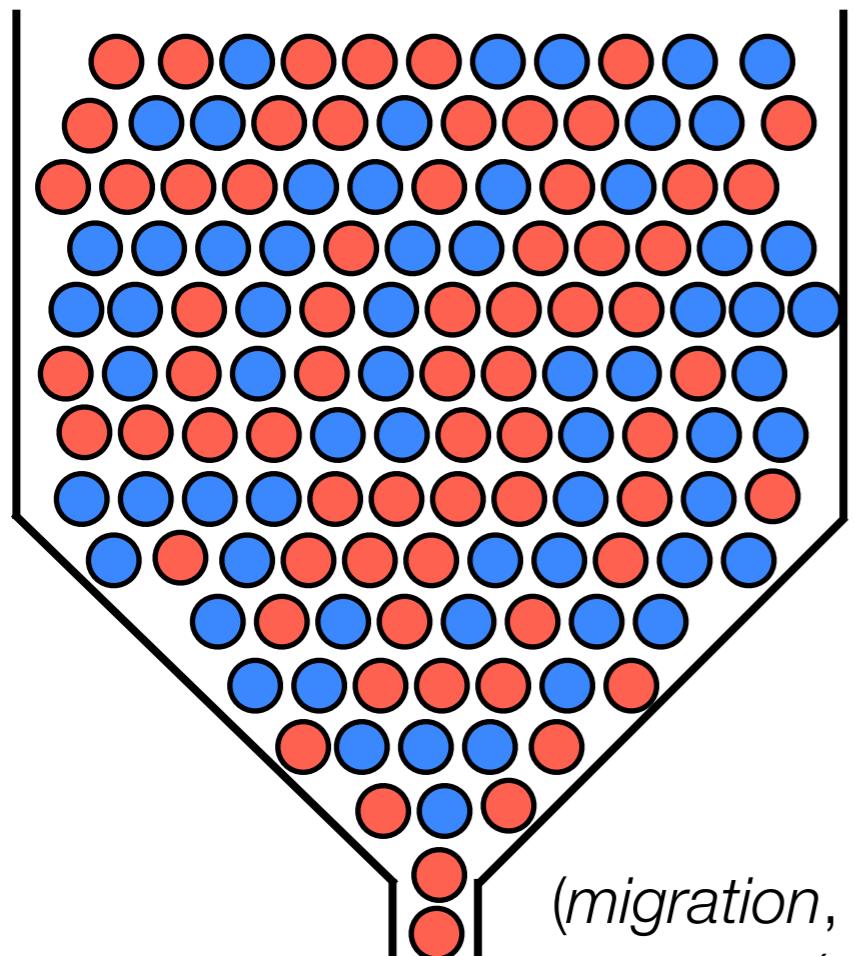
Genetic bottleneck

“Mixed” population



Genetic bottleneck

“Mixed” population

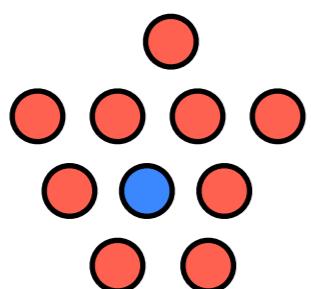


Bottleneck event

(migration, war, disease, nearly extinct animals)

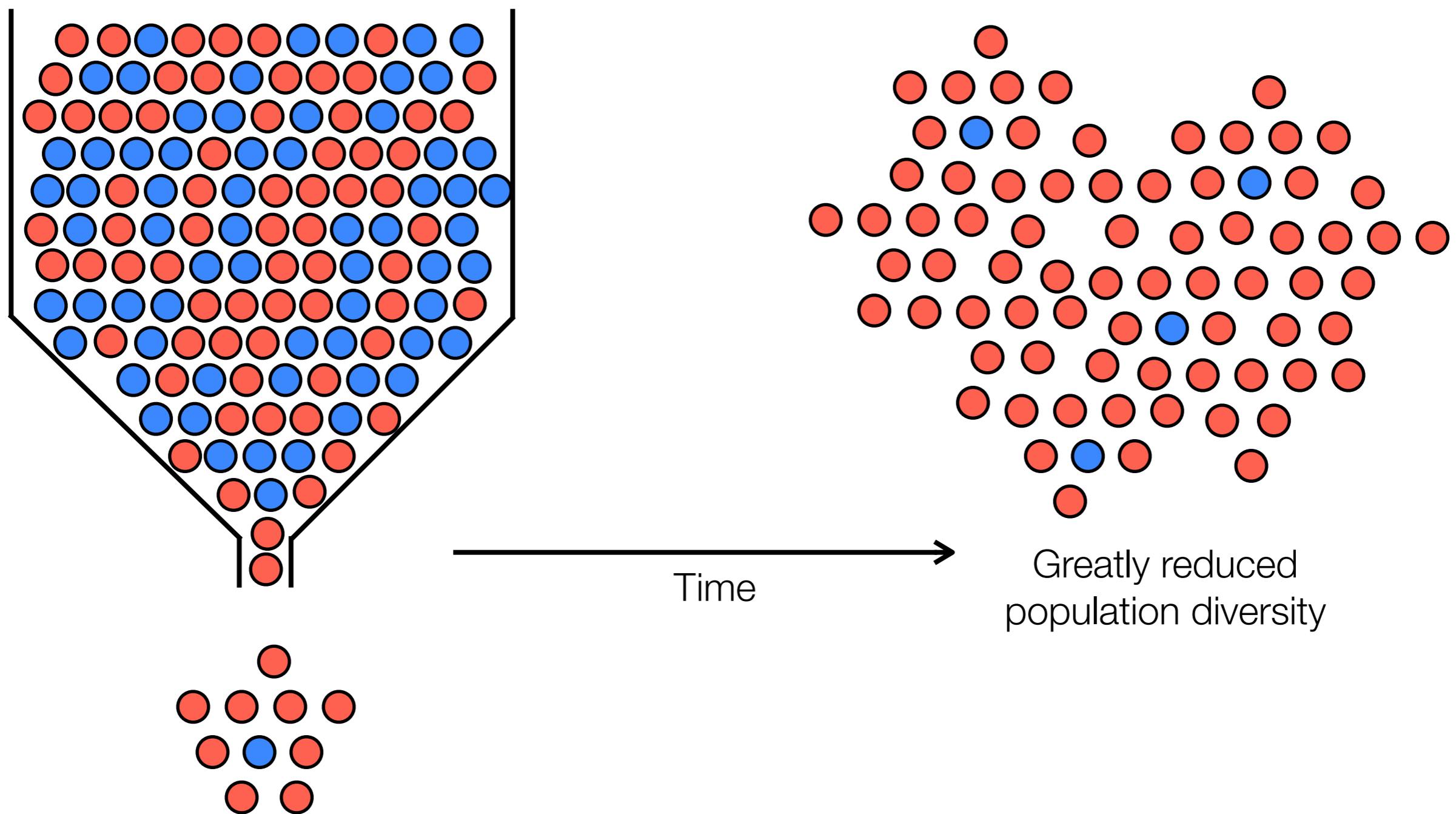
(e.g., elephant seals in Baja)

greatly reduces population size



Genetic bottleneck

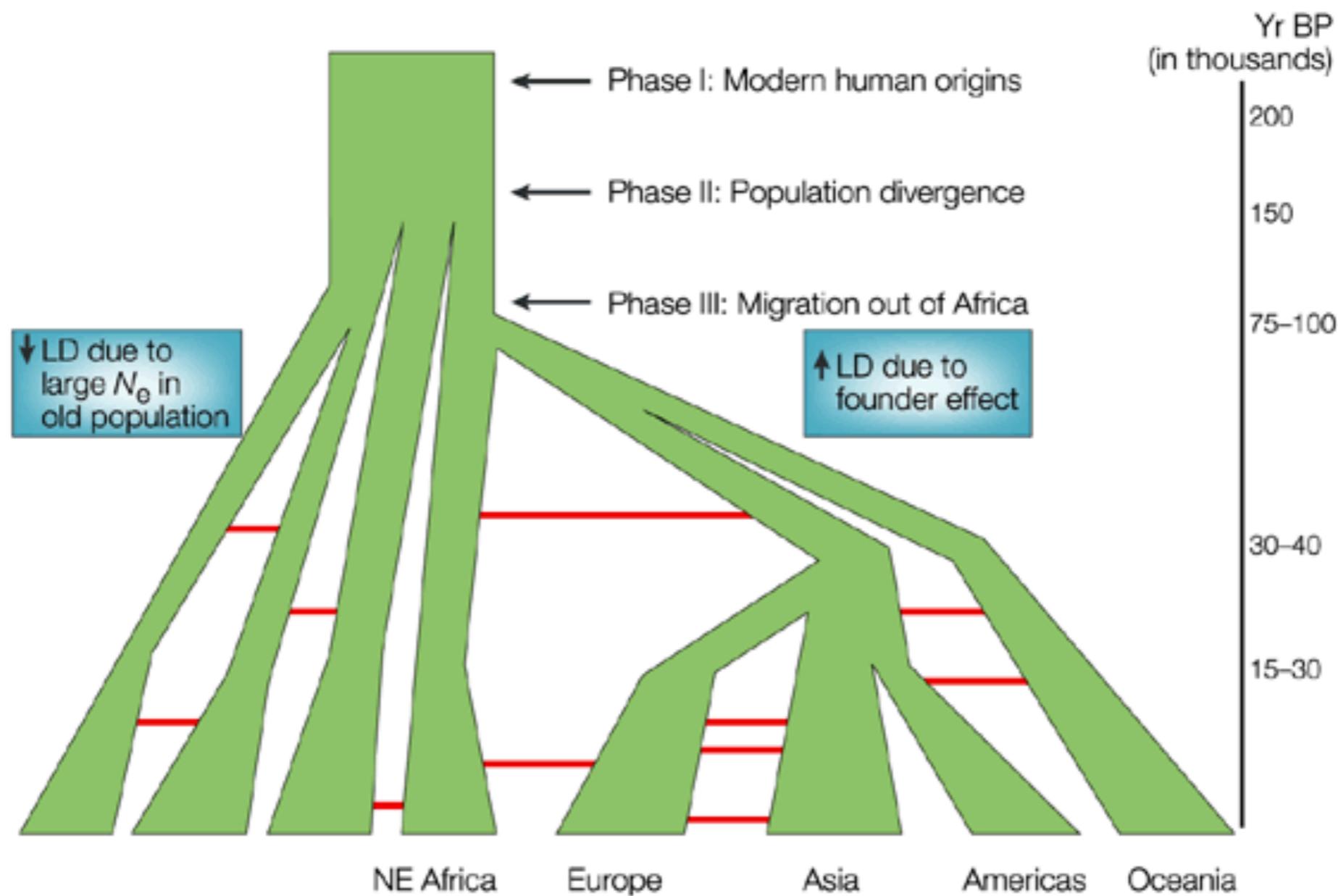
“Mixed” population



Greatly reduced
population diversity

Out of Africa Theory

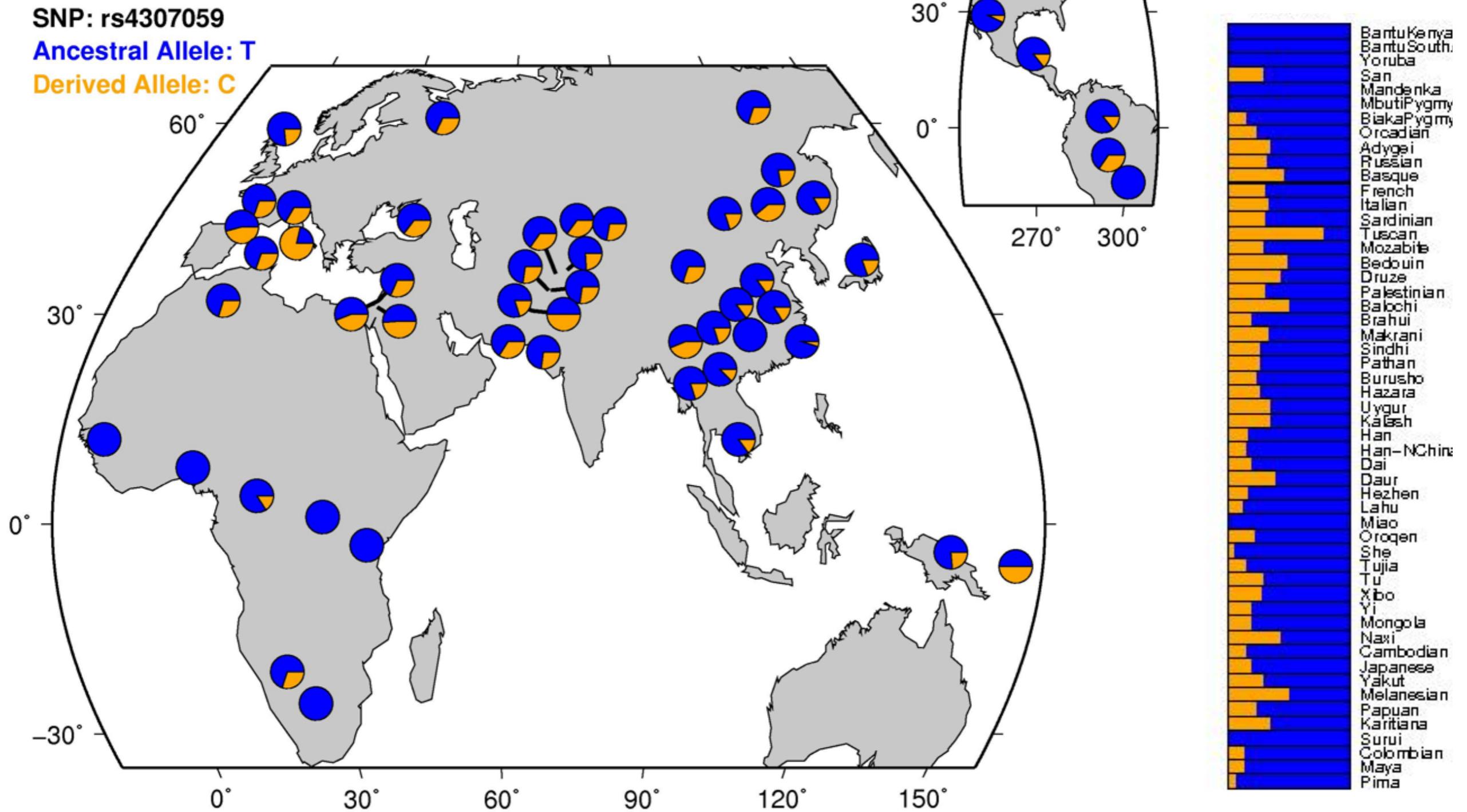
bottlenecks & reduced diversity



A founder effect.



