
Big Data Analysis

Introduction to Next-Generation Sequencing



Genome Size Does Matter!

Species	<i>T2 phage</i>	<i>Escherichia coli</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Paris japonica</i>
Genome Size	170,000 bp	4.6 million bp	130 million bp	3.2 billion bp	150 billion bp
Common Name	 Virus	 Bacteria	 Fruit fly	 Human	 Canopy Plant

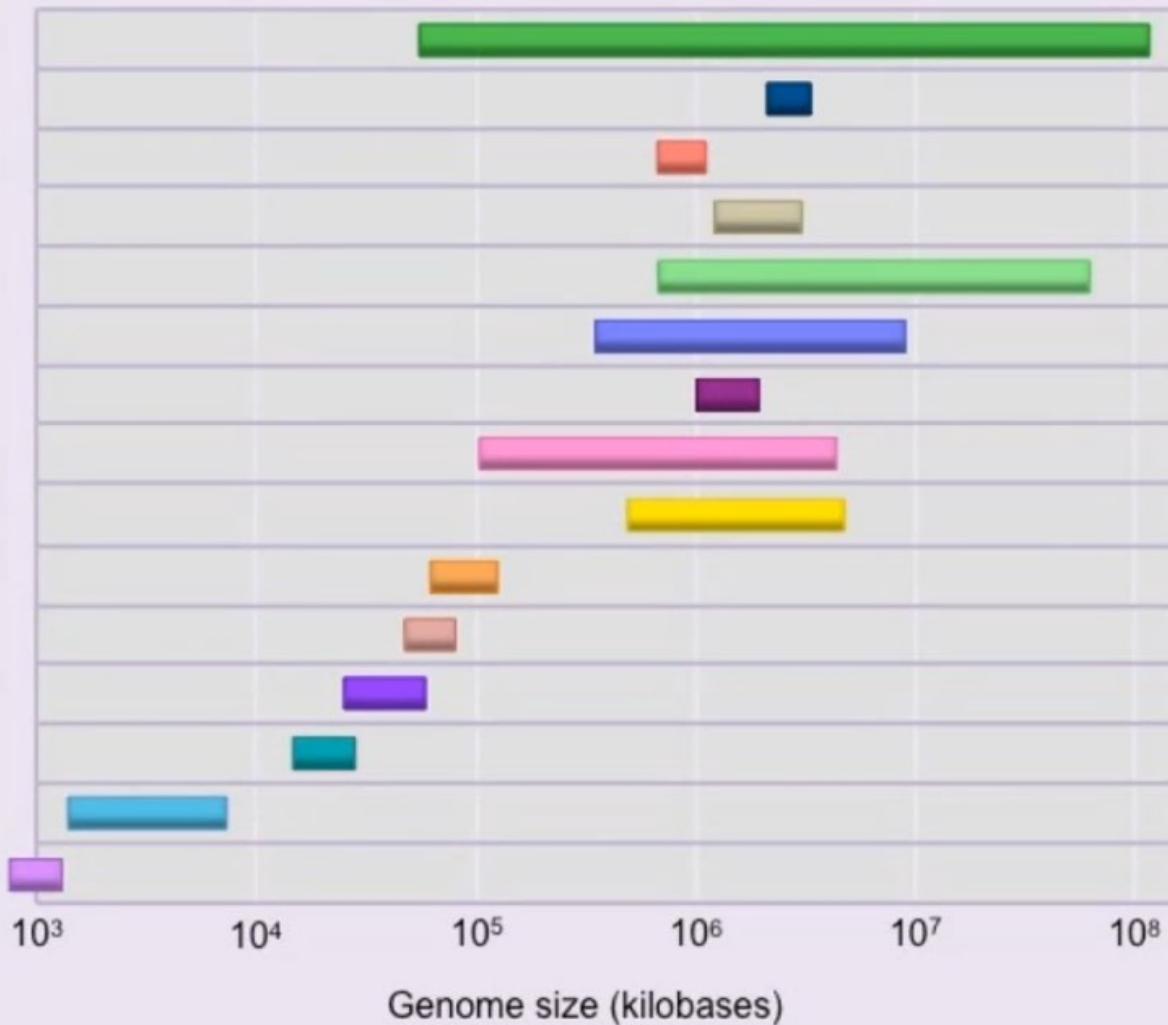
Largest known genome (*Paris japonica*) – 150 billion bp

Smallest known genome (*Carsonella ruddi*) – 160 kb

<https://ib.bioninja.com.au>

Are there any correlation with genome size and phenotypic complexity?

- Flowering Plants
- Mammals
- Birds
- Reptiles
- Amphibians
- Fish
- Crustaceans
- Insects
- Molluscs
- Worms
- Molds
- Algae
- Fungi
- Bacteria
- Viruses





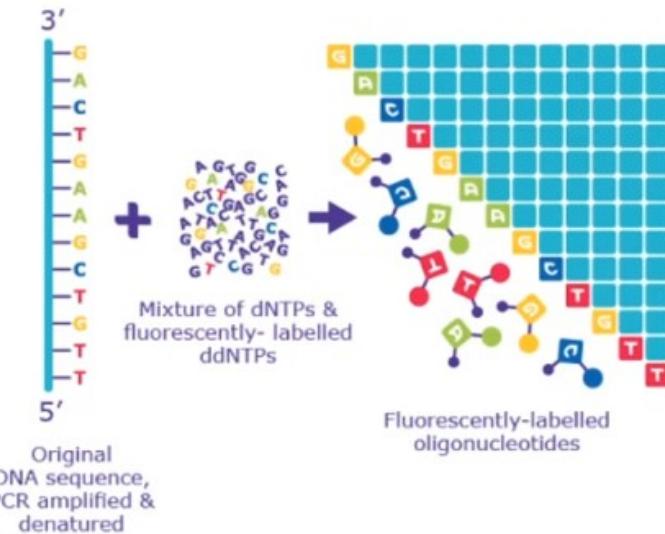
Sanger Sequencing Overview (1st Generation)

The Nobel Prize in Chemistry 1958 - Protein Sequencing
The Nobel Prize in Chemistry 1980 - Nucleic acid sequencing

Frederick Sanger

1

PCR with fluorescent,
chain-terminating ddNTPs



2

Size separation by capillary
gel electrophoresis

Large fragments

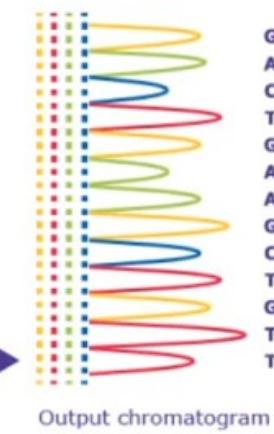
Small fragments

Laser beam

Photomultiplier

3

Laser excitation & detection
by sequencing machine



<https://www.youtube.com/watch?v=KTstRrDTmWI>



High-throughput sequencing (Next-Generation Sequencing)

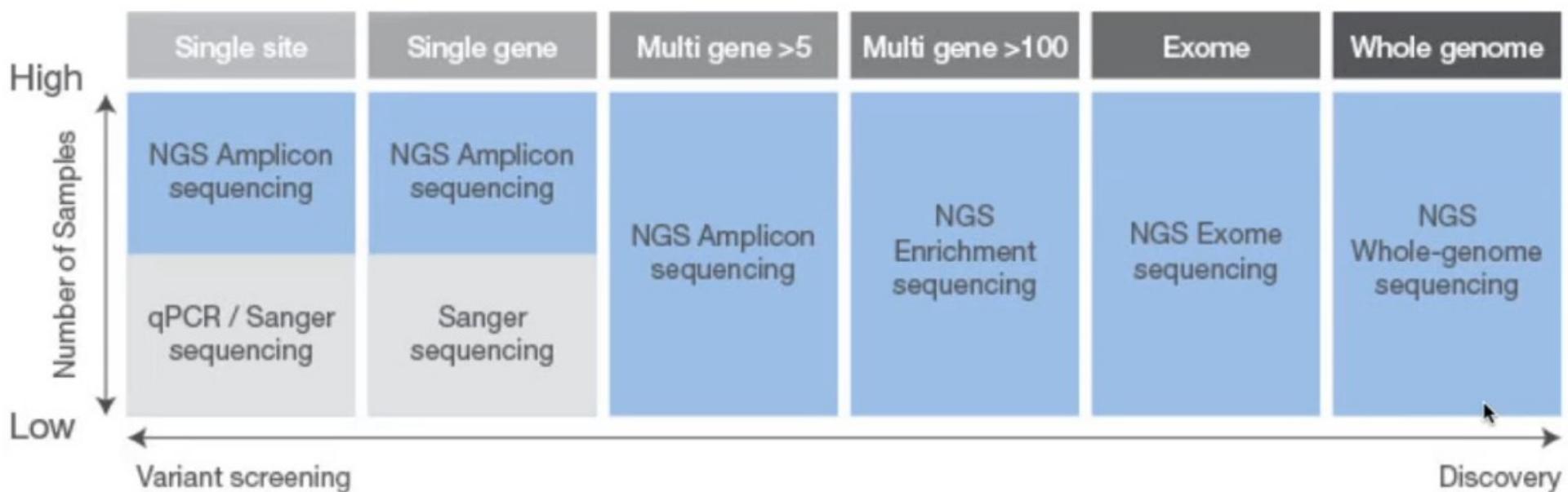
The critical difference between Sanger (dideoxy) sequencing and NGS is **sequencing volume**.

While the Sanger method only sequences a **single DNA fragment** at a time, NGS is **massively parallel**, sequencing millions of fragments simultaneously per run.

This high-throughput process translates into sequencing **hundreds to thousands of genes** at one time.

NGS also offers **greater discovery power** to detect novel or rare variants with deep sequencing.

High throughput sequencing (Next-Generation Sequencing)



High-throughput sequencing (Next-Generation Sequencing)

Advantages

- Higher sensitivity to detect low-frequency variants
- Faster turnaround time for high sample volumes
- Comprehensive genomic coverage
- Lower limit of detection
- Higher throughput with sample multiplexing
- Ability to sequence hundreds to thousands of genes or gene regions simultaneously

High-throughput sequencing (Next-Generation Sequencing)

“It is remarkable to reflect on the fact that the first human genome, famously co-published in Science and Nature in 2001, required 15 years to sequence and cost nearly three billion dollars.

In contrast, the HiSeq X® Ten System, released in 2014, can sequence over 45 human genomes in a single day for approximately \$1000 each.”

High-throughput sequencing (Next-Generation Sequencing)



Illumina sequencing (2nd Generation)

The “sequencing-by-synthesis” technology is used by Illumina. It was originally developed by Shankar Balasubramanian and David Klenerman at the University of Cambridge.

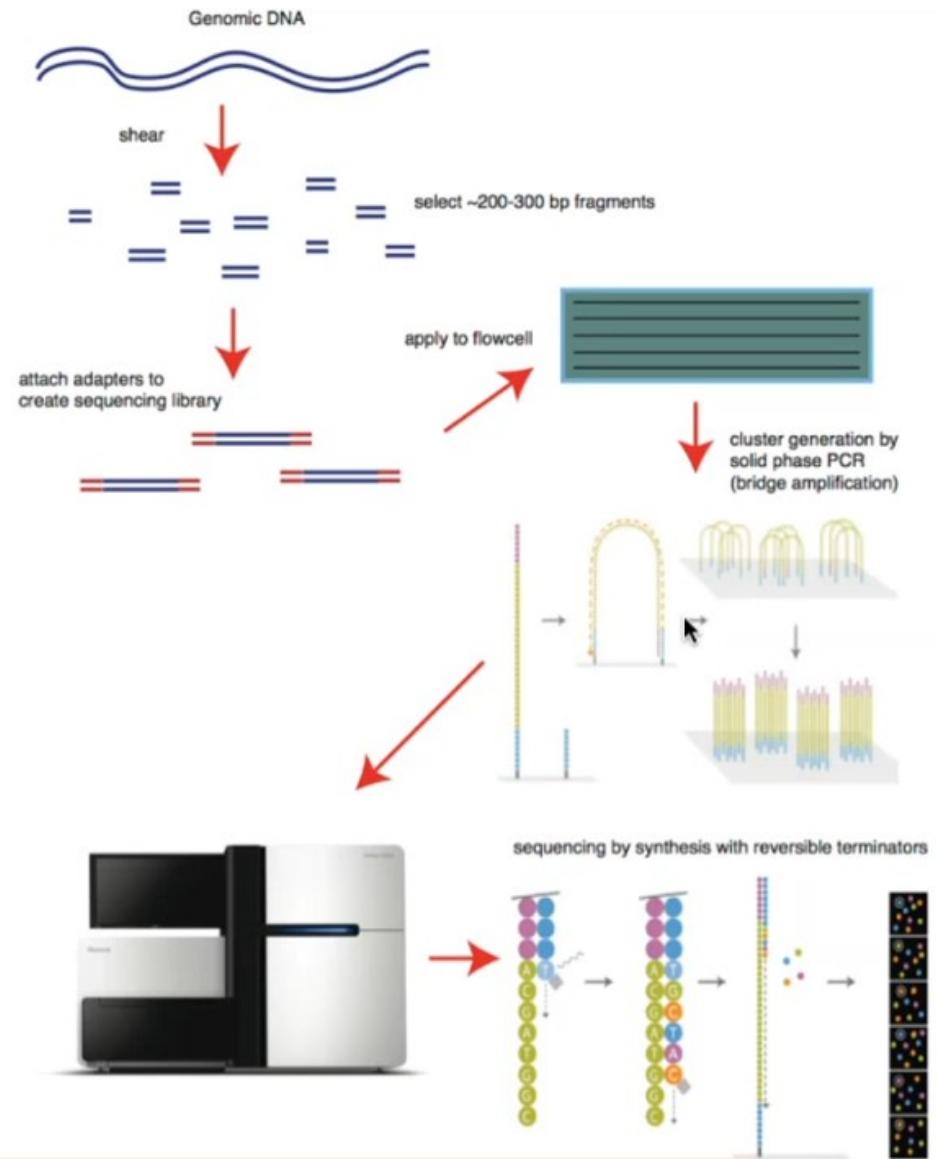
DNA or cDNA samples are randomly fragmented, usually into segments of 200 to 600 base pairs. These fragments are then ligated to adaptors and made single-stranded

Each fragment is amplified on the flow cell, and unlabeled nucleotides and polymerization enzymes are added. These additions, called “Bridge amplification,” connect and lengthen the fragments of DNA on the flow cell

Illumina’s “sequencing by synthesis” involves a proprietary method whereby four labeled reversible dNTP terminators, primers and DNA polymerase are added to the templates on the flow cell.

When excited by a laser, fluorescence from each cluster can be detected, which identifies the first base.

Base calls are made from signal intensity measurements during each cycle, reducing error rates further



<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina sequencing (2nd Generation)

The “sequencing-by-synthesis” technology is used by Illumina. It was originally developed by Shankar Balasubramanian and David Klenerman at the University of Cambridge.

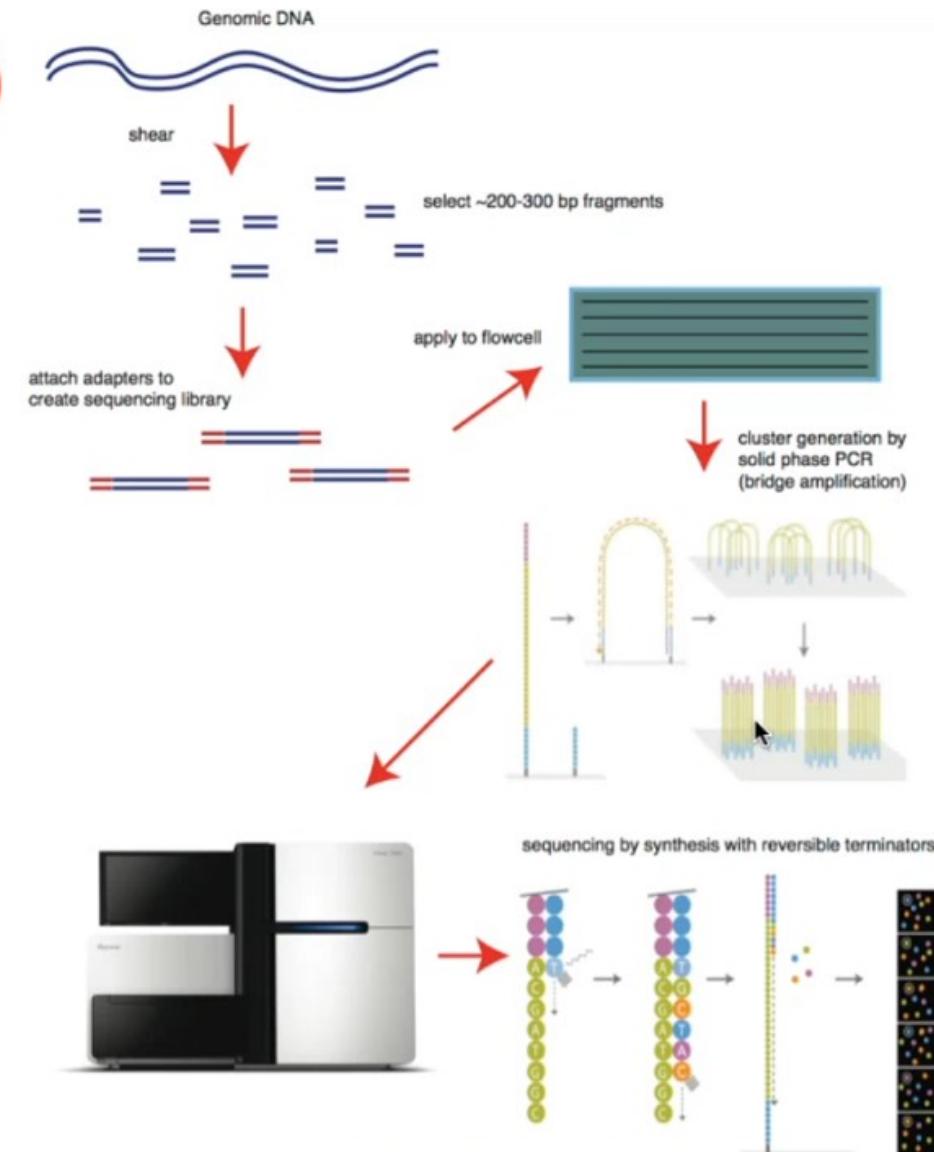
DNA or cDNA samples are randomly fragmented, usually into segments of 200 to 600 base pairs. These fragments are then ligated to adaptors and made single-stranded

Each fragment is amplified on the flow cell, and unlabeled nucleotides and polymerization enzymes are added. These additions, called “Bridge amplification,” connect and lengthen the fragments of DNA on the flow cell

Illumina’s “sequencing by synthesis” involves a proprietary method whereby four labeled reversible dNTP terminators, primers and DNA polymerase are added to the templates on the flow cell.

When excited by a laser, fluorescence from each cluster can be detected, which identifies the first base.

Base calls are made from signal intensity measurements during each cycle, reducing error rates further



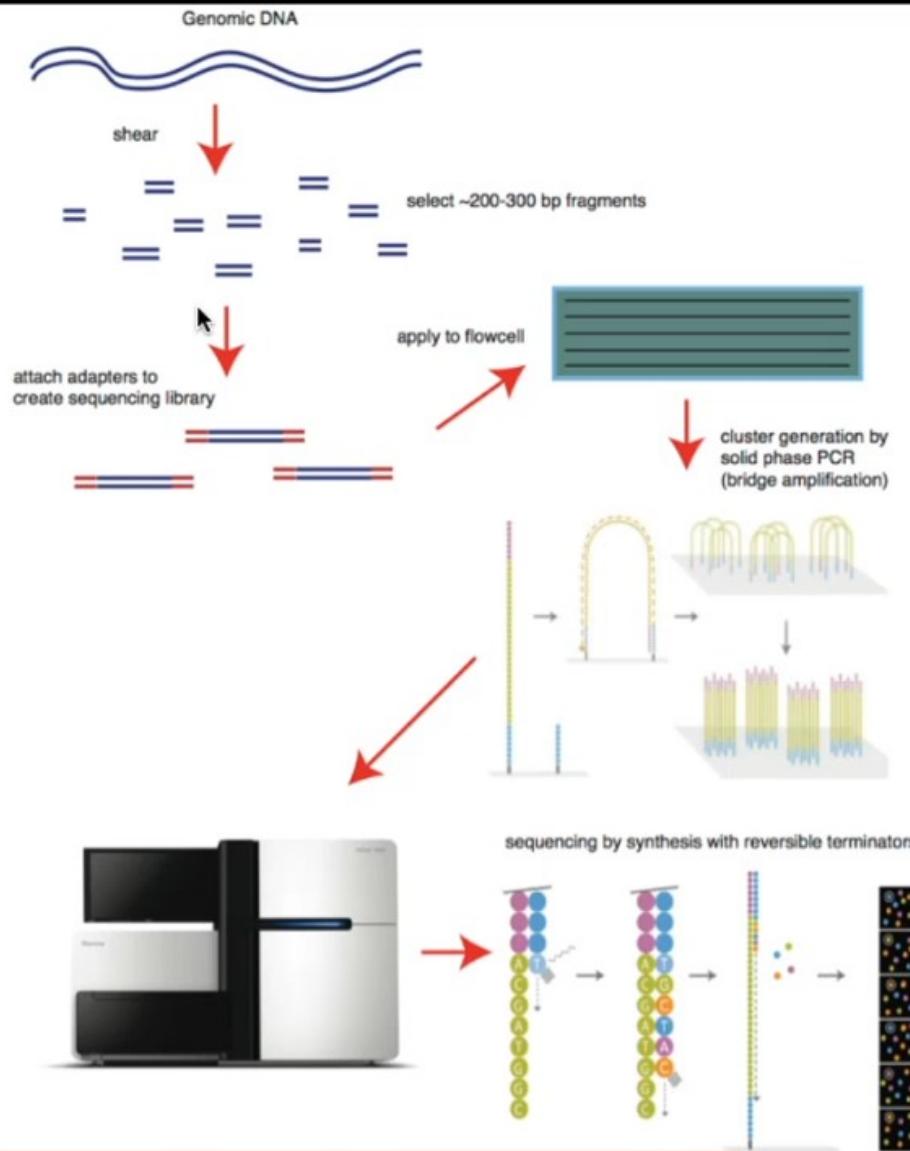
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina sequencing

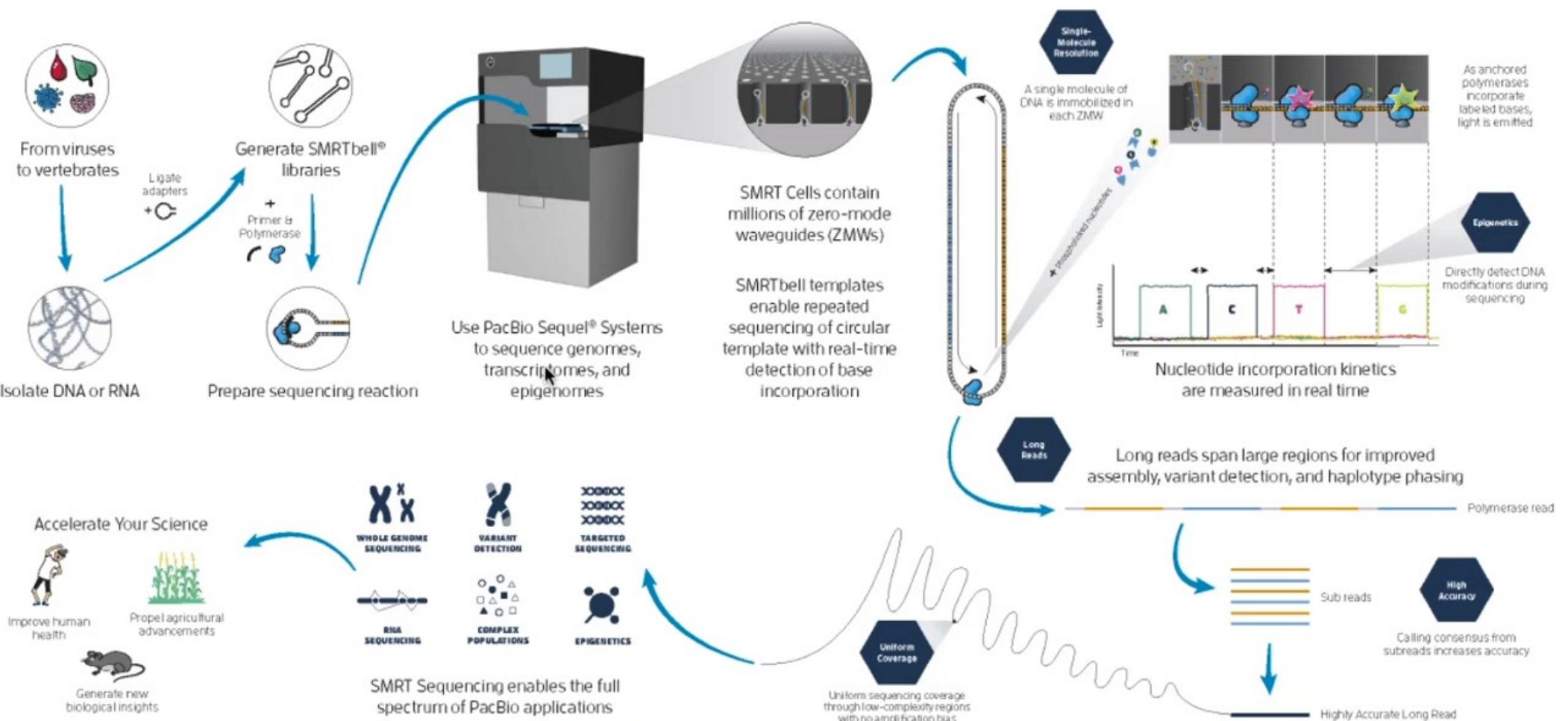
Disadvantages

The short read lengths hinder

1. assigning reads to complex parts of the genome
2. phasing of variants
3. resolving repeat regions
4. introduce gaps and ambiguous regions in *de novo* assemblies



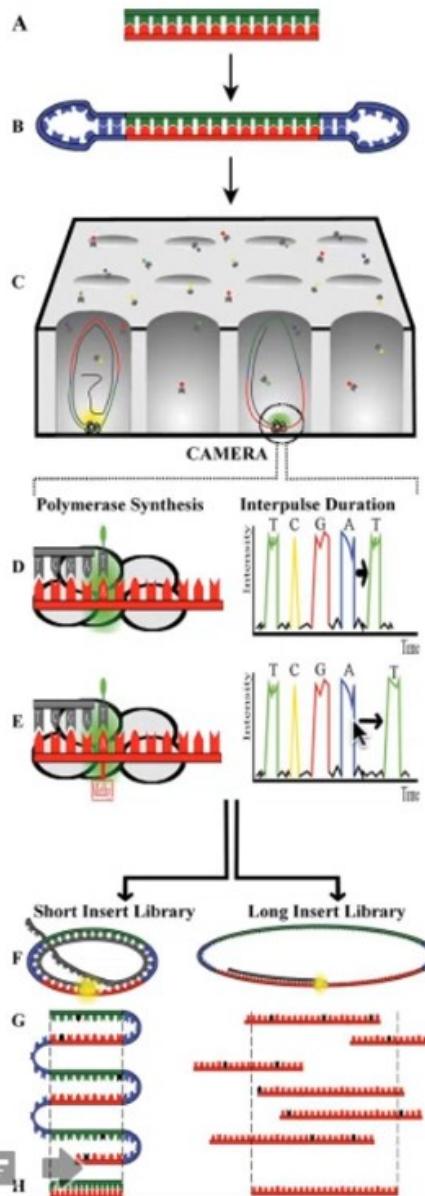
Pacbio sequencing (3rd Generation)



<https://www.youtube.com/watch?v=v8p4ph2MAvI>

www.pacb.com

PN: PS100-032919



Pacbio sequencing

Comparison of PacBio sequencing platforms to two current industry standards

Platform	Read length	Number reads	Error rate	Run time
PacBio RSII (per SMRT cell)	Average 10–16 kb	~55 000	13–15%	0.5–6 hours
PacBio Sequel (per SMRT cell)	Average 10–14 kb	~365 000	13–15%	0.5–10 hours
Illumina HiSeq 4000	2 × 150 bp	5 billion	~0.1%	<1–3.5 days
Illumina MiSeq	2 × 300 bp	25 million	~0.1%	4–55 hours

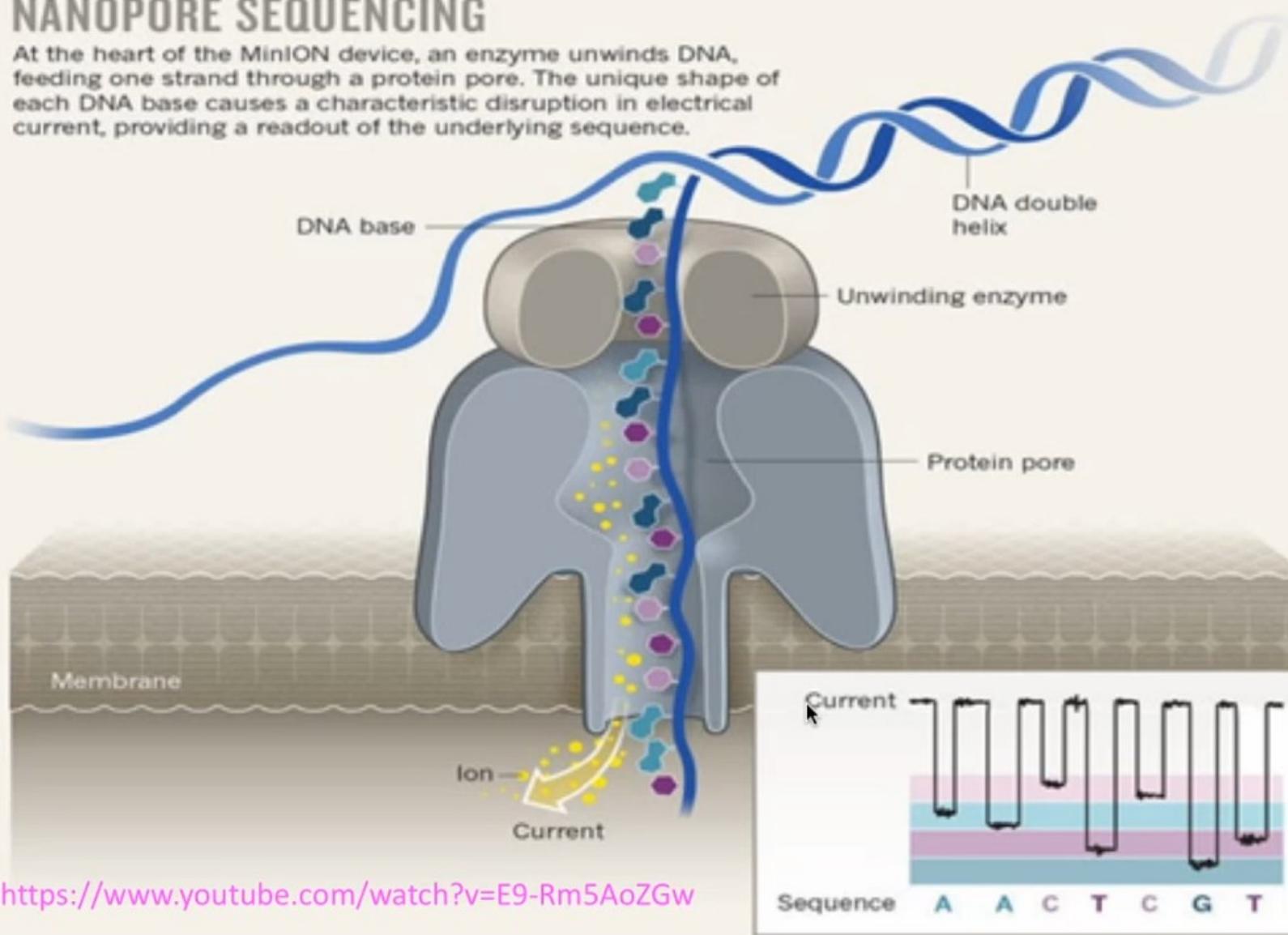
Advantages

1. Long Reads
2. Uniform Coverage
3. Epigenetics
4. Single-molecule resolution



NANOPORE SEQUENCING

At the heart of the MinION device, an enzyme unwinds DNA, feeding one strand through a protein pore. The unique shape of each DNA base causes a characteristic disruption in electrical current, providing a readout of the underlying sequence.



Oxford Nanopore Sequencing (4th Generation)

Advantages

Long Reads (> kb to Mb)

Direct DNA/RNA Sequencing

Epigenetics

Economical

Genome Assembly



Genomic
DNA

Next-generation
DNA sequencing

... CATTCA...
... AGCCATTAG...
... GGTAGTTAG...
... GGTAAACTAG...
... TATAATTAG...
... CGTACCTAG...

millions-billions of *reads*
~30-1000 nucleotides

Resequencing



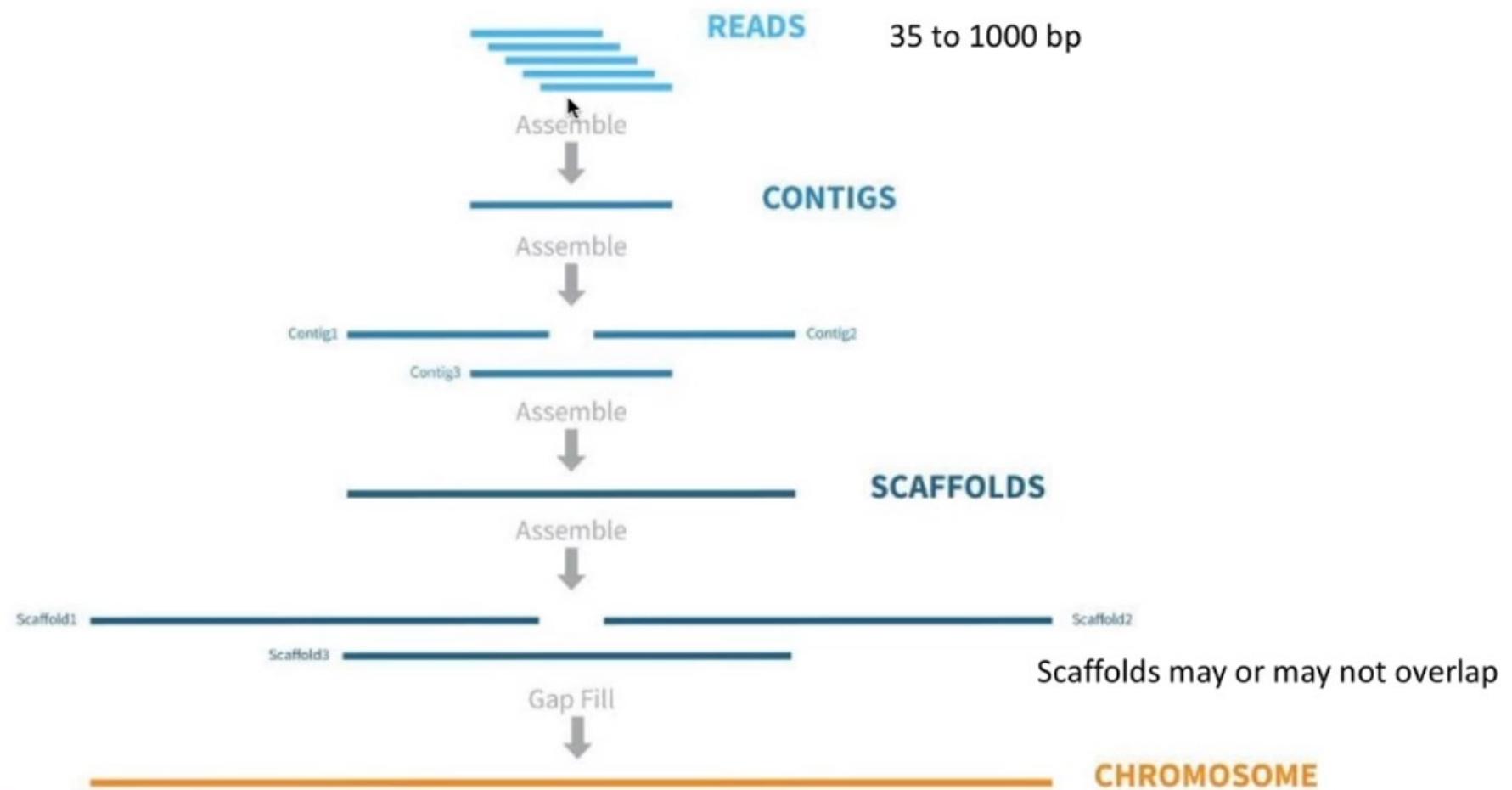
Align reads to *reference genome* and identify variants

***De novo* assembly**



Construct genome sequence from overlaps between reads

De novo Genome Assembly

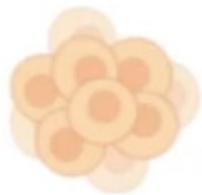


Big Data Analysis

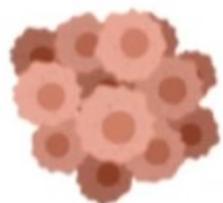
Transcriptomics



Why the normal cell turns into a tumour cell?