Population stratification



What do allele frequencies tell us?

- Sequencing many individuals allows us to look at genetic variation and their frequencies at a moment in time.
- But, **BUT!**, the frequency of a variant is in the population is just a snapshot in time
 - e.g., a frame in a movie of selection, genetic drift.
- For example, a variant may be **rare** b/c it just occurred, b/c of purifying selection, or b/c of genetic drift.
- Or, a variant may be **common** because it is beneficial or because of genetic drift.

A map of human genetic variation



Sequencing is dirt cheap.
Let's catalog ~all genetic variation!

ARTICLE

doi:10.1038/nature09534

A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium*

The 1000 Genomes Project Consortium*

population-scale sequencing

More on this in the afternoon...

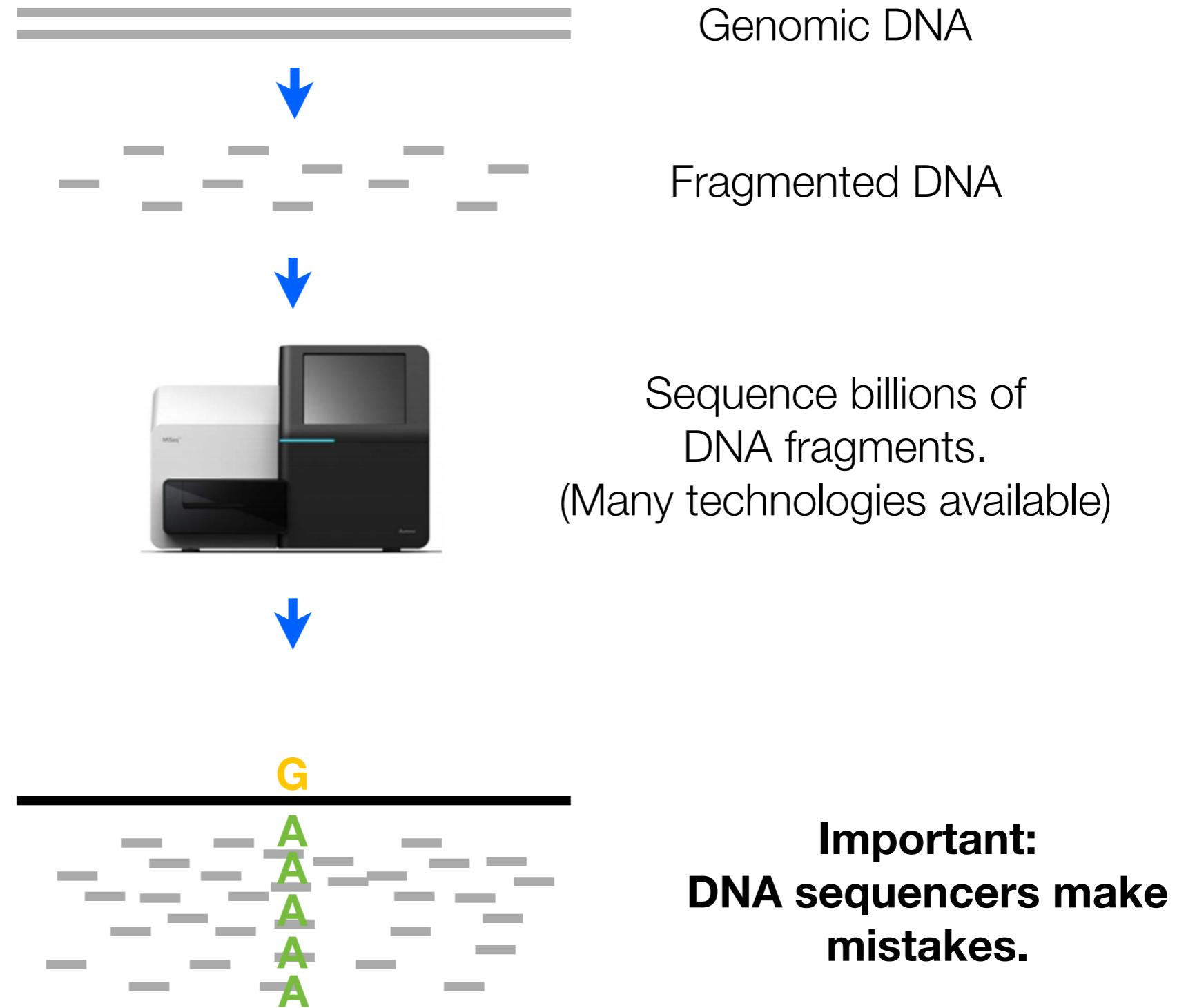
Detecting genetic variation

CGCAAATTGCCGGATTCCTTGCCTGCATGTAGTTAACGAGATTGCCAGCACCGGGTATCATTACCACTTTCTTTCTTAACCTGCCGTCA
GCCCTTTCTTGACCTCTTCTTCTGTTCATGTGTATTGCTGTCTCTAGCCAGACTTCCCCTGCCAGGTCTGACTTCCAGCAACTGCTGGCCTGTGCCAGG
GTGCAAGCTGAGCACTGGAGTGGAGTTCTGTGGAGAGGGAGCCATGCCTAGAGTGGATGGCCATTGTTCATCTTCTGGCCCCCTG
TTGTCTGCATGTAACCTAACCAACCAGGCATAAGGGAAAGATTGGAGGAAAGATGAGTGAGAGCATCAACTCTCACAACCT
AGGCCAGTAAGTAGTGCTTGCTCATCTCCTGGCTGTGATACGTGGCCGCCCTCGCTCCAGCAGCTGGACCCCTACCTGCCGTCT
GCTGCCATCGGAGCCAAAGCCGGCTGTGACTGCTCAGACCAGCCGGCTGGAGGGAGGGCTCAGCAGGTCTGGCTTGGCCCTGG
AGAGCAGGTGGAAGATCAGGCAGGCCATCGCTGCCACAGAACCCAGTGGATTGCCCTAGGTGGATCTCTGAGCTCAACAAGCCCTCT
CTGGGTGGTAGGTGCAGAGACGGGAGGGCAGAGCCGAGGCACAGCCAAGAGGGCTGAAGAAATGGTAGAACGGAGCAGCTGGTGT
GTGTGGGCCACCAGGCCAGGCTCCTGTCTCCCCCAGGTGTGGATGCCAGGCATGCCCTCCCCAGCATCAGGTCTCCAGAG
CTGCAGAAGACGACGCCACTGGATCACACTCTTGTGAGTGCTCCAGTGTGCAGAGGTGAGAGGAGAGTAGACAGTGAGTGGGA
GTGGCGTCGCCCTAGGGCTCTACGGGGCCGGCTCCTGTCTCCTGGAGAGGGCTTCGATGCCCTCCACACCCTTTGATCTTCCC
TGTGATGTCATCTGGAGCCCTGCTGCTTGCCTGGCCTATAAACGCTCCTAGTCTGGCTCCAAGGCCTGGCAGAGTCTTCCCAGGG
AAGCTACAAGCAGCAAACAGTCTGCATGGTCATCCCCCTCACTCCCAGCTCAGAGCCCAGGCCAGGGCCCCAAGAAAGGCTCTGG
TGGAGAACCTGTGCATGAAGGCTGTCAACCAGTCATAGGCAAGCCTGGCTGCCTCCAGCTGGTCAGACAGACAGGGCTGGAGAAGG
GGAGAAGAGGAAAGTGAGGTTGCCCTGCTCCTACCTGAGGCTGAGGAAGGAGAAGGGATGCACTGTTGGGAGGCAGCTGTA
ACTCAAAGCCTCTGCTTCCACGAAGGCAGGCCATCAGGCACCAAGGGATTCTGCCAGCATAGTGCTCCTGGACCAGTGAT
ACACCCGGACCCTGCTGGACACGCTGTTGGCTGGATCTGAGCCCTGGAGGTCAAAGCCACCTTGGTCTGCCATTGCTGC
TGTGTGGAAGTTCACTCCTGCCCTTCCCTAGAGCCTCCACCACCCGAGATCACATTCTCACTGCCCTTGTCTGCCAGT
TTCACCAGAAGTAGGCCTTCCCTGACAGGCAGCTGCACCACTGCCTGGCGCTGTGCCCTTGTCTGCCAGTGGAGACGGTG
TTTGTGATGGCCTGGCTGCAGGGATCCTGCTACAAAGGTGAAACCCAGGAGAGTGAGTGGAGTCCAGAGTGTTGCCAGGACCCAGGCA
CAGGCATTAGTGCCGTTGGAGAAAACAGGGGAATCCGAAGAAATGGTGGCTGGCCATCCGTGAGATCTTCCCAGGTG
TTTCTCTGGAAGCCTTTAAGAACACAGTGGCGCAGGCTGGTGGAGCCCTGGAGACAGGCAAGACAGAAGTCCCCGCC
CAGCTGTGGCCTCAAGCCAGCCTCCGCTCCTGAAGCTGGTCTCCACACAGTGCTGGTCCACCCCTCCAAAGGAAGTAGG
TCTGAGCAGCTTGTCTGGCTGTCCATGTCAGAGCAACGGCCAAGTCTGGTCTGGGGGGAGGTGTCATGGAGCCCTACGA
TTCCCAAGTCGTCCTCGCCTCTGCCTGTGGCTGCGGTGGCGCAGAGGAGGGATGGAGTCTGACACGCCAAAGGCTCCT
CCGGGCCCTCACAGCCCCAGGTCTTCCCAGAGATGCCCTGGAGGGAAAAGGCTGAGTGAGGGTGGTGGAAACCCCTGG
CCCCAGCCCCGGAGACTAAATACAGGAAGAAAAGGCAGGACAGAATTACAAGGTGCTGGCCAGGGCGGGCAGCGCCCTGCC
CTACCCCTGCGCCTCATGACCGGAGCCATAGCCCAGGCAGGGAGGGCTGAGGACCTCTGGTGGCGGCCAGGGCTCCAGCATGTGCC
TAGGGGAAGCAGGGGCCAGCTGGCAAGAGCAGGGGGTGGCAGAAAGCACCCGGTGGACTCAGGGCTGGAGGGAGGAGGCGATCTTG
CCCAAGGCCCTCCGACTGCAAGCTCCAGGGCCGCTCACCTGCTCCTGCTCCTGCTGCTGCTTCTCCAGCTTGCCTCC
GCTGCGCAGCTTGGCCTGCGATGCCCTAGCTGGCGATGGACTCTAGCAGAGTGCCAGCCACCAGGGCAACCACCTCCC

It all starts with DNA sequencing

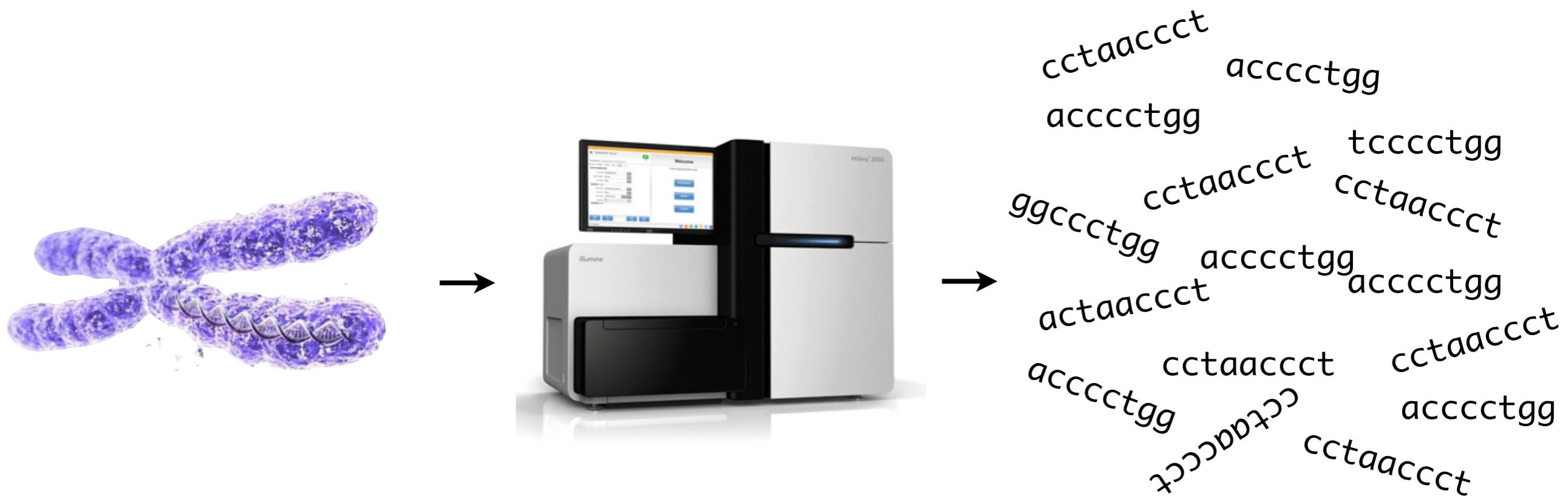
*Sequence alignment is non-trivial
and the strategy depends
upon the type of genetic variation
you are trying to detect*

Align DNA to a reference genome. Comparing sample DNA to reference reveals genetic differences.



Sequence mapping: where does the read belong in the genome?

- d. the genome is big and repetitive.
- e. DNA sequencing yields short sequences (<100 characters)

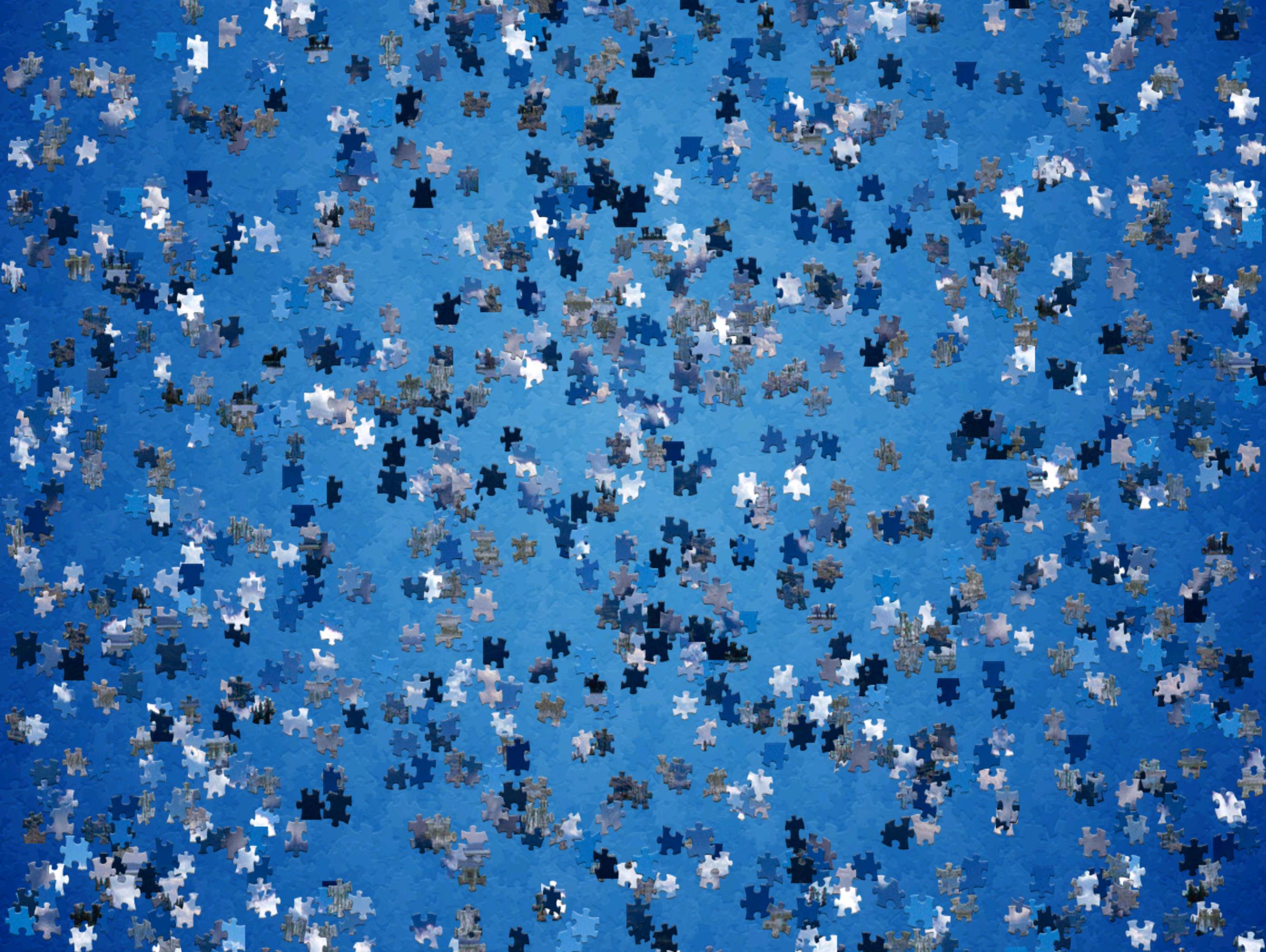


Where the small subsequences belong?

How do we “stitch” them together to make a complete genome?

“The reference genome”





Sequence alignment

Smith-Waterman, Needleman-Wunsch

Reference

cgggtatccaa

Read

ccctaggtcccc

What is the best alignment?

Reference	cgggtatccaa
Read	ccctaggtcccc
Reference	cggta--t-ccaa
Read	ccc-taggcccc-a

Reference	cgggtatccaa
Read	ccctaggtcccc
Reference	cggta--t-ccaa
Read	ccc-taggcccc-a
Reference	cggta---tccaa
Read	cc---ctaggcccc

Reference cgggatatccaa

Read cccttaggtcccc

Reference cggta---t-ccaa

Read ccc-taggcccc-a

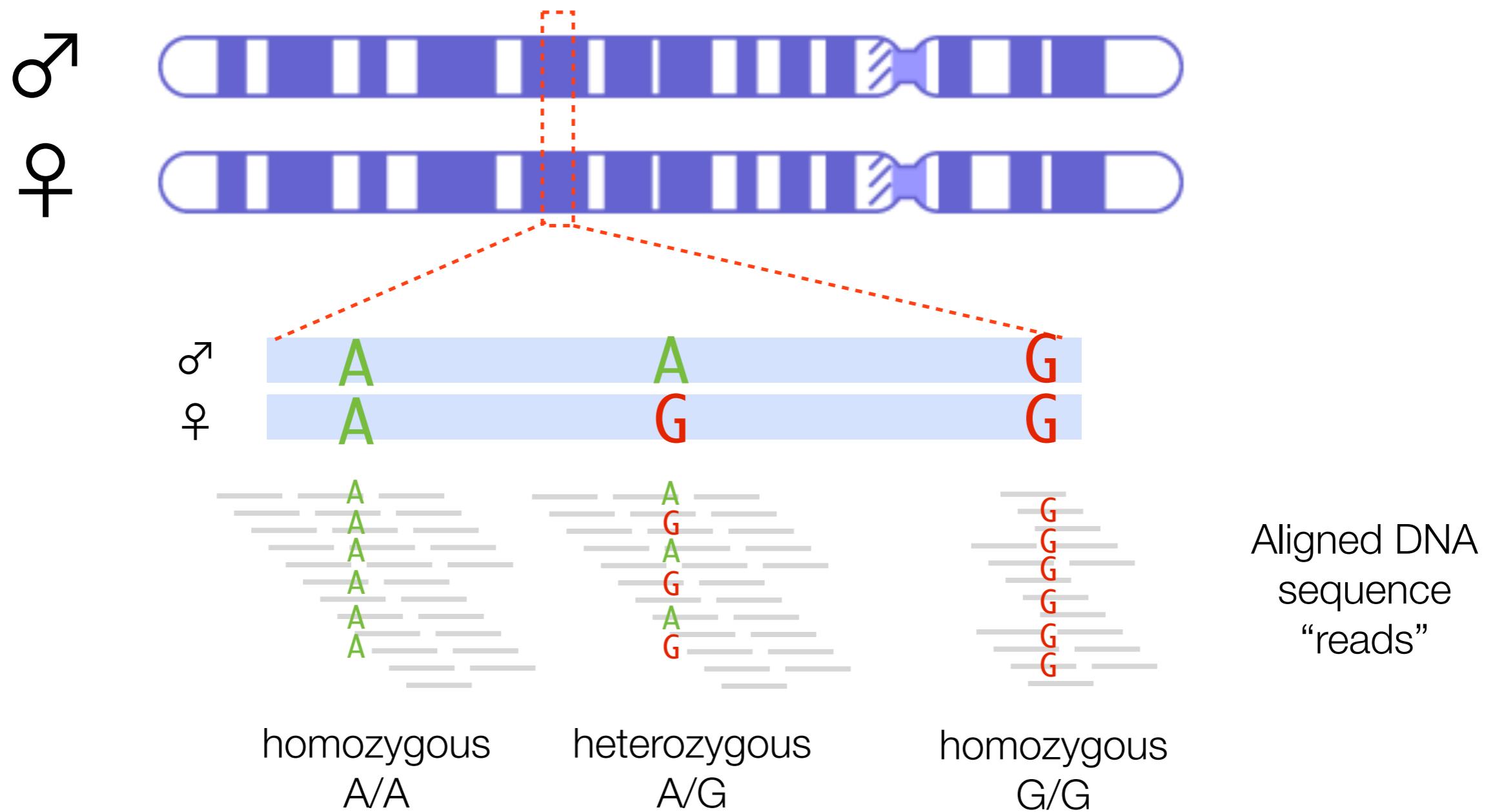
Reference cggta---tccaa

Read cc---ctaggcccc

Reference c-gggta---tccaa

Read cc---ctaggcccc

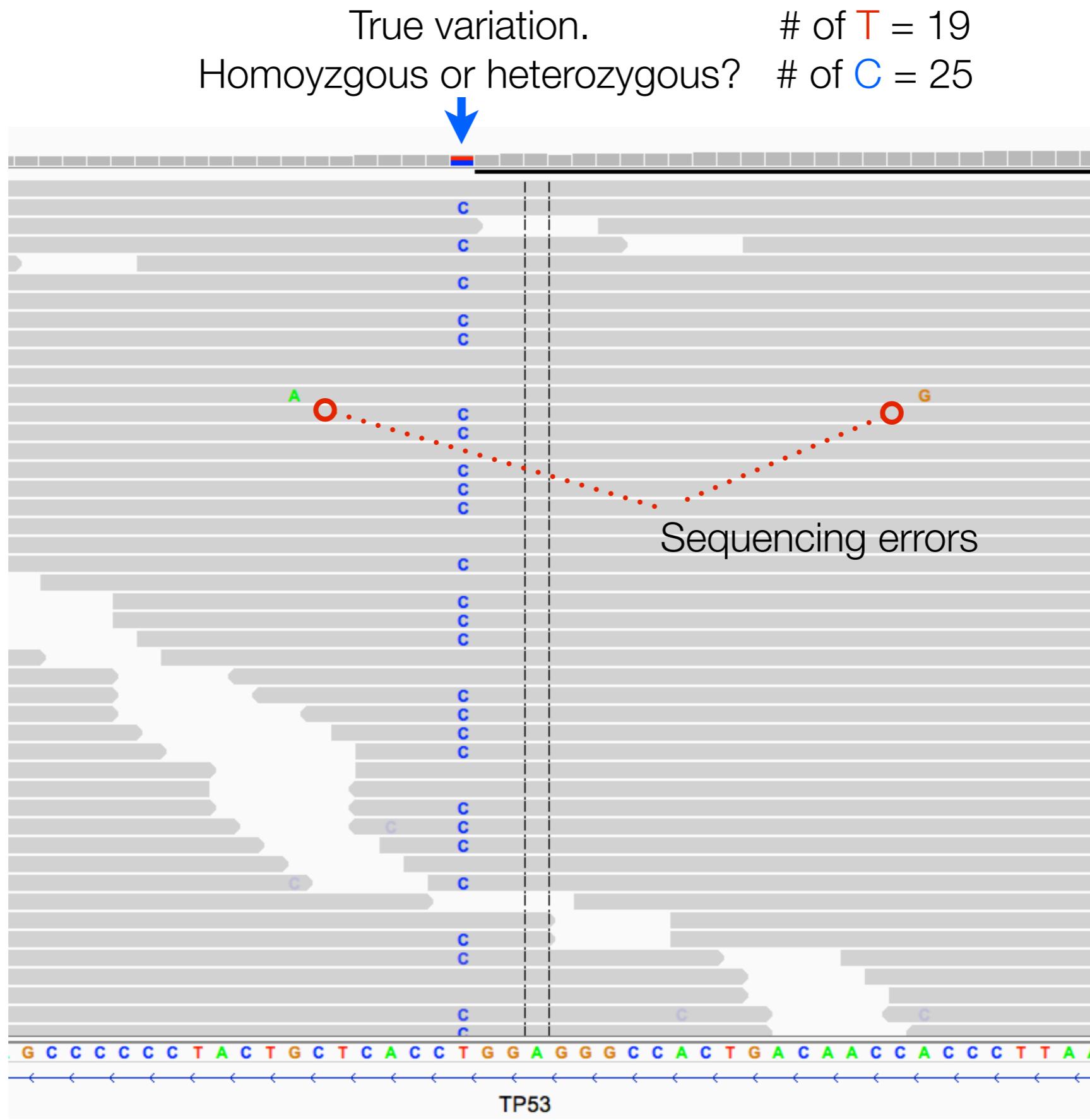
Recall: we are diploid.