Normal Cell

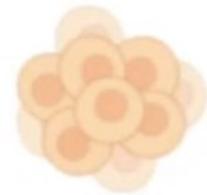


Tumour Cell

There is an underlying **difference** in the genetic mechanism



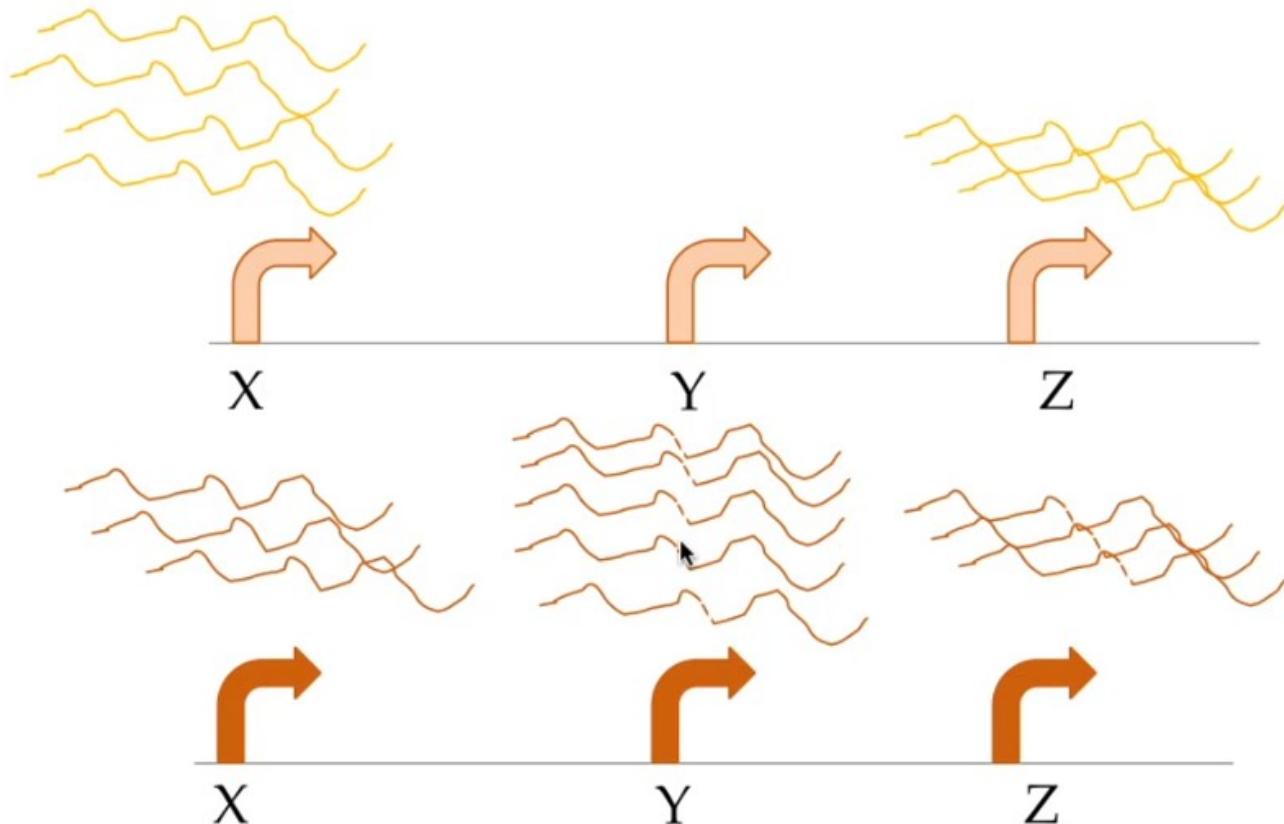
What is the **difference in the expression of genes** between normal and tumour cells?



Normal Cell



Tumour Cell



What are the genes that are upregulated (active) or downregulated (inactive) ?

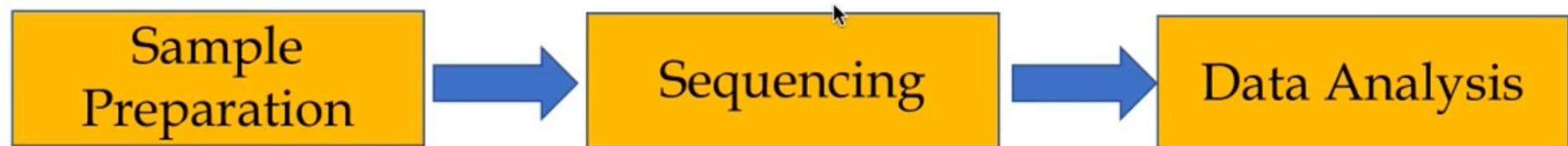
How much the genes are expressed?

Next-Gen Sequencing - RNA-seq

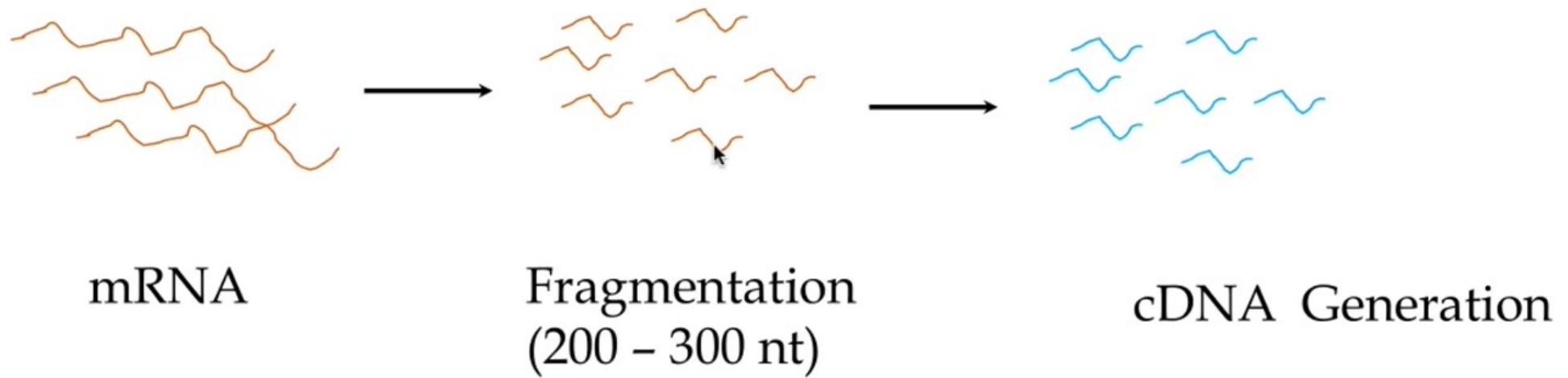
Steps in RNA-seq



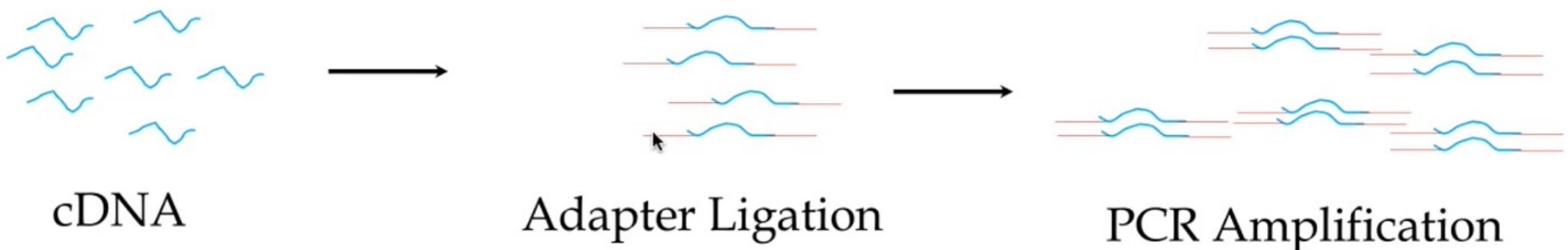
Steps in RNA-seq



Sample (Library) Preparation



Sample (Library) Preparation



Quality Check

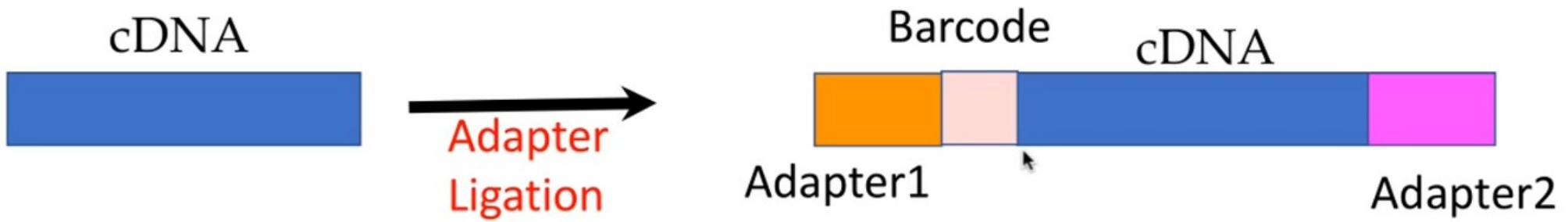


Check for optimum concentration of sample

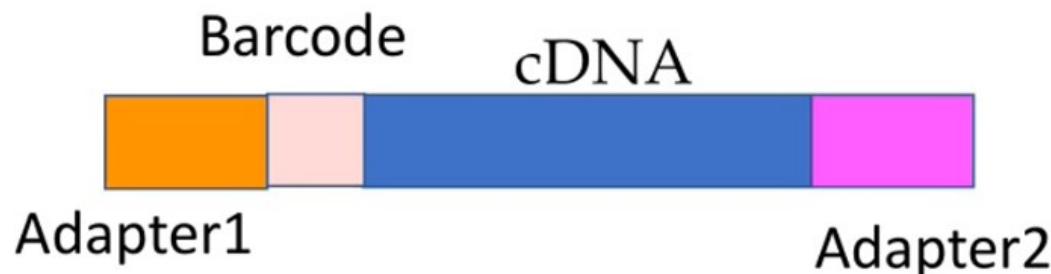
Check for uniform length of the fragments (don't want too short and too long fragments)



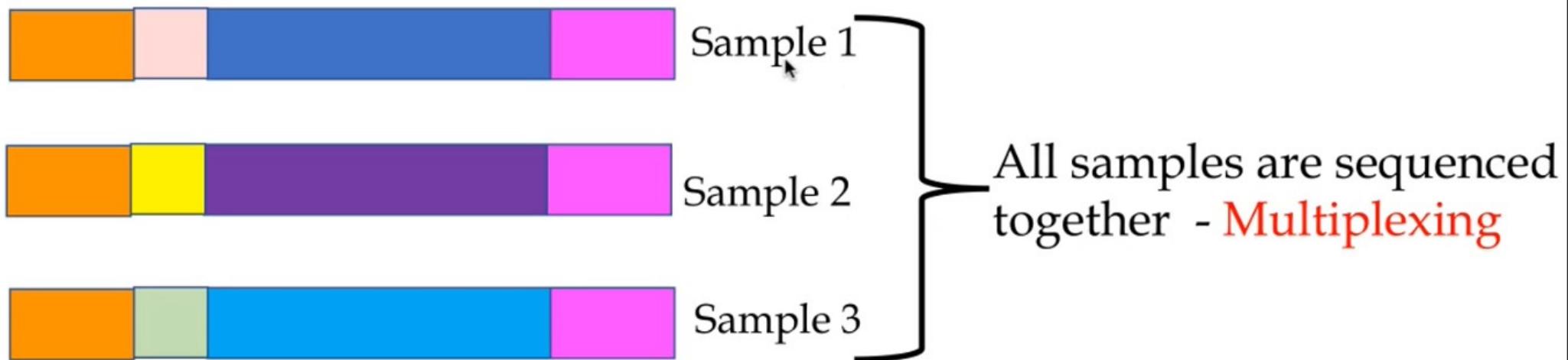
Sample (Library) Preparation



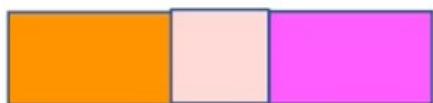
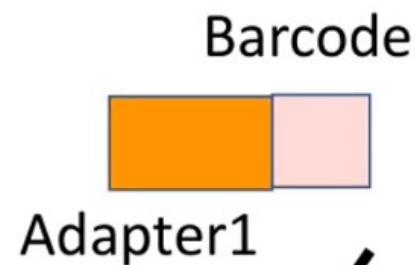
Sample (Library) Preparation



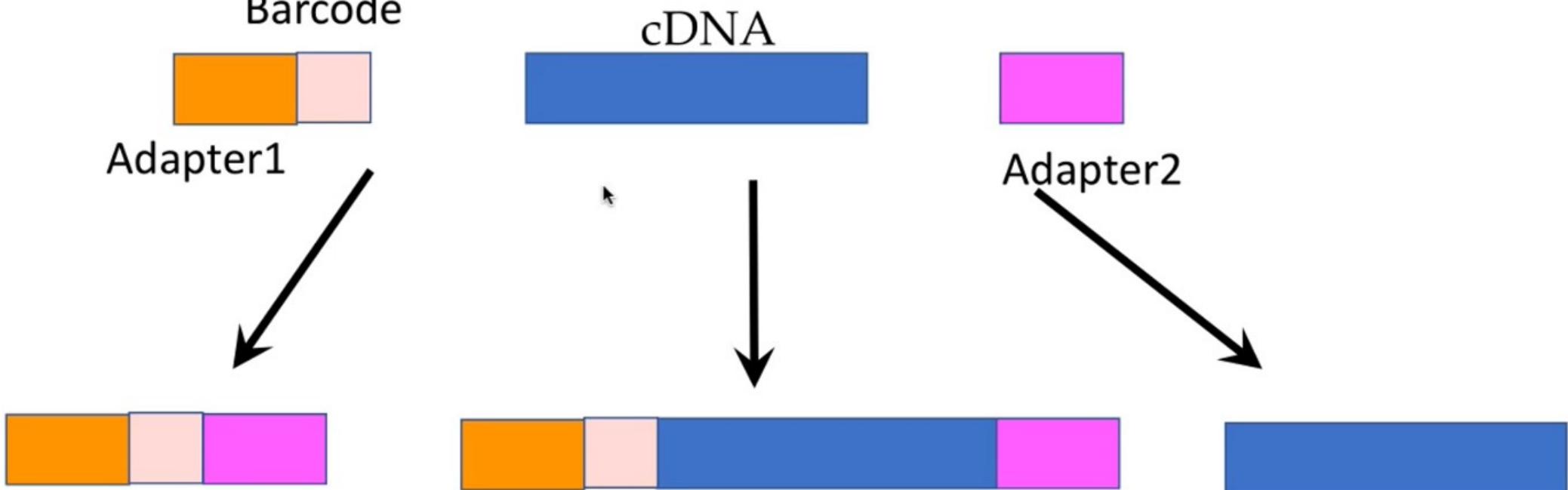
Barcode - A short unique sequence



Sample (Library) Preparation Artefact(s)



Wrong Ligation



Stranded vs Unstranded Library Preparation

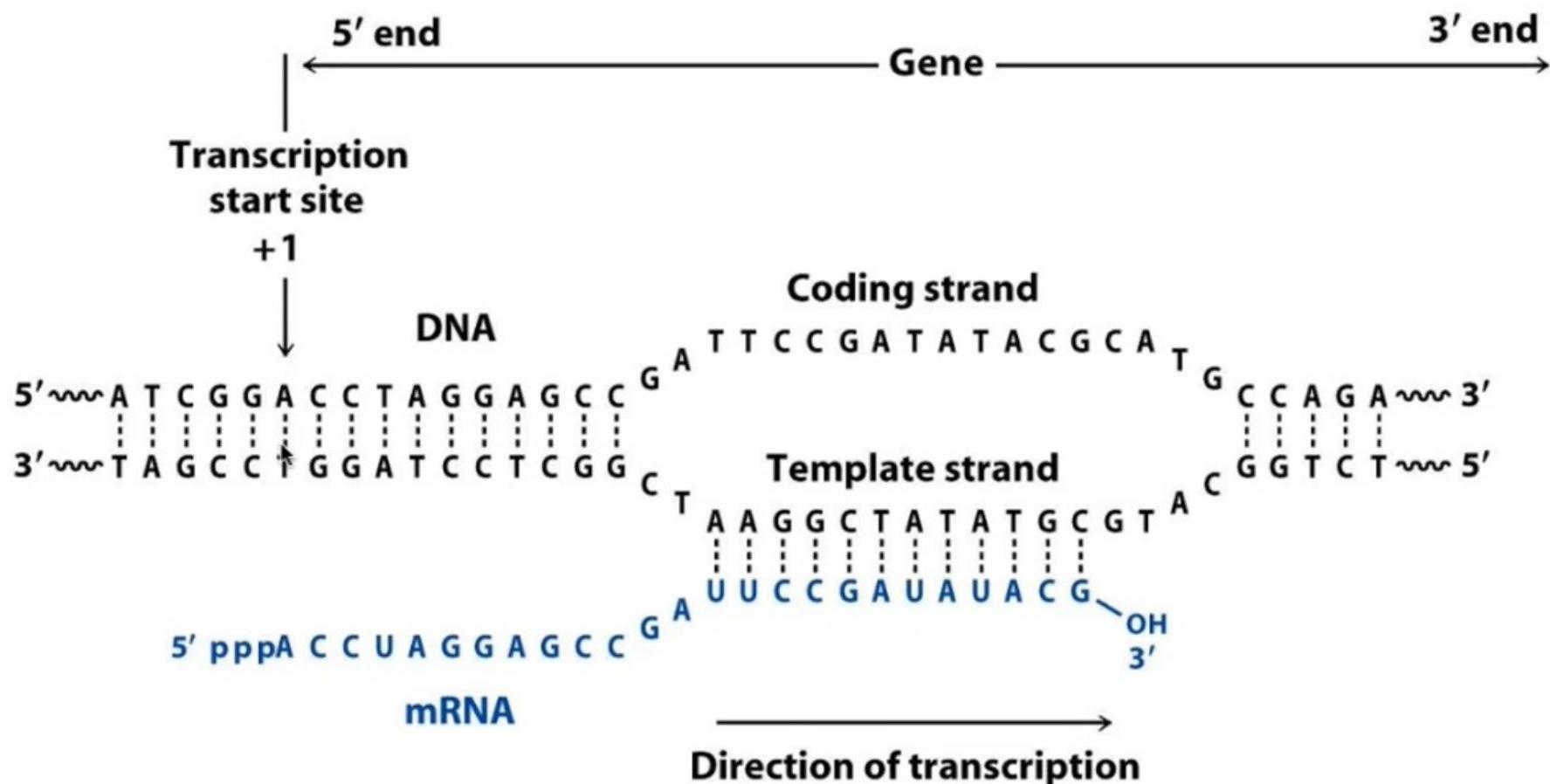
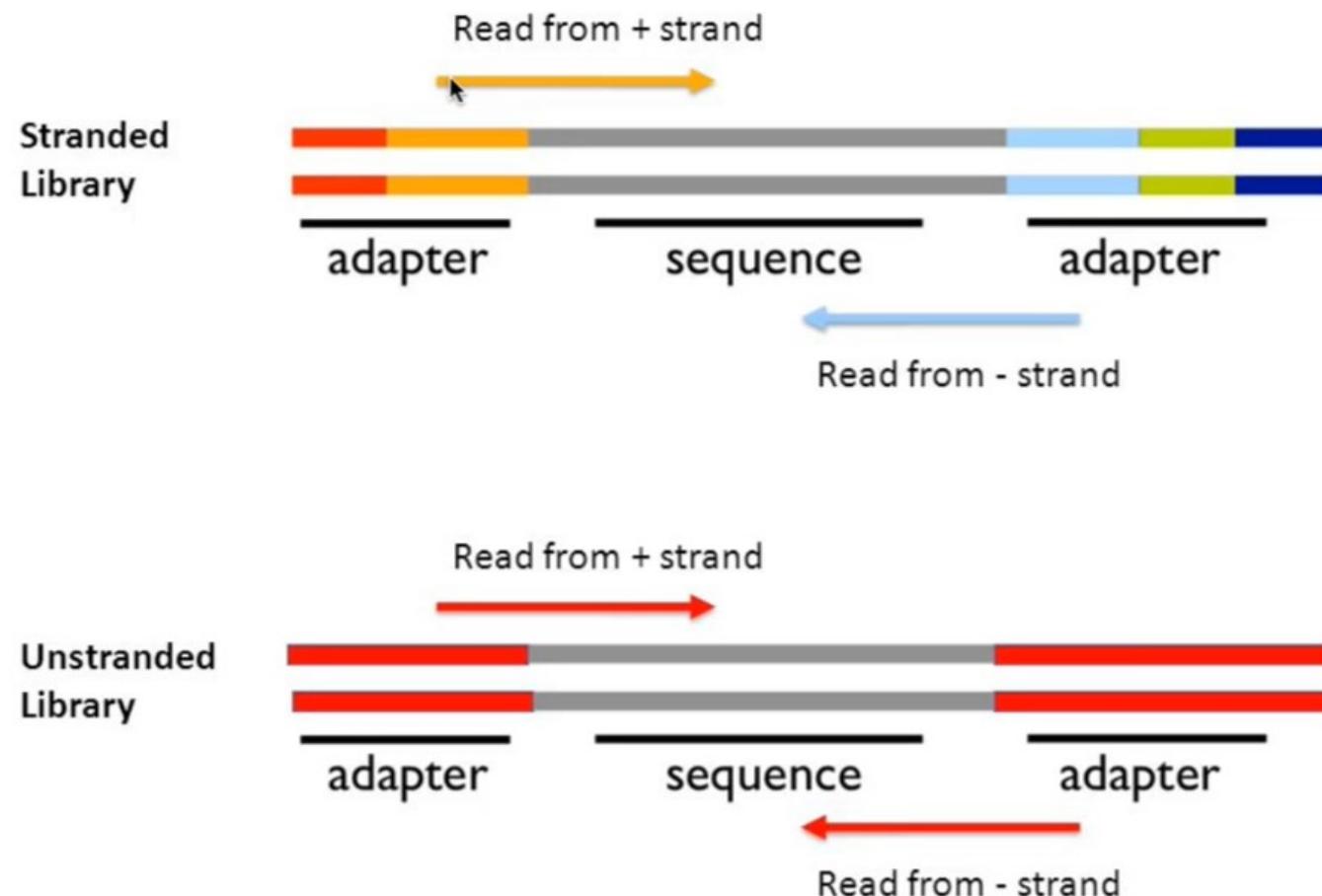


Figure 21-5 Principles of Biochemistry, 4/e

Stranded vs Unstranded Library Preparation

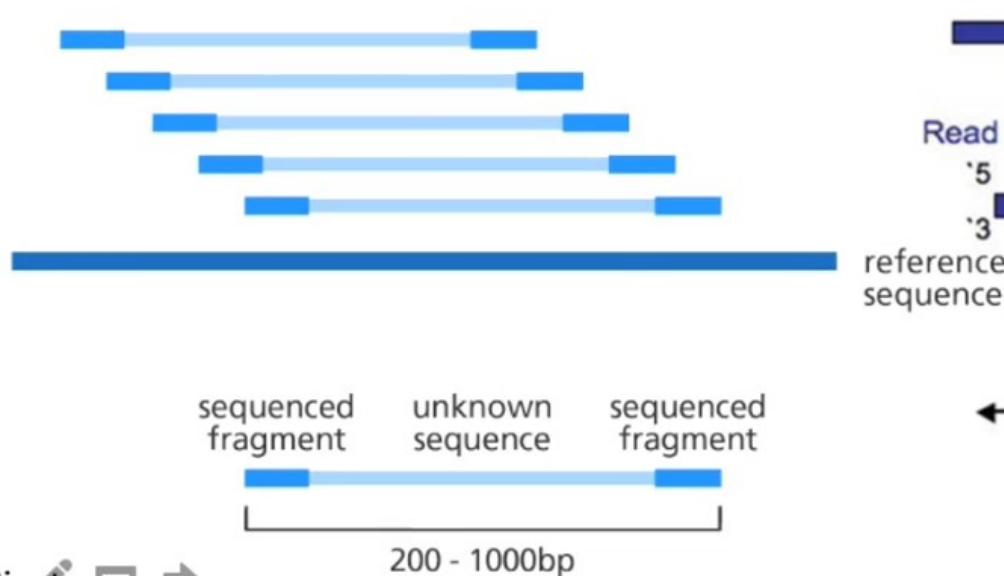
Stranded libraries provide information on the strand that codes for the RNA. This allows to disentangle transcription of overlapping genes.



Single-end reads

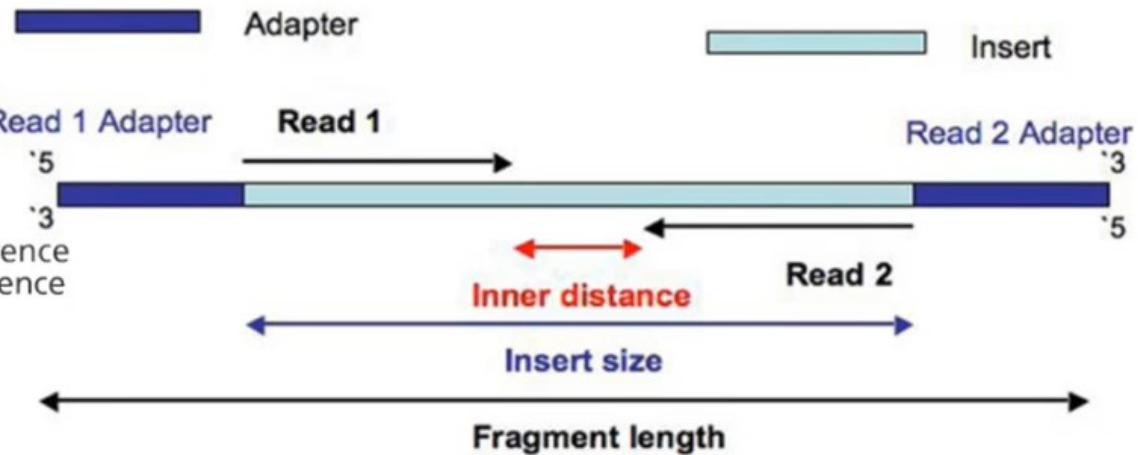


Paired-end reads

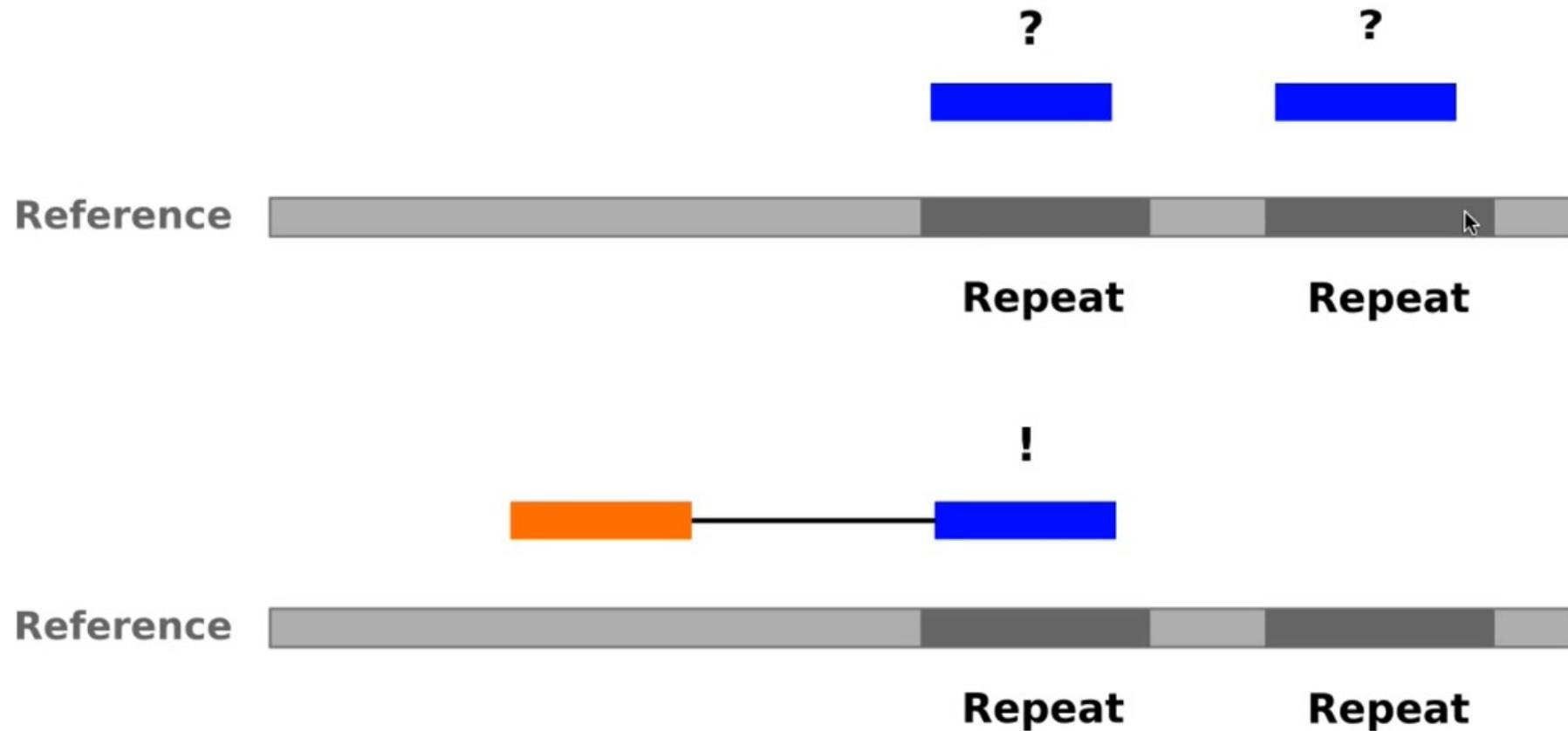


Single Vs. Paired End Read

- Paired-end (mate-pair) sequencing enables sequencing both ends of the same fragment.
- This increases the coverage.
- Helps to resolve ambiguous alignment to the reference sequence



Single Vs. Paired End Read



Illumina sequencing (2nd Generation)

The “sequencing-by-synthesis” technology is used by Illumina. It was originally developed by Shankar Balasubramanian and David Klenerman at the University of Cambridge.

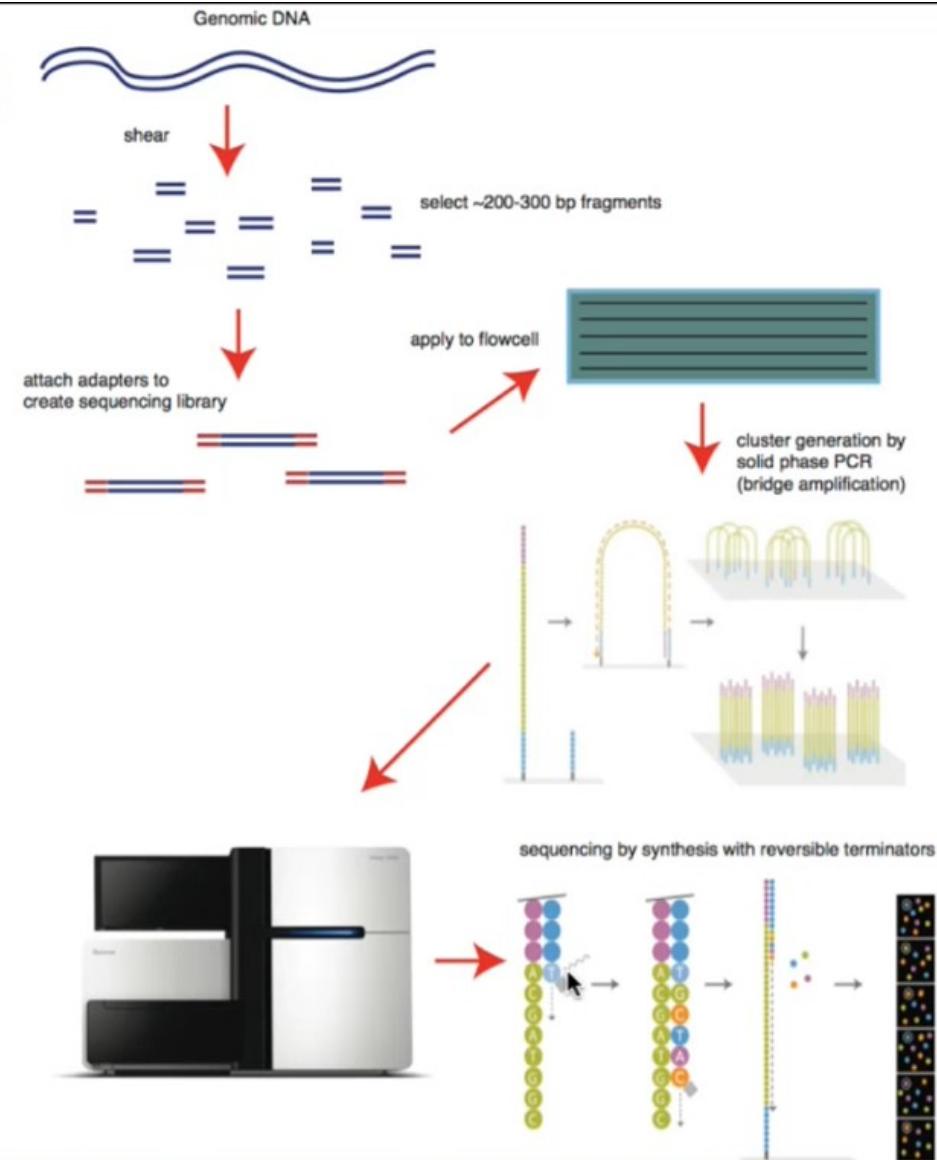
DNA or cDNA samples are randomly fragmented, usually into segments of 200 to 600 base pairs. These fragments are then ligated to adaptors and made single-stranded

Each fragment is amplified on the flow cell, and unlabeled nucleotides and polymerization enzymes are added. These additions, called “Bridge amplification,” connect and lengthen the fragments of DNA on the flow cell

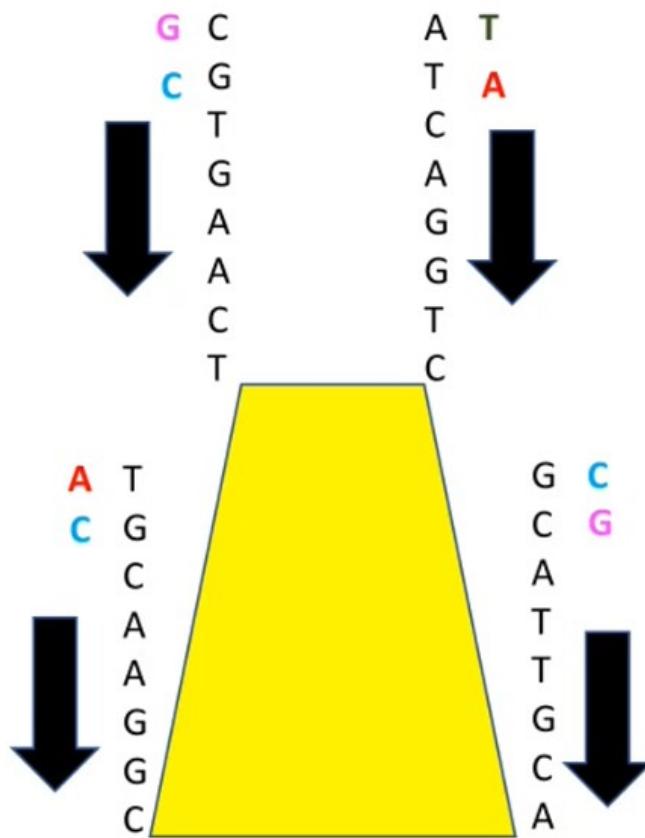
Illumina’s “sequencing by synthesis” involves a proprietary method whereby four labeled reversible dNTP terminators, primers and DNA polymerase are added to the templates on the flow cell.