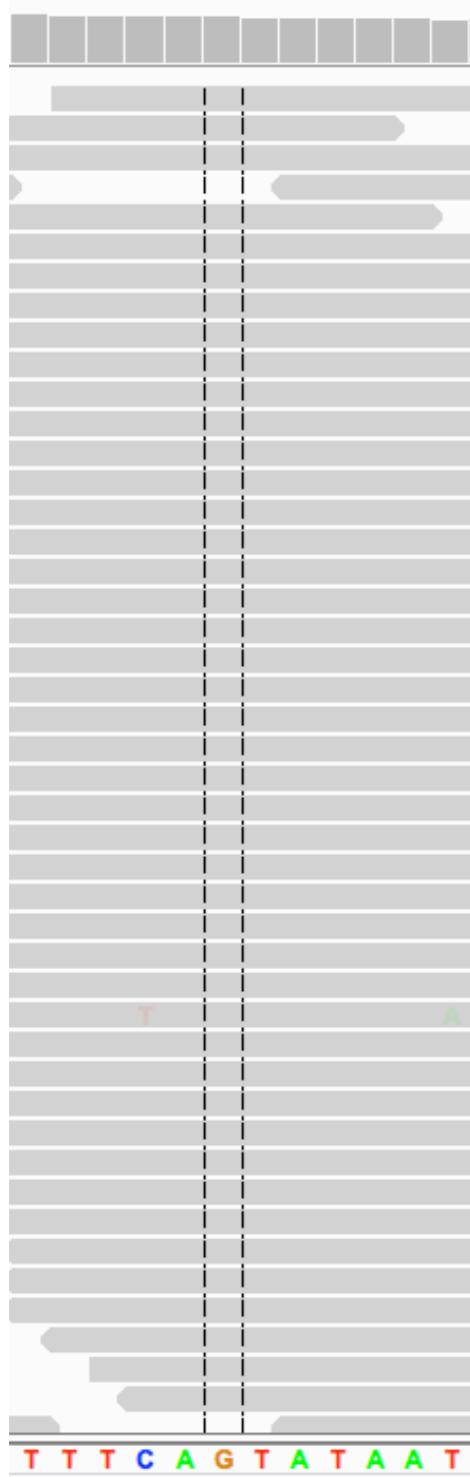
Single-nucleotide polymorphisms (SNPs)



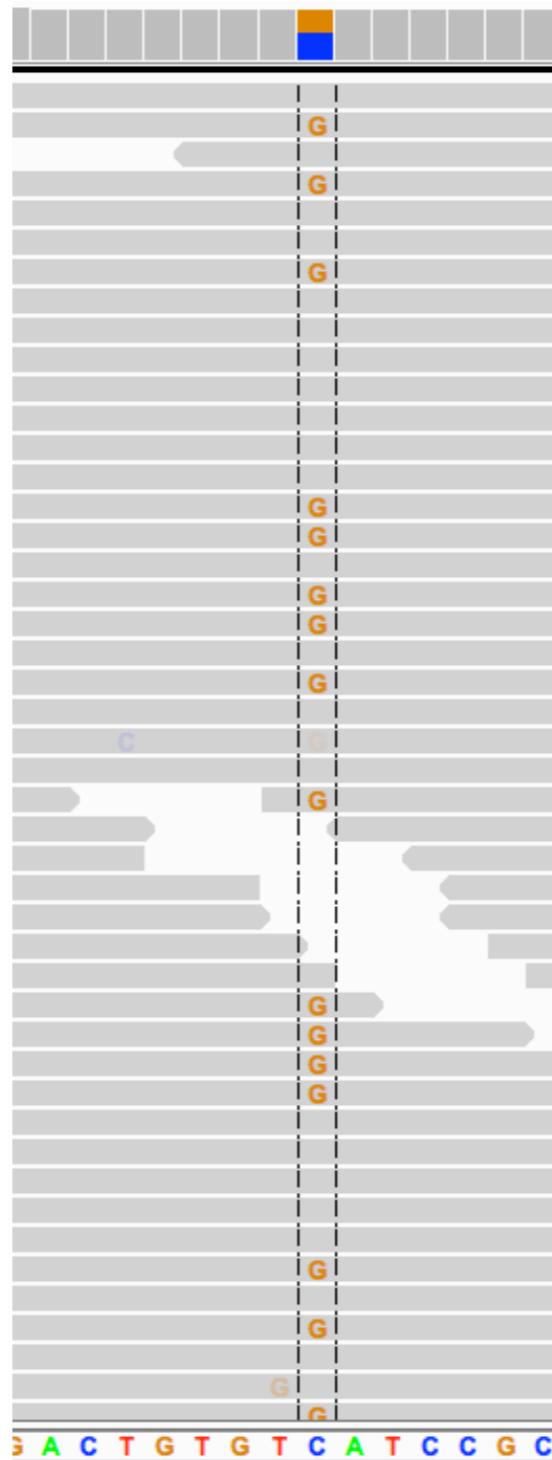
Single-nucleotide polymorphisms (SNPs)



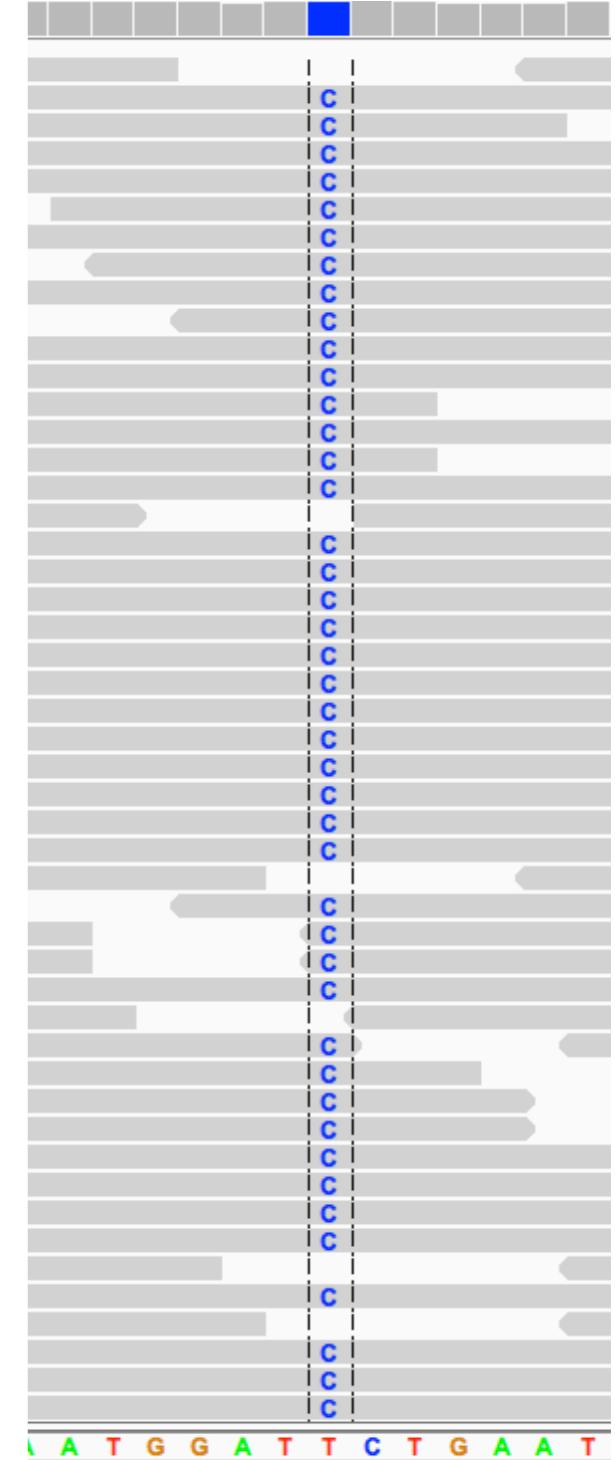
Different SNP genotypes



Homozygous for reference
(i.e., both chroms same as ref)

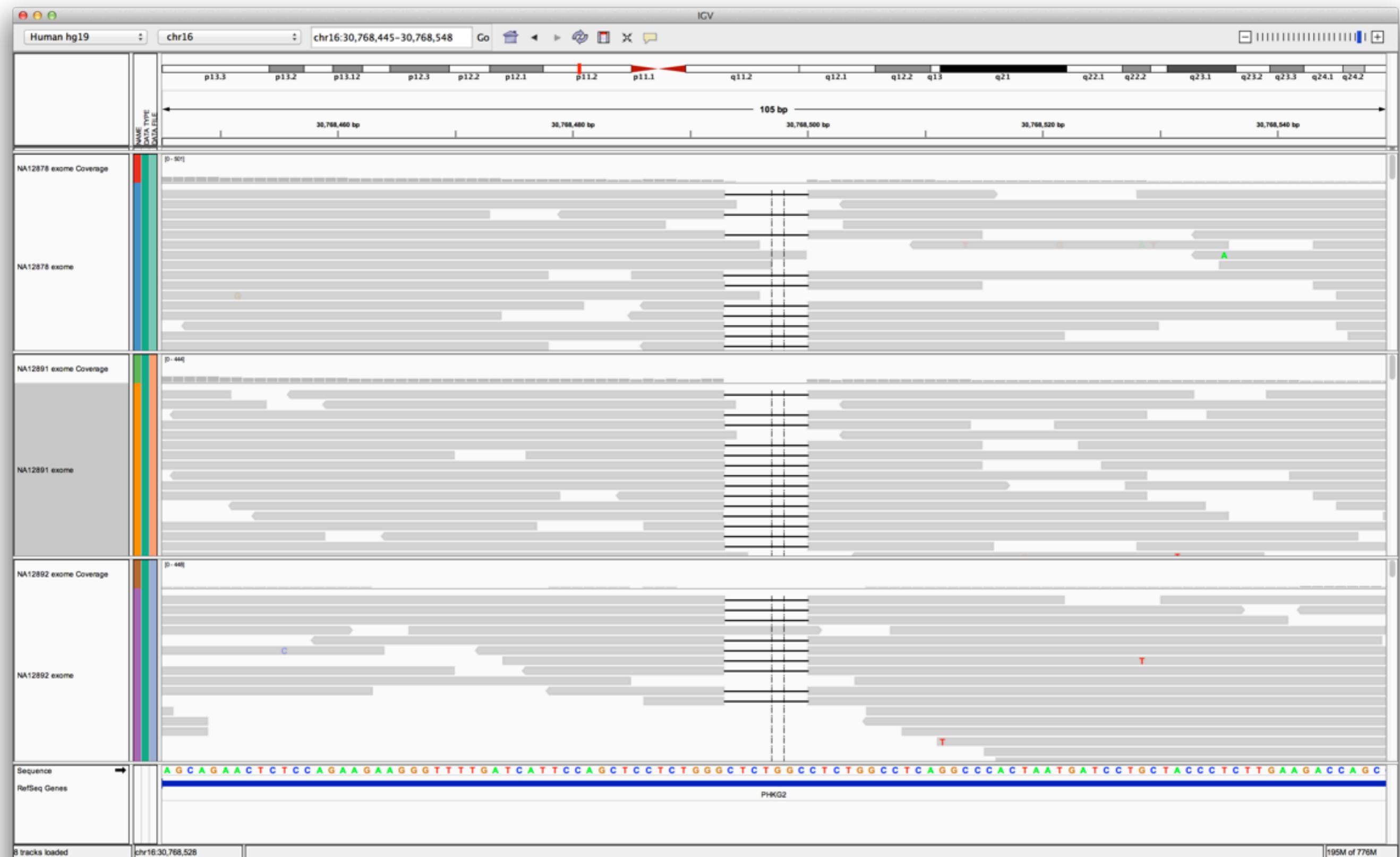


Heterozygous
(i.e., 1 chrom same as ref, 1 diff.)

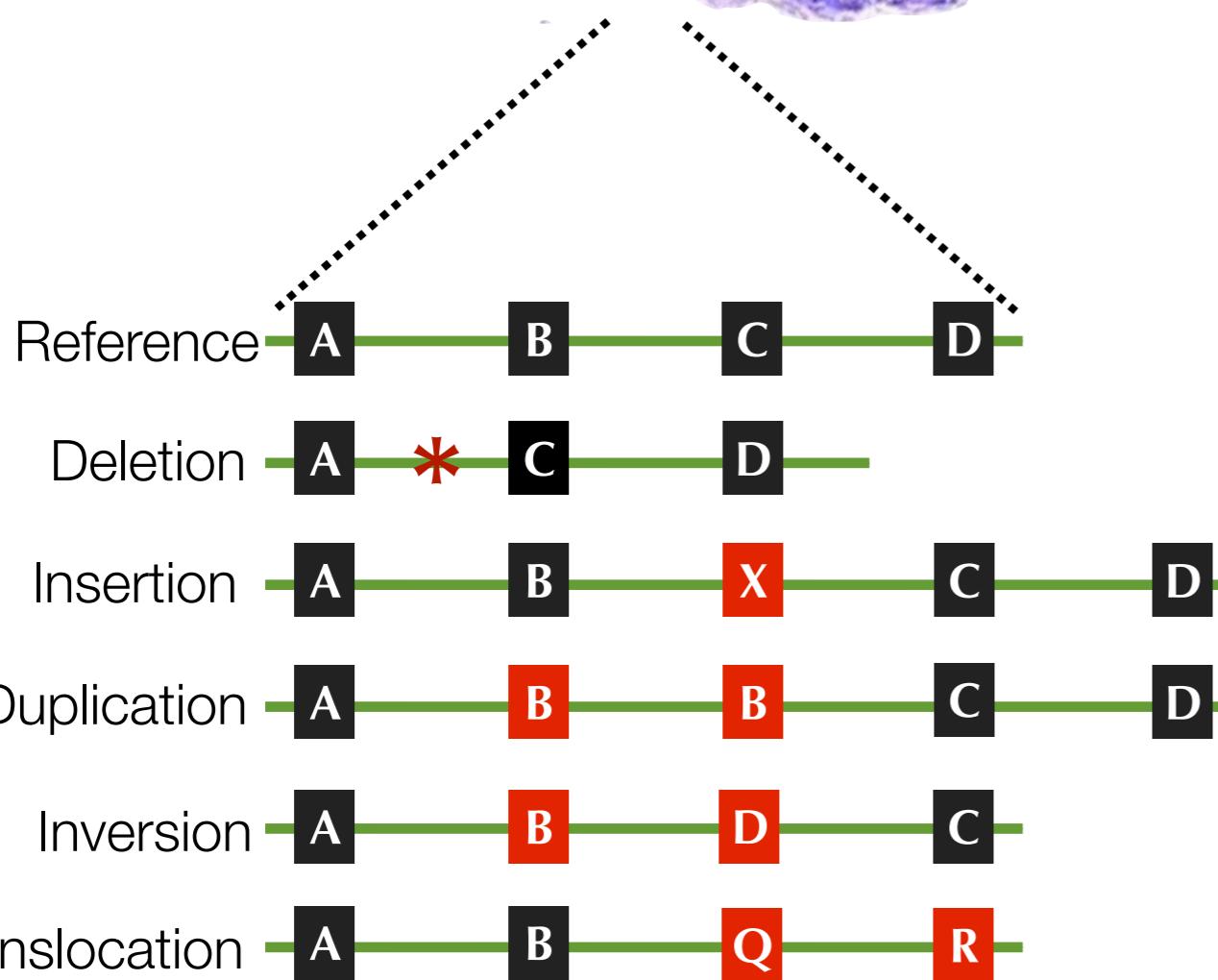
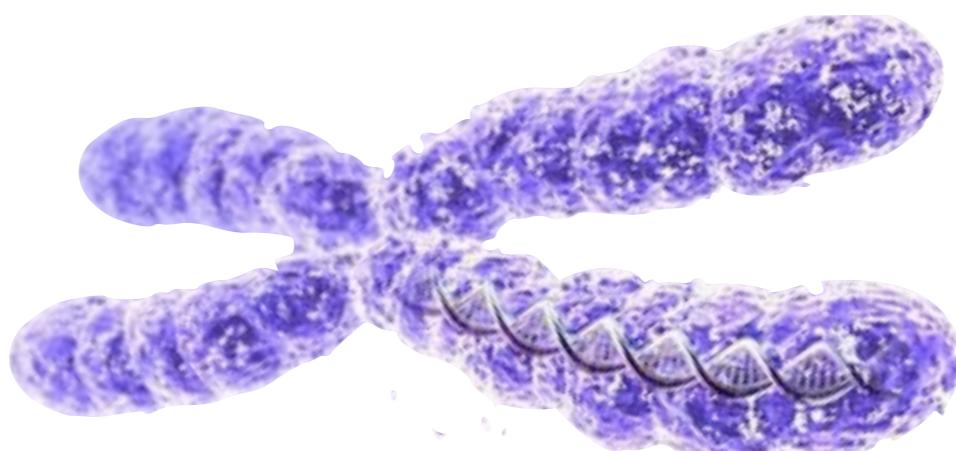


Homozygous for reference
(i.e., both chroms diff than ref)

Insertion-deletion polymorphisms (INDELS)

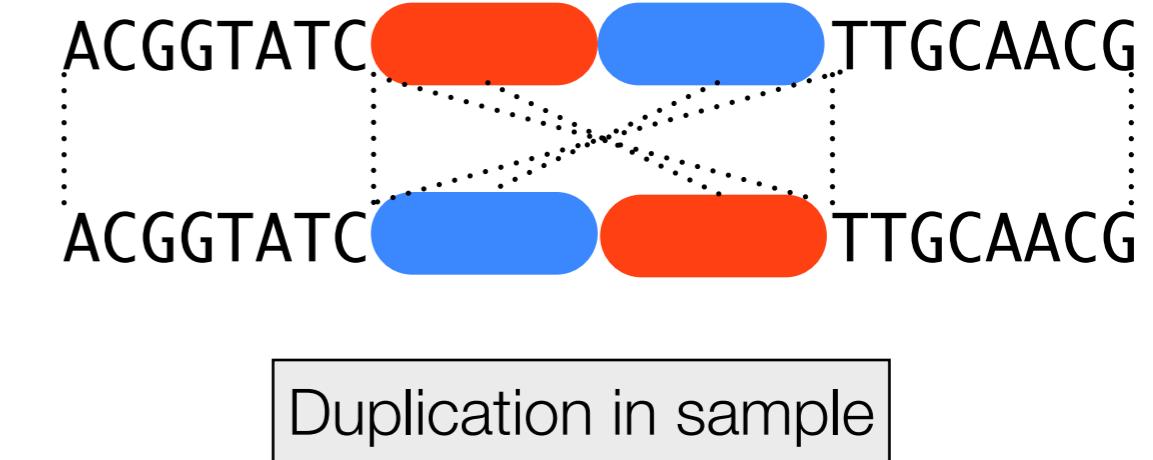
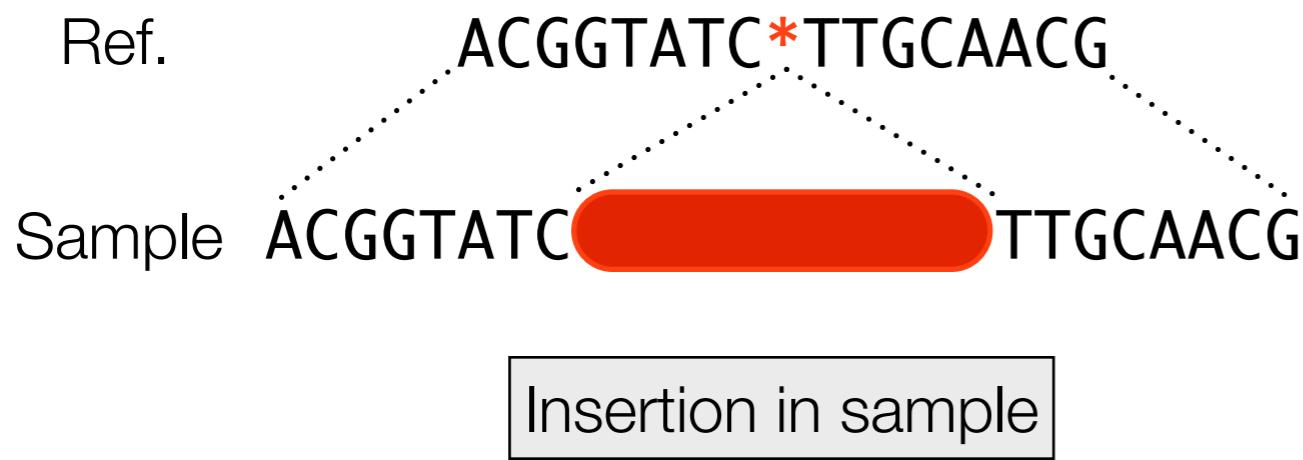
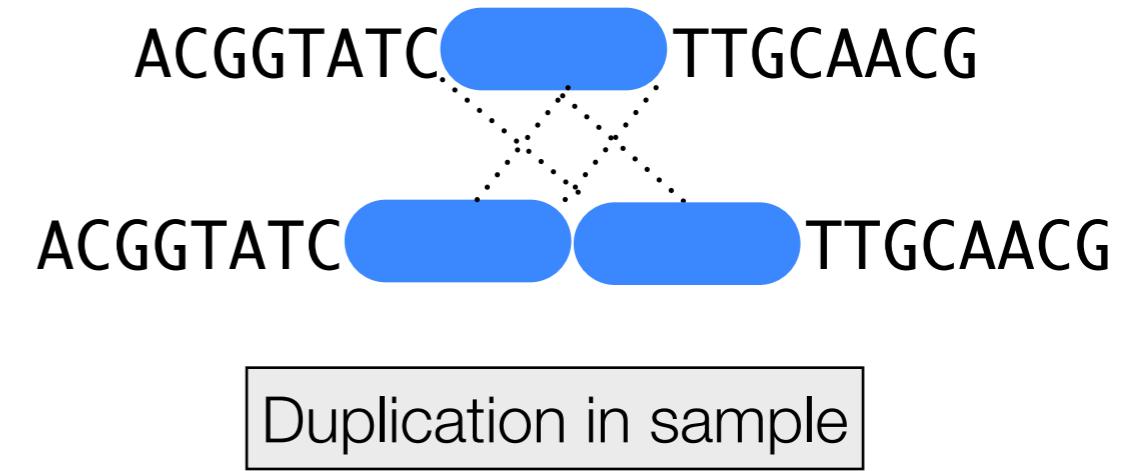
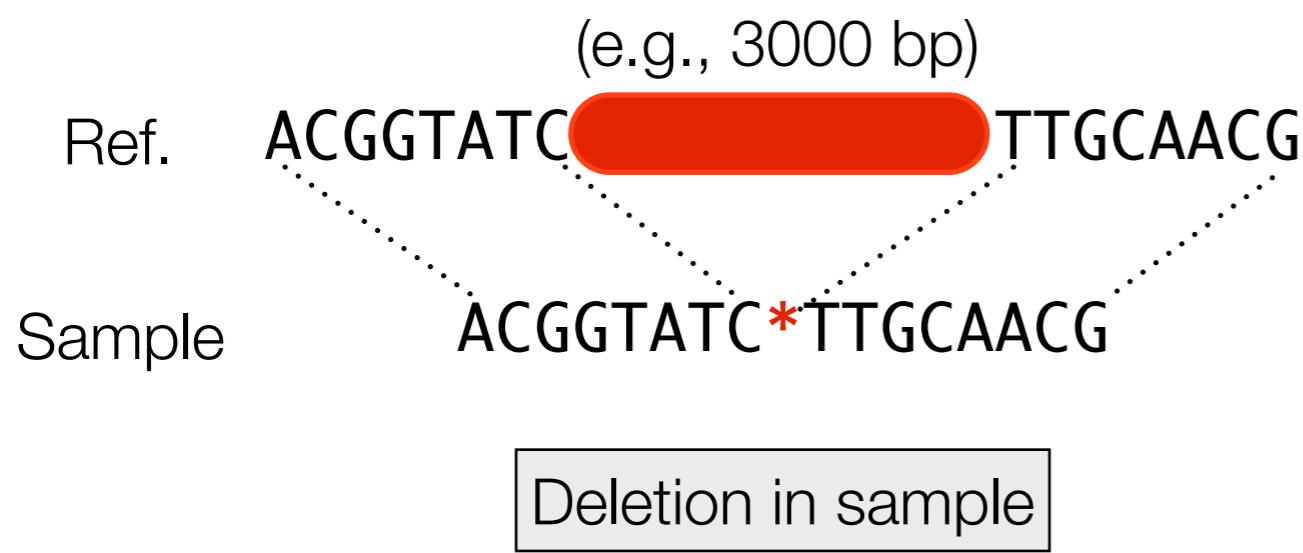


Structural Variation

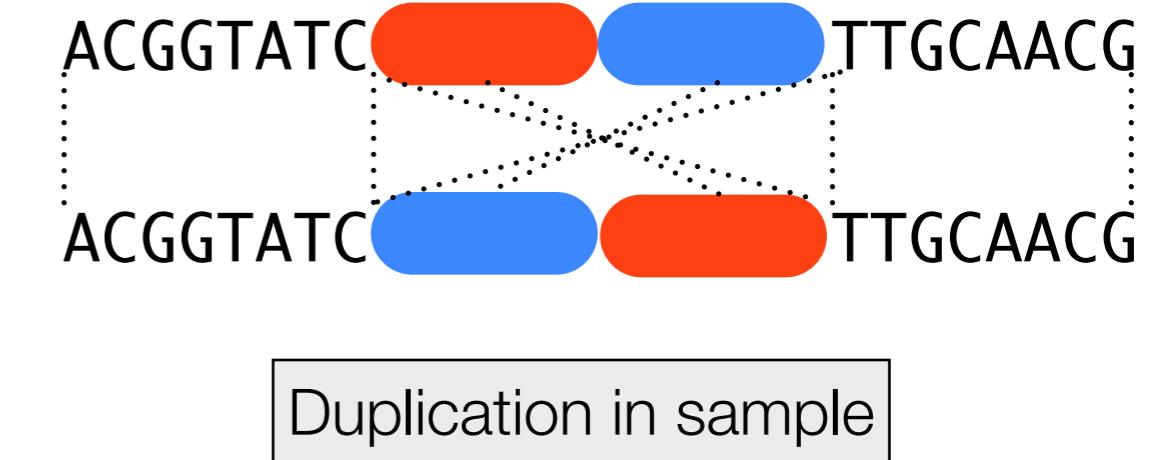
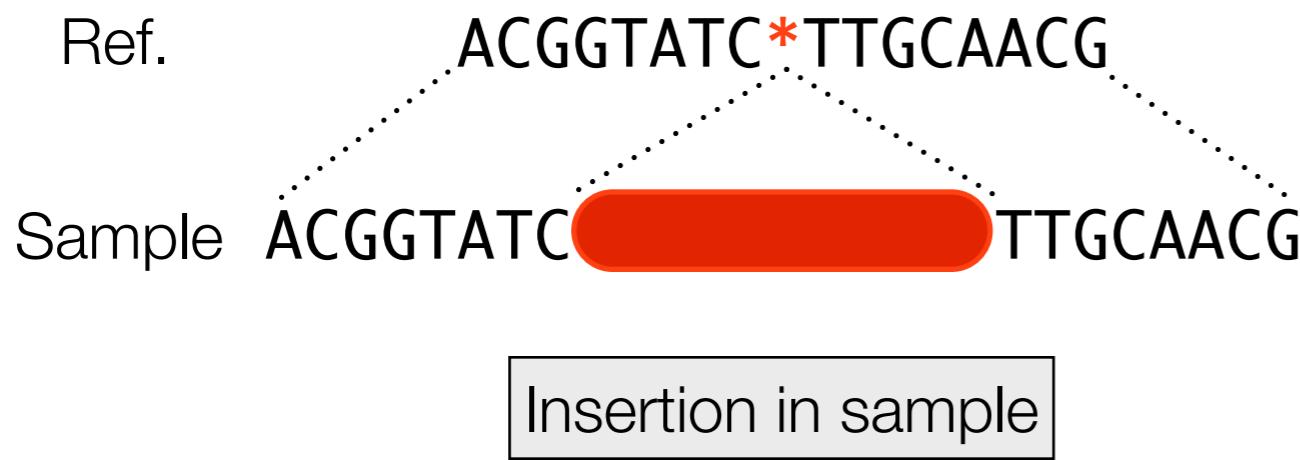
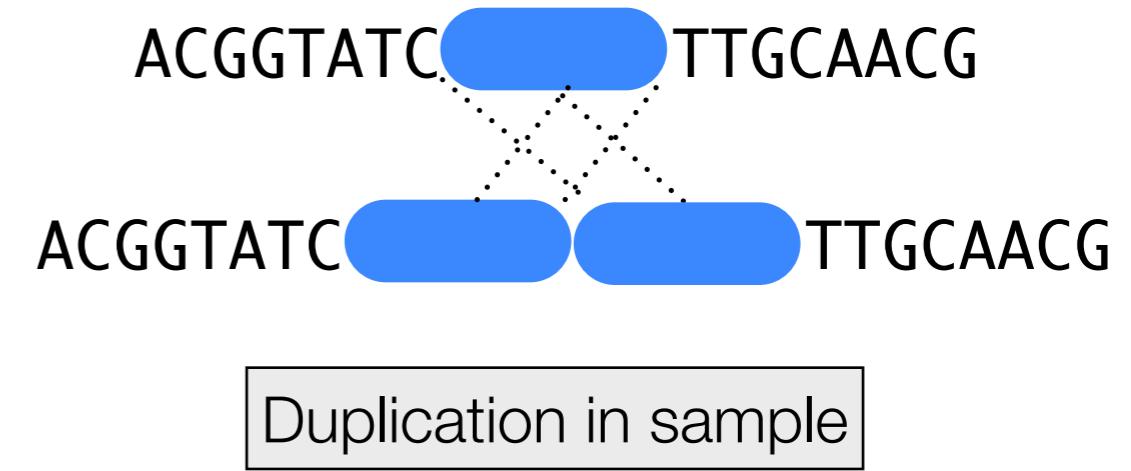
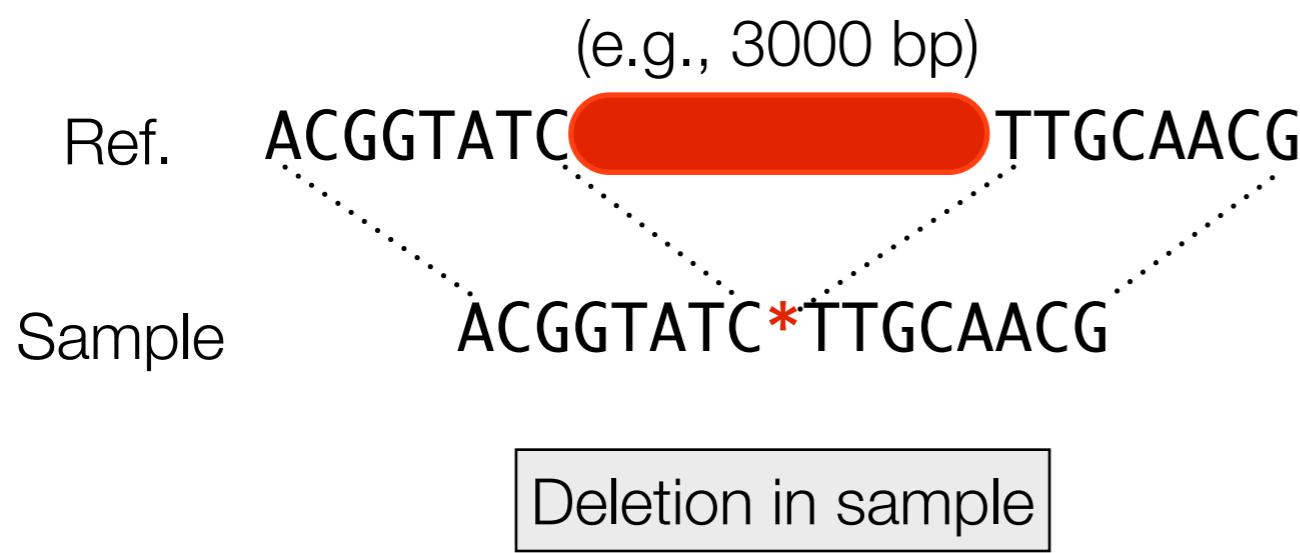