When excited by a laser, fluorescence from each cluster can be detected, which identifies the first base.

Base calls are made from signal intensity measurements during each cycle, reducing error rates further

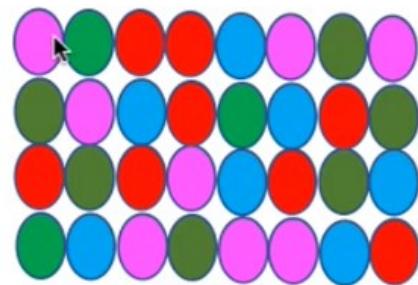


<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

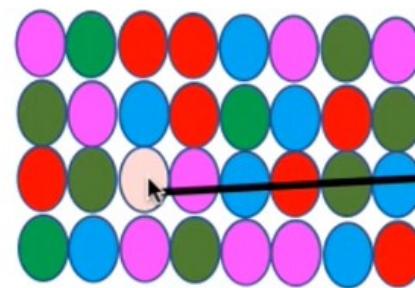
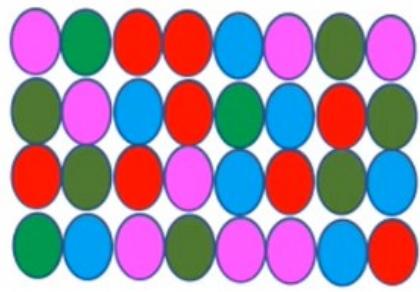


View from top of the flow cell

G A T C



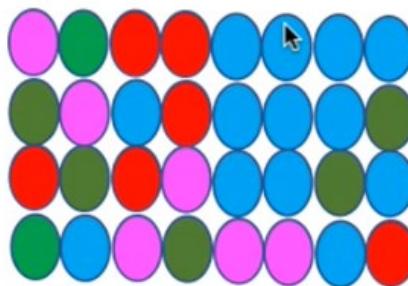
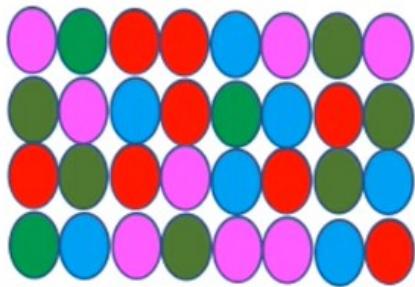
G A T C



Such bases will be called out with low quality scores

Probe is dull.
Machine finds it difficult to call the correct base

G A T C



Same bases occur very closely
– “Low complexity”.
This leads to ambiguity in
finding the correct position of
bases due to blurring of
colours.

Such bases will be called out with low quality scores

Steps in NGS Data Analysis



FASTQ Format

@SEQ_ID

[Header Information Line starts with @ and followed by sequence identifier]

GATTGGGGTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGT_{TT} [Sequence]

+ [+ followed by sequence identifier or description]

!"*((***+))%%%%++)(%%%%%).1***-+*")**55CCF>>>>CCCCCCC65 [Quality score identifier]

Symbols = Low quality

Numbers = better

Alphabets = Best



Data Processing and Analysis Steps

1. Remove the low quality reads and adapters
2. Align the reads to the genome of interest
3. Count the number of reads per gene

Data Processing and Analysis Steps

1. Remove the low quality reads and adapters
2. Align the reads to the genome of interest
3. Count the number of reads per gene

Quality Scores

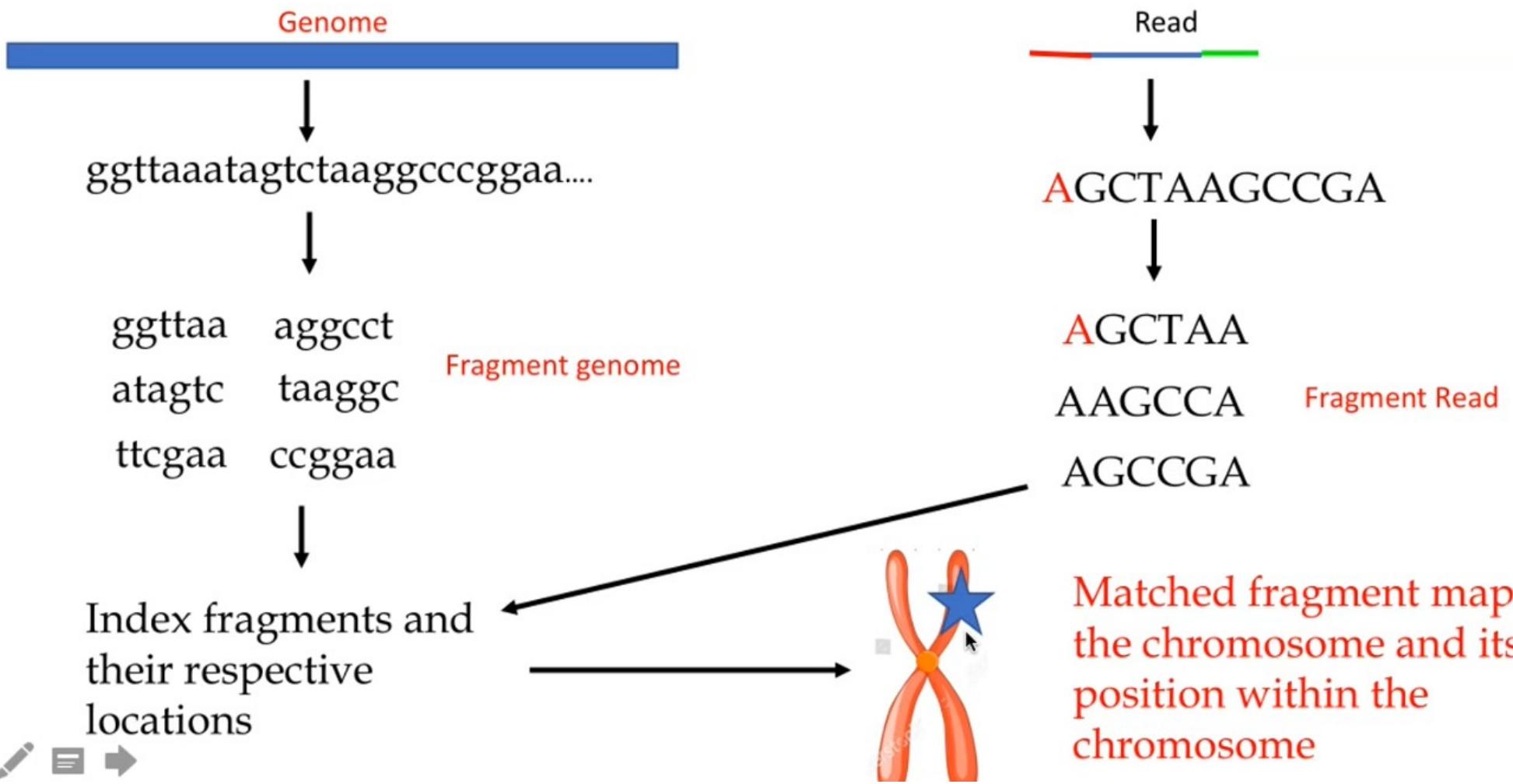
$$Q = -10 \log_{10}(P)$$

Q - Phred Quality Score

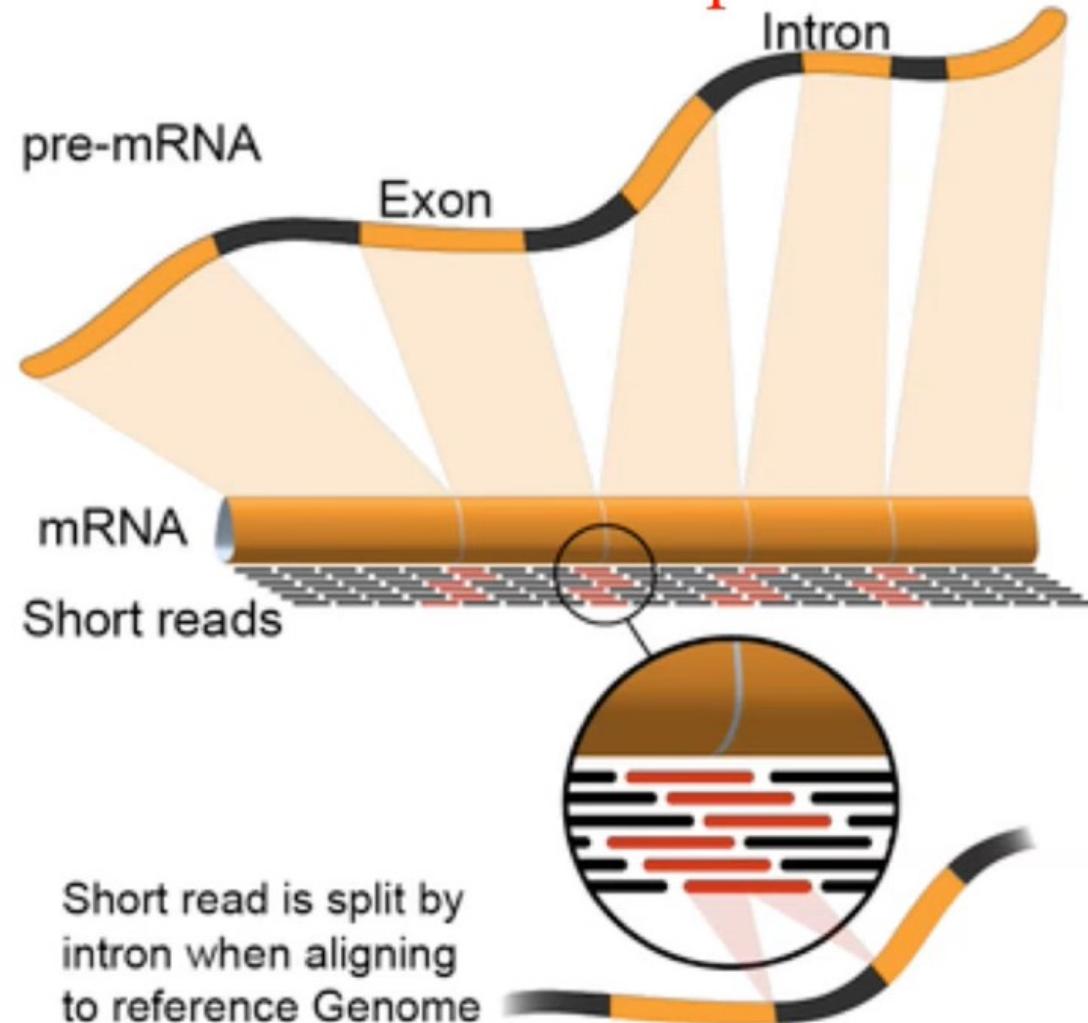
P - Probability of calling an incorrect base

Q	P	Base Call Accuracy
10	0.1	90%
20	0.01	99%
30	0.001	99.9%
40	0.0001	99.99%
50	0.00001	99.999%

Genome Alignment



Alignment to Reference Sequence



Alignment involves
the knowledge of
splice junctions

Sequenced reads have to be normalised for

1. Sequencing depth

2. Length of genes

→ no. of times a particular gene gets sequenced

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Sample 1, 2 and 3 are replicates – experiments repeated multiple times under similar experimental conditions.

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Sample 3 has more reads than others irrespective of gene sizes.

It represents that Sample 3 has more depth than others.

1. Sample 1 and 2 may have more low quality reads
2. Sample 3 has slightly higher concentration than sample 1 & 2 on the flow cell

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Gene C has twice as many reads as Gene B regardless of samples.

It represents that Sample 3 has more depth than others.

Reads per kilo base of transcript per Million mapped Reads (RPkM)

$$\text{RPkM} = \frac{\text{\# of Reads} \times 10^6}{\text{Total \# of Reads} \times \text{Gene size}}$$



Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1
Total	35	45	106
Scaling Factor	3.5	4.5	10.6

Scaling factor for Read counts = $\frac{\text{Total # counts}}{10}$

In Reality, we have to divide by 10^6



Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	1.43	1.78	1.42
B (2 kb)	2.86	2.67	2.83
C (4 kb)	5.71	5.56	5.56
D (10 kb)	0	0	0.09

Scaled for Read count

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Raw Count



RPkM

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	1.43	1.78	1.42
B (2 kb)	1.43	1.39	1.42
C (4 kb)	1.43	1.78	1.42
D (10 kb)	0	0	0.009

Scaled for Read count and Gene length

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Raw Count

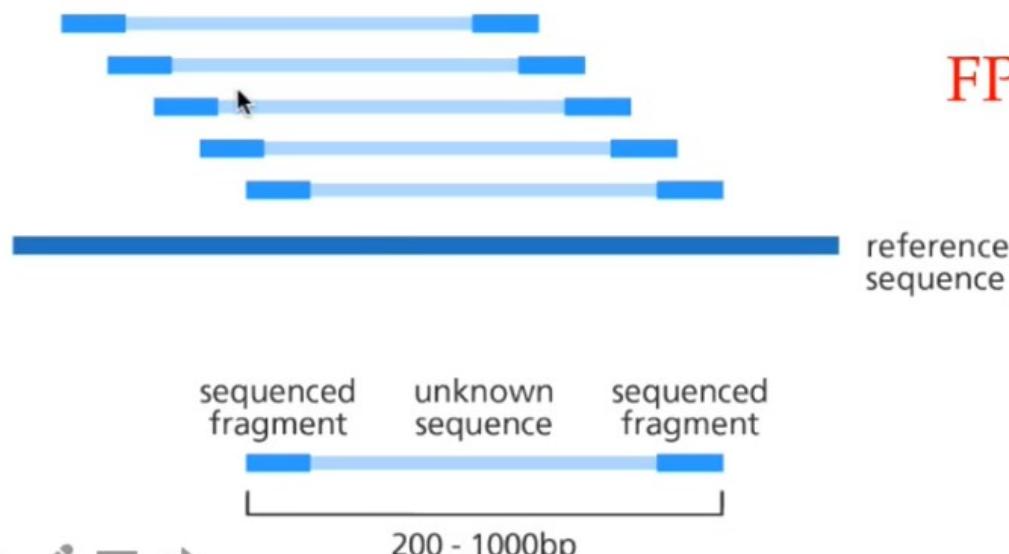


Single-end reads



RPkM = Reads Per kilobase Million

Paired-end reads



FPkM = Fragments Per kilobase Million

Both reads mapped to the same fragment are NOT counted twice

Transcripts Per Million (TPM)

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5 ↗	8	15
B (2 kb)	5	6	15
C (4 kb)	5	6.25	15
D (10 kb)	0	0	0.1
Total	15	20.25	45.1

Scaled for
Gene length -
 Divide each
 read count by
 respective
 gene length

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	10	12	30
C (4 kb)	20	25	60
D (10 kb)	0	0	1

Raw Count

Transcripts Per Million (TPM)

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	5	8	15
B (2 kb)	5	6	15
C (4 kb)	5	6.25	15
D (10 kb)	0	0	0.1
Scaling Factor	1.5	2.025	4.51

Scaled for Gene length

Scaling factor for Gene Size = Total # counts after gene length scaling
10

In Reality, we have to divide by 10^6

Transcripts Per Million (TPM)

Gene name	Sample 1	Sample 2	Sample 3	
A (1 kb)	3.33	3.95	3.326	Scaled for Gene length and Read count - Divide
B (2 kb)	3.33	2.96	3.326	each read count by respective scaling factor
C (4 kb)	3.33	3.09	3.326	
D (10 kb)	0	0	0.02	

Gene name	Sample 1	Sample 2	Sample 3	
A (1 kb)	5	8	15	Raw Count
B (2 kb)	10	12	30	
C (4 kb)	20	25	60	
D (10 kb)	0	0	1	



RPkM Vs. TPM

RPkM

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	1.43	1.78	1.42
B (2 kb)	1.43	1.39	1.42
C (4 kb)	1.43	1.78	1.42
D (10 kb)	0	0	0.009
Total	4.29	4.5	4.25

Total Reads
are different.

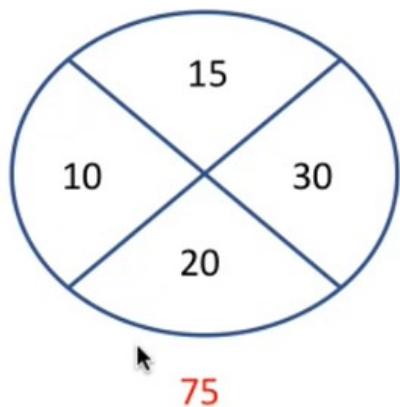
TPM

Gene name	Sample 1	Sample 2	Sample 3
A (1 kb)	3.33	3.95	3.326
B (2 kb)	3.33	2.96	3.326
C (4 kb)	3.33	3.09	3.326
D (10 kb)	0	0	0.02
Total	10	10	10

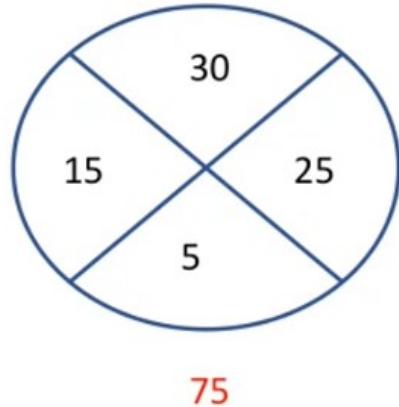
Total Reads
are same.
Comparisons
across
samples make
sense.

Sample 1

TPM



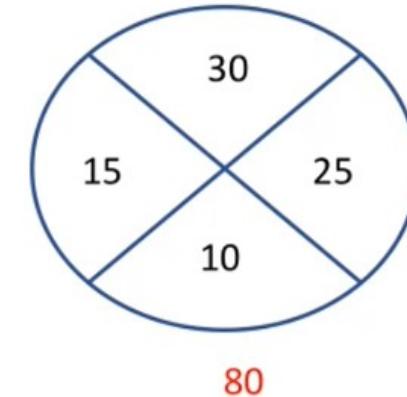
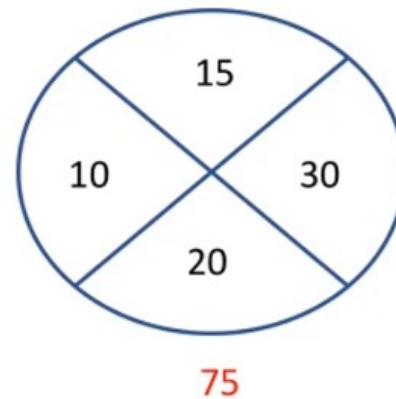
Sample 2



Total Reads are the same. Read counts
can be compared across samples



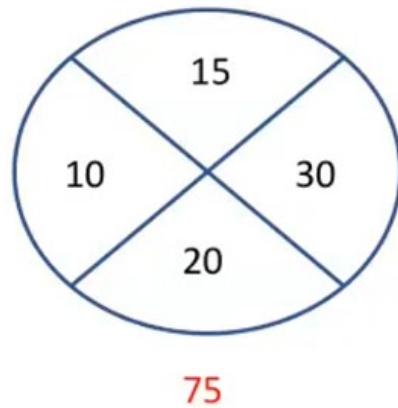
RPKM/FPKM



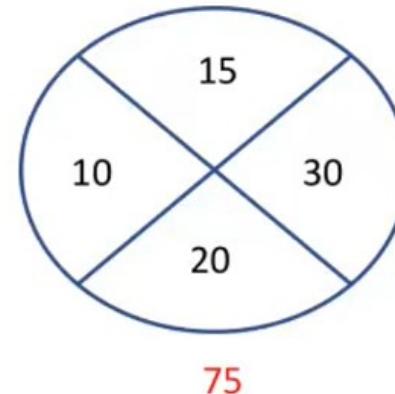
Total Reads are **not** the same. Read counts
can **not** be compared across samples

Sample 1

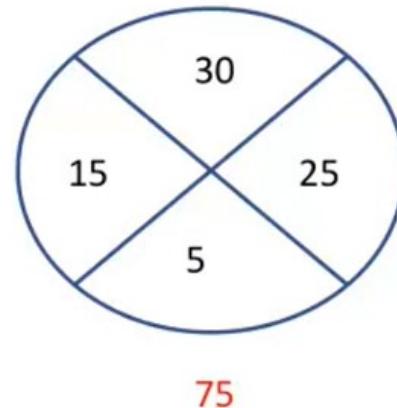
TPM



RPKM/FPKM



Sample 2



Total Reads are the same. Read counts can be compared across samples

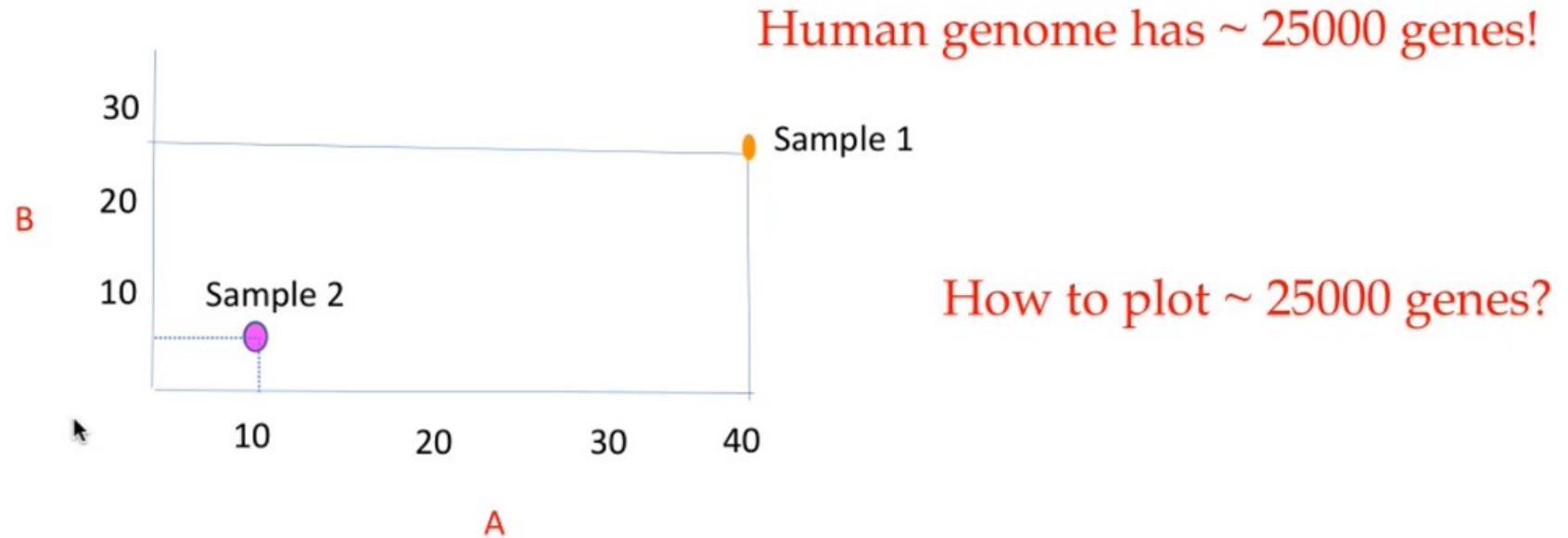


Total Reads are **not** the same. Read counts can **not** be compared across samples

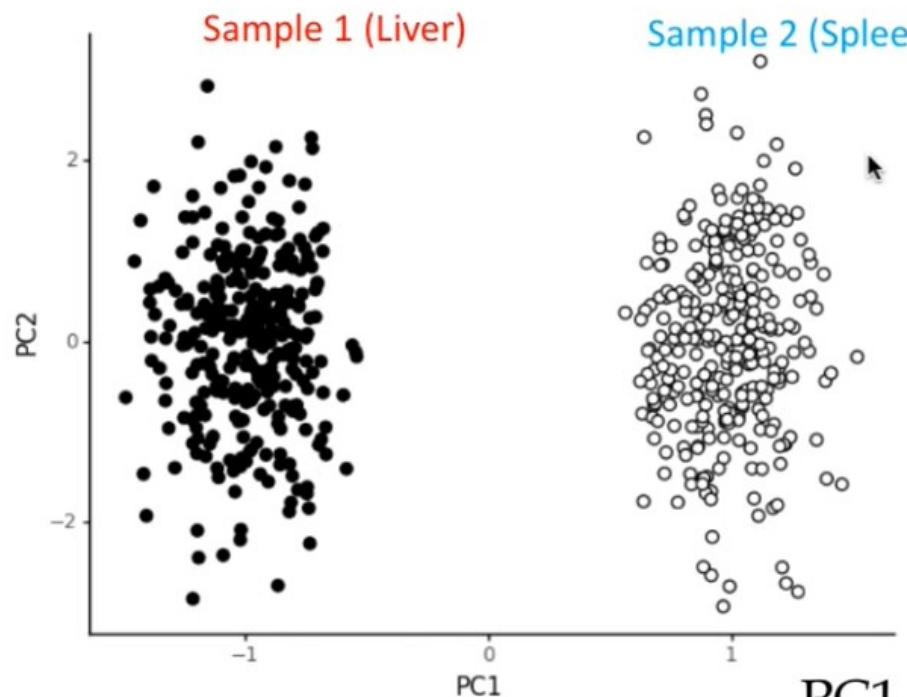


Plot the Data

Gene	Sample 1	Sample 2
A	40	10
B	25	5



Principal Component Analysis (PCA)



→ PCA reduces the dimensionality of the data

PC1 captures most of the variations in the data

PC2 captures second most variations in the data

Differential Gene Analysis

Problems in Library Normalisation

1. Differences in Library Size (Sequencing depth)
2. Differences in Library Composition

Problem 2 is not addressed by RPKM/FPKM/TPM

Differential Gene Analysis

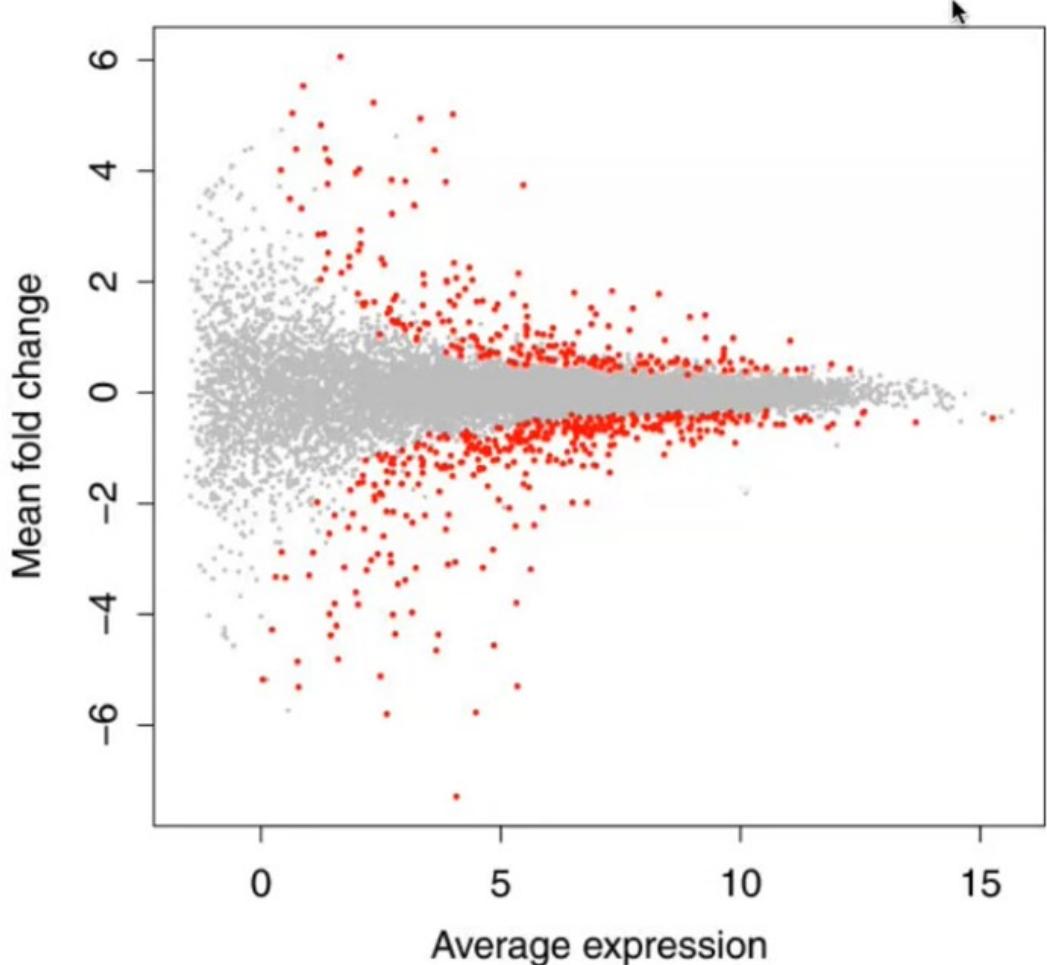
Differences in library composition

Gene name	Liver	Spleen
A	30	235
B	24	188
C	0	0
D	563	0
E	5	39
F	13	102
Total	635	635

← D is the only differentially expressed gene

For other genes in Spleen sample, the reads got distributed

Differential Gene Analysis (MA plot)



Mean Average

Gene expressions that are different
between WT and Mutant Sample

No significant change between WT
and Mutant sample

Horizontal : Left to Right

→ Low to High Expression

Pic Source: Fernandes et al, Microbiome, 2014



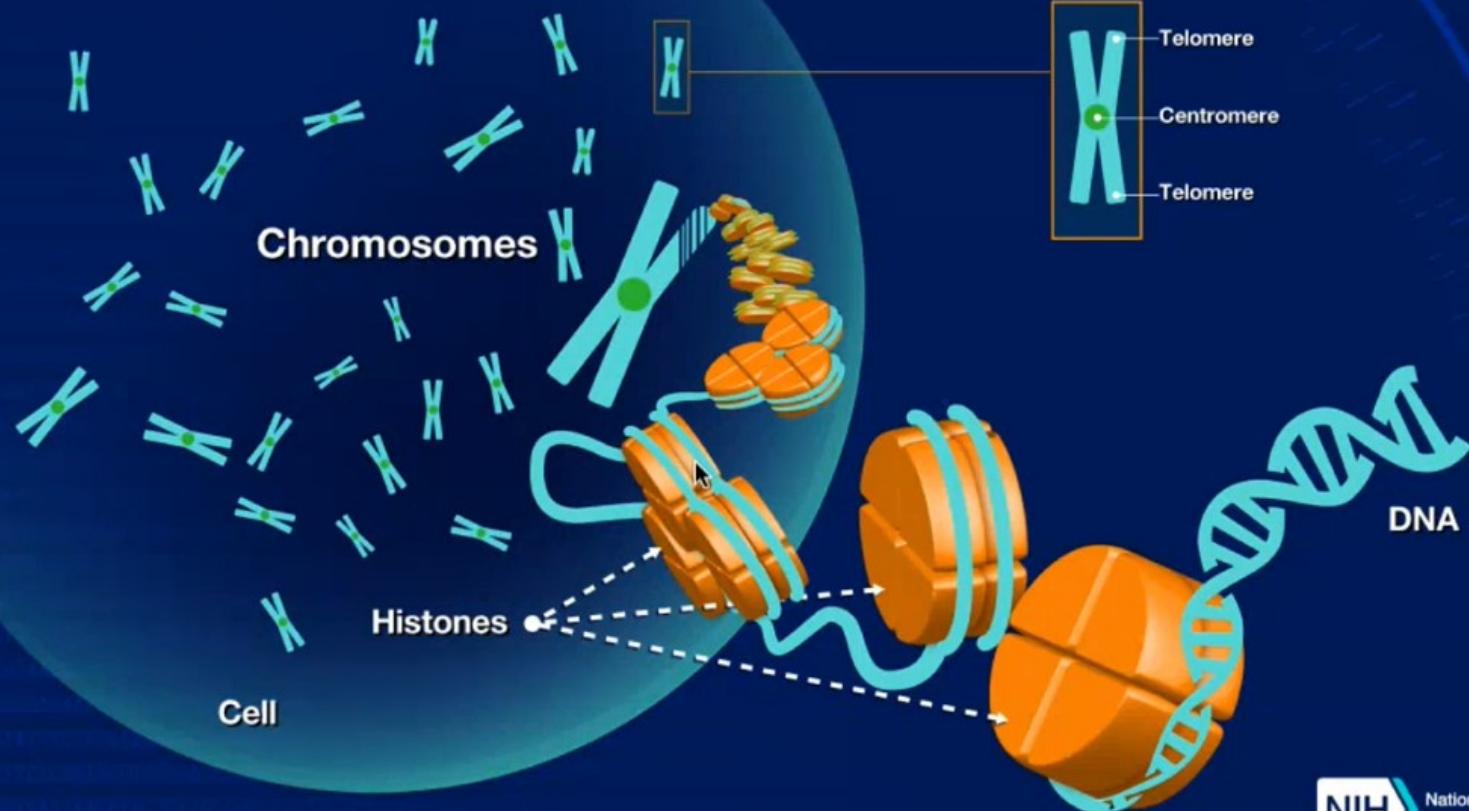
Big Data Analysis

CHIP-seq



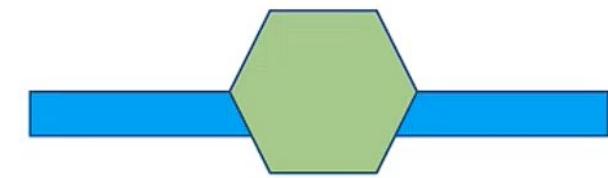
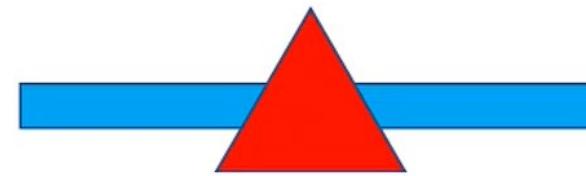
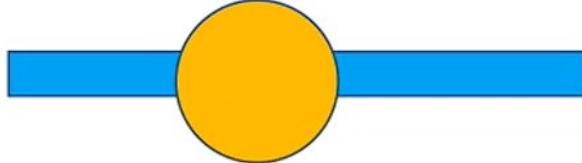
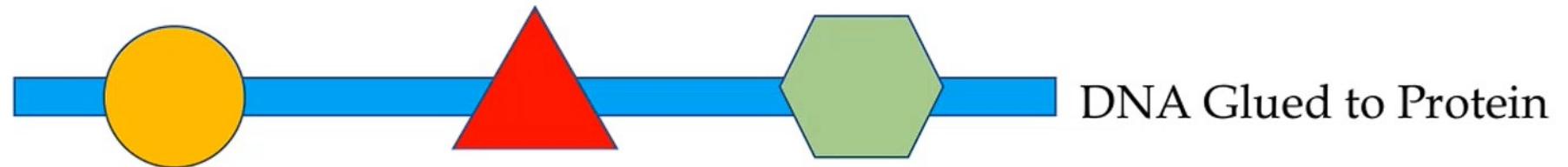
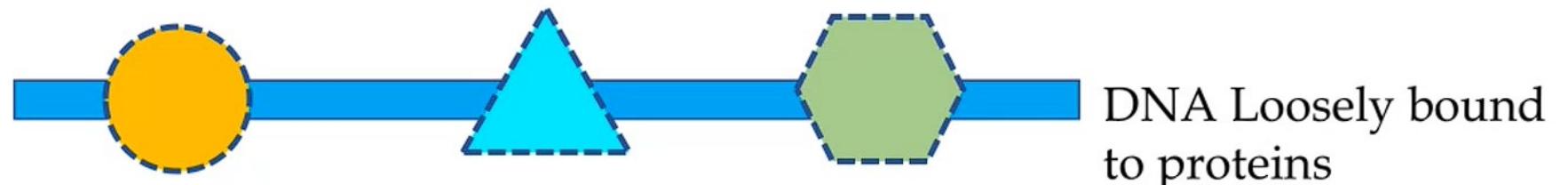
Chromosome

NHGRI FACT SHEETS
genome.gov



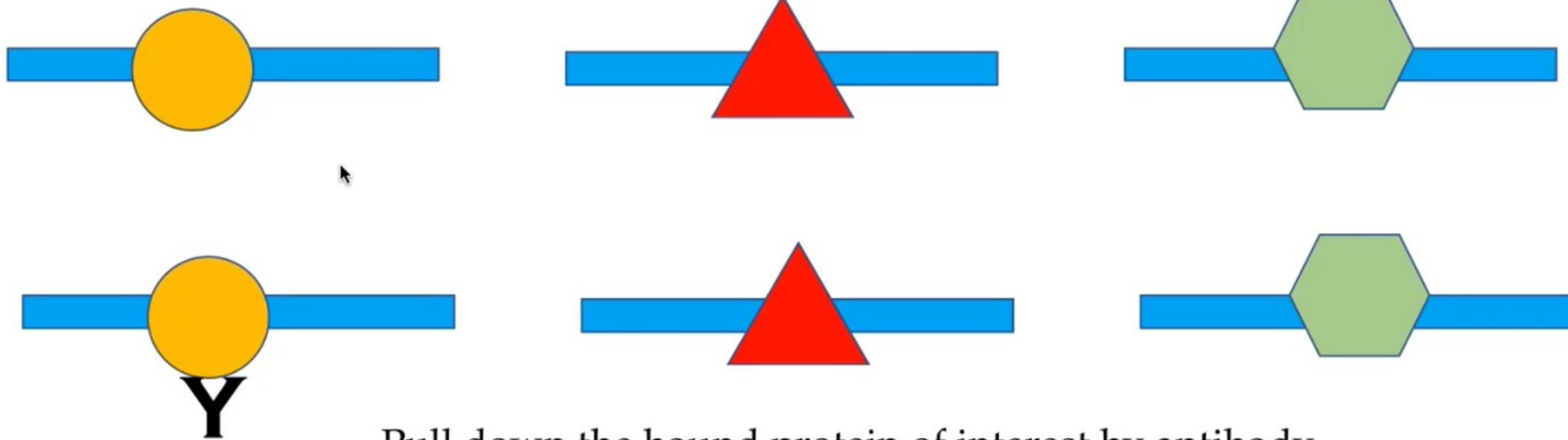
Chromatin = DNA + Histones + Auxiliary proteins

How do we find the region bound by protein of interest?

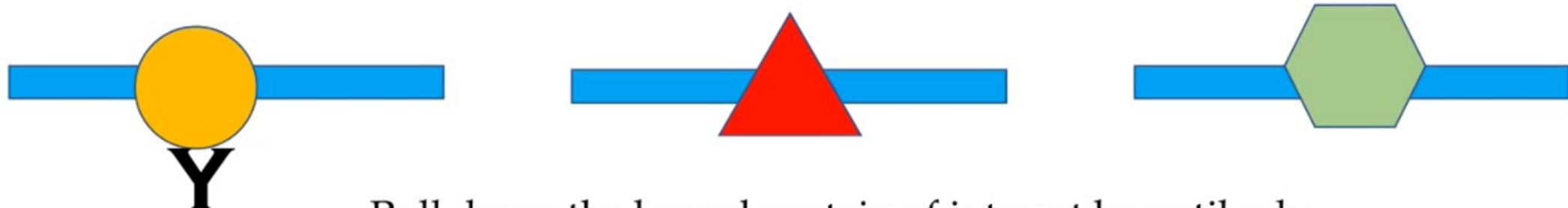


Fragmentation of DNA (300 bp)

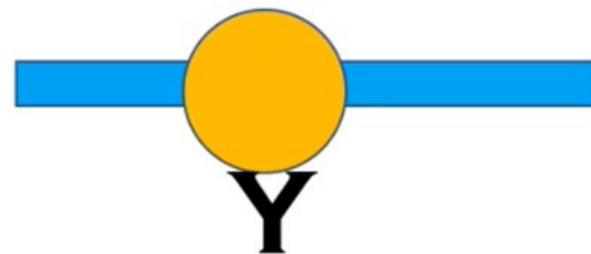
CHIP-seq - Chromatin Immunoprecipitation followed by high throughput sequencing



ChIP-seq - Chromatin Immunoprecipitation followed by high throughput sequencing



Pull down the bound protein of interest by antibody



Remove the glue, the protein of interest and the antibody



DNA fragment bound by protein of interest



Adapter ligation at the end and PCR amplification (Library Preparation)

Illumina sequencing (2nd Generation)

The “sequencing-by-synthesis” technology is used by Illumina. It was originally developed by Shankar Balasubramanian and David Klenerman at the University of Cambridge.

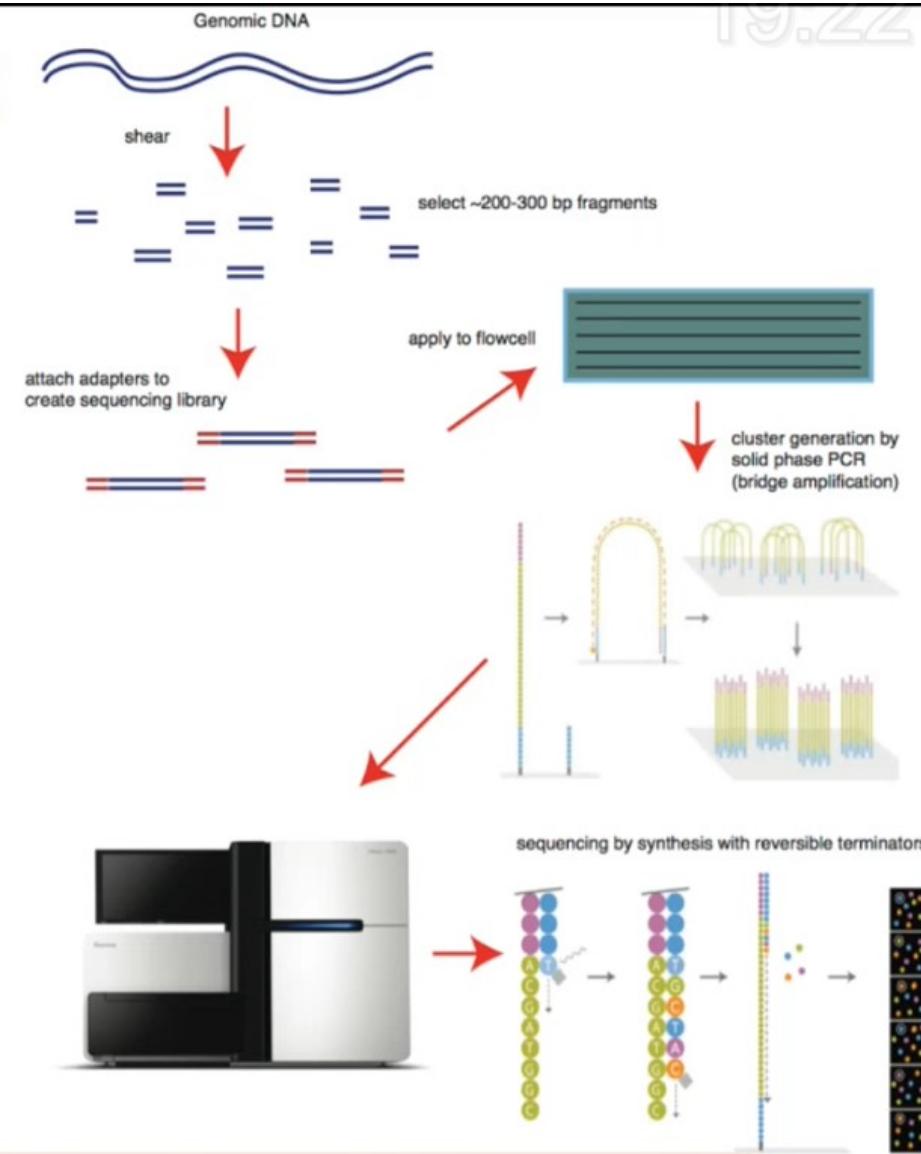
DNA or cDNA samples are randomly fragmented, usually into segments of 200 to 600 base pairs. These fragments are then ligated to adaptors and made single-stranded

Each fragment is amplified on the flow cell, and unlabeled nucleotides and polymerization enzymes are added. These additions, called “Bridge amplification,” connect and lengthen the fragments of DNA on the flow cell

Illumina’s “sequencing by synthesis” involves a proprietary method whereby four labeled reversible dNTP terminators, primers and DNA polymerase are added to the templates on the flow cell.

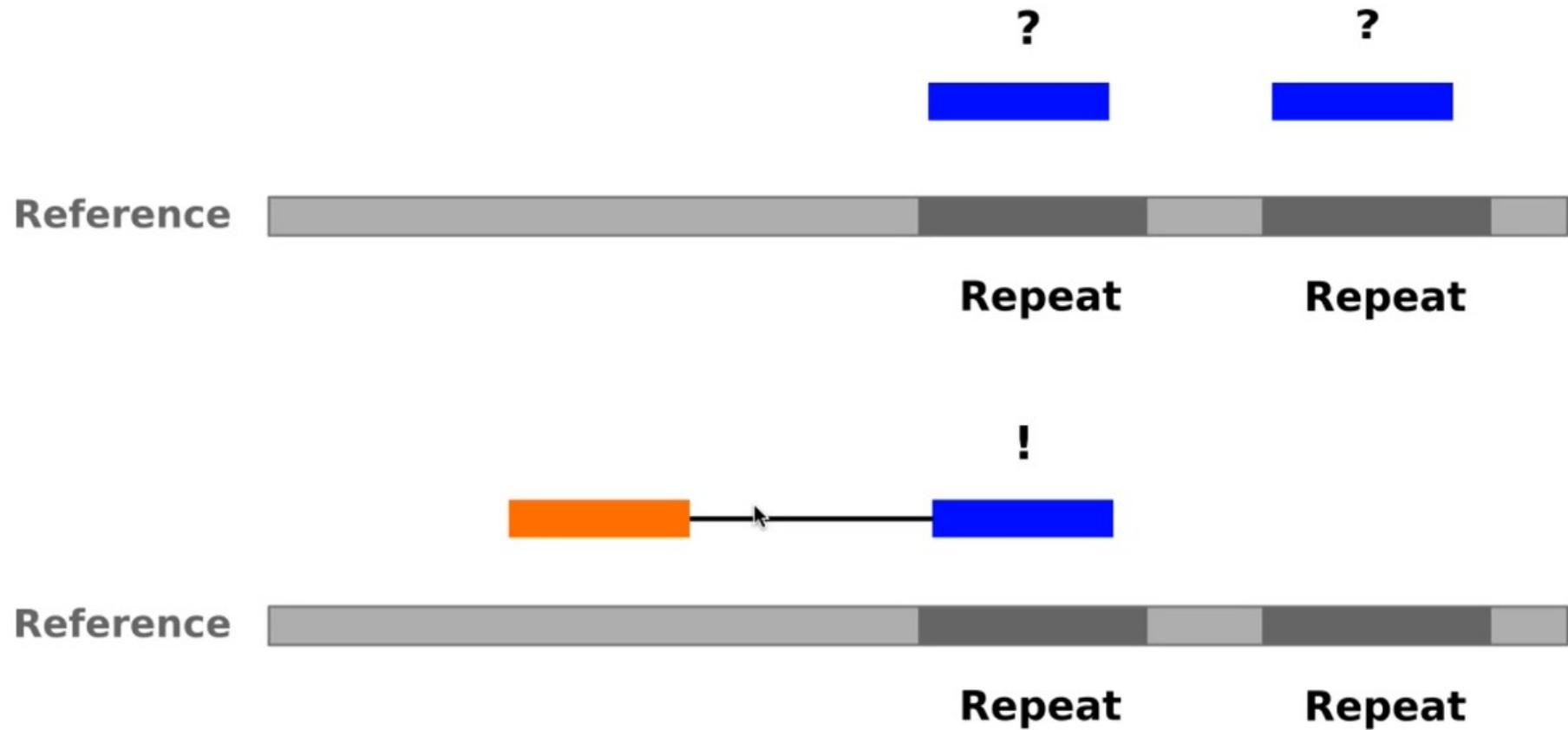
When excited by a laser, fluorescence from each cluster can be detected, which identifies the first base.

Base calls are made from signal intensity measurements during each cycle, reducing error rates further



<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Single Vs. Paired End Read



Steps in NGS Data Analysis



FASTQ Format

@SEQ_ID [Header Information Line starts with @ and followed by sequence identifier]

GATTGGGGTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT [Sequence]

+ [+ followed by sequence identifier or description]

!"*(((***+))%%%%++)(%%%%%).1***-+*'')**55CCF>>>>CCCCCCC65 [Quality score identifier]

Symbols = Low quality

Numbers = better

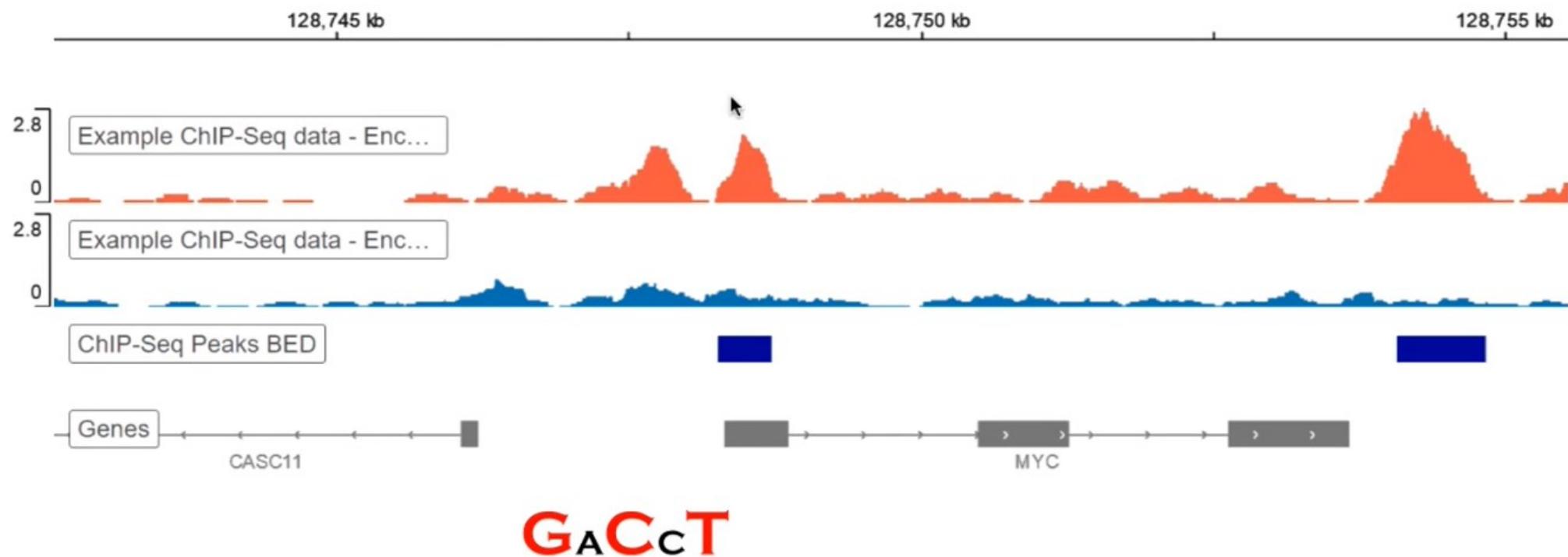
Alphabets = Best

Data Processing and Analysis Steps

1. Remove the low quality reads and adapters
2. Align the reads to the genome of interest (SAM/BAM format)
3. Count the number of reads per gene (Treated vs. Control)



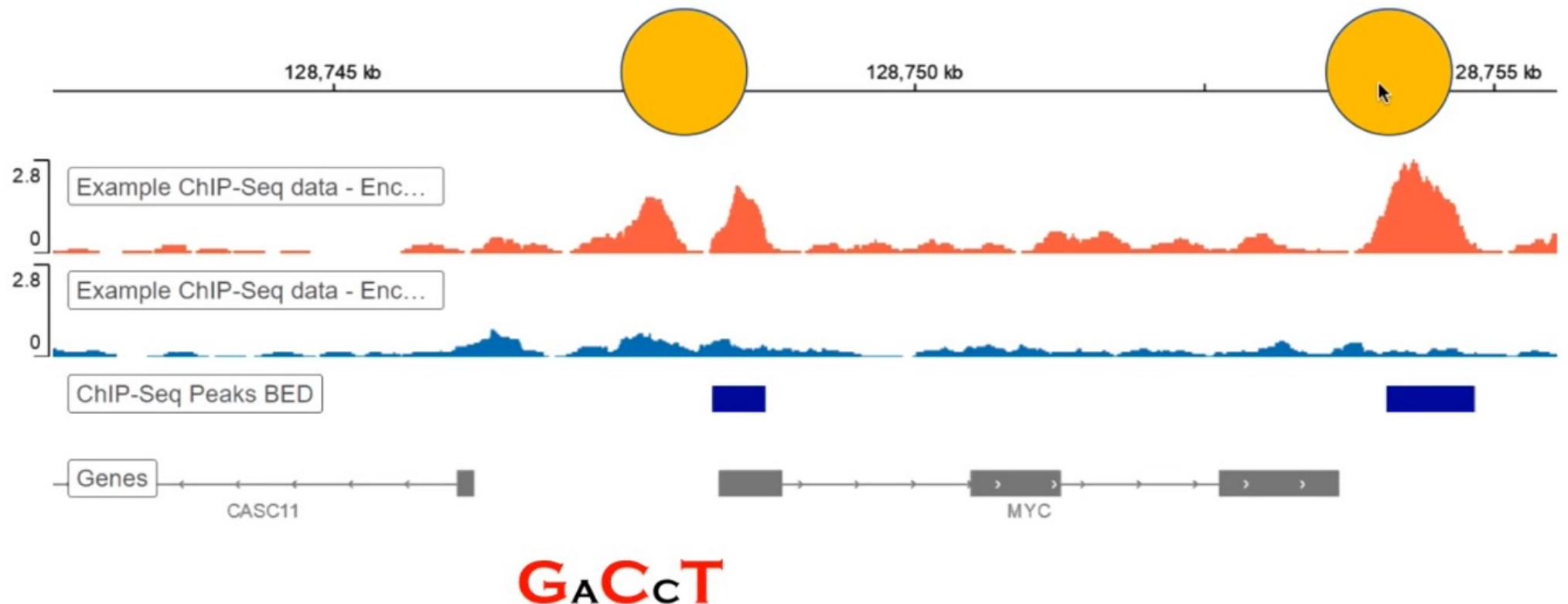
VISUALISING THE READS MAPPED ONTO THE REFERENCE GENOME

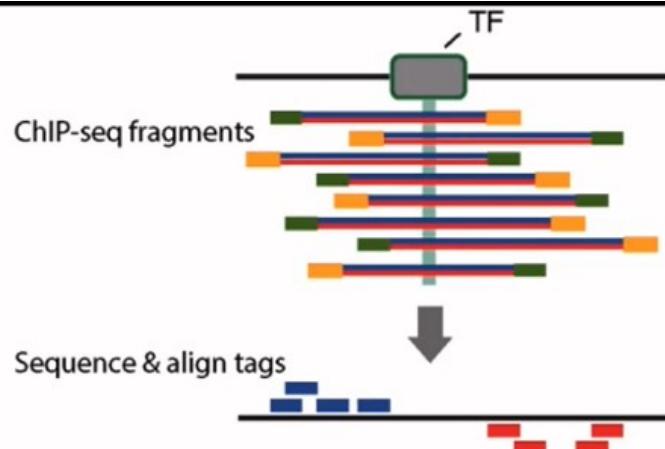


Basepairtech.com



VISUALISING THE READS MAPPED ONTO THE REFERENCE GENOME





PEAK CALLING

Peak-finding

- 1) Shift or extend tags
- 2) Build tag density landscape
- 3) Find max. locations

Peak-pairing

- 1) Build stranded tag density landscapes
- 2) Find max. locations on each strand
- 3) Pair opposite strand nearby peaks

Probabilistic binding detection

- 1) Probabilistically assign tags to binding events
- 2) Update binding event locations and model
- 3) Iterate between 1) & 2) until algorithm converges



Phylogenetic Analysis





What is Phylogenetics?

- The taxonomical classification of organisms on the basis of their degree of evolutionary relatedness.
 - originated by Willi Henning, a German zoologist. *Grundzüge einer Theorie der phylogenetischen Systematik* (1950),
 - *A Dictionary of Earth Sciences*, © Oxford University Press 1999



Tree of Life

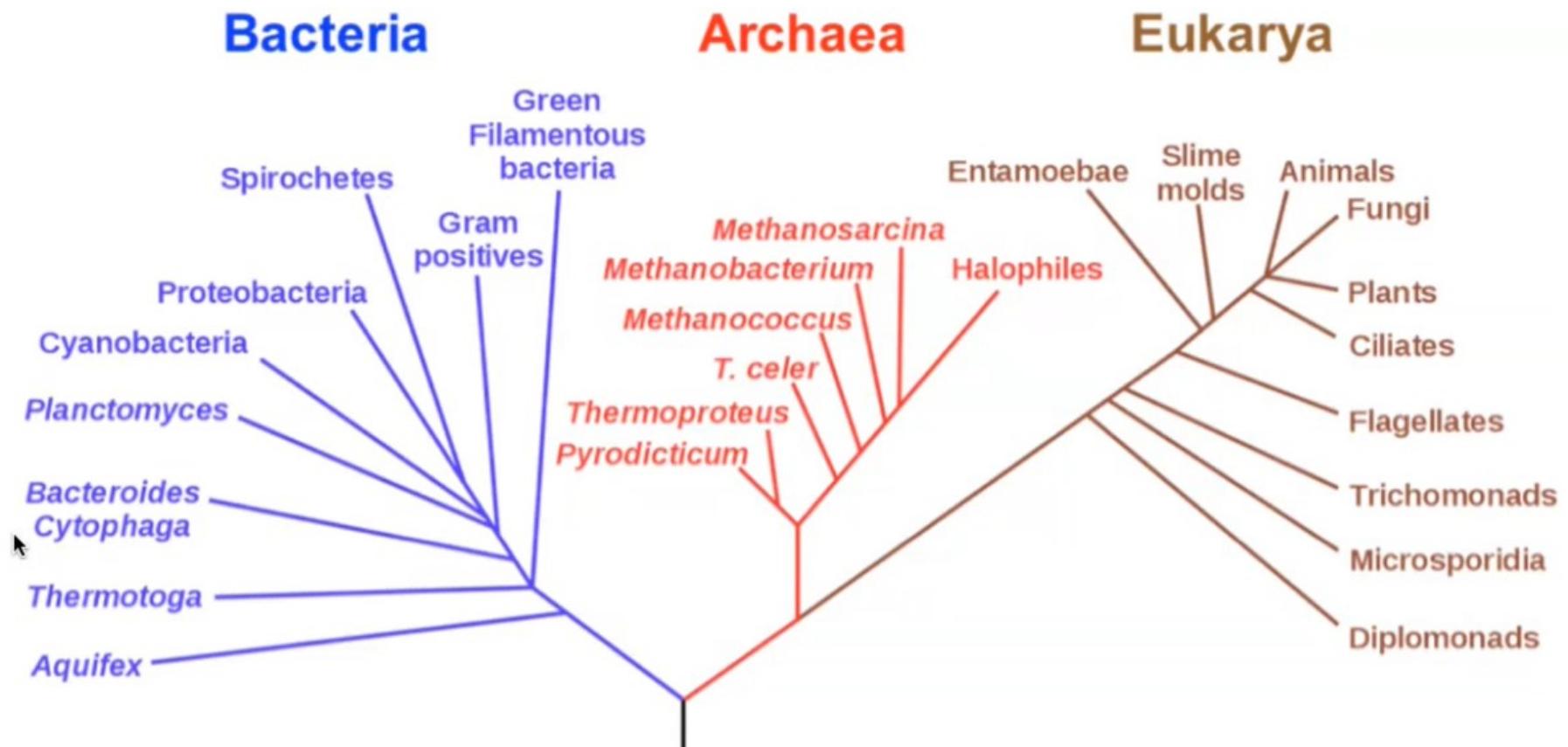


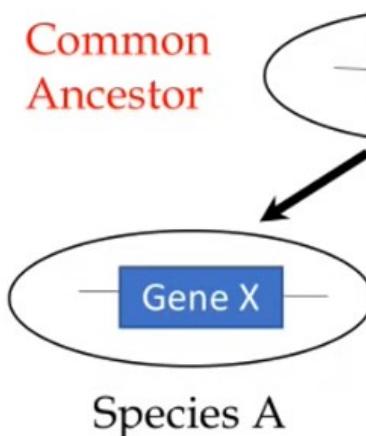
Image credit: NASA Astrobiology Institute & Eric Gaba

Evolutionary Relationships

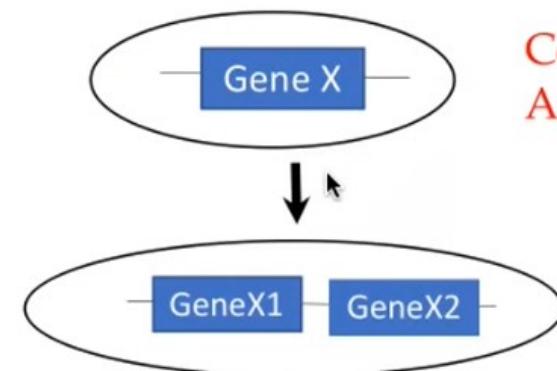
PAST
↓
PRESENT



HOMOLOGY:
Descent from a common ancestor



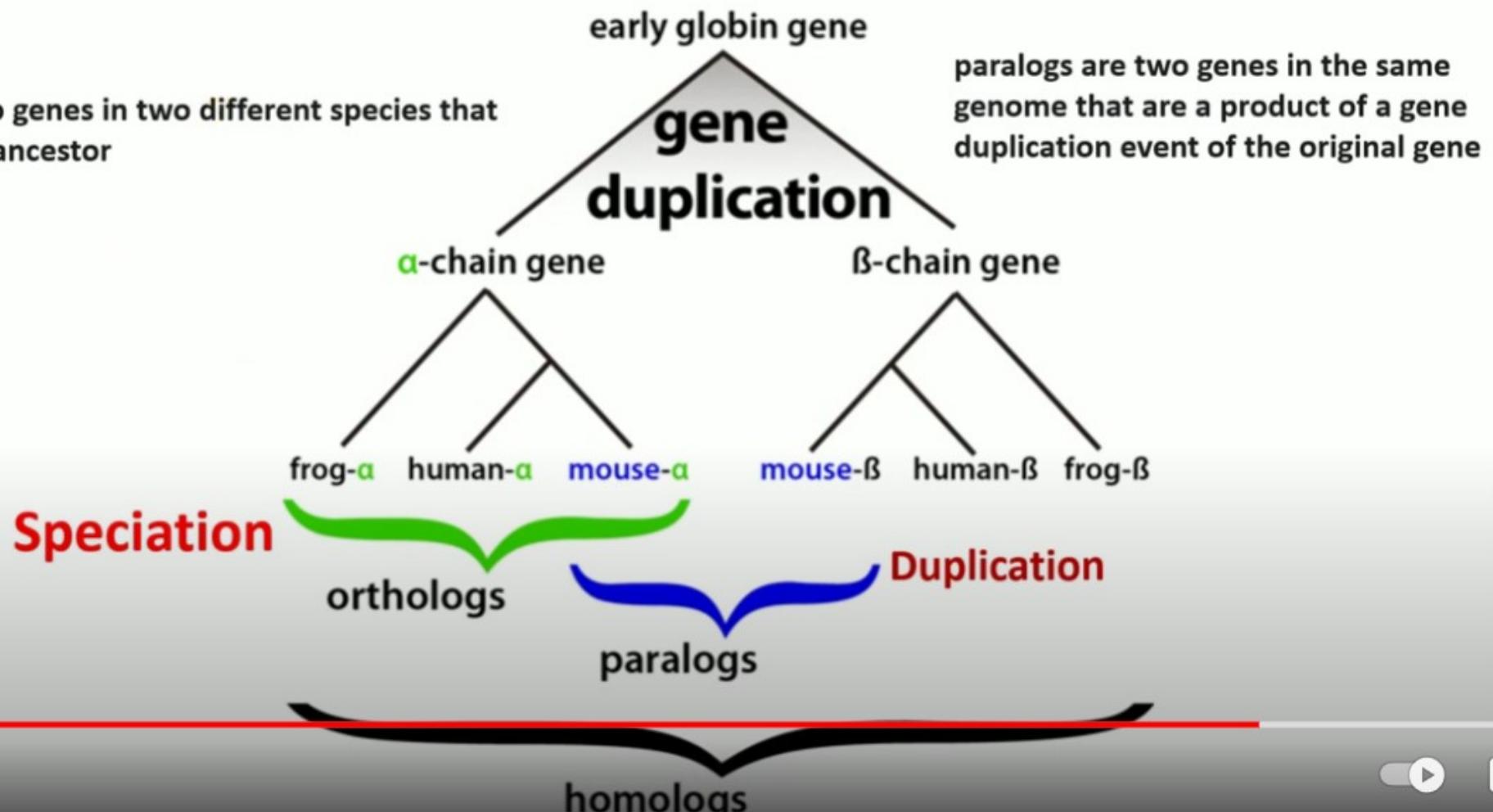
ORTHOLOGY: Descent from a common ancestor via speciation



PARALOGY: Descent from a common ancestor via gene duplication

Homologs are the genes or proteins that are similar due to their shared ancestor or common origin

Orthologs are two genes in two different species that share a common ancestor



Similarity: Homology vs Analogy

Homology: Similarity in characteristics resulting from shared ancestry.

Analogy: The similarity of characteristics between two species that are not closely related; attributable to convergent evolution.

Similar due to inheritance

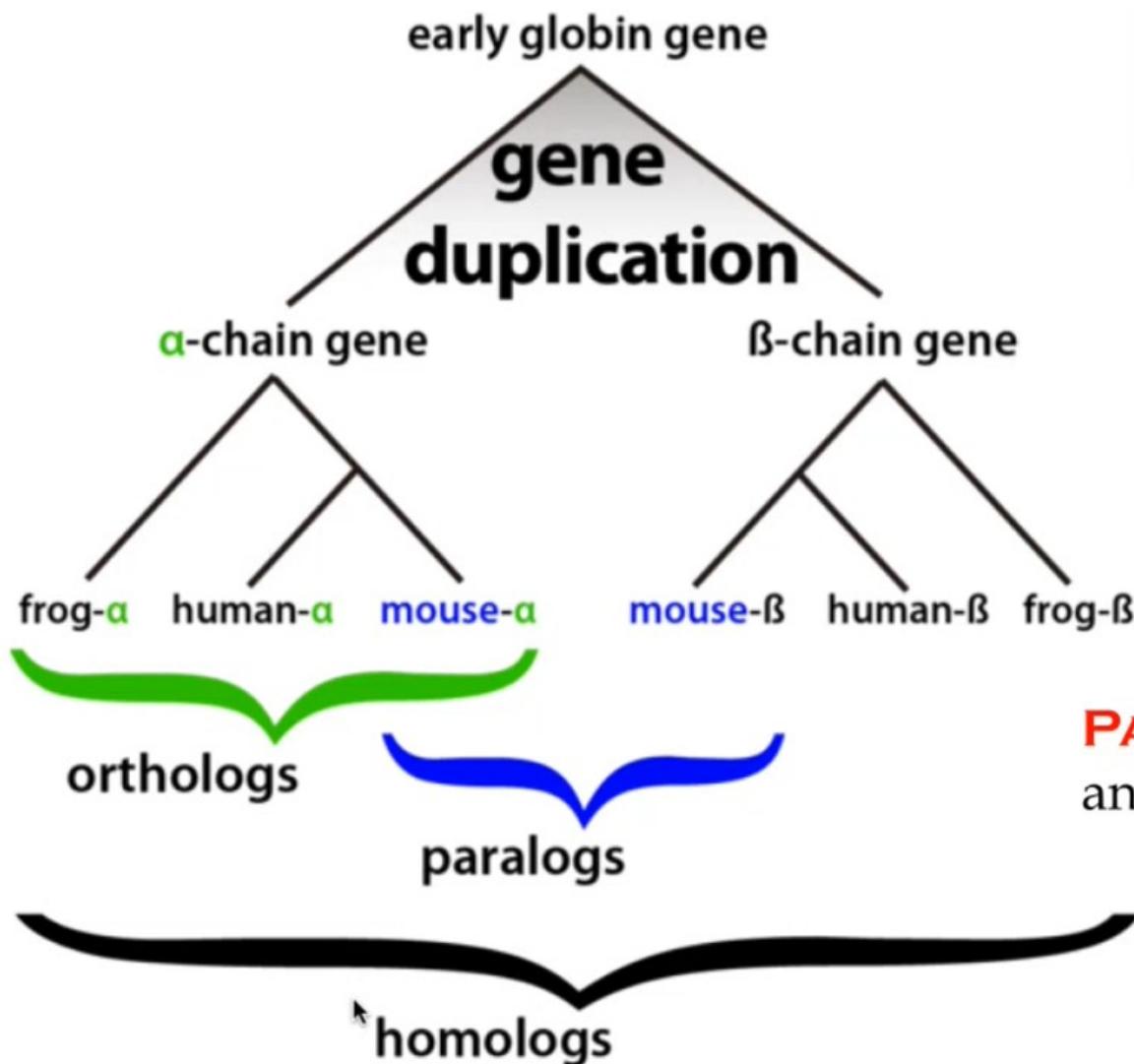


Two sisters: homologs

*Similar due to...
uh...other factors*



Two "Elvis": analogs



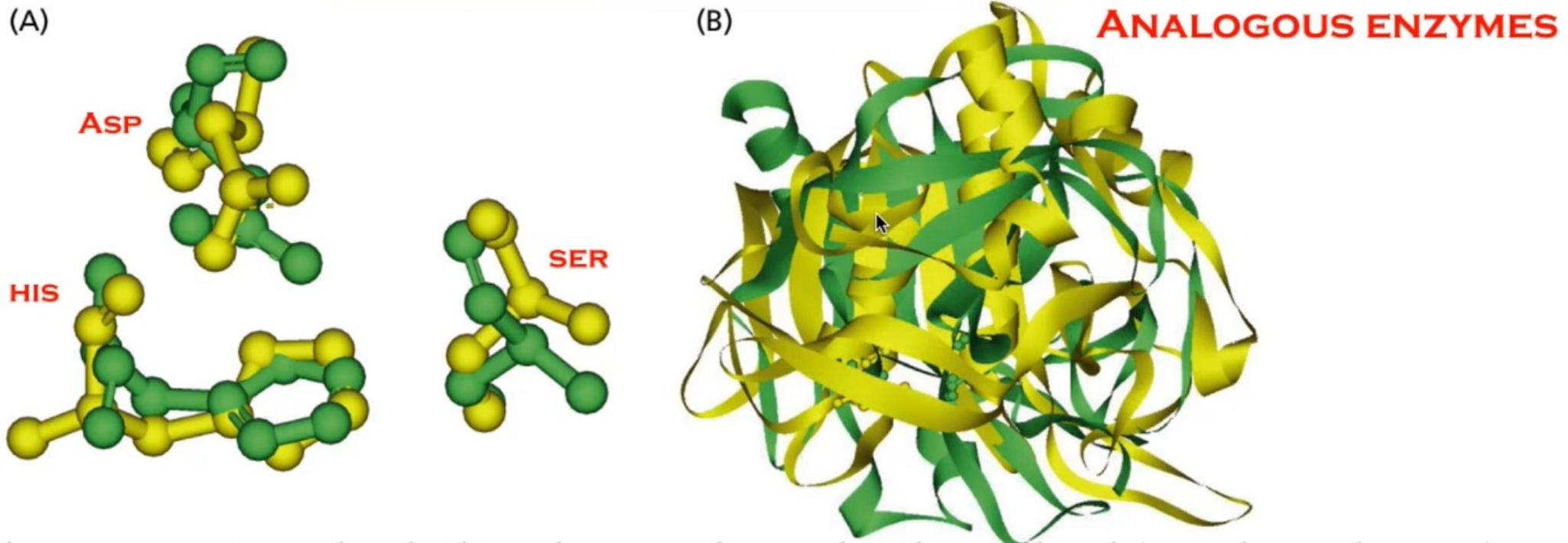
Homology
(Divergent Evolution)

ORTHOLOG: Descent from a common ancestor via speciation

PARALOG: Descent from a common ancestor via gene duplication

Image Credit: Popo H. Liao

Homoplasy (Convergent Evolution)



Chymotrypsin and subtilisin have independently evolved (non-homologous). They have different structure that doesn't overlay with each other. But, they have catalytic residues that independently evolved to retain the same catalytic triad.

GOAL OF PHYLOGENETIC ANALYSIS

TO RECONSTRUCT THE EVOLUTIONARY HISTORY



A phylogenetic tree is a diagram that proposes an hypothesis for the reconstructed evolutionary relationships between a set of objects.