- large (>1kb) differences that affect the copy number, orientation, or location of genomic segments
- Common in mammalian genomes (thousands between two people)
- A hallmark of cancer
- A major cause of spontaneous disease
- more are functional than SNPs
- very challenging to identify

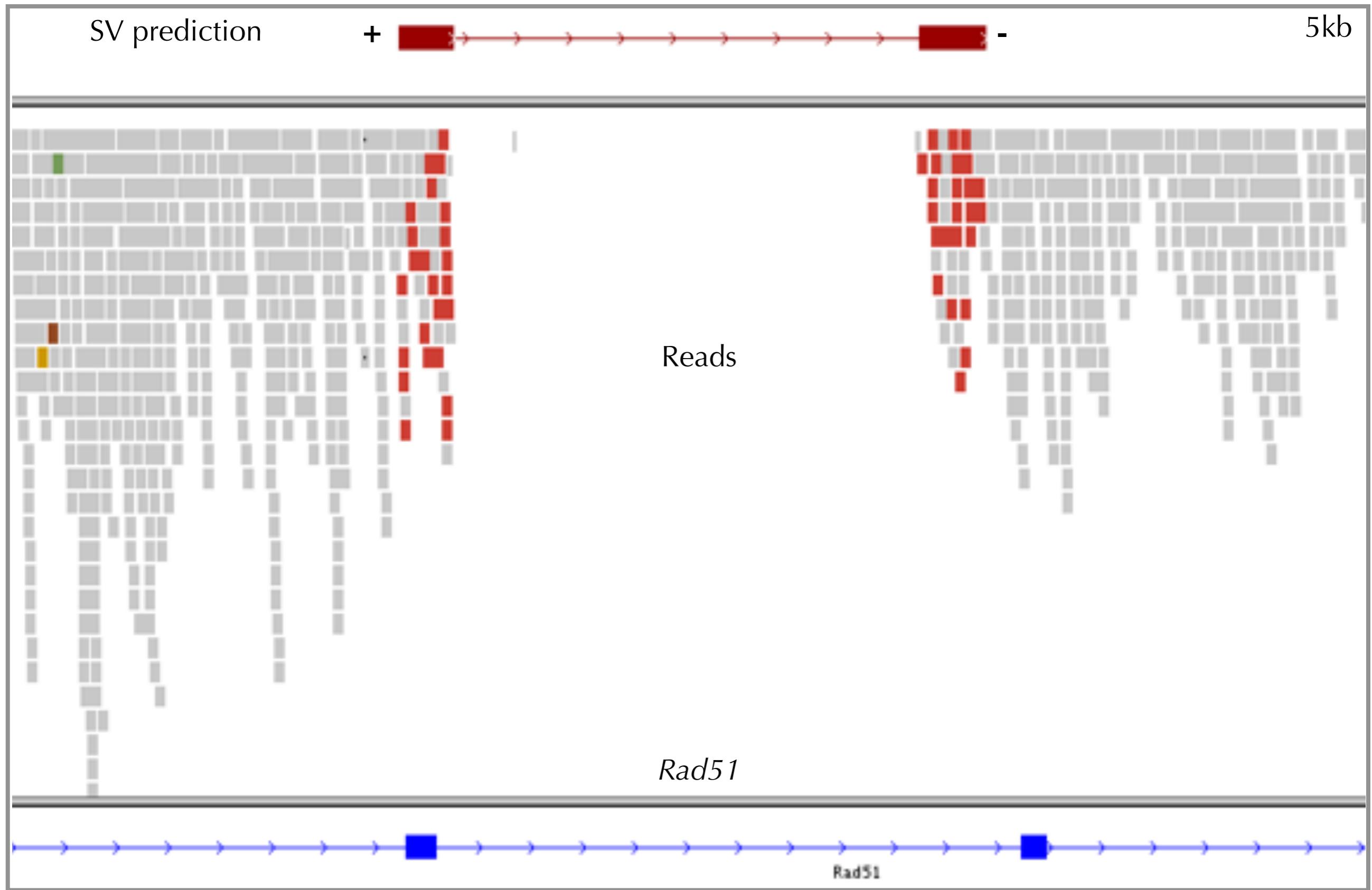
SV Alignment structures



SV patterns



A real deletion detected with sequencing



Functional consequences

Impact sometimes hard to predict.

synonymous (silent)

	L	Q	T
Normal	ctg	cag	act
Mutated	ctg	caa	act
	L	Q	T

non-synonymous (missense)

	L	Q	T
Normal	ctg	cag	act
Mutated	ctg	cgg	act
	L	R	T

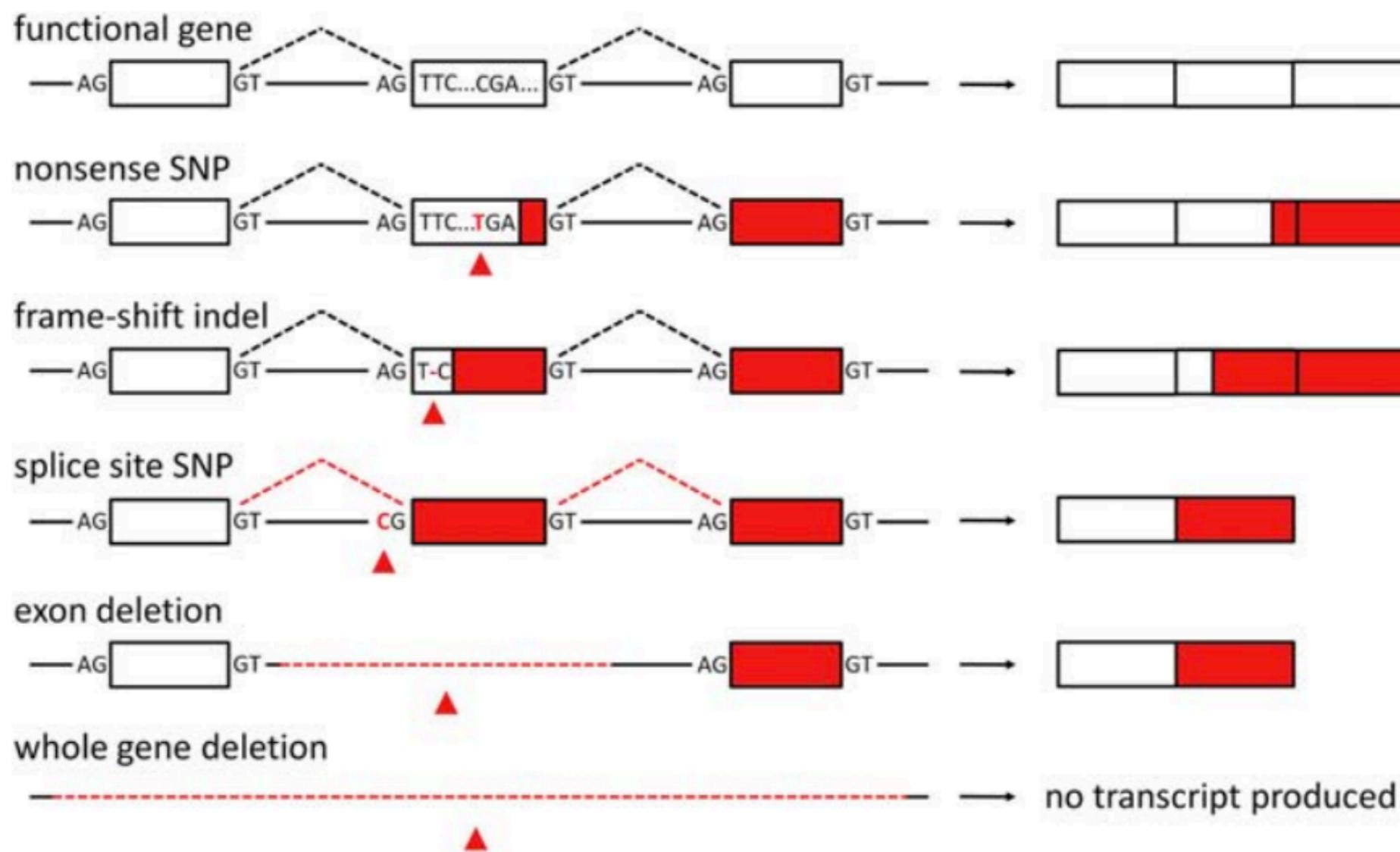
stop-gain (nonsense)