Xenology

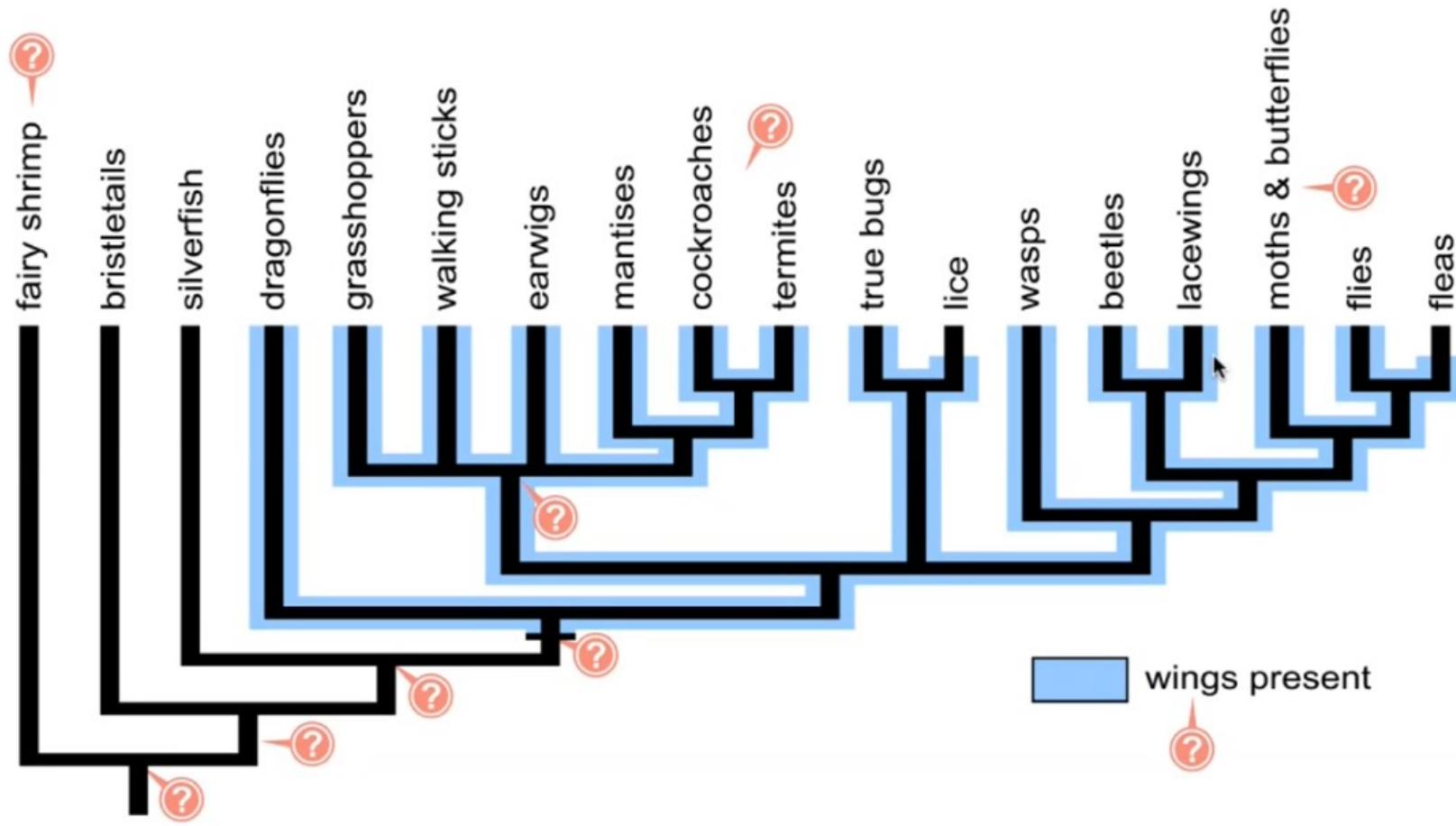
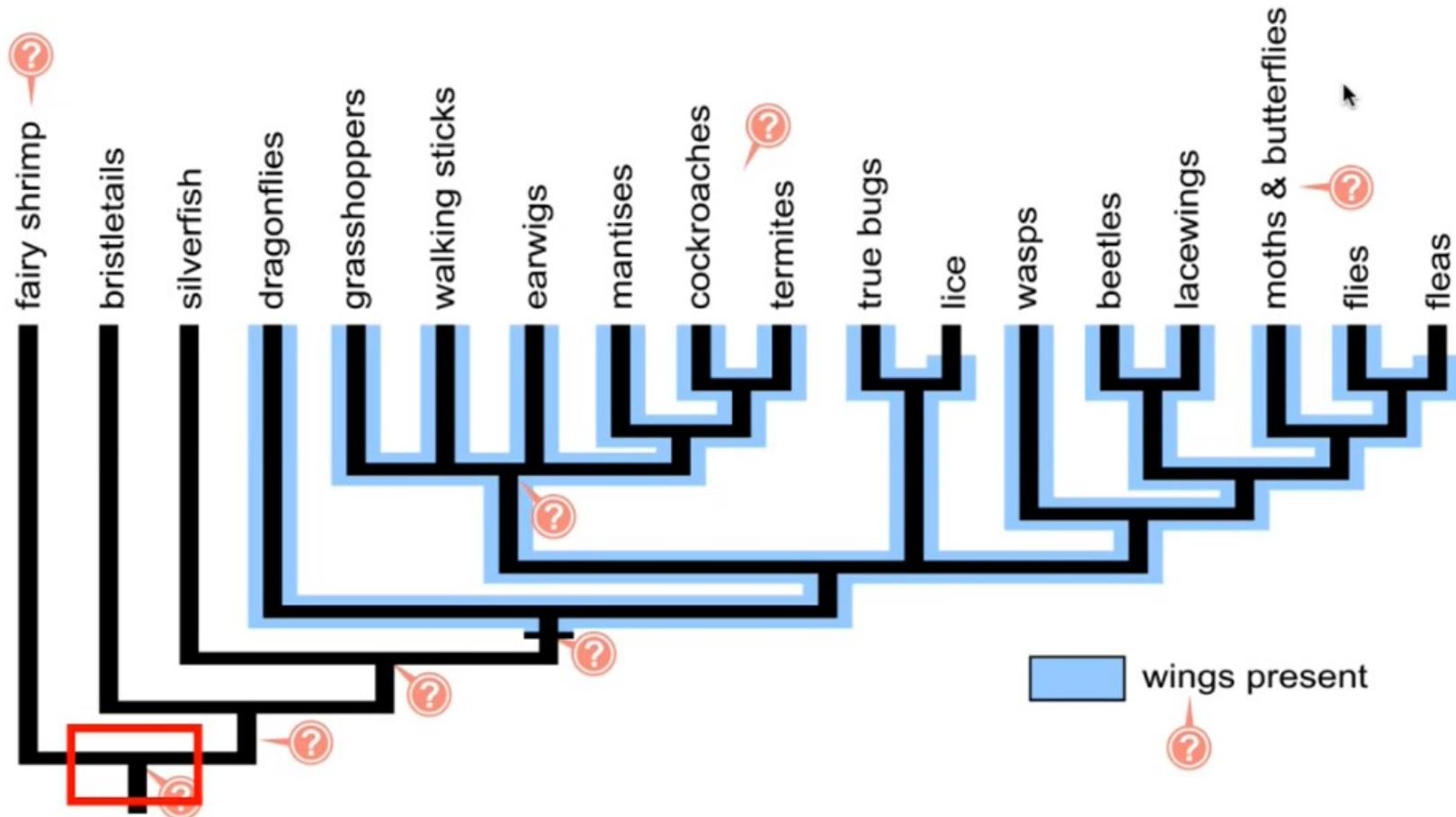


Image Credit: evolution.berkeley.edu

**TREES
DEPICT
EVOLUTIONARY RELATIONSHIPS
NOT
EVOLUTIONARY PROGRESS**





ROOT - the common ancestor of all the species shown here



Image Credit: evolution.berkeley.edu

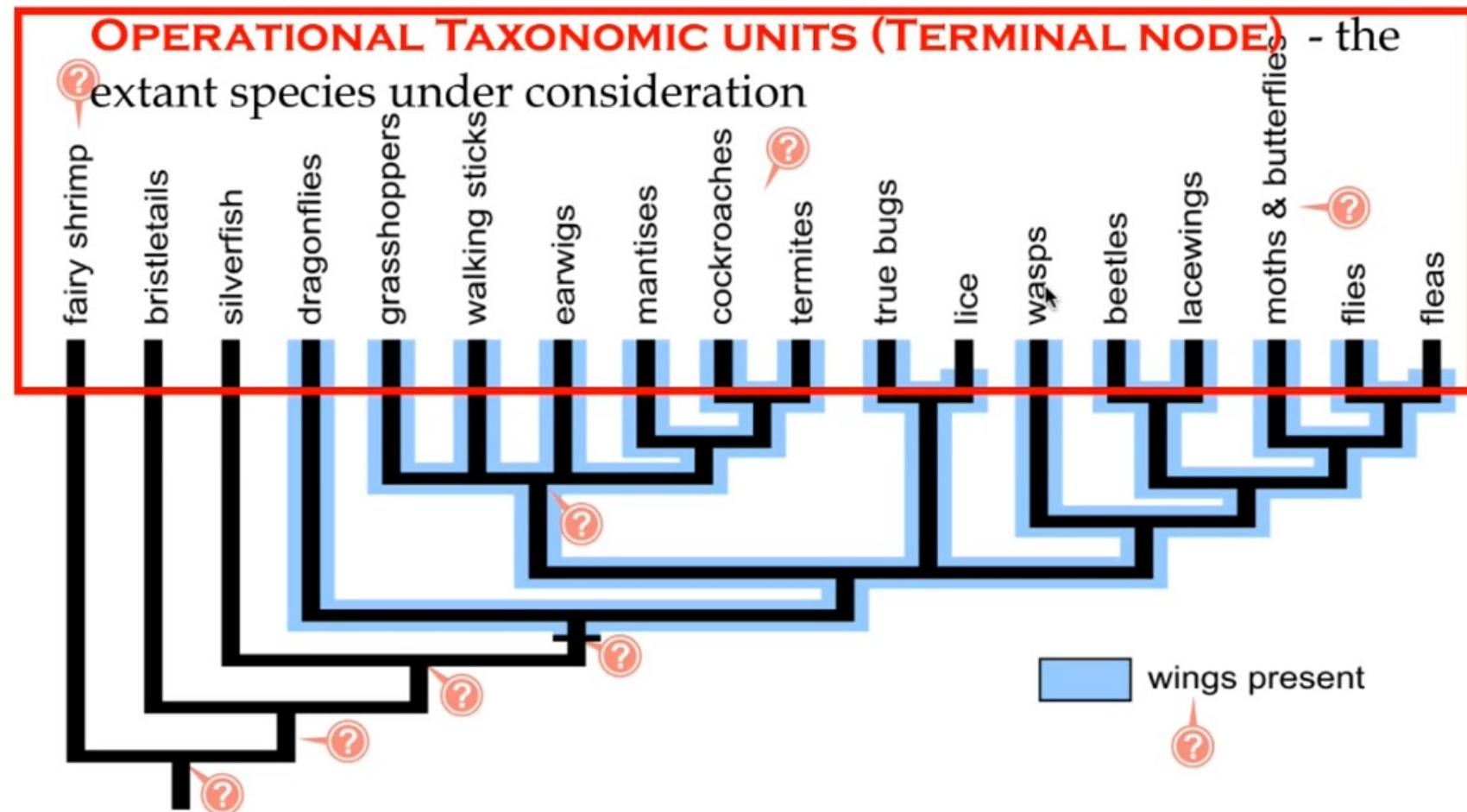


Image Credit: evolution.berkeley.edu

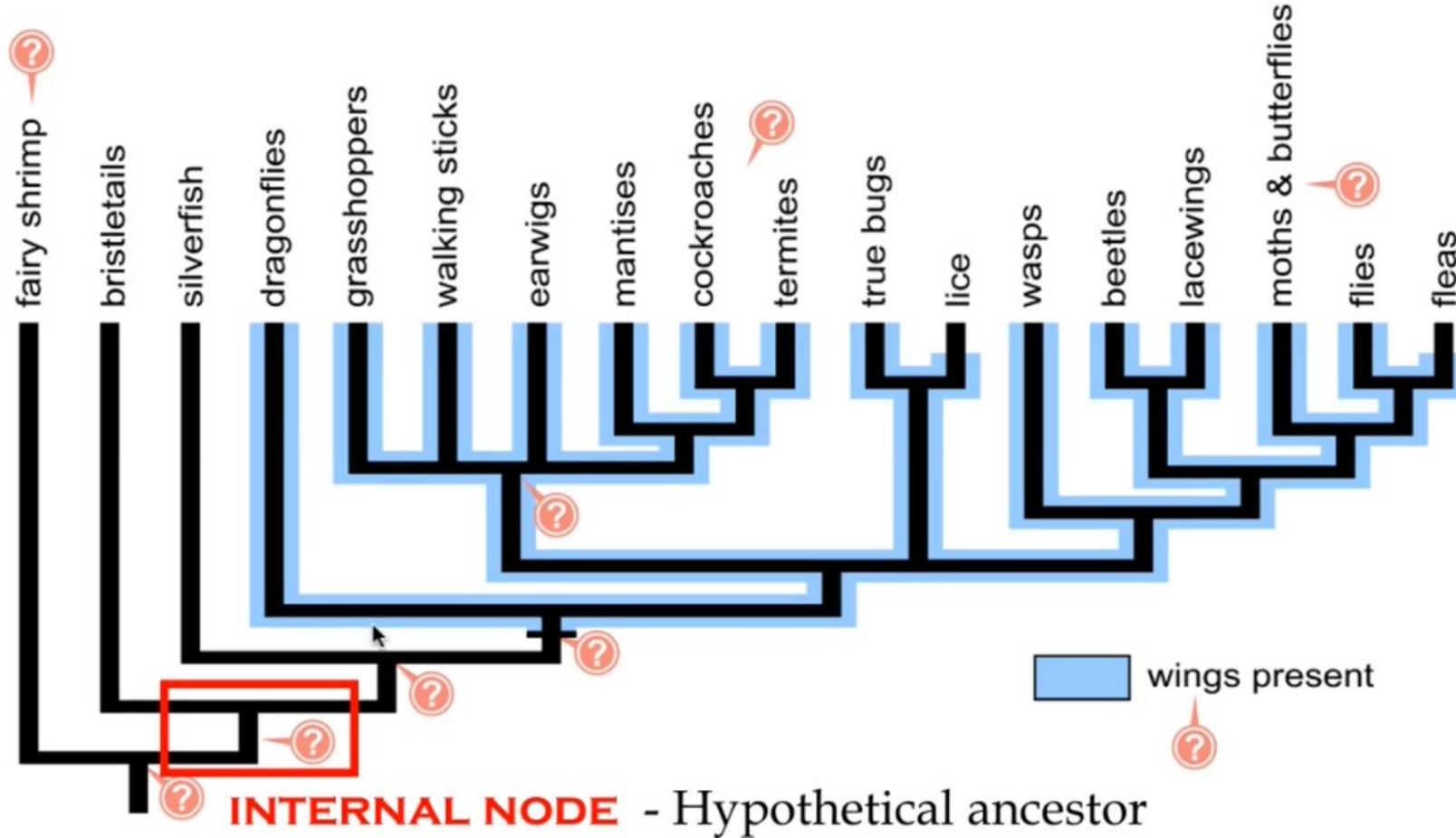
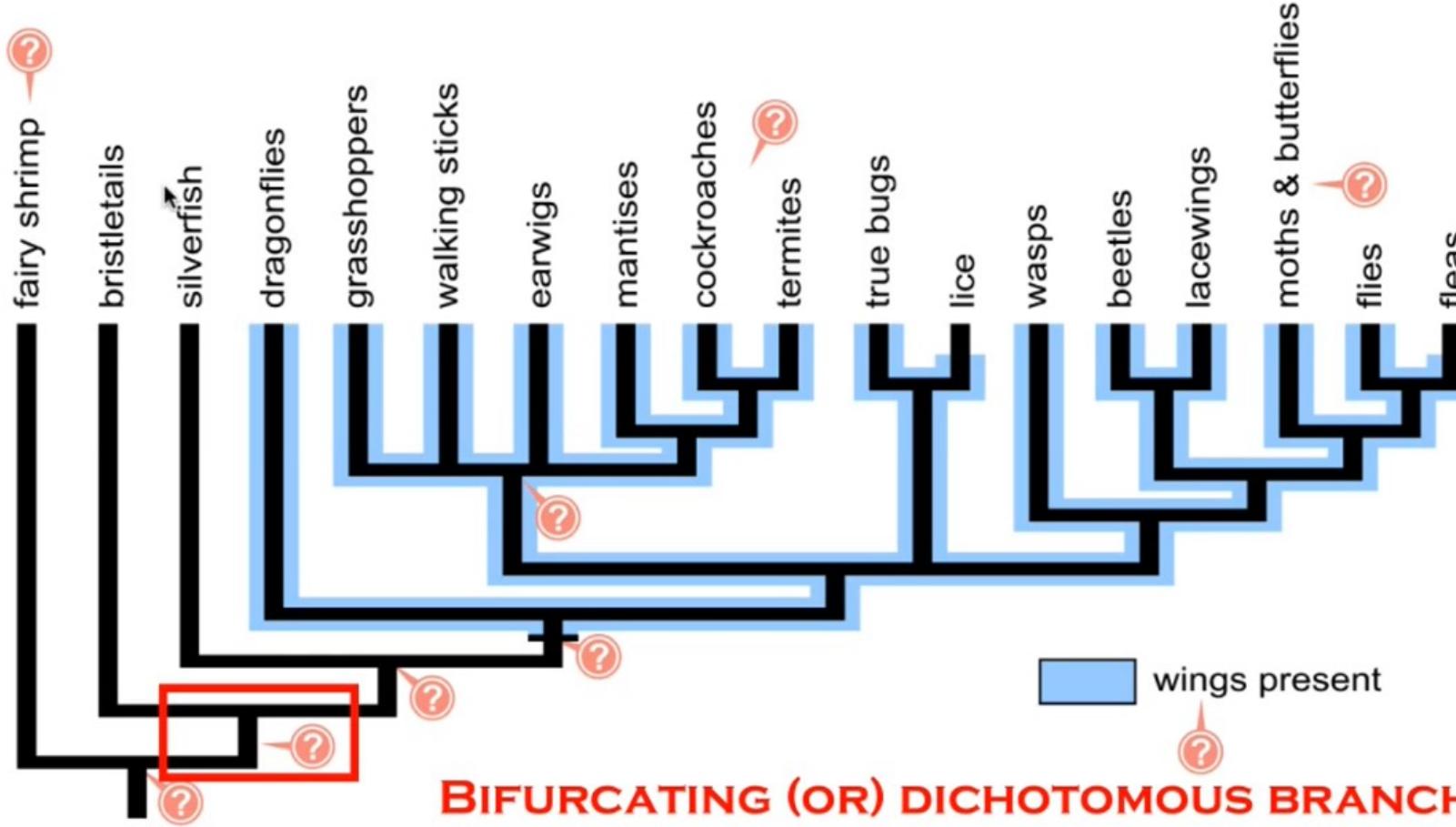
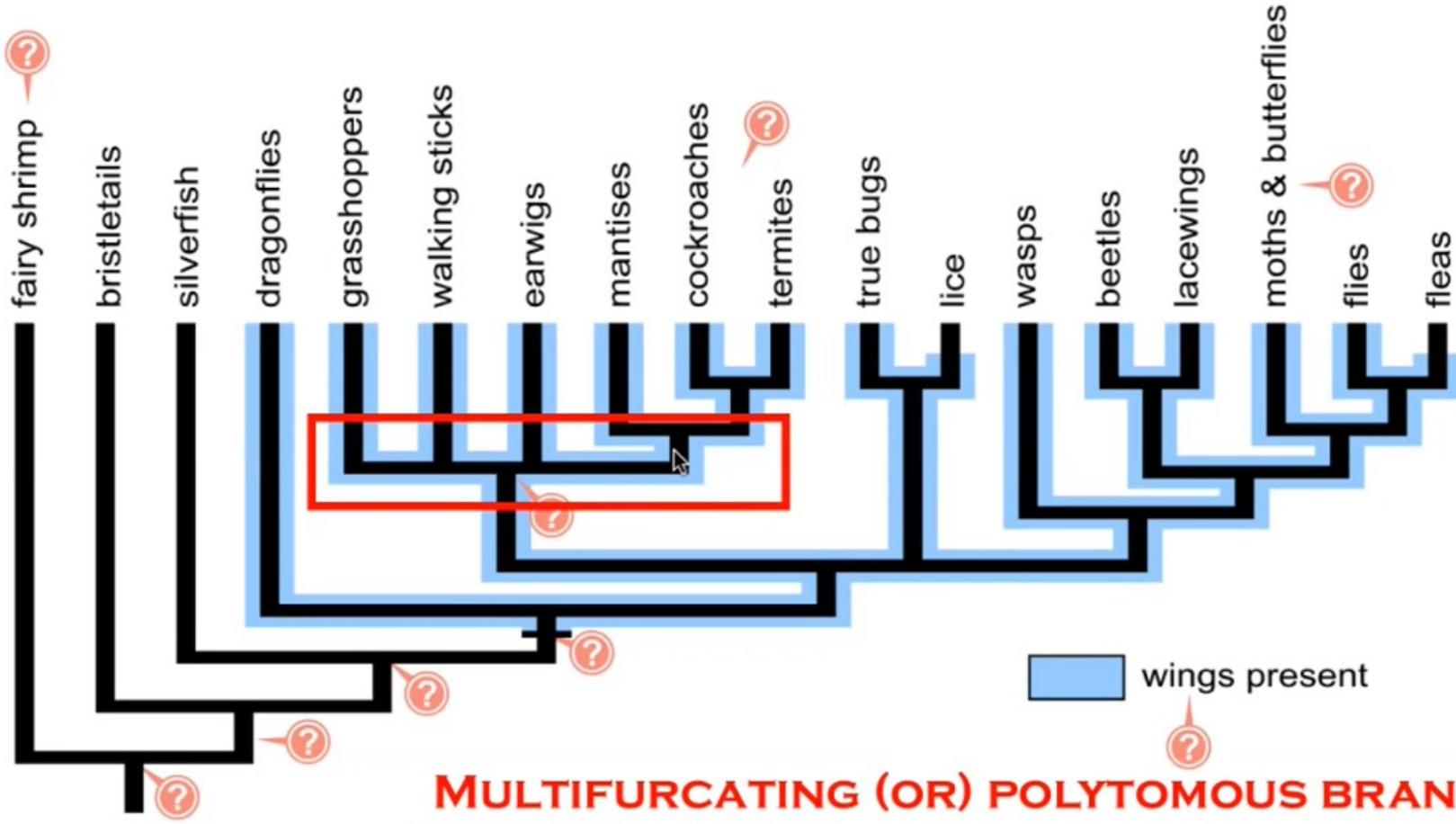


Image Credit: evolution.berkeley.edu



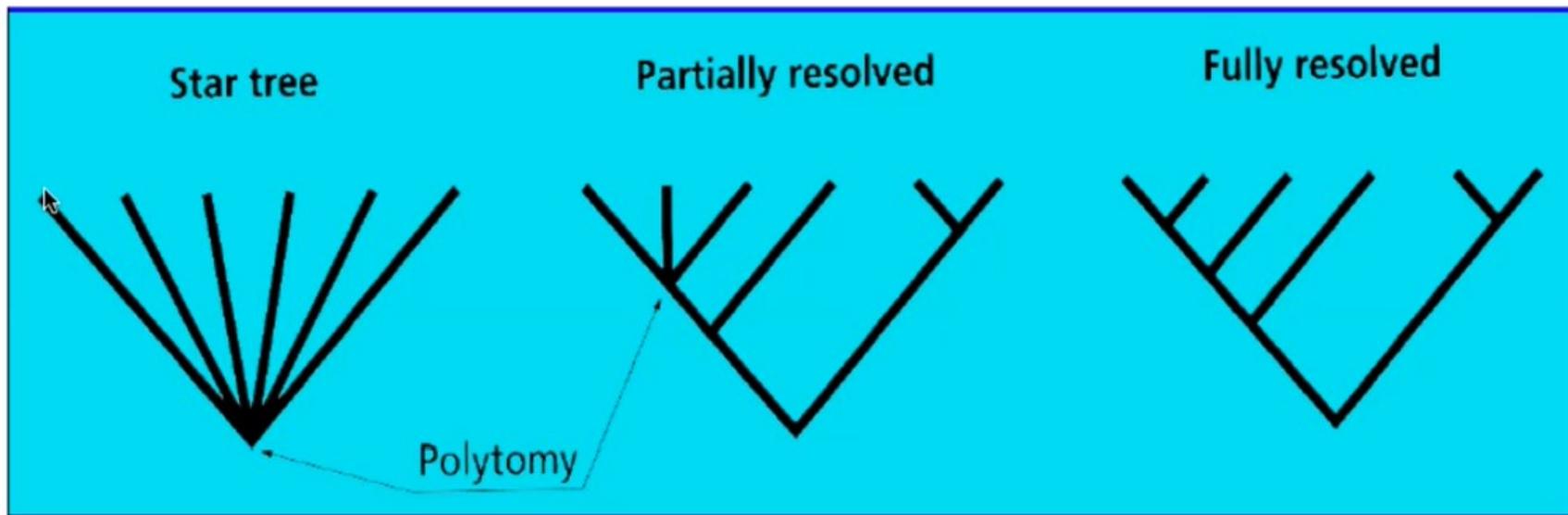
BIFURCATING (OR) DICHOTOMOUS BRANCHING PATTERN - Each ancestral node giving rise to **TWO** descendant nodes only

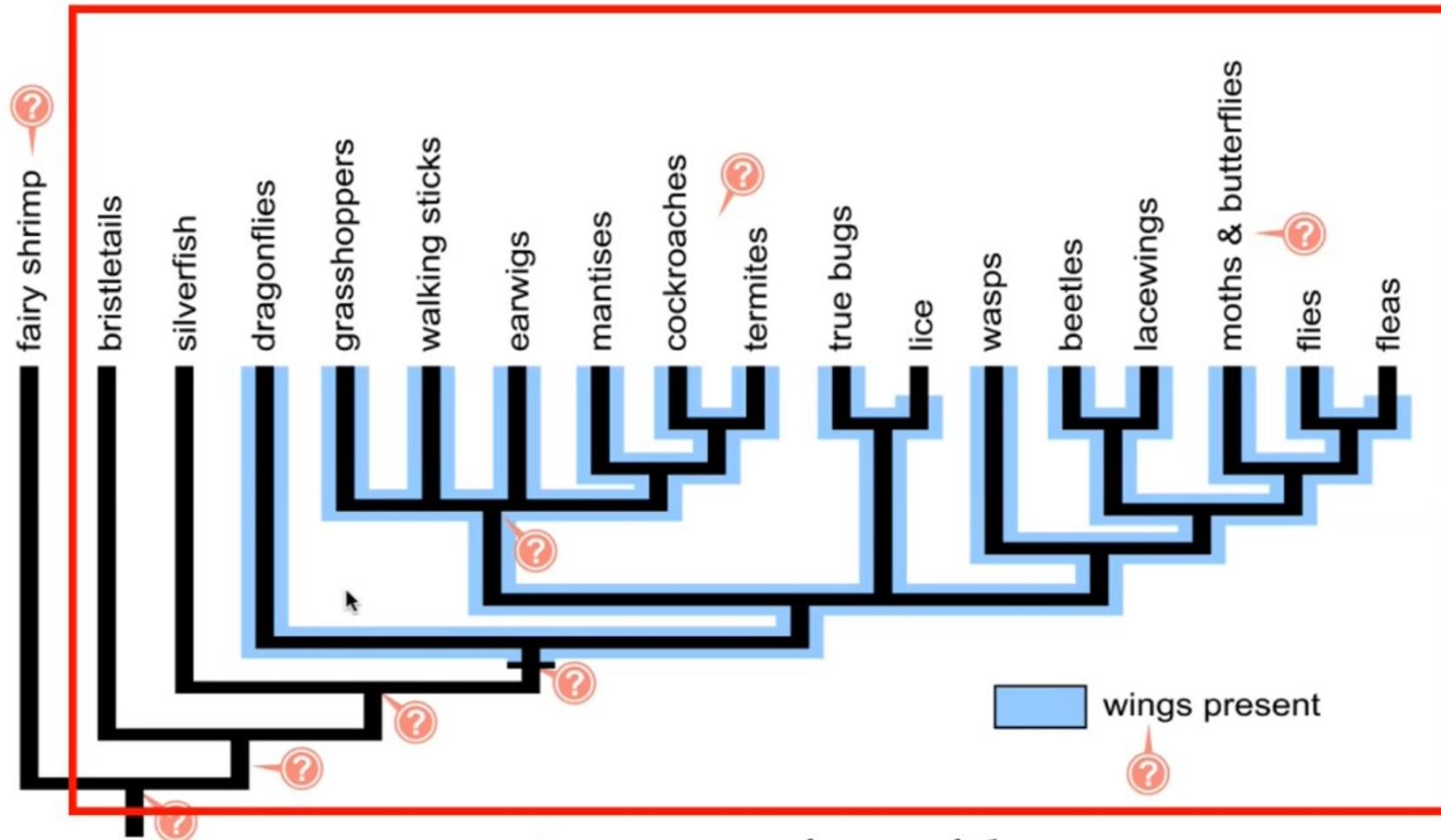
Image Credit: evolution.berkeley.edu



MULTIFURCATING (OR) POLYTOMOUS BRANCHING PATTERN – Each ancestral node giving rise to **MORE THAN TWO** descendant nodes

Image Credit: evolution.berkeley.edu

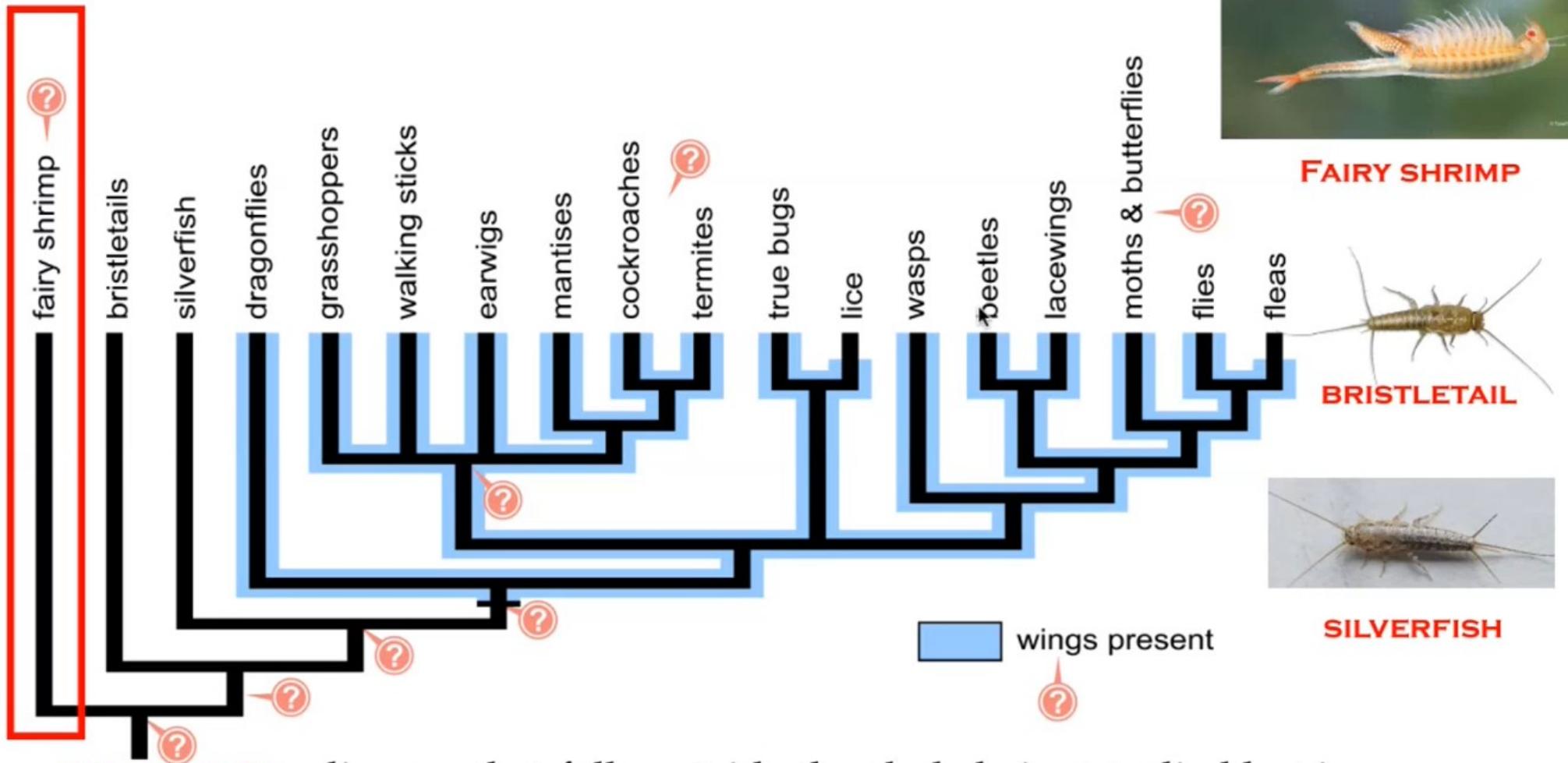




INGROUP - focus of the tree



Image Credit: evolution.berkeley.edu

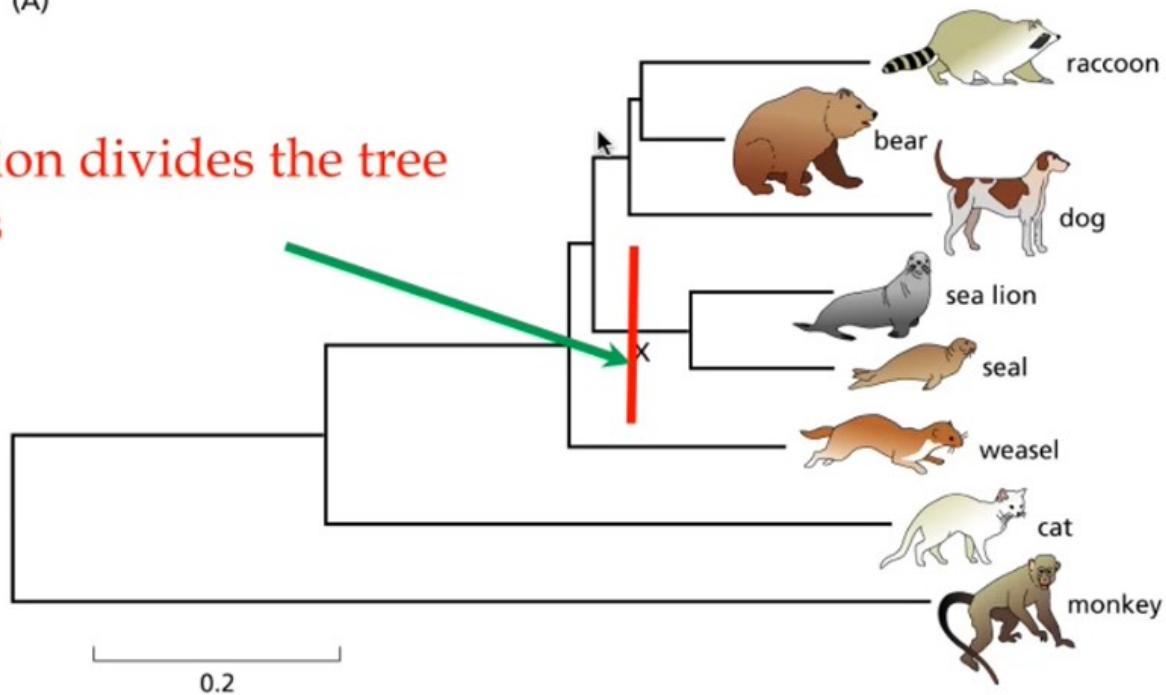


OUTGROUP - lineage that falls outside the clade being studied but is closely related to that clade

Image Credit: evolution.berkeley.edu

(A)

Split or Partition divides the tree into two parts

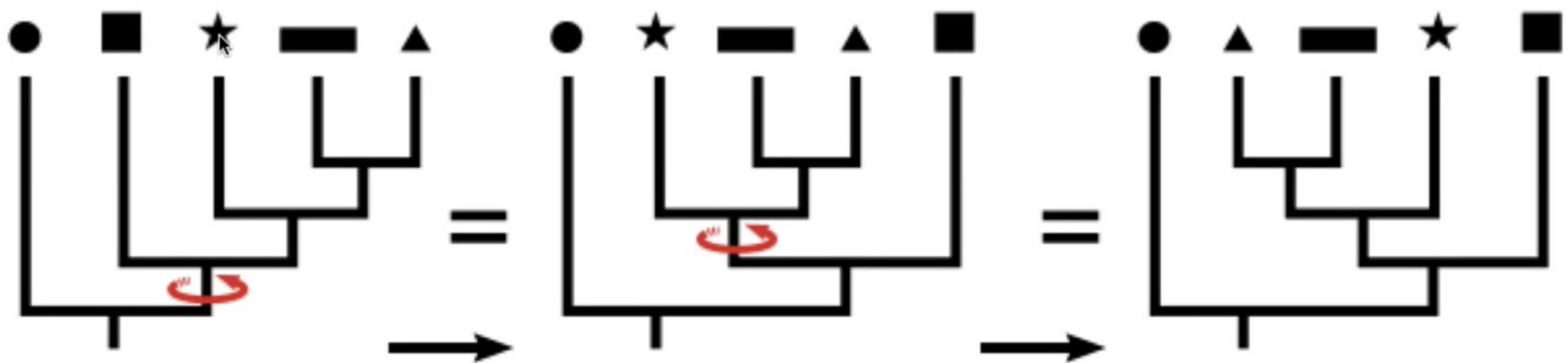


(B)

raccoon	bear	dog	sea lion	seal	weasel	cat	monkey
*	*						
*	*	*		*	*		
*	*	*	*	*	*		
*	*	*	*	*	*	*	

Every split produces two groups. Hence one group is suffice to be defined.

Rotations about the branches change the order of the taxa but NOT the evolutionary relationships



Rotations about the branches change the order of the taxa but NOT the evolutionary relationships

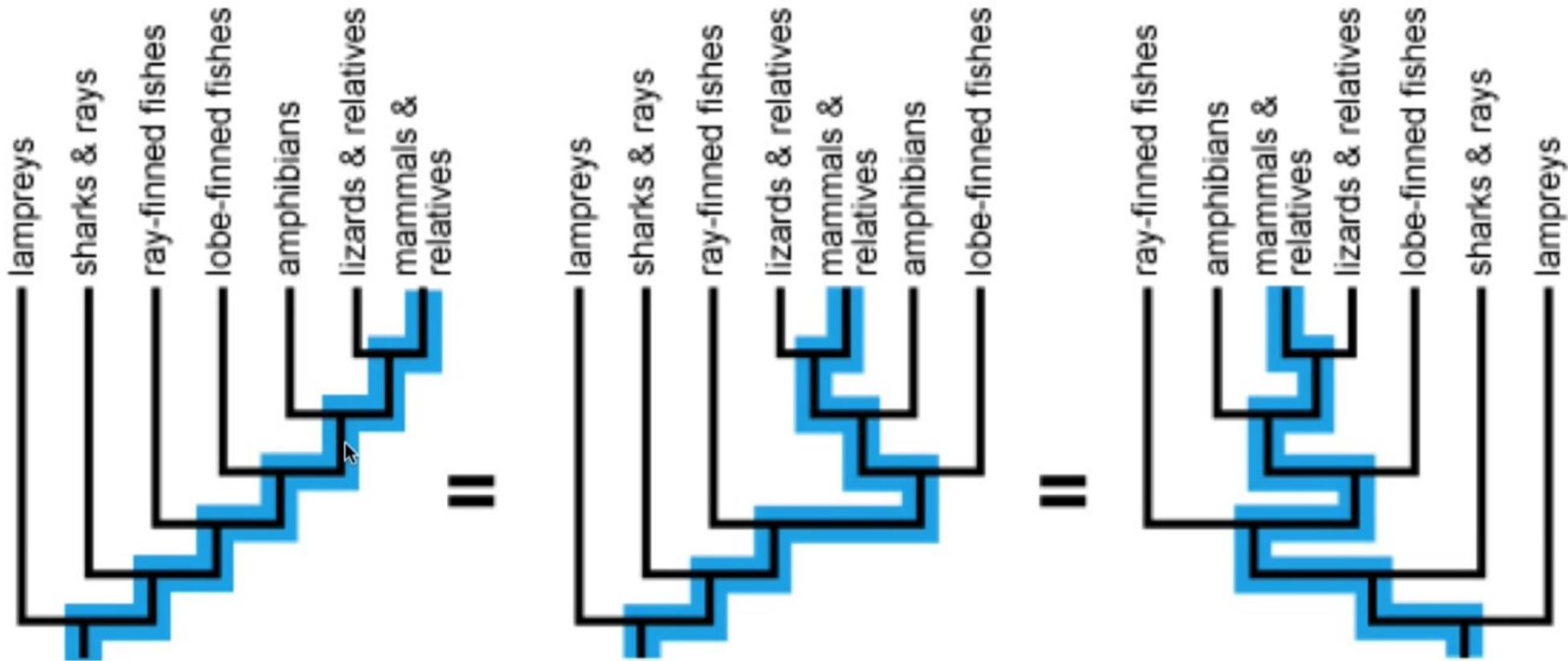
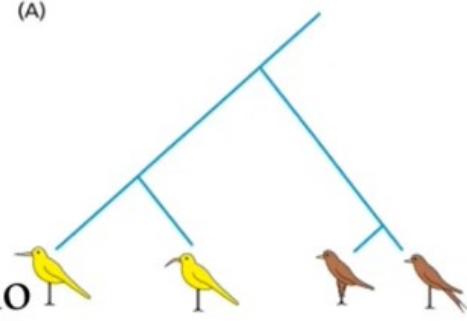


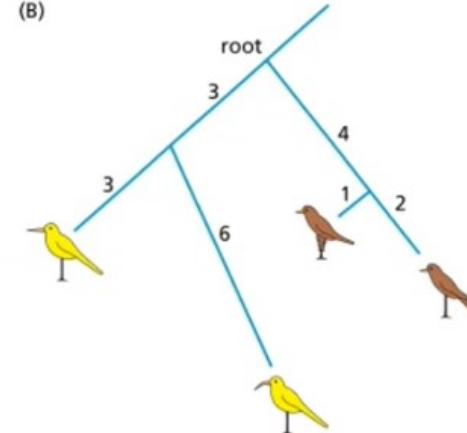
Image Credit: evolution.berkeley.edu

(A)



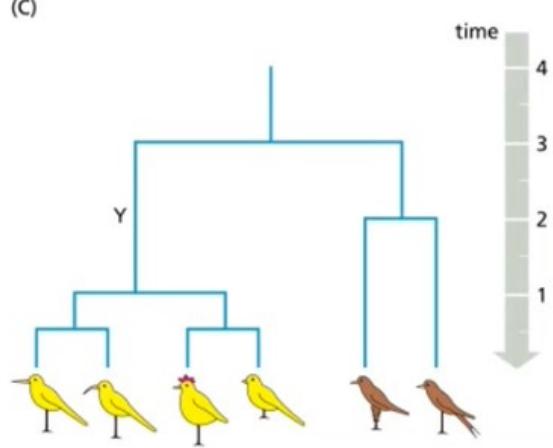
A ROOTED CLADOGRAM –
Genealogy but no timing of divergence

(B)



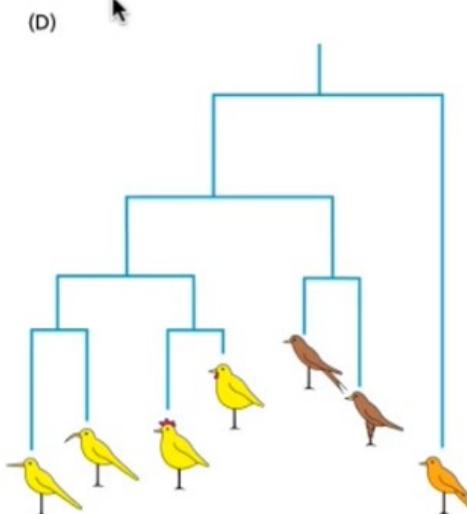
ADDITIVE TREE –
branch length
Proportional to
number of mutations

(C)



ULTRAMETRIC TREE – constant rate of
Mutation assumed along all branches
(molecular clock)

(D)

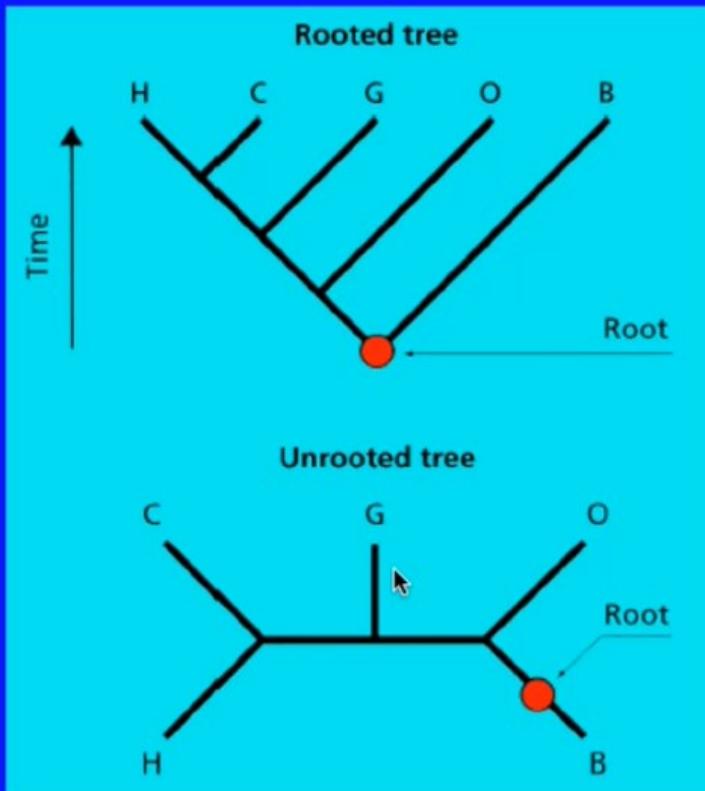


ADDITIVE TREE WITH OUTGROUPS

Rooted and Unrooted Trees

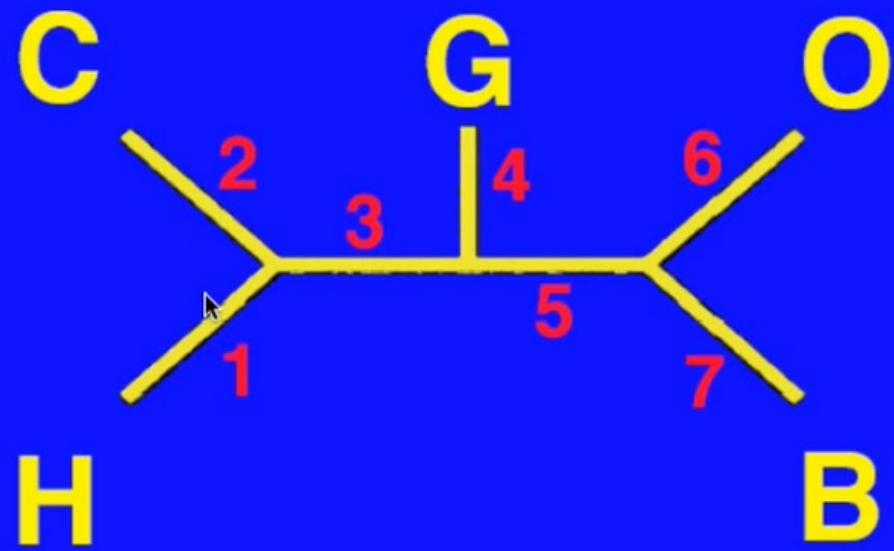
Rooted trees.
Have direction
corresponding
to evolutionary
time.

Unrooted trees.
we cannot talk
of ancestors and
descendants

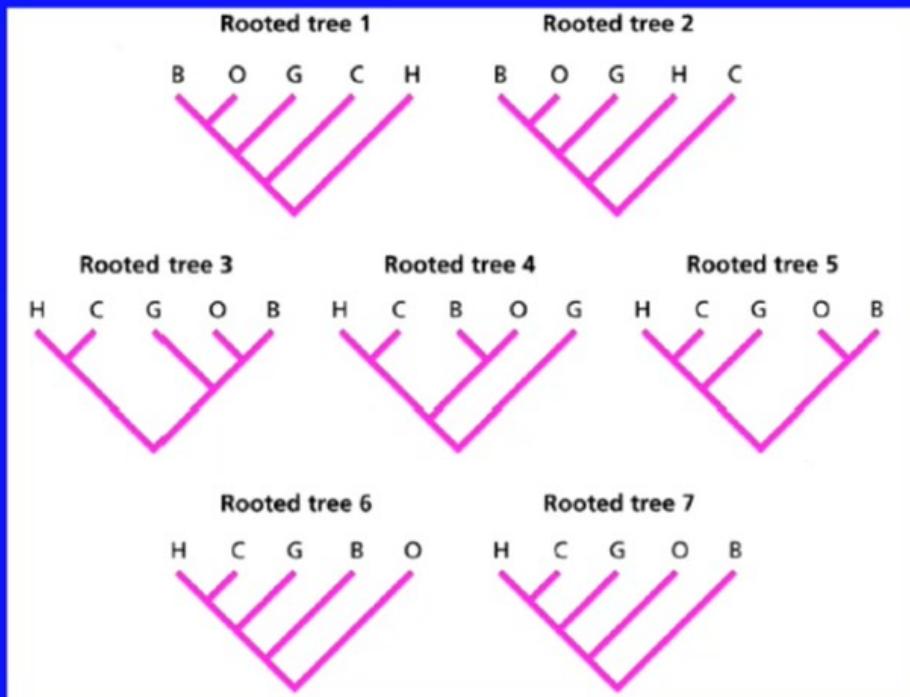


H = human
C = chimp
G = gorilla
O = orang
B = gibbon

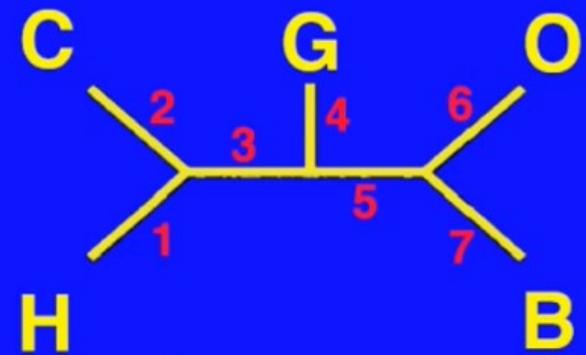
For this unrooted tree....



There are 7 different corresponding rooted trees



UNROOTED TREE



N.B. All correspond to same unrooted network shown previously but have the "root" placed on one of the 7 different branches of the original unrooted network

Number of possible trees increases exponentially with increasing number of sequences

Number of sequences	Number of unrooted trees	Number of rooted trees
2	1	1
3	1	3
4	3	15
5	15	105
6	105	945
7	945	10395
8	10395	135135
9	135135	2027025
10	2027025	34459425

$$\text{Rooted Trees} = \frac{(2n-3)!}{2^{n-2} (n-2)!}$$

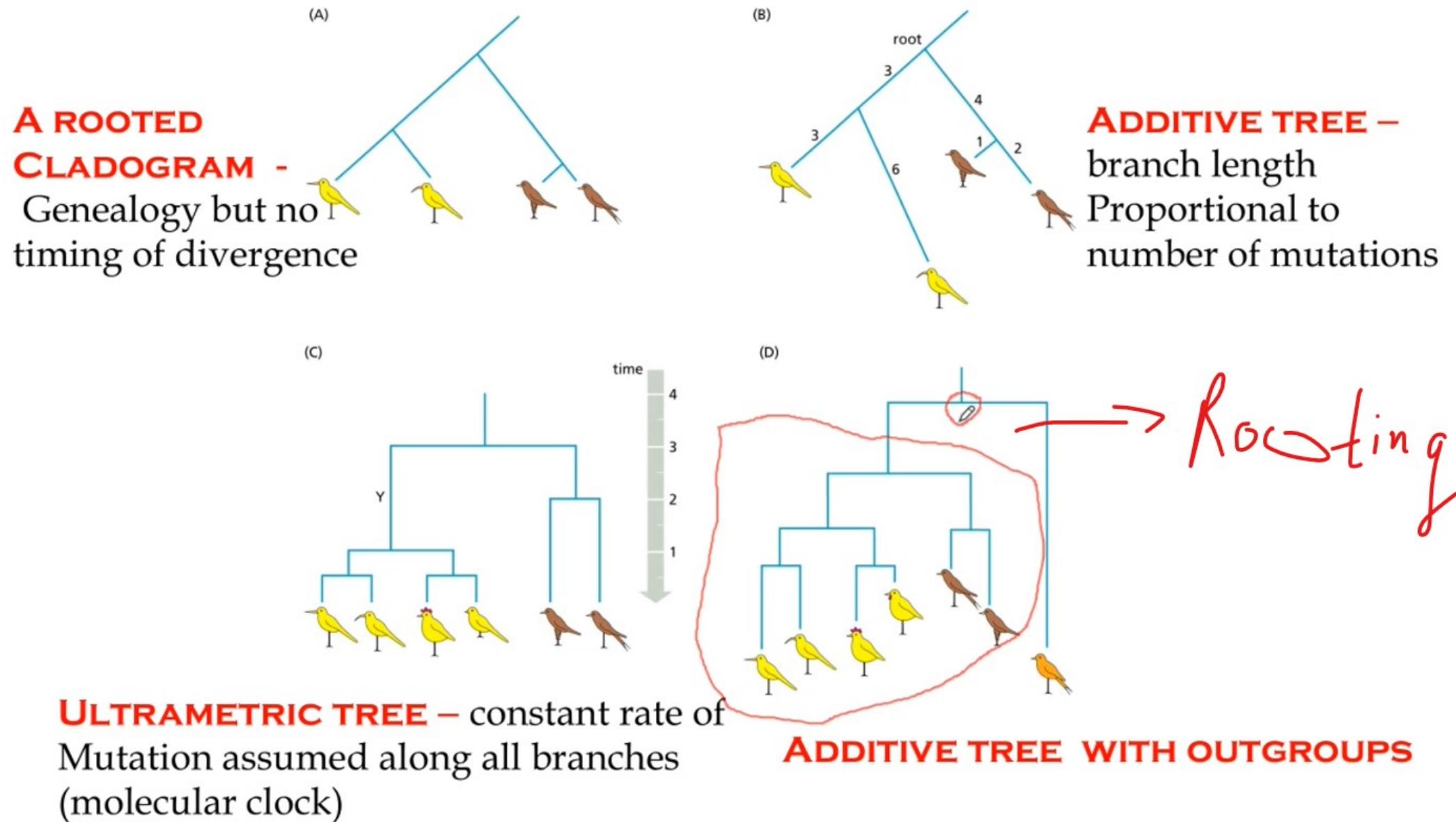
$n > 2$

$$\text{Unrooted Trees} = \frac{(2n-5)!}{2^{n-3} (n-3)!}$$

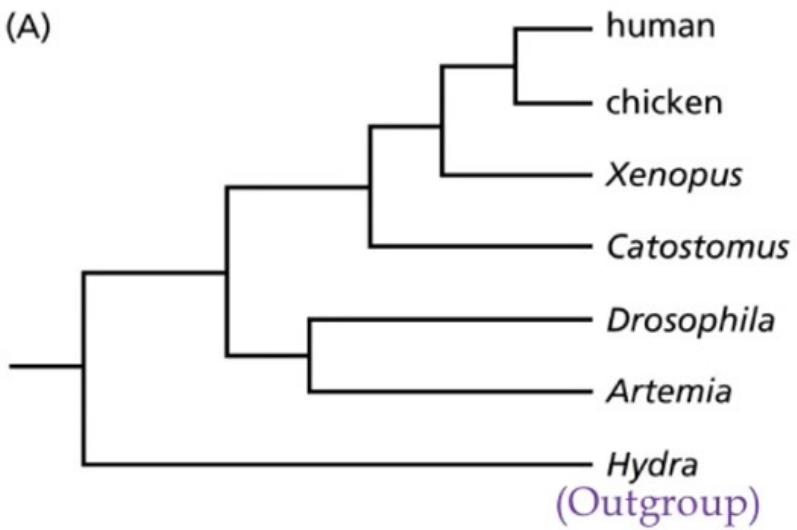
$n > 3$

For ten sequences there are more than 34 million rooted trees.

- **For 20 sequences** $8,200,794,532,637,891,559,000$ **trees.**
- **For a recent study of 135 human mtDNA sequences** 2.113×10^{267} **trees.**
- **Larger than number of particles in the known universe.**
(actually, larger than the volume of the universe measured in Plank Constant Units)

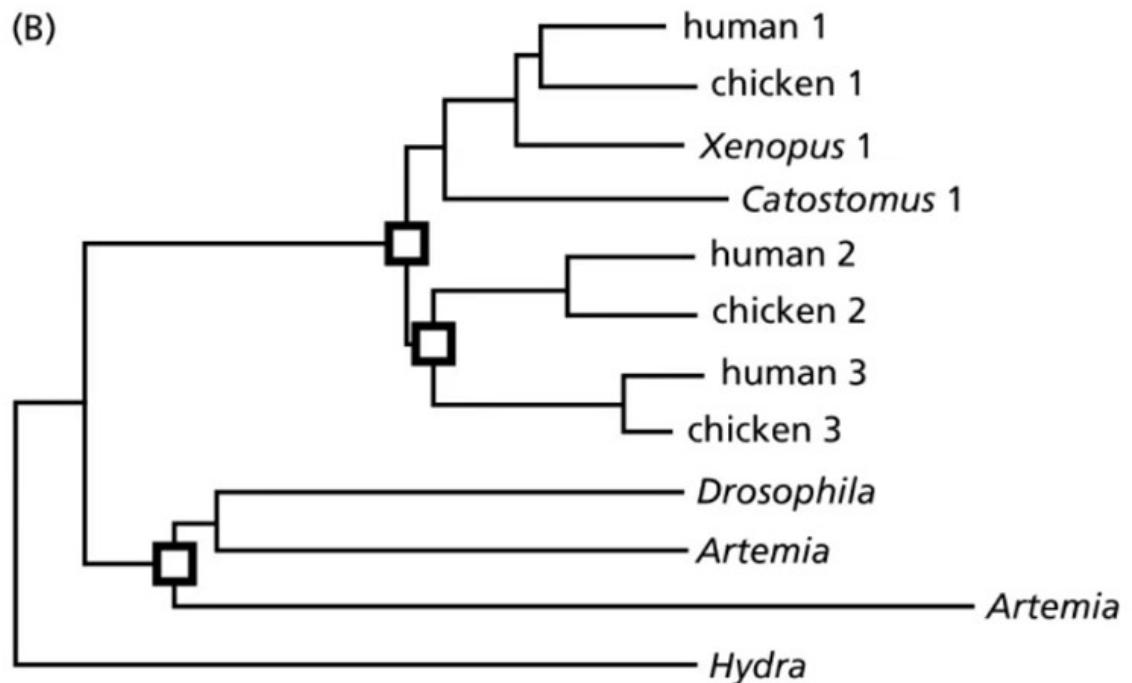


(A)



SPECIES TREE

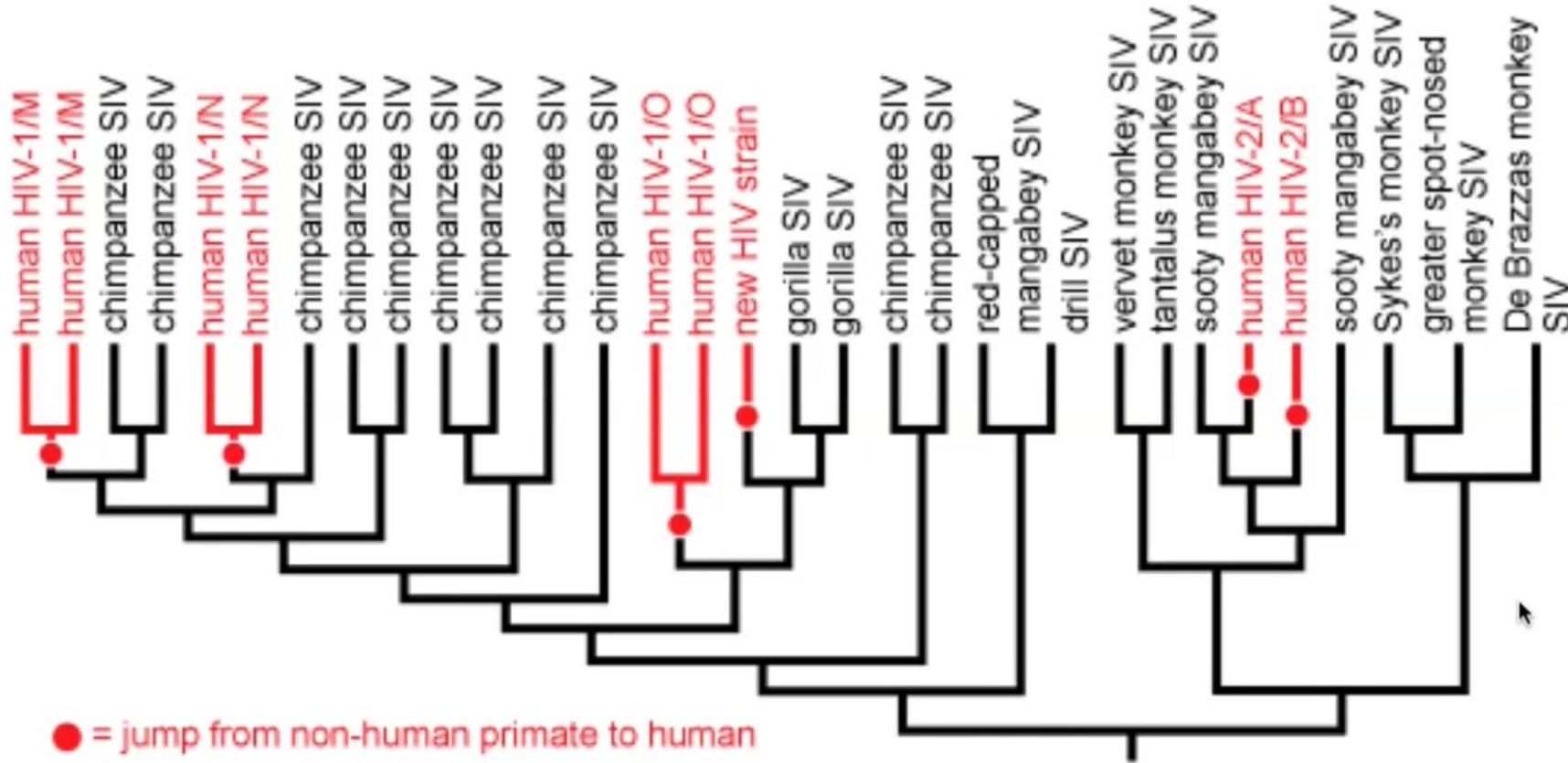
(B)



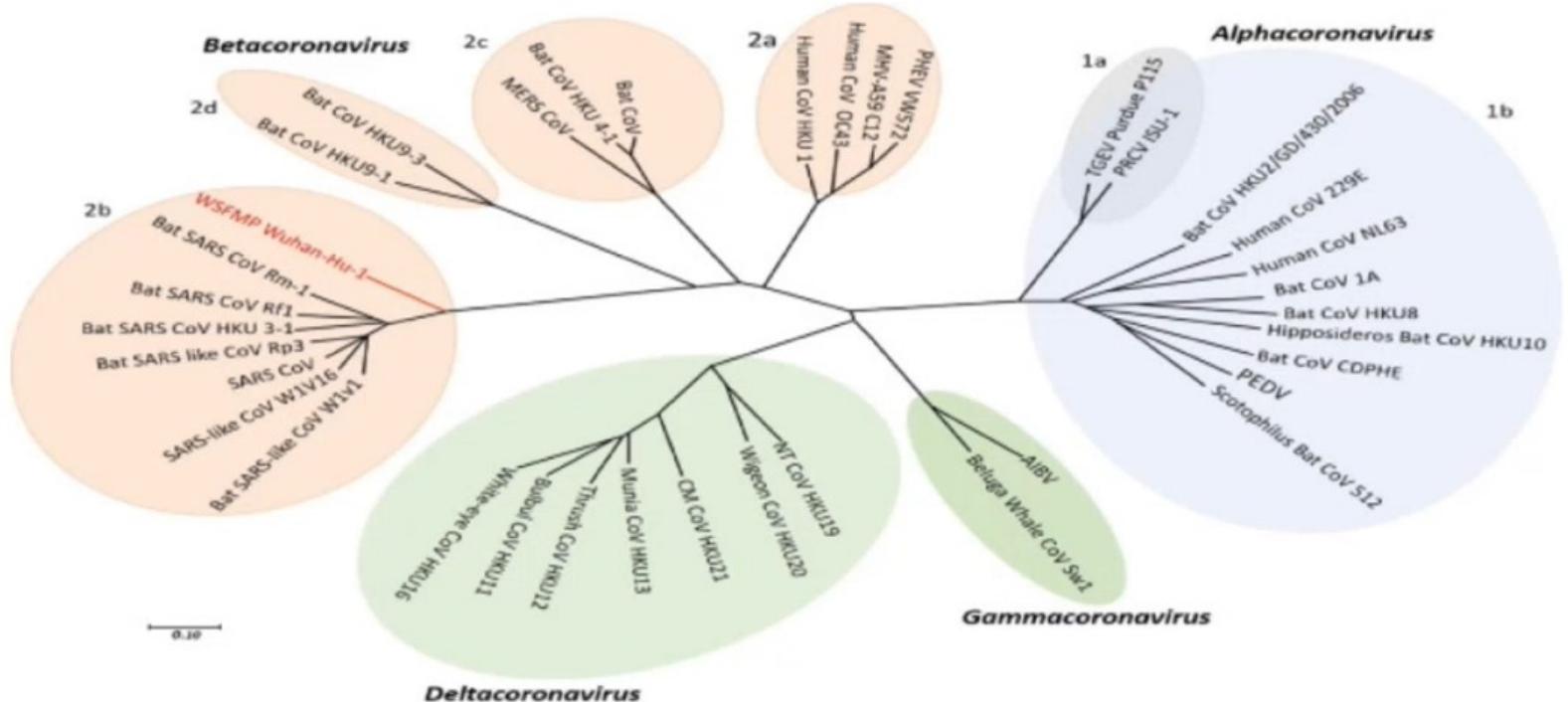
GENE TREE

Na^+-K^+ Pump Membrane Protein Family

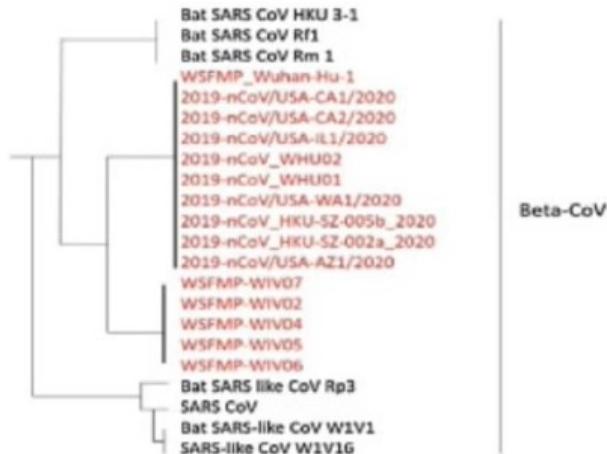
Phylogenetics reveals that HIV has evolved from Simian Immunodeficiency virus (SIV) several times



Keele et al, Nature 313:523-526 (2006)



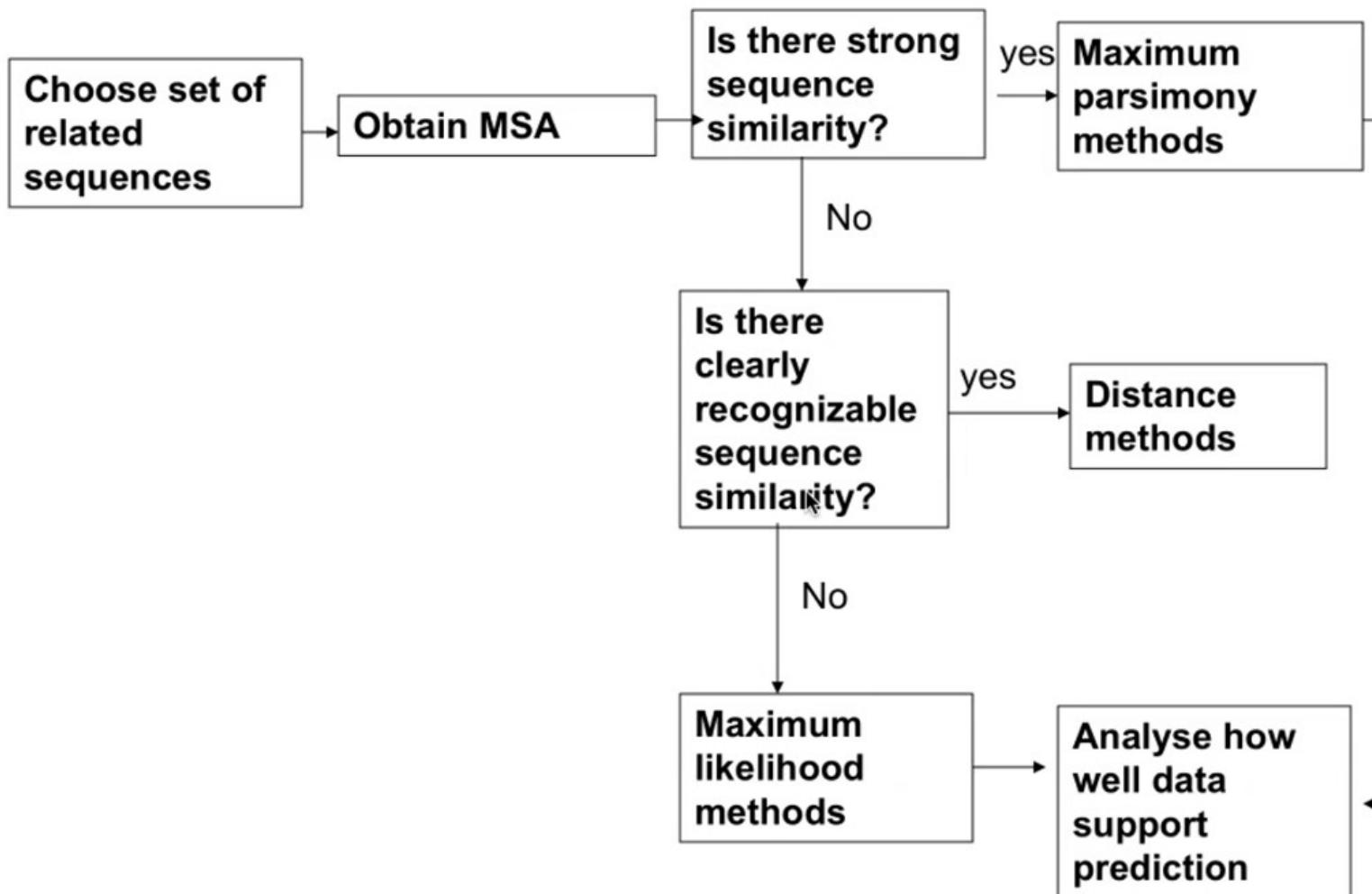
Phylogenetic analysis predicted
that the SARS-CoV-2 got
transmitted to humans from bats



Phylogenetic Analysis

Evolutionary Trees

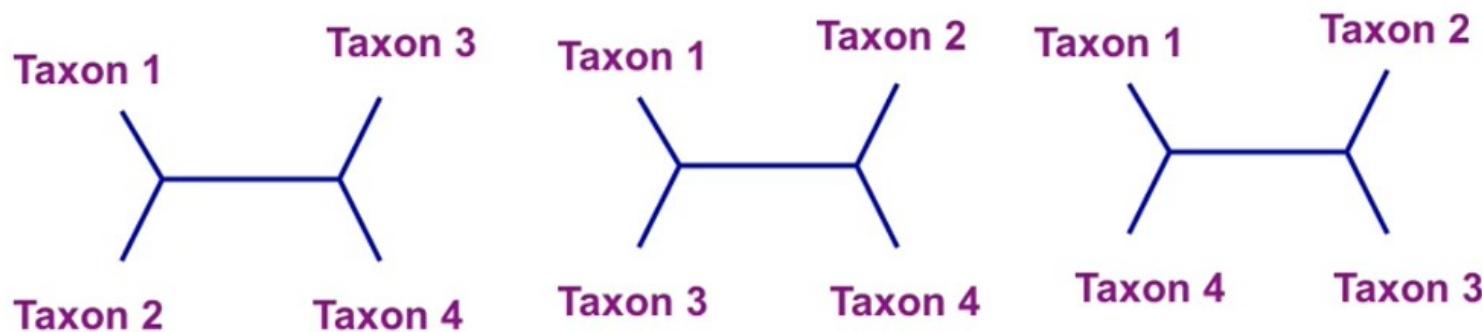
- Maximum Parsimony
- Distance based
- Maximum likelihood



Maximum Parsimony Method

- Predicts the evolutionary tree that minimizes the number of steps required to generate the observed variation in the sequences
- MSA is required to predict which sequence positions are likely to correspond
- Trees that produce the smallest number of changes overall for all sequence positions are identified
- All possible trees relating a group of sequences are examined

Sequence	1	2	3	4	5	6	7	8	9
1	A	A	G	A	G	T	G	C	A
2	A	G	C	C	G	T	G	C	G
3	A	G	A	T	A	T	C	C	A
4	A	G	A	G	A	T	C	C	G



All possible unrooted trees are considered

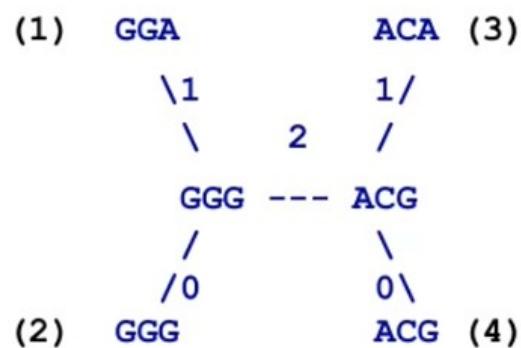
Sequence	1	2	3	4	5	6	7	8	9
1	A	A	G	A	G	T	G	C	A
2	A	G	C	C	G	T	G	C	G
3	A	G	A	T	A	T	C	C	A
4	A	G	A	G	A	T	C	C	G

Informative and non-informative sites

A site is informative only when there are at least two different kinds of nucleotides at the site, each of which is represented in at least two of the sequences under study.