	L	Q	T
Normal	ctg	cag	act
Mutated	ctg	tag	act
	L	STOP	T

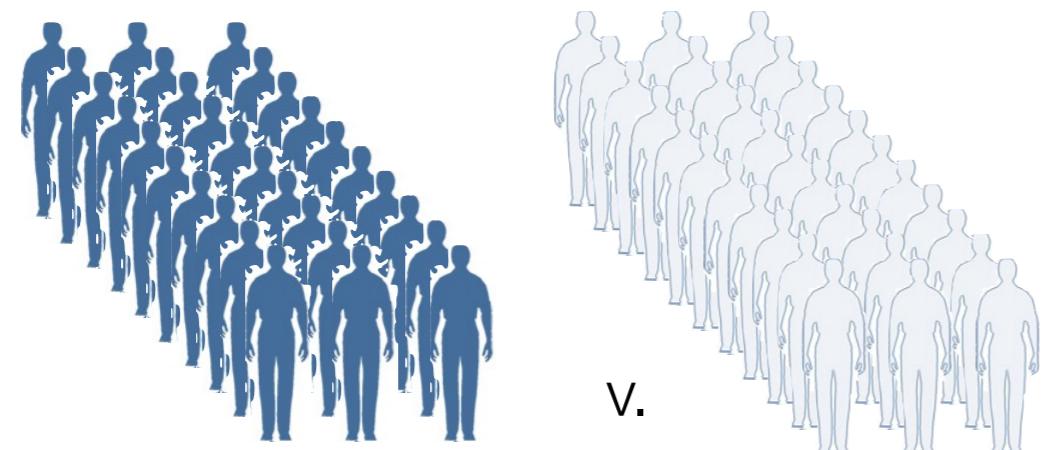
stop-loss

	L	STOP	T
Normal	ctg	tag	act
Mutated	ctg	cag	act
	L	Q	T

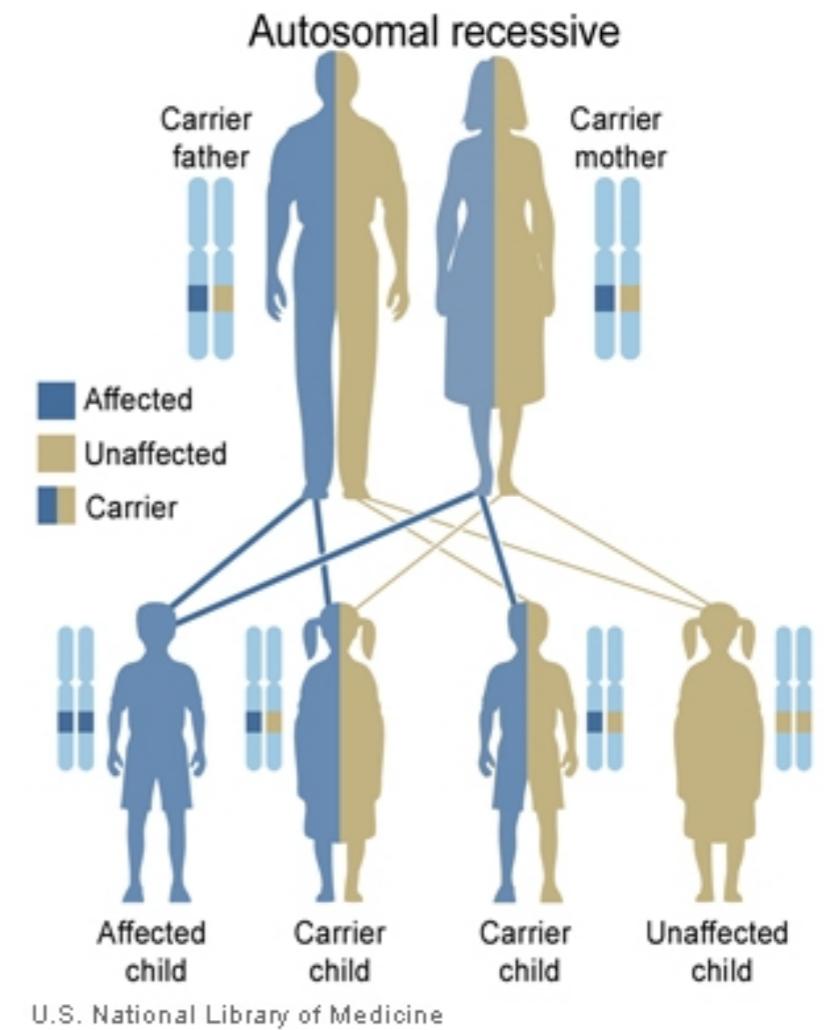
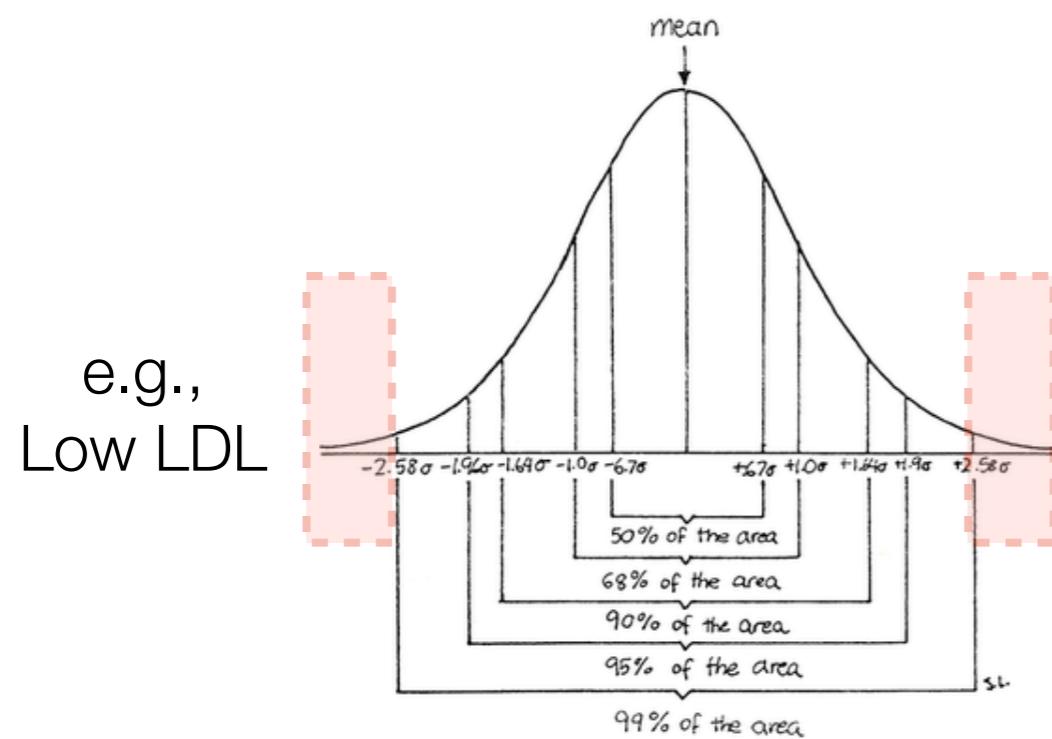
Loss of function mutations



Interpreting relevance to disease.



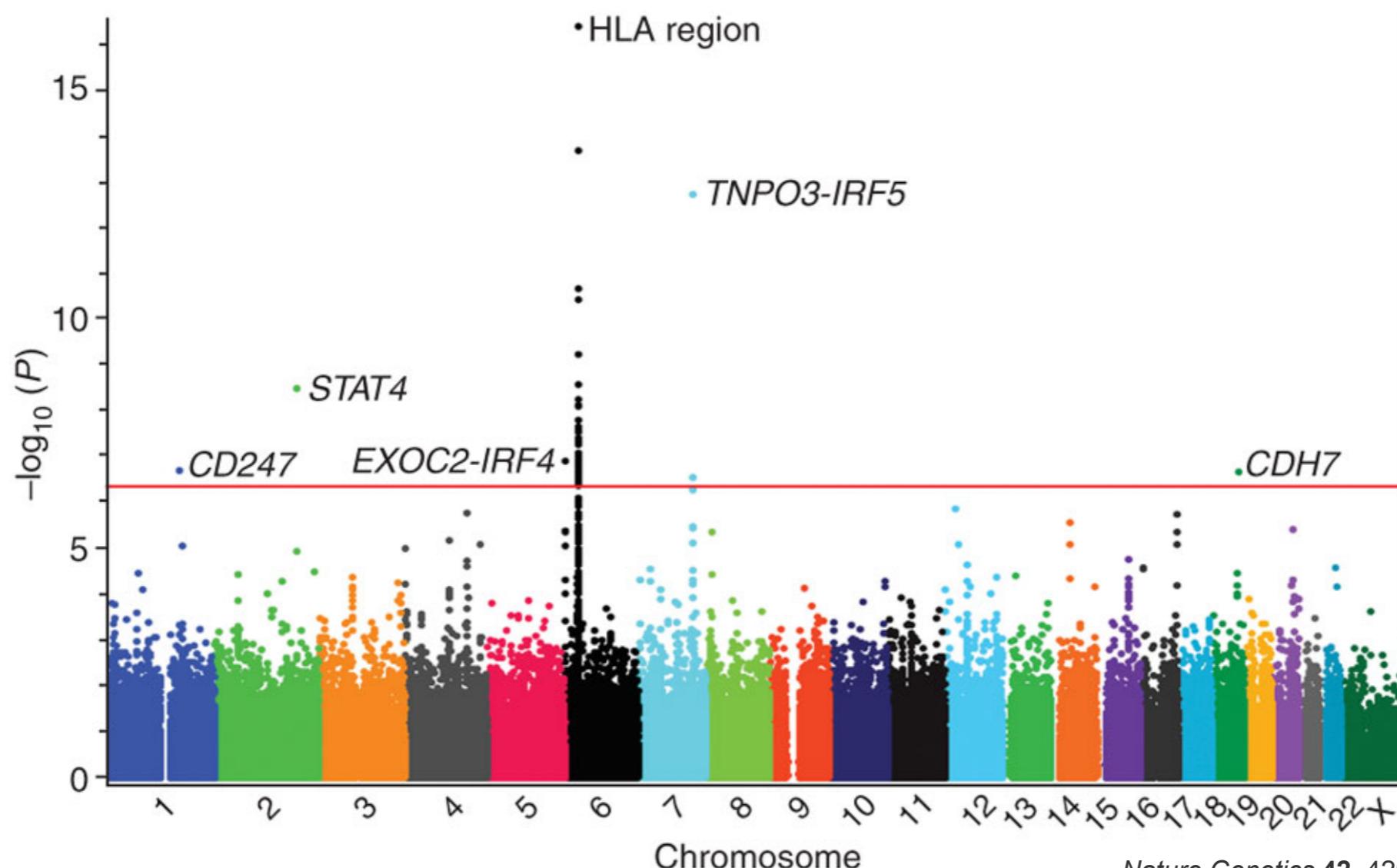
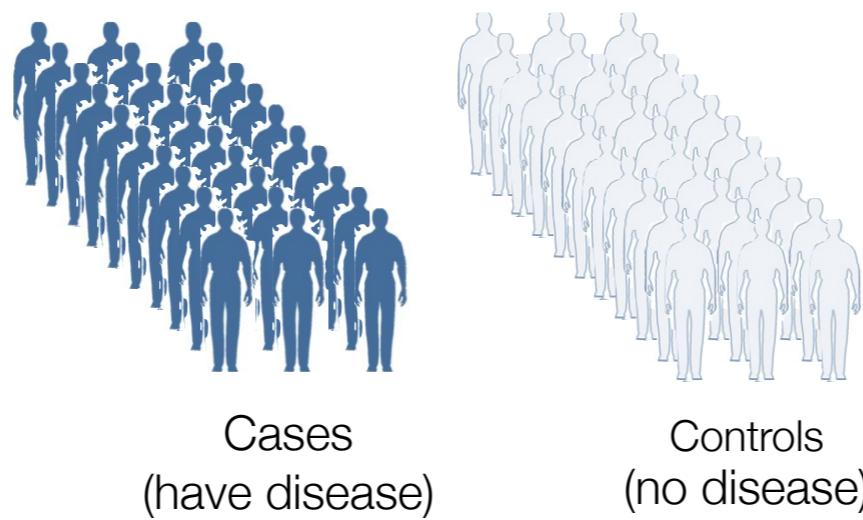
Cases
(have disease)
Controls
(no disease)
Complex diseases
(multiple genes contribute to risk)



e.g.,
High LDL

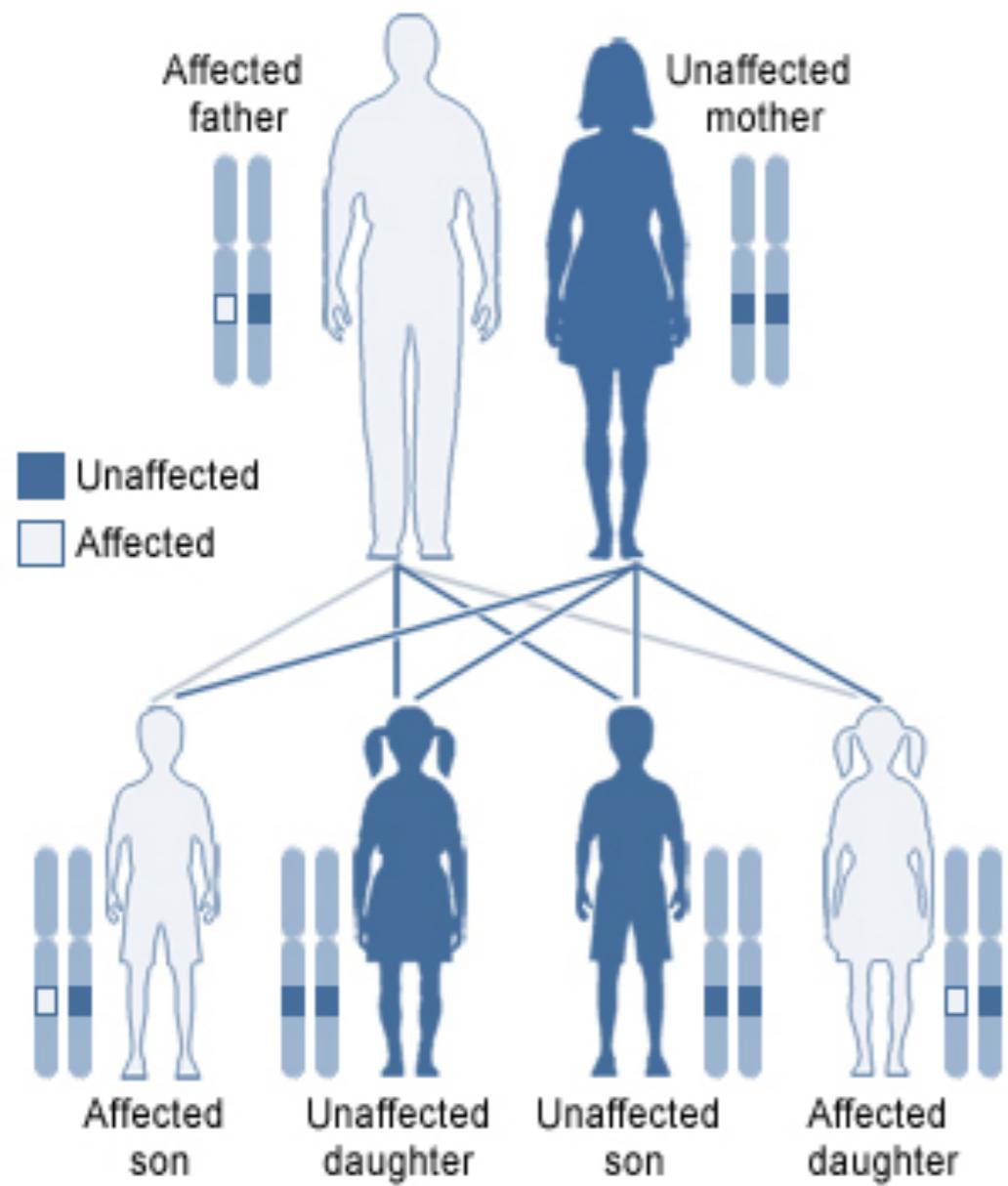
Sampling individuals at extremes of trait distribution

Genome-wide association studies



Single gene diseases (autosomal)