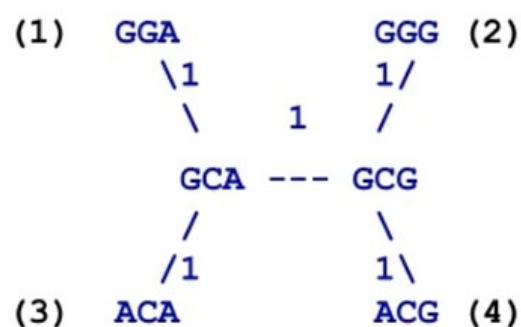
1 GGA
2 GGG
3 ACA
4 ACG



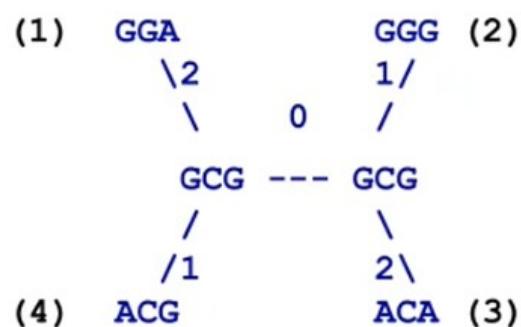
Number of mutations

Tree I: 4

Tree I is chosen because it requires the smallest number of changes (4) at the informative sites.



Tree II: 5



Tree III: 6



Maximum Parsimony Method

Disadvantages

Does not use all the information (only informative sites are used)

Does not provide information on branch length

Advantages

Evaluates different trees

Provides information about ancestral sequences

Distance Based Method

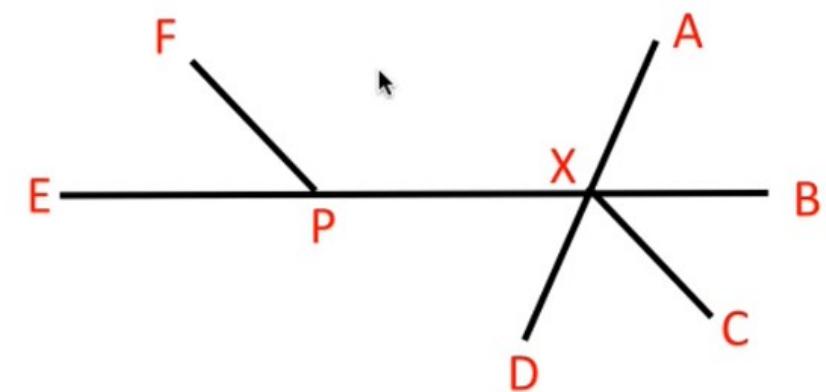
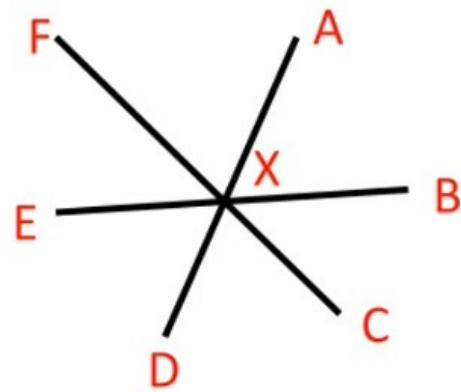
Sequence	1	2	3	4	5	6	7	8	9
P	A	A	G	A	G	T	G	C	A
Q	A	G	C	C	G	T	G	C	G
R	A	G	A	T	A	T	C	C	A
S	A	G	A	G	A	T	C	C	G



	P	Q	R	S
P		4	5	6
Q			5	4
R				2
S				

Distance Based Method

Neighbour Joining Method



	A	B	C	D	E
A		5	4	9	8
B			5	10	9
C				7	6
D					7
E					

$$U_i = \sum_{j=1}^N d_{ij} \quad \dots\dots(1) \quad U_A = 5 + 4 + 9 + 8 = 26$$

$$\delta_{ij} = d_{ij} - \left(\frac{U_i + U_j}{N - 2} \right) \quad \dots\dots(2) \quad \delta_{AB} = 5 - \left(\frac{26+29}{3} \right) = \frac{-40}{3}$$

	A	B	C	D	E	U
A		5	4	9	8	26
B			5	10	9	29
C				7	6	22
D					7	33
E						30

N = 5

$3 \delta_{ij}$

	A	B	C	D	E
A		-40	-36	-32	-32
B			-36	-32	-32
C				-34	-34
D					-42
E					

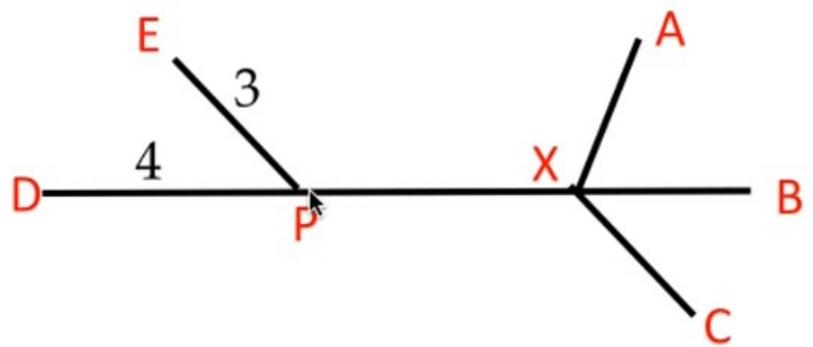
	A	B	C	D	E
A		5	4	9	8
B			5	10	9
C				7	6
D					7
E					

U
26
29
22
33
30

N = 5

$3 \delta_{ij}$

	A	B	C	D	E
A		-40	-36	-32	-32
B			-36	-32	-32
C				-34	-34
D					-42
E					



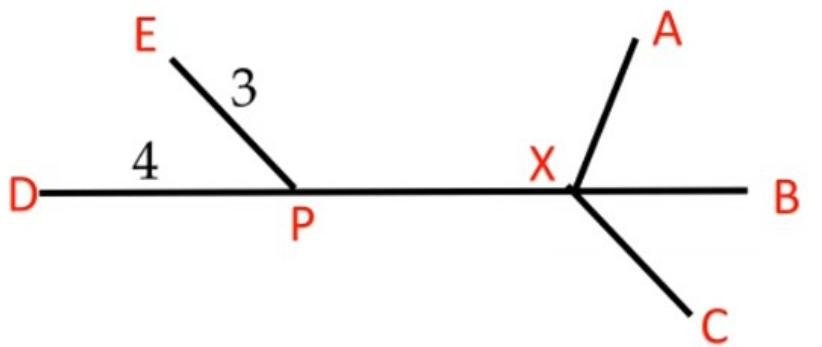
$$b_{DP} = \frac{1}{2} \left(d_{DE} + \left(\frac{U_D - U_E}{N-2} \right) \right)$$

$$b_{DP} = \frac{1}{2} \left(7 + \left(\frac{33 - 30}{3} \right) \right) = 4$$

$$b_{iY} = \frac{1}{2} \left(d_{ij} + \left(\frac{U_i - U_j}{N-2} \right) \right) \dots\dots(3)$$

$$b_{jY} = d_{ij} - b_{iY} \dots\dots(4)$$

$$b_{EP} = 7 - 4 = 3$$



$$b_{Yk} = \frac{1}{2} (d_{ik} + d_{jk} - dij) \dots\dots(5)$$

$$b_{PA} = \frac{1}{2} (d_{AD} + d_{AE} - dDE) = \frac{(9+8-7)}{2} = 5$$

	A	B	C	P	U
A		5	4	5	14
B			5	6	16
C				3	12
P					14

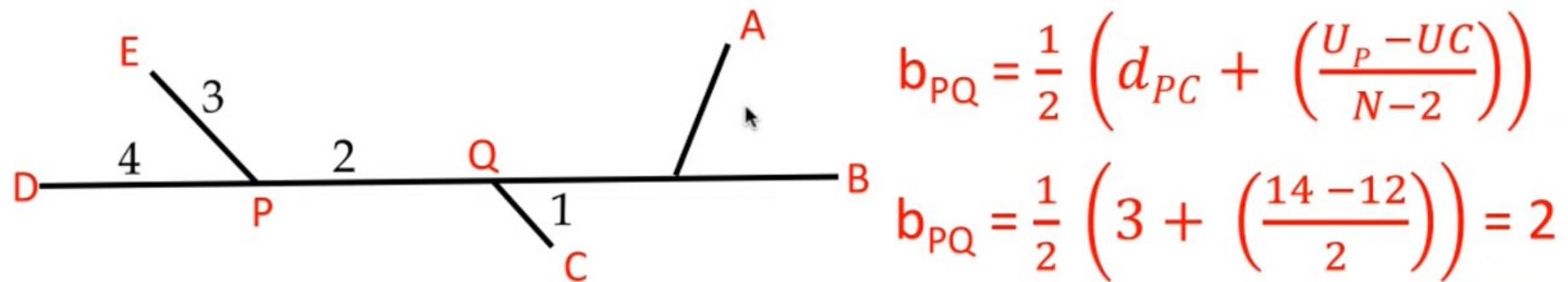
N = 4

2 δ_{ij}

	A	B	C	P
A		-20	-18	-18
B			-18	-18
C				-20
P				

$$U_i = \sum_{j=1}^N dij$$

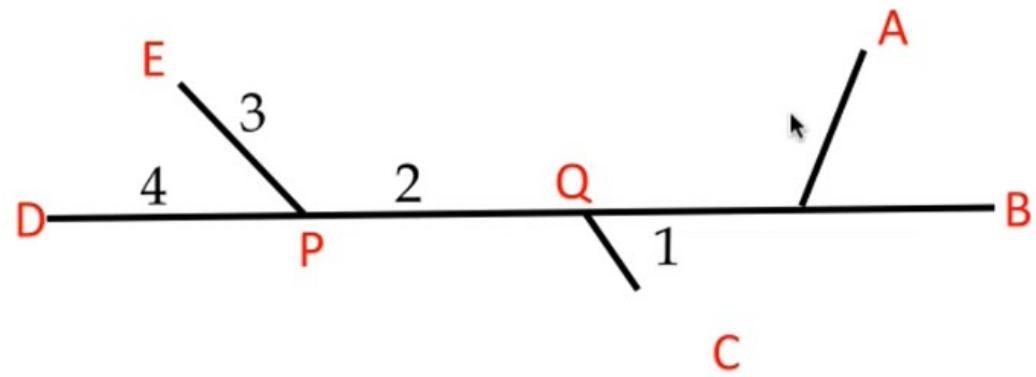
$$\delta_{ij} = dij - \left(\frac{Ui + Uj}{N - 2} \right)$$



$$b_{iY} = \frac{1}{2} \left(d_{ij} + \left(\frac{U_i - Uj}{N-2} \right) \right)$$

$$b_{jY} = d_{ij} - b_{iY}$$

$$b_{QC} = b_{PQ} - b_{PC} = 3 - 2 = 1$$



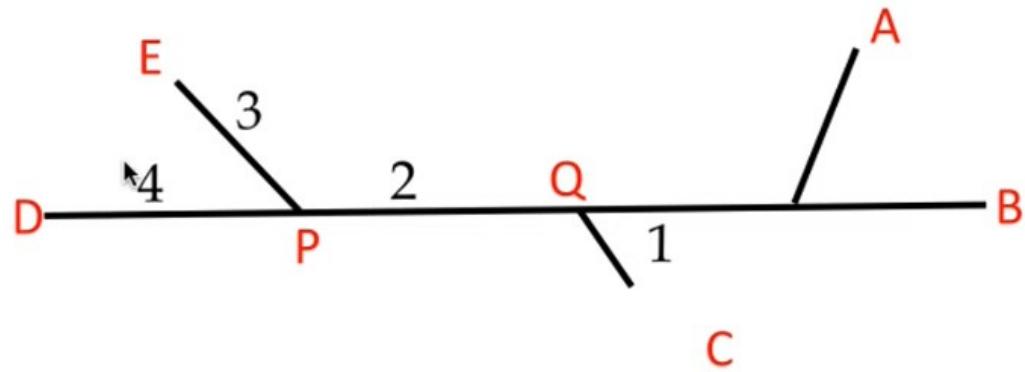
$$b_{Yk} = \frac{1}{2} (d_{ik} + d_{jk} - dij) \dots\dots(5)$$

	A	B	Q
A		5	3
B			4
Q			

U
8
9
7

$$U_i = \sum_{j=1}^N d_{ij}$$





$$b_{Yk} = \frac{1}{2} (d_{ik} + d_{jk} - dij) \dots\dots(5)$$

$$b_{QA} = \frac{1}{2} (d_{PA} + d_{CA} - dPC) = \frac{(5+4-3)}{2} = 3$$

	A	B	Q
A		5	3
B			4
Q			

U
8
9
7

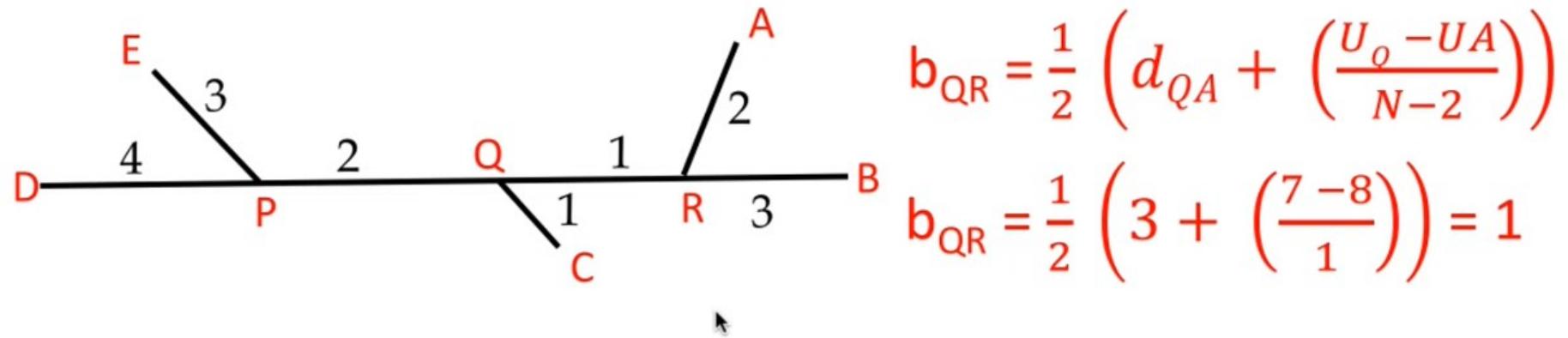
N = 3

δ_{ij}

	A	B	Q
A		-12	-12
B			-12
Q			

$$U_i = \sum_{j=1}^N d_{ij}$$

$$\delta_{ij} = dij - \left(\frac{U_i + U_j}{N - 2} \right)$$



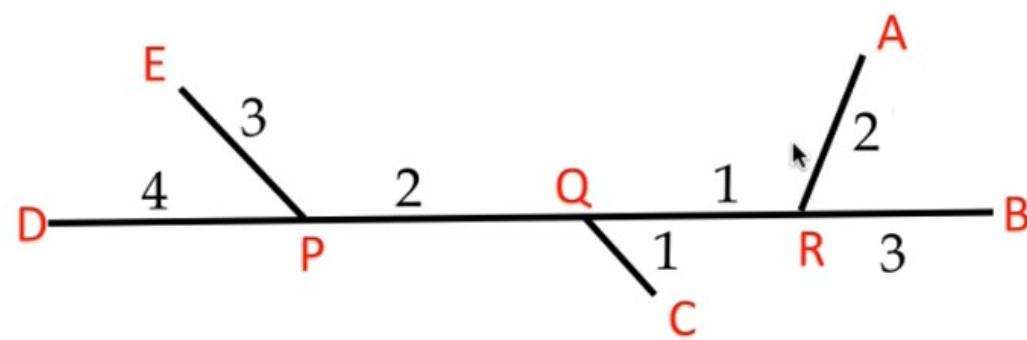
$$b_{QR} = \frac{1}{2} \left(d_{QA} + \left(\frac{U_Q - U_A}{N-2} \right) \right)$$

$$b_{QR} = \frac{1}{2} \left(3 + \left(\frac{7 - 8}{1} \right) \right) = 1$$

$$b_{iY} = \frac{1}{2} \left(d_{ij} + \left(\frac{U_i - U_j}{N-2} \right) \right)$$

$$b_{jY} = d_{ij} - b_{iY}$$

$$b_{RB} = b_{QB} - b_{QR} = 4 - 1 = 3$$



- Unrooted Tree
- Additive Tree

Neighbour-joining Method

Advantages

Robust enough to recover the correct tree topology

Provides information on branch length - Avoids negative branch length

Disadvantages

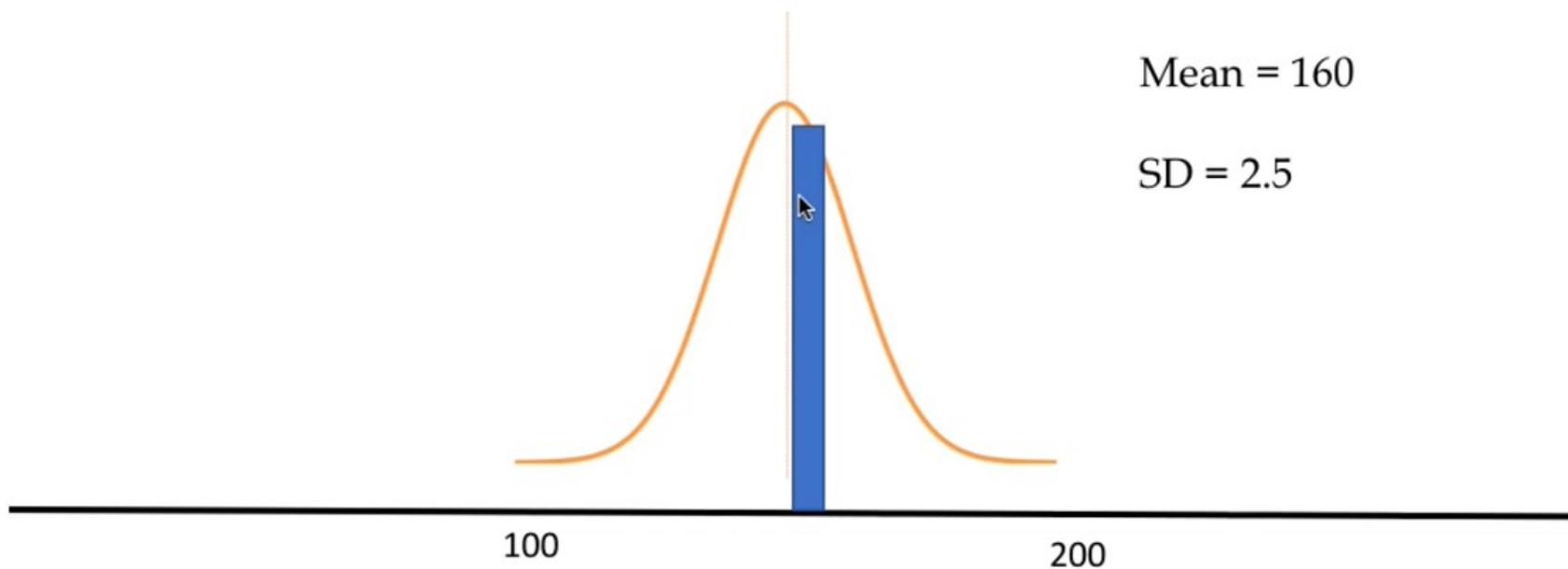
Sequence information is lost

Does not provide information about ancestral sequences

Maximum Likelihood

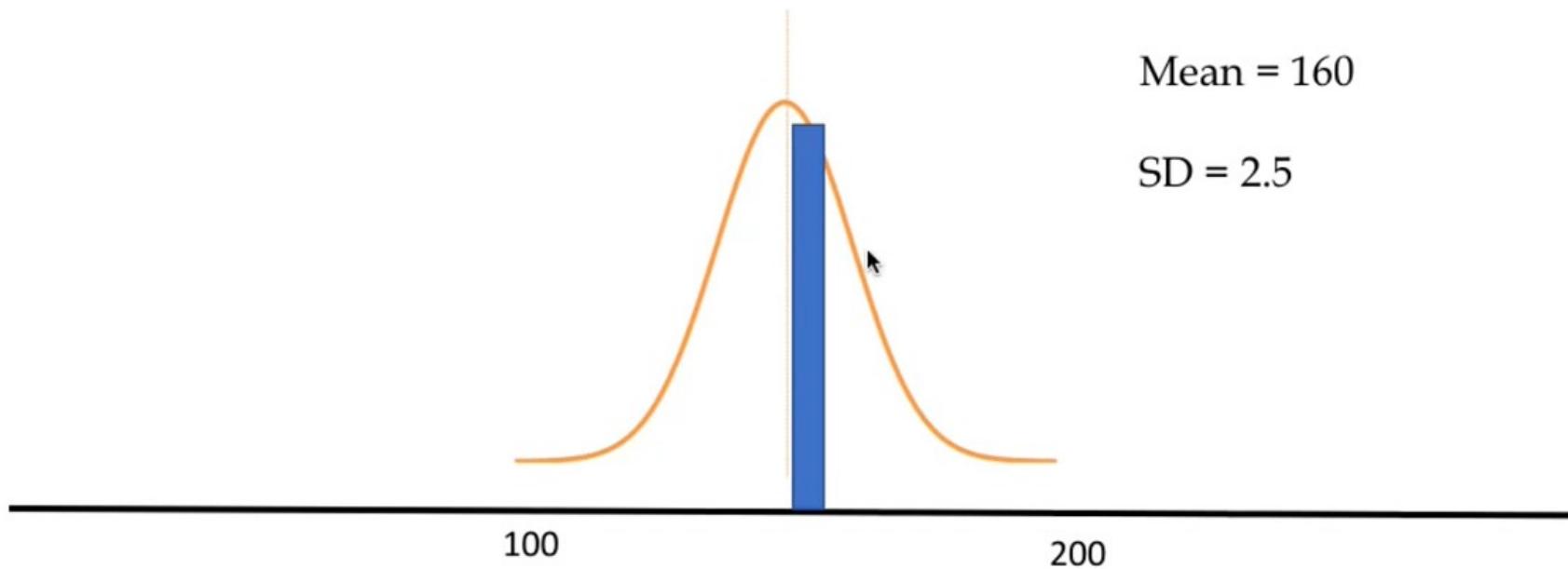


Probability versus Likelihood



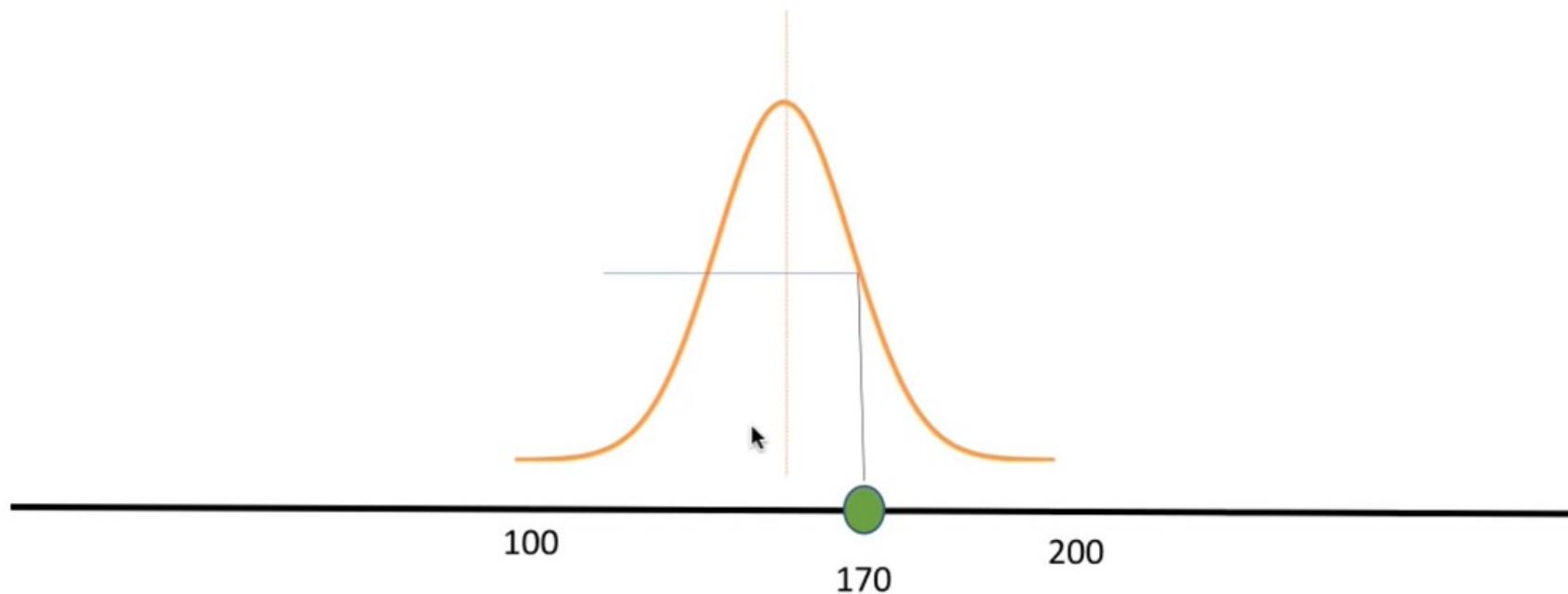
$P(X \geq 150 \text{ and } X < 160) = \text{Area under the curve}$

Probability versus Likelihood



$P(X \geq 150 \text{ and } X < 160 \mid \text{mean} = 160 \text{ and } \text{SD} = 2.5) = \text{Area under the curve}$

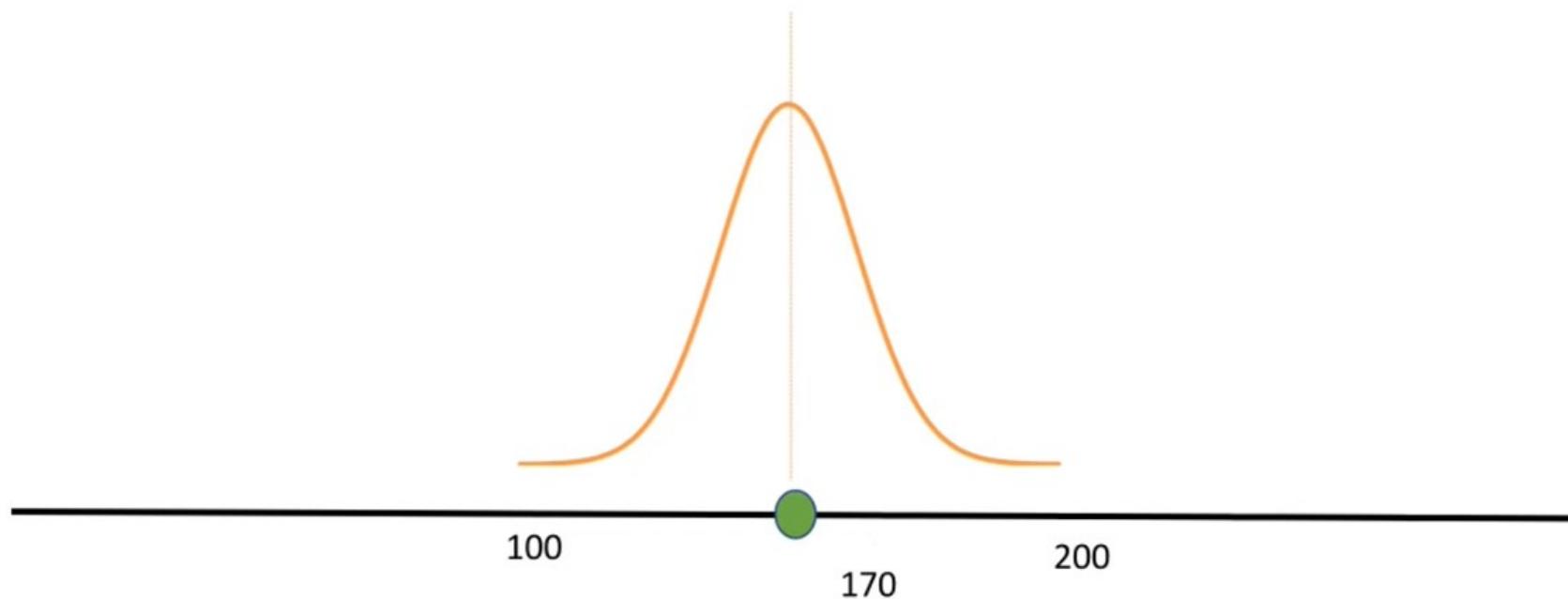
Probability versus Likelihood



$L(\text{mean} = 160 \text{ and } \text{SD} = 2.5 \mid X = 170 \mid) = 0.15$



Probability versus Likelihood

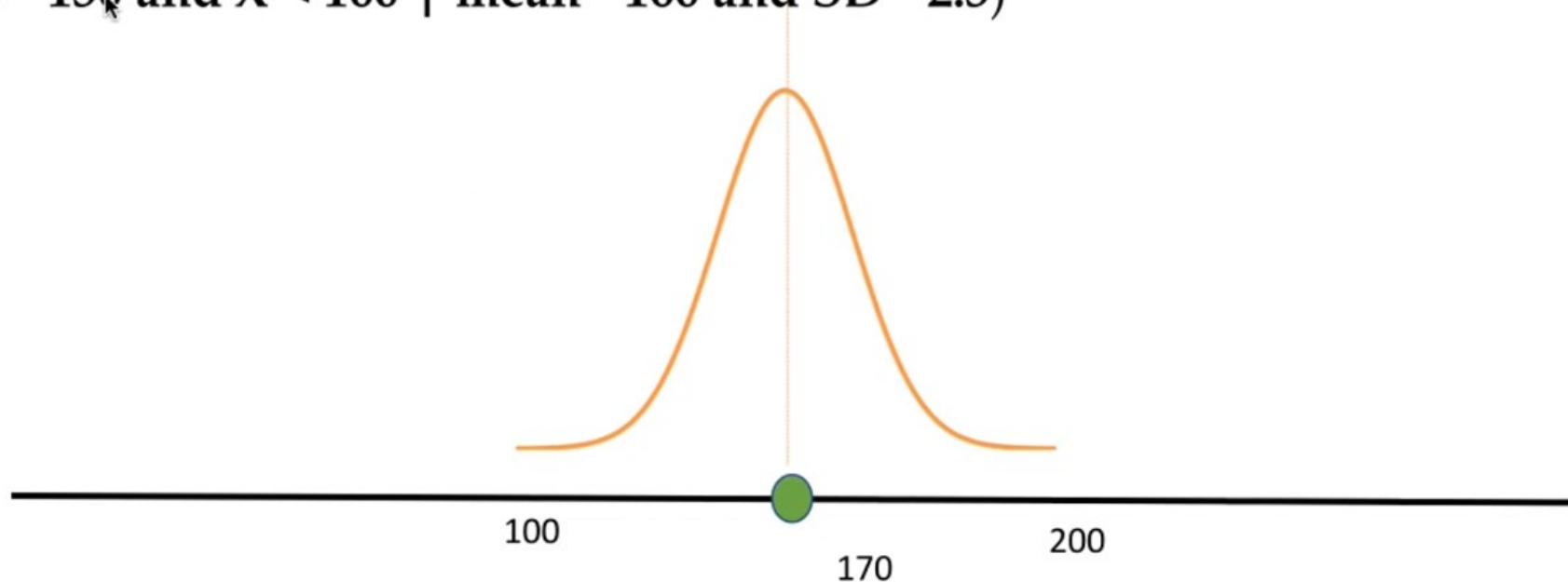


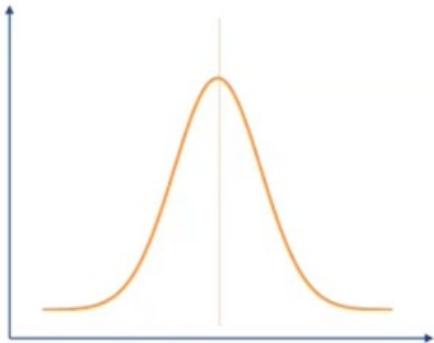
$$L(\text{mean} = 160 \text{ and } \text{SD} = 2.5 \mid X = 160 \mid)^\star = 0.25$$

Probability versus Likelihood

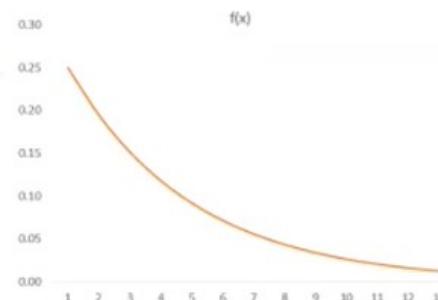
$L(\text{mean} = 160 \text{ and } \text{SD} = 2.5 \mid X = 160)$

$P(X >= 150 \text{ and } X < 160 \mid \text{mean} = 160 \text{ and } \text{SD} = 2.5)$