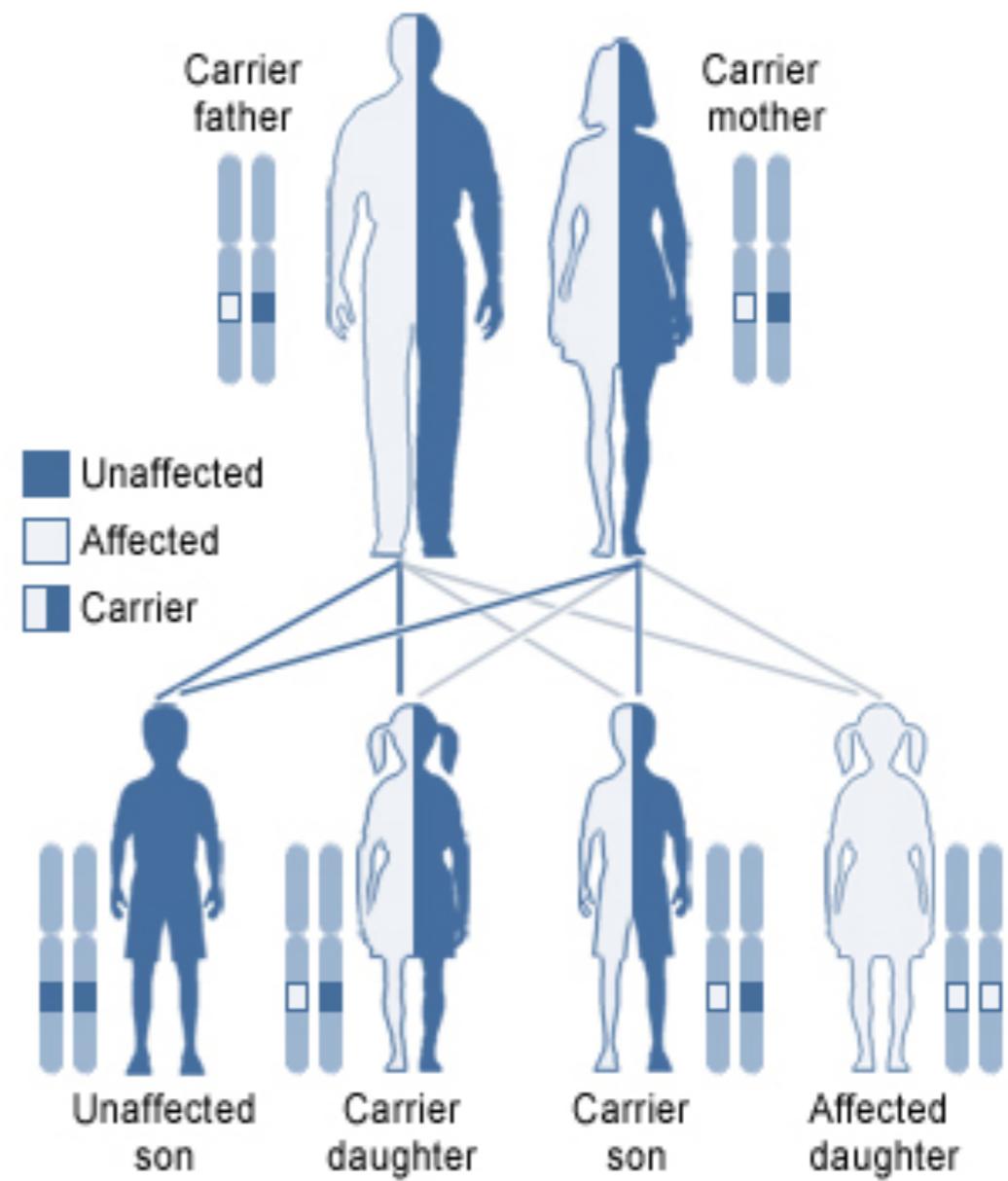
Autosomal dominant



U.S. National Library of Medicine

e.g., Huntington's disease

Autosomal recessive

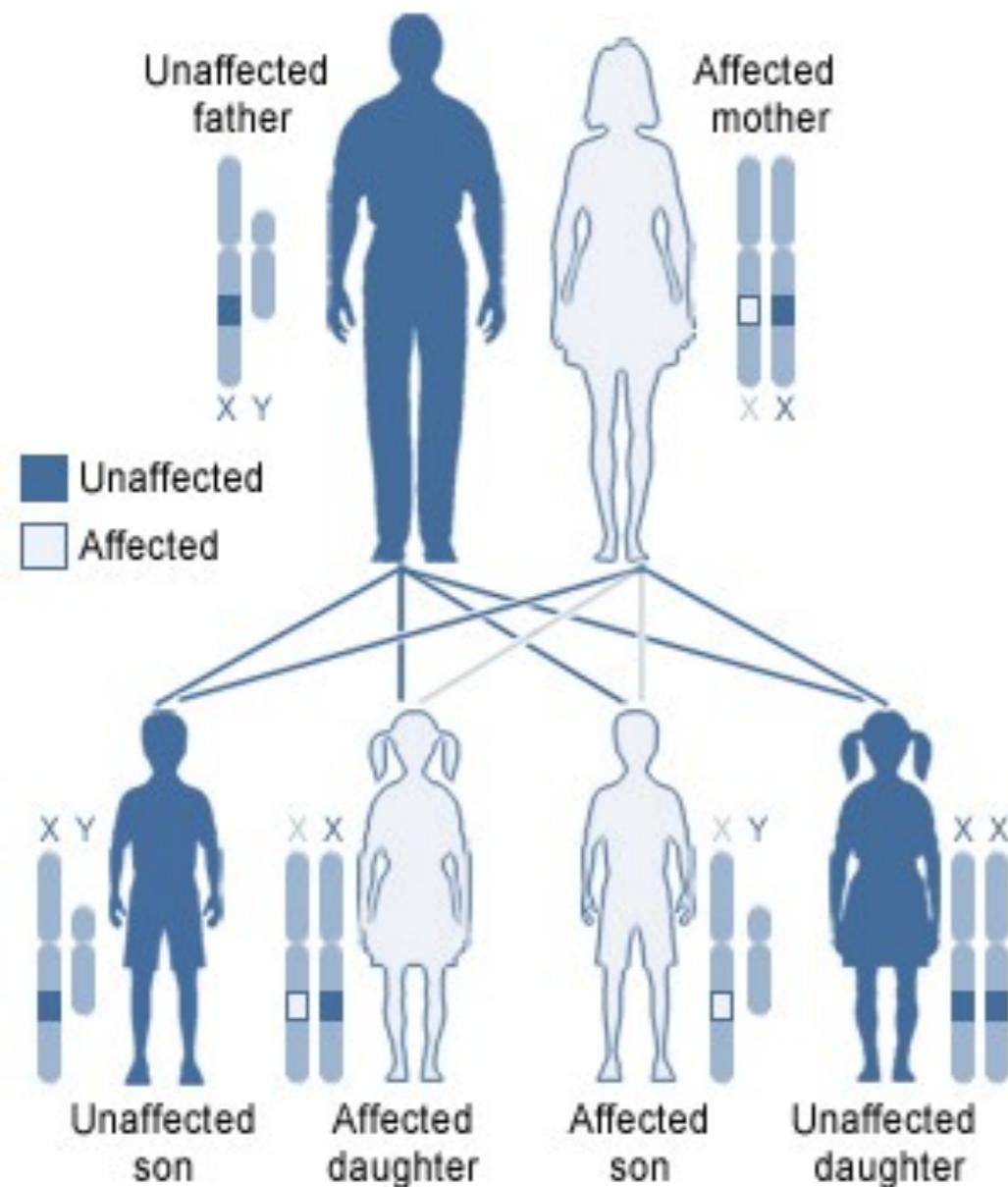


U.S. National Library of Medicine

e.g., cystic fibrosis

Single gene diseases (X-linked dominant)

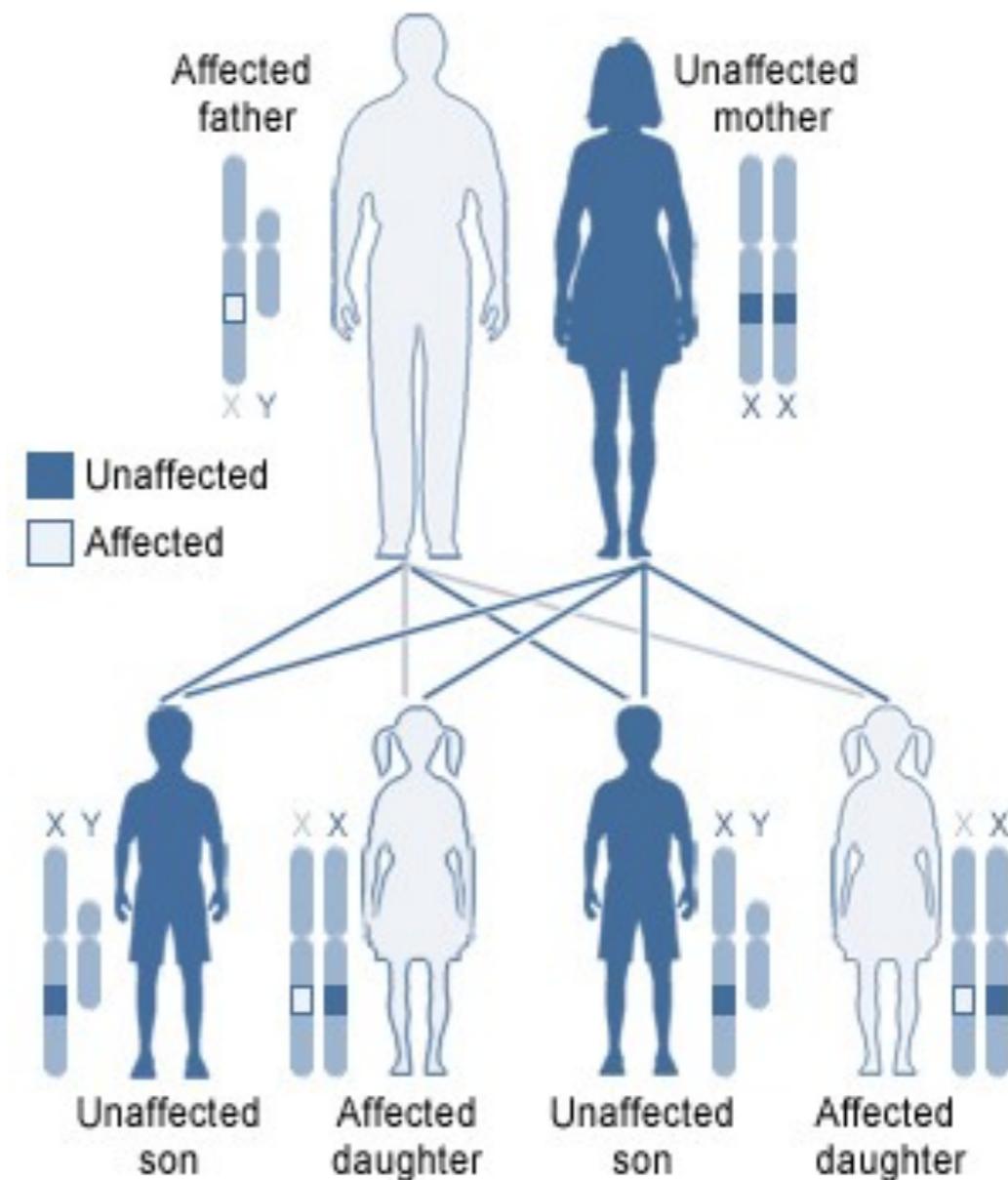
X-linked dominant, affected mother



U.S. National Library of Medicine

Both sons and daughters
can be affected.

X-linked dominant, affected father

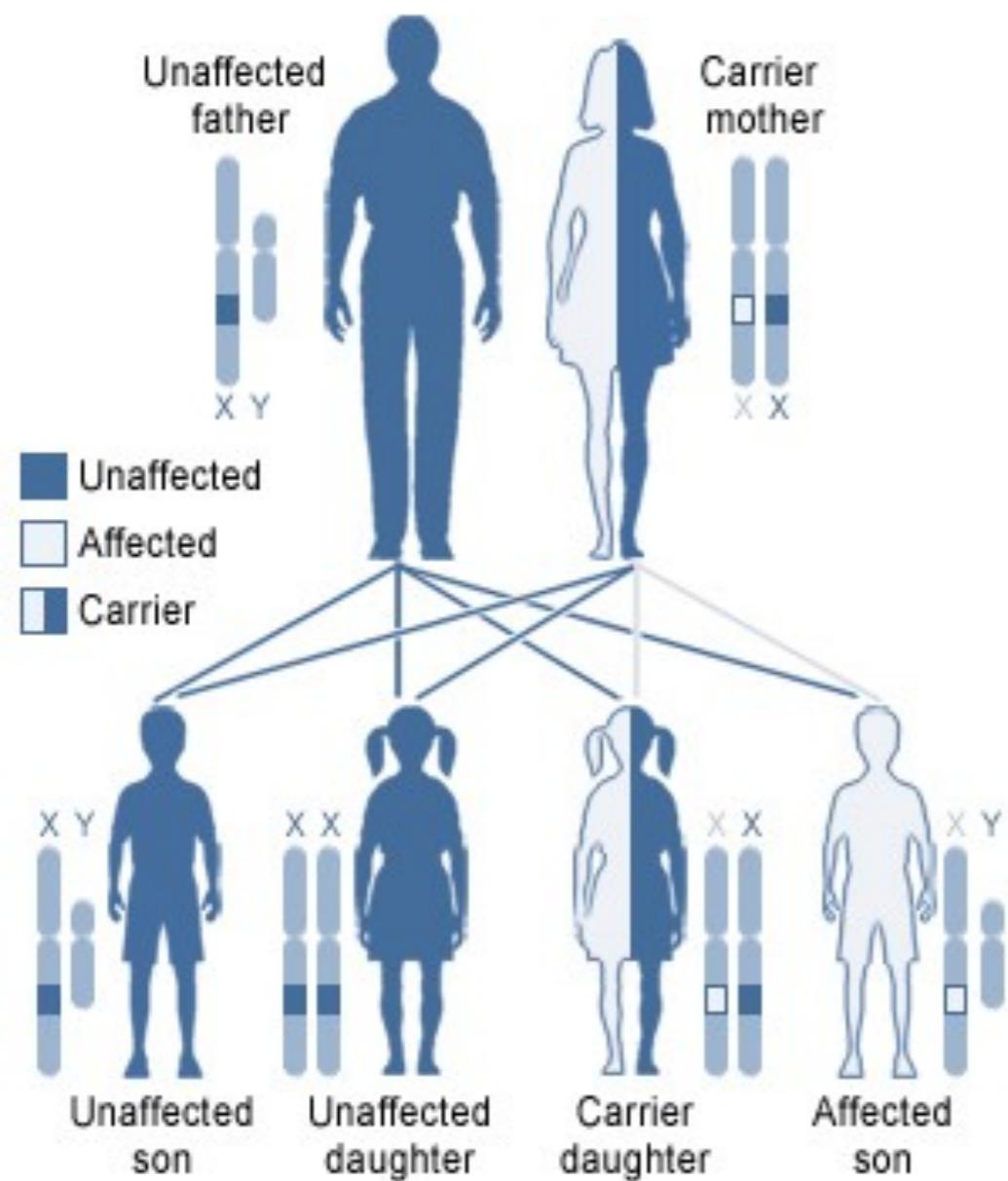


U.S. National Library of Medicine

Affected dads cannot pass on
disease to sons

Single gene diseases (X-linked recessive)

X-linked recessive, carrier mother



X-linked recessive, affected father