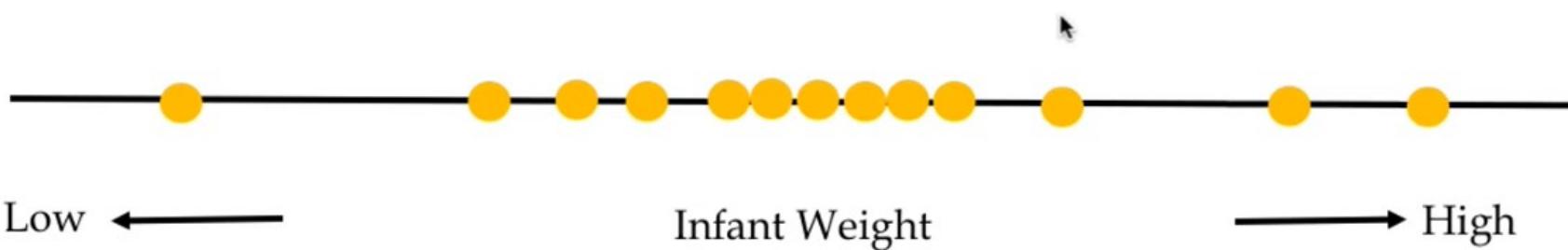




Normal Distribution



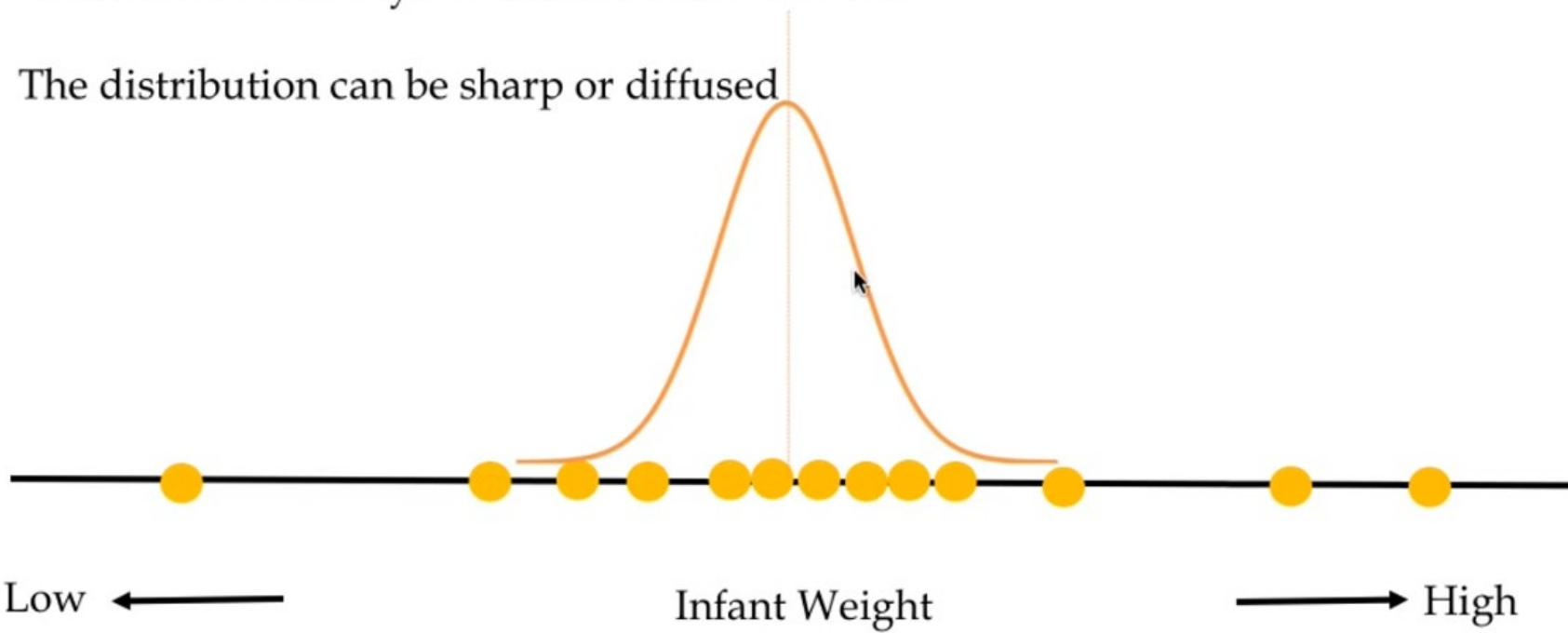
Exponential Distribution



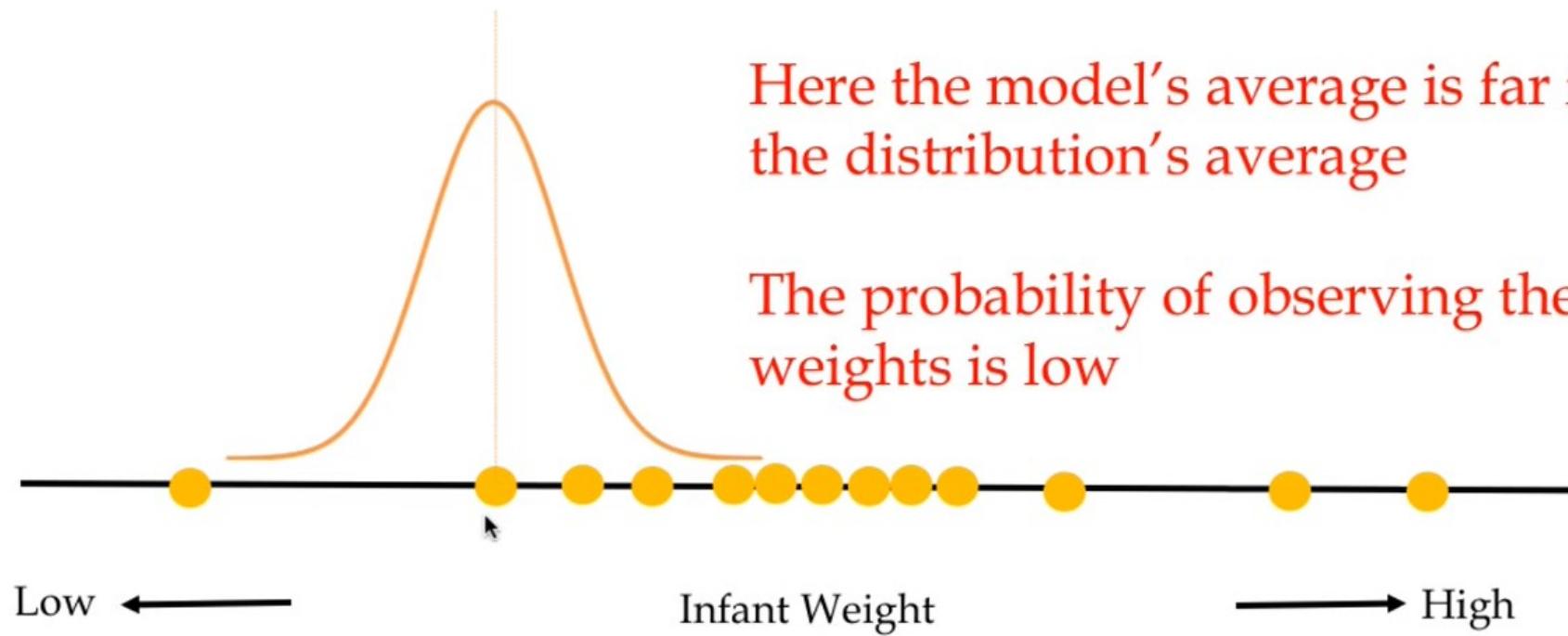
Infant Weight

Normal Distribution

1. Most of the measurements are close to the mean of the distribution
2. Measurements are symmetrical about the mean
3. The distribution can be sharp or diffused



Normal Distribution



Normal Distribution

Here the model's average is close to the distribution's average

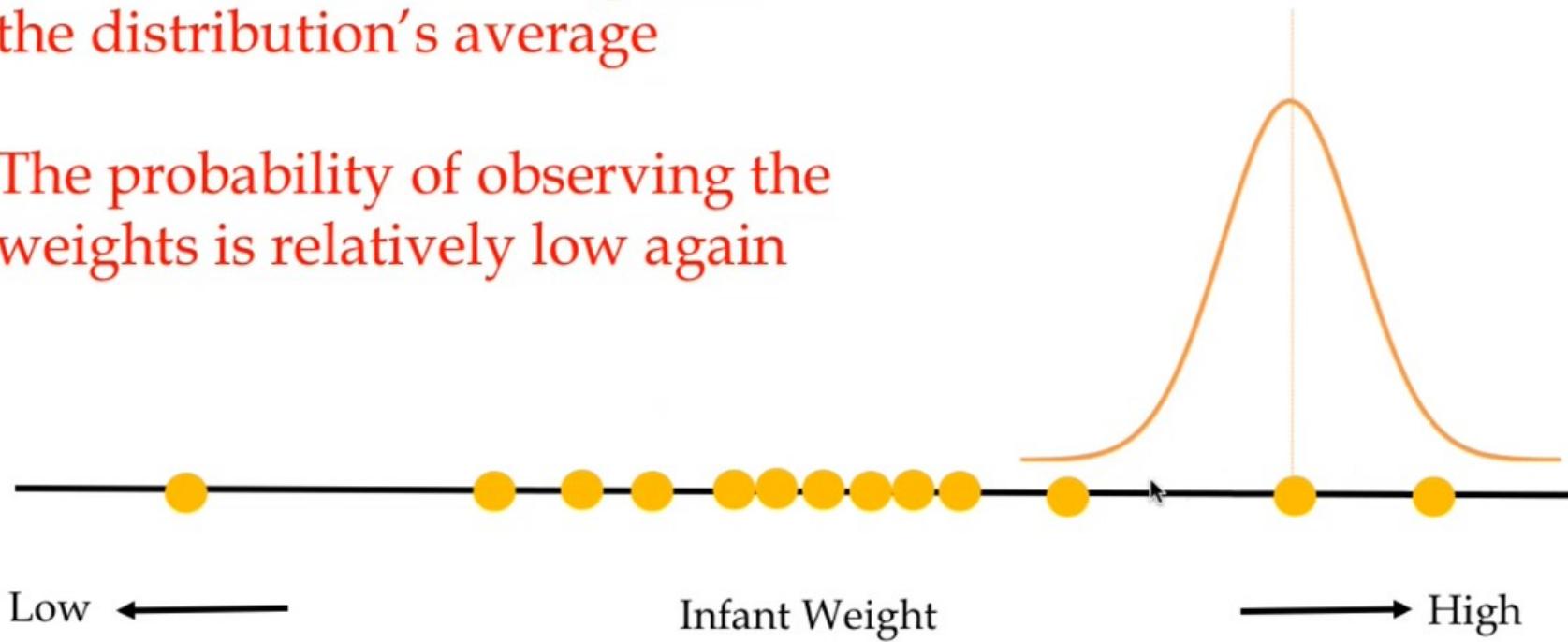
The probability of observing the weights is relatively high

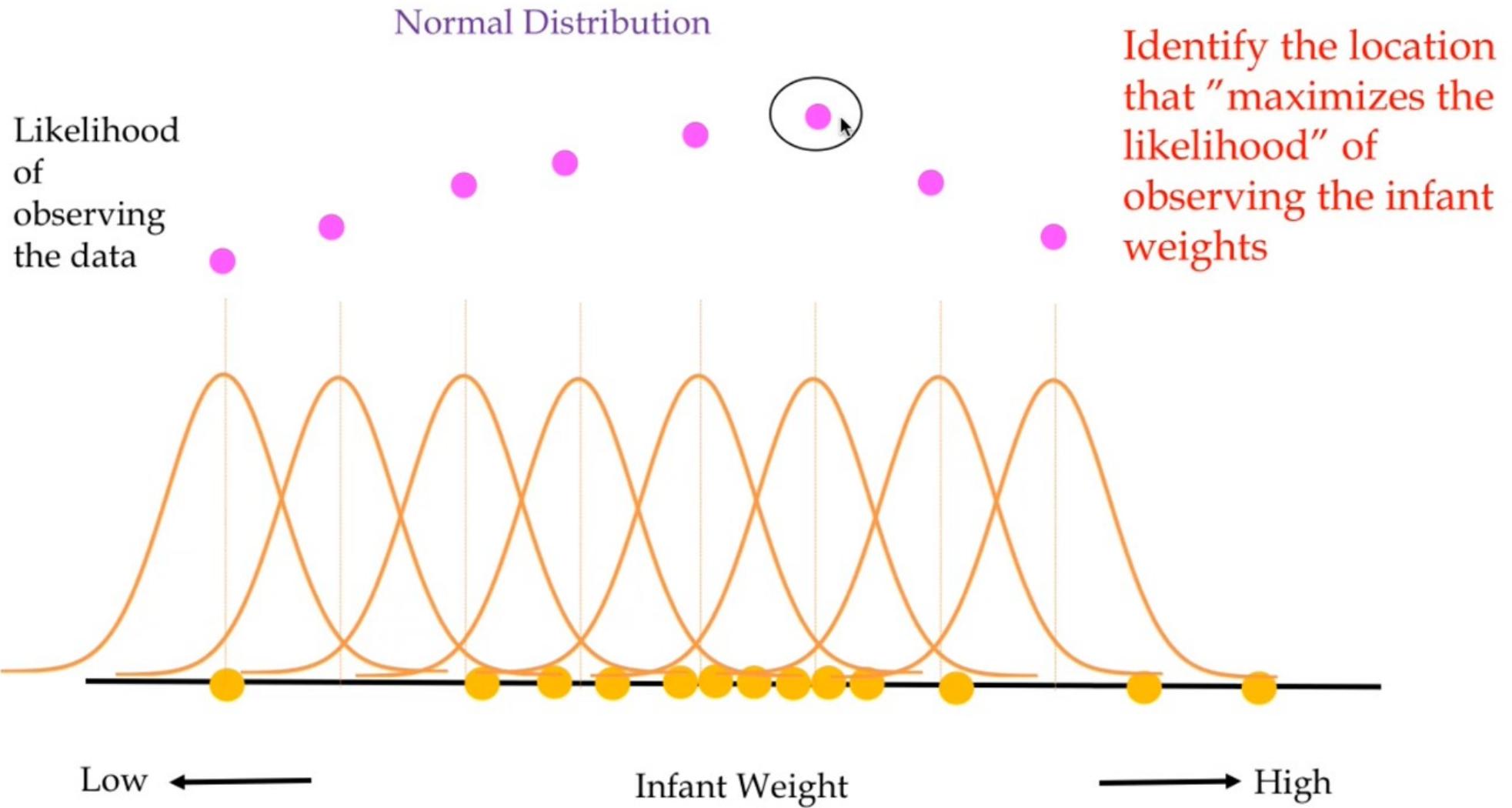


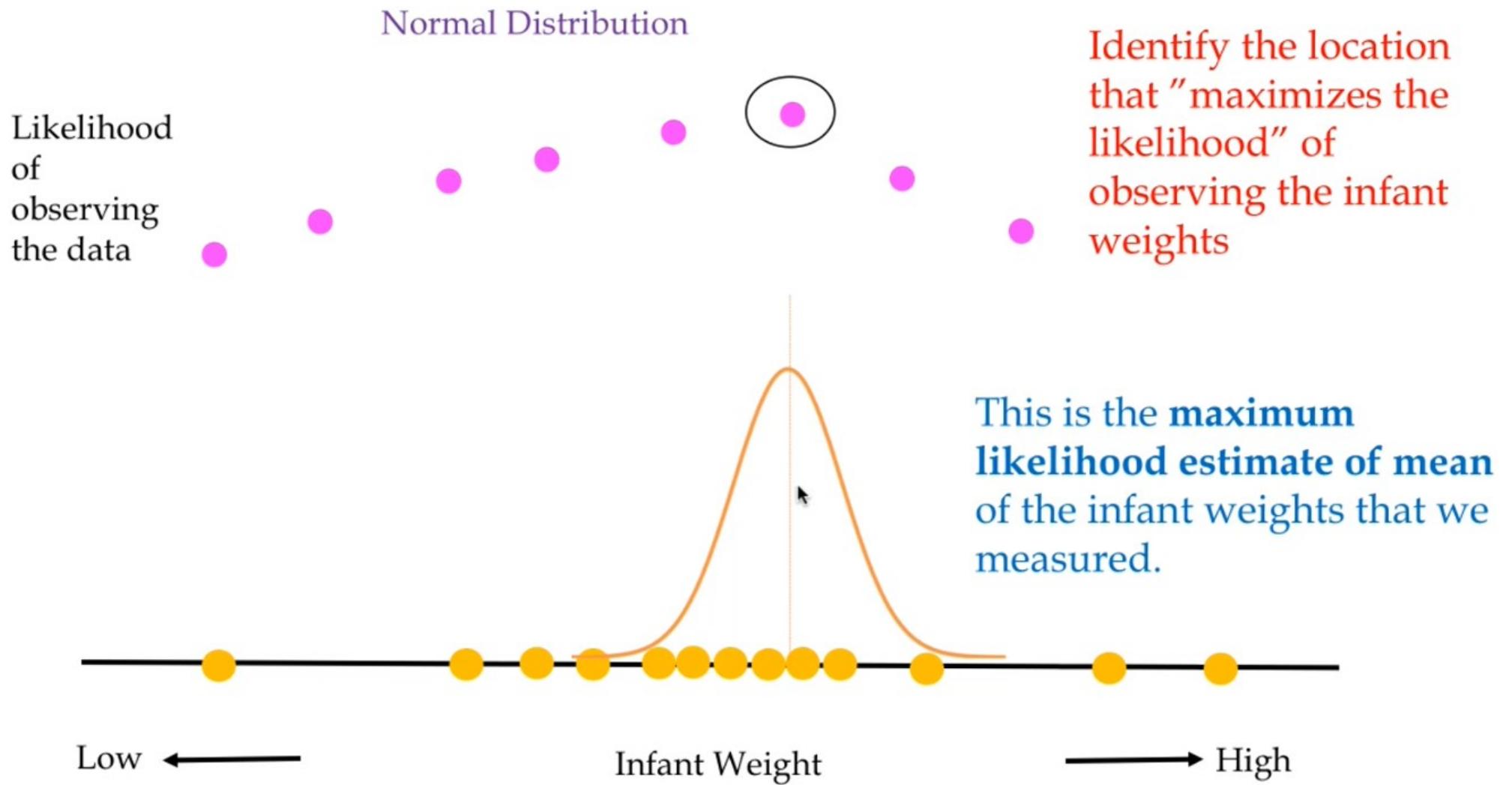
Normal Distribution

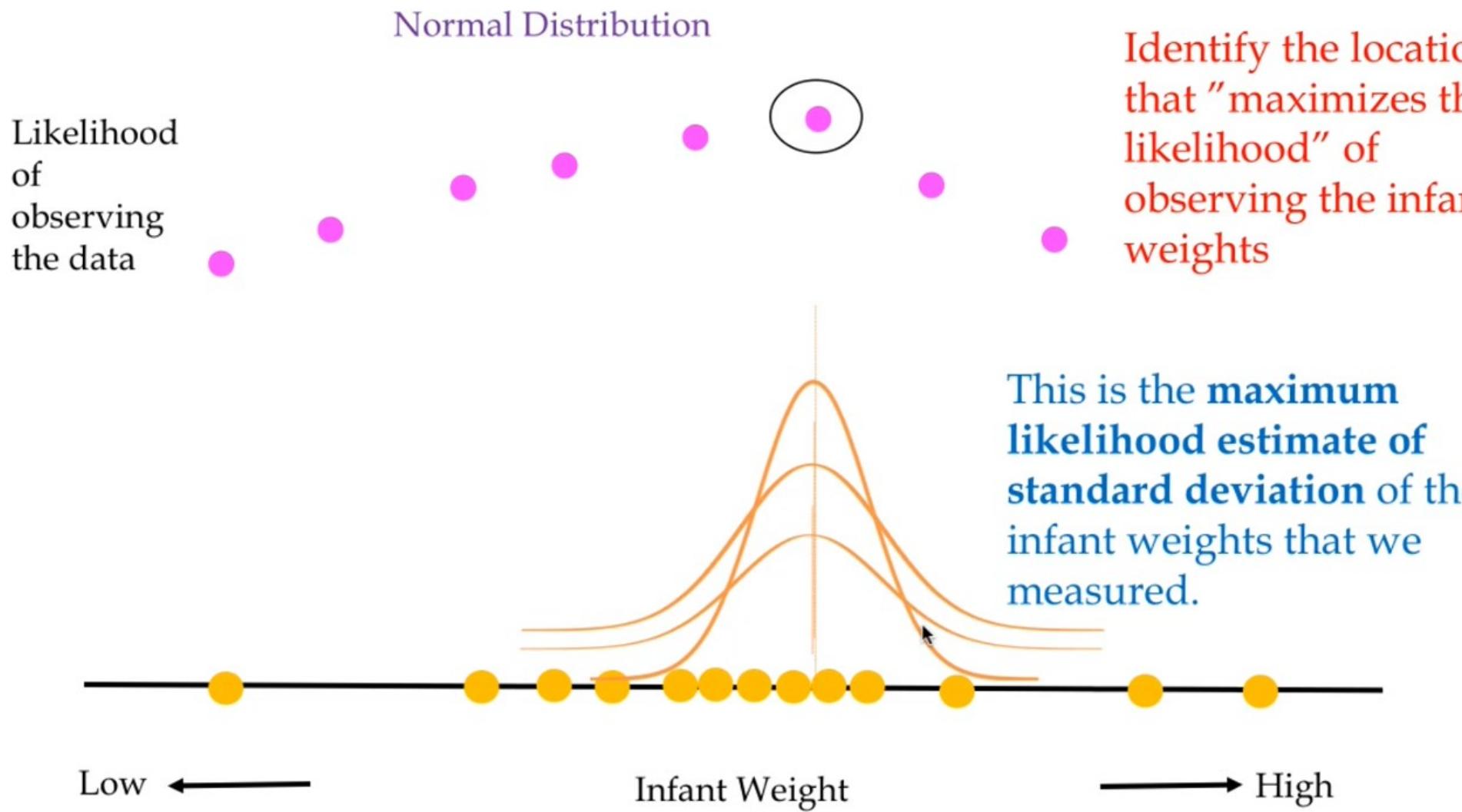
Here the model's average is far from the distribution's average

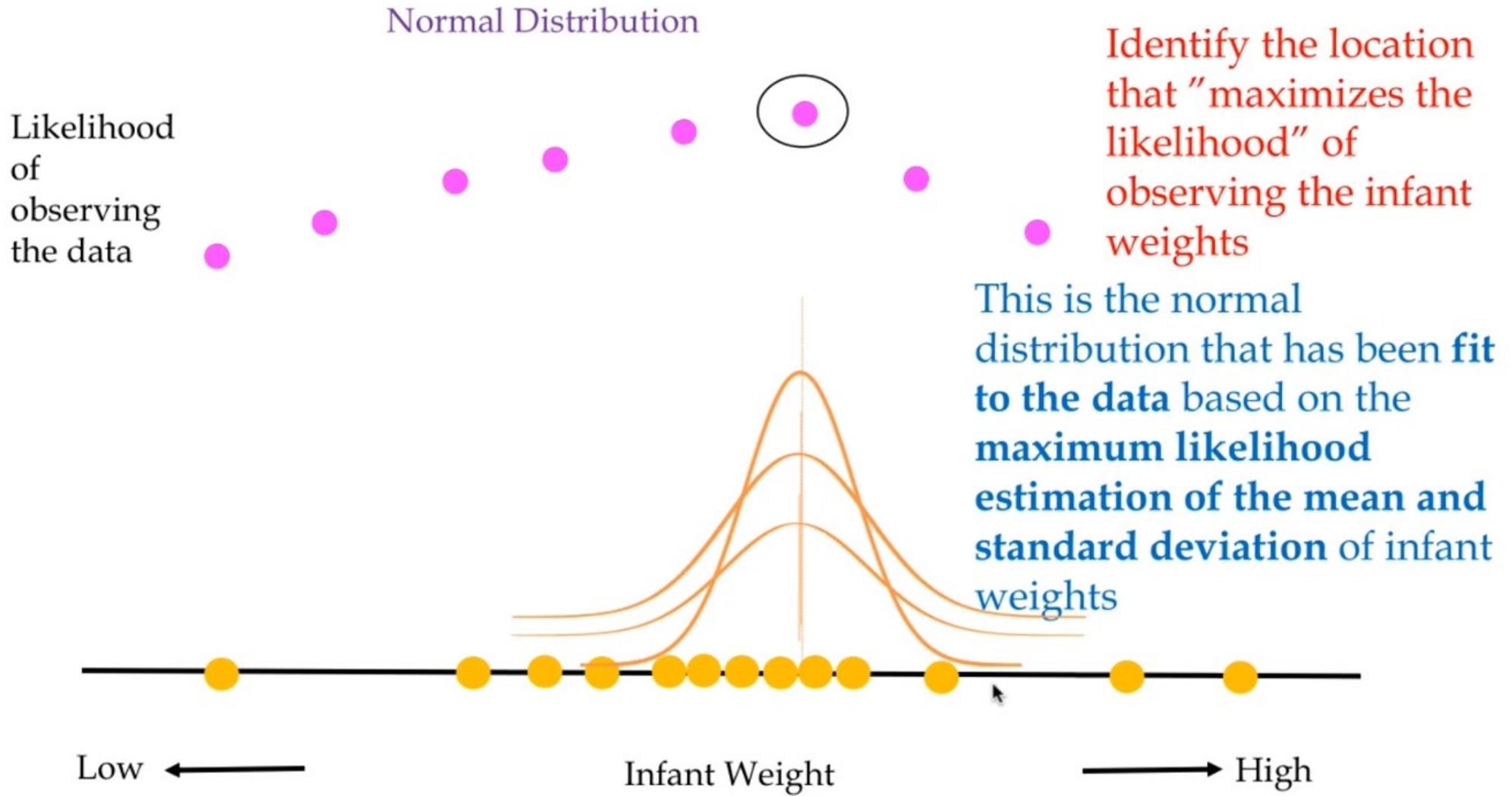
The probability of observing the weights is relatively low again











Maximum Likelihood Method

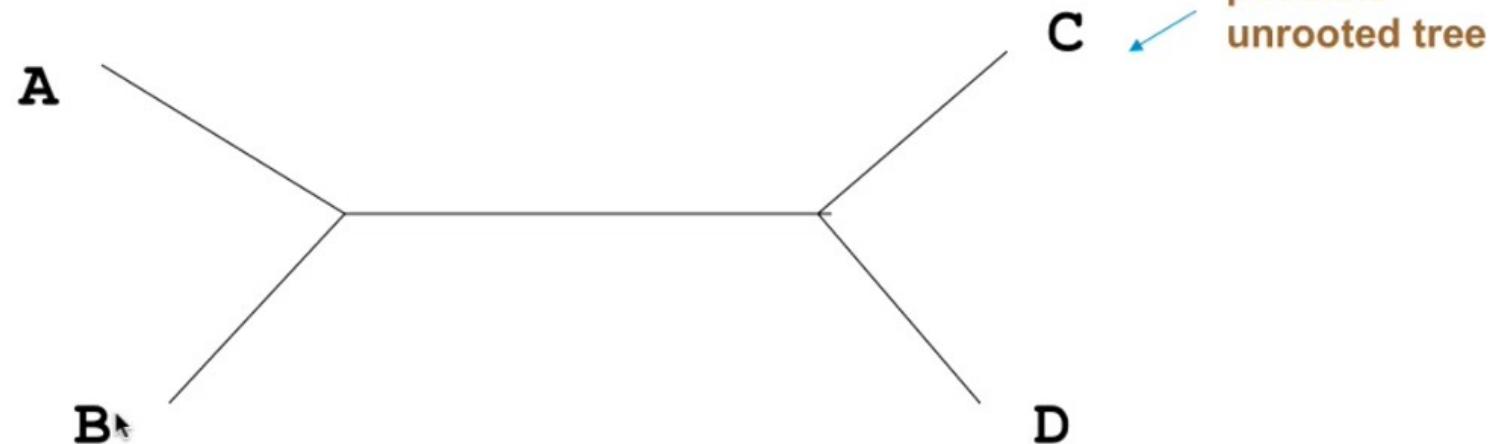
- Uses probability calculations to find a tree that best accounts for the variation in a set of sequences
- All possible trees are considered
- No. of sequence changes/mutations that might have occurred is considered for each tree
- Method can be used to explore relationships between the most diverse sequences
- Computationally intense

Seq. A: **ACGCGTTGGG**

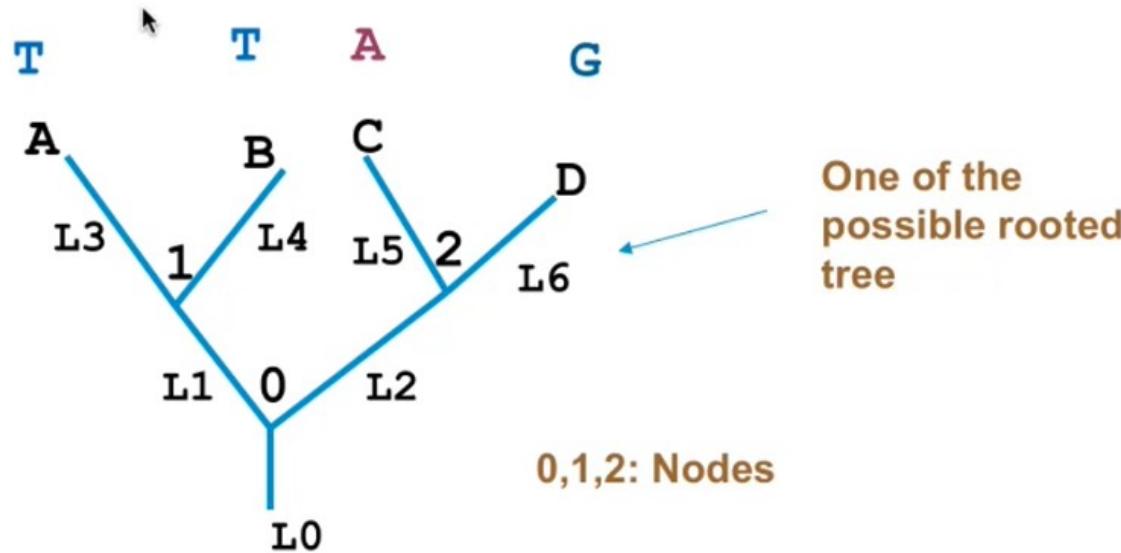
Seq. B: **ACGCGTTGGG**

Seq. C: **ACGCAAATGAA**

Seq. D: **ACACAGGGAA**



One of the
possible
unrooted tree



Four possible bases in three nodes (0,1,2)

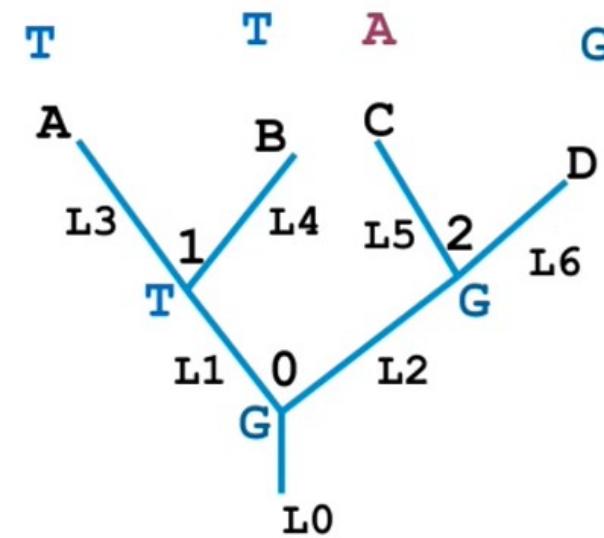
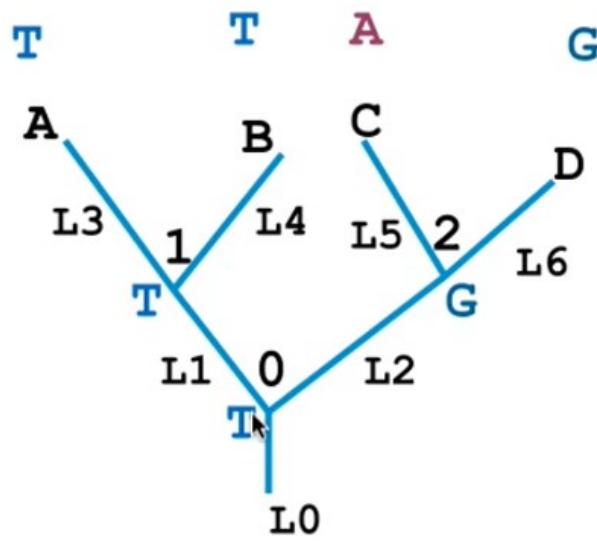
$4 \times 4 \times 4 = 64$ possible combinations

L1-L5: Likelihood values - probability of base change per site

Likelihood of the tree (L) =

$$L_0 * L_1 * L_2 * L_3 * L_4 * L_5 * L_6$$

Transversion: Purine \leftrightarrow Pyrimidine
Transition: Pyrimidine \leftrightarrow Pyrimidine
Purine \leftrightarrow Purine



$$L(\text{Tree}) = L(\text{Tree1}) + L(\text{Tree2}) + \dots + L(\text{Tree64})$$

Calculation is repeated for all other columns

Maximum Likelihood Method

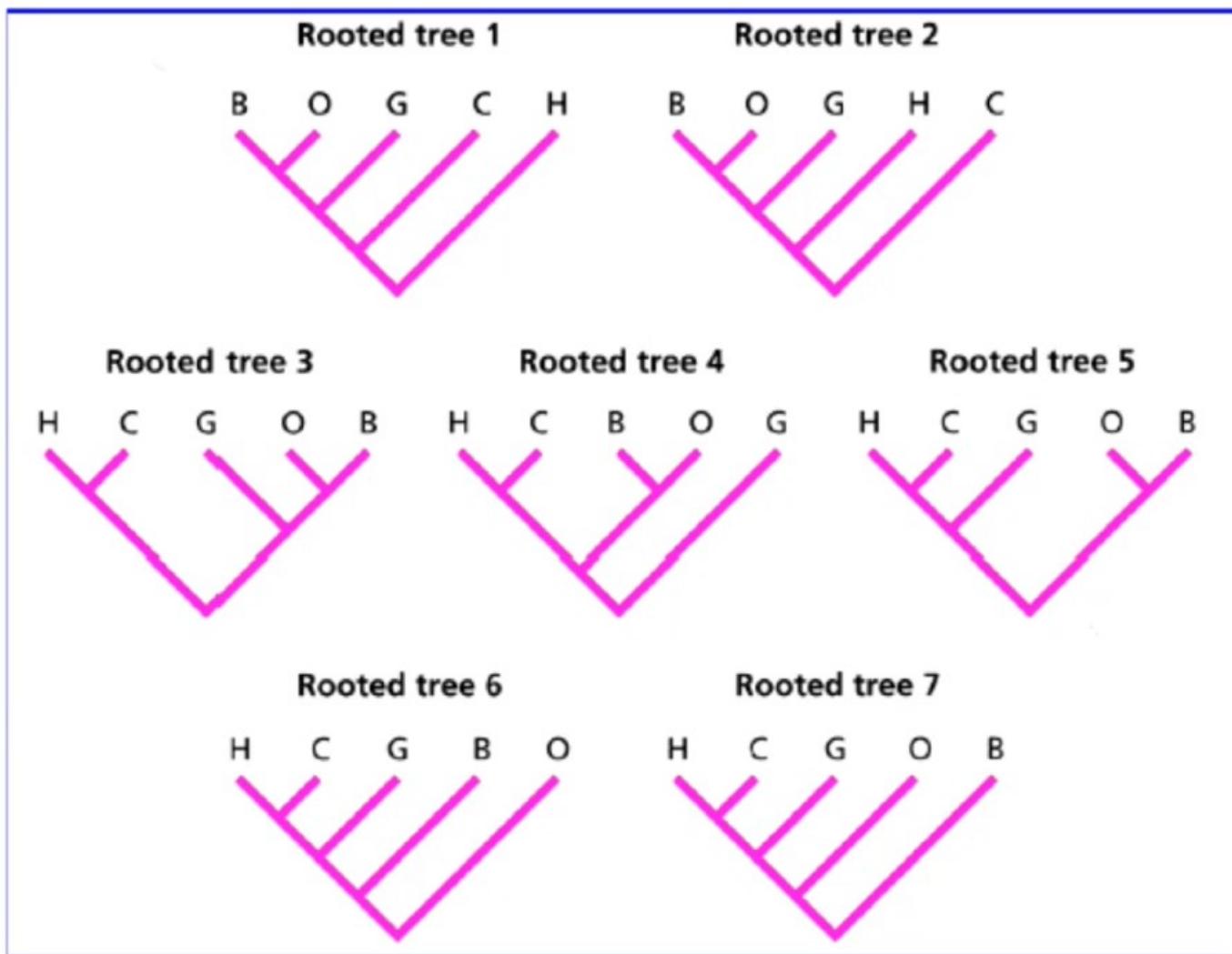
→

Find the likelihood of the trees for all columns

Each of the three possible trees is evaluated

One with the highest likelihood score is identified

Assessing Tree Topologies



Comparing Tree Topologies

- To resolve ambiguities in tree topology
- Useful when the trees produced by several methods have been applied to a single data set or by the same method applied to different data sets.
- All trees are treated equally in the analysis
- The analysis can identify support across a range of techniques or data for a given topology
- To estimate support for tree topology for a given data set constructed by a particular tree construction method bootstrap analysis is employed

Boot strapping

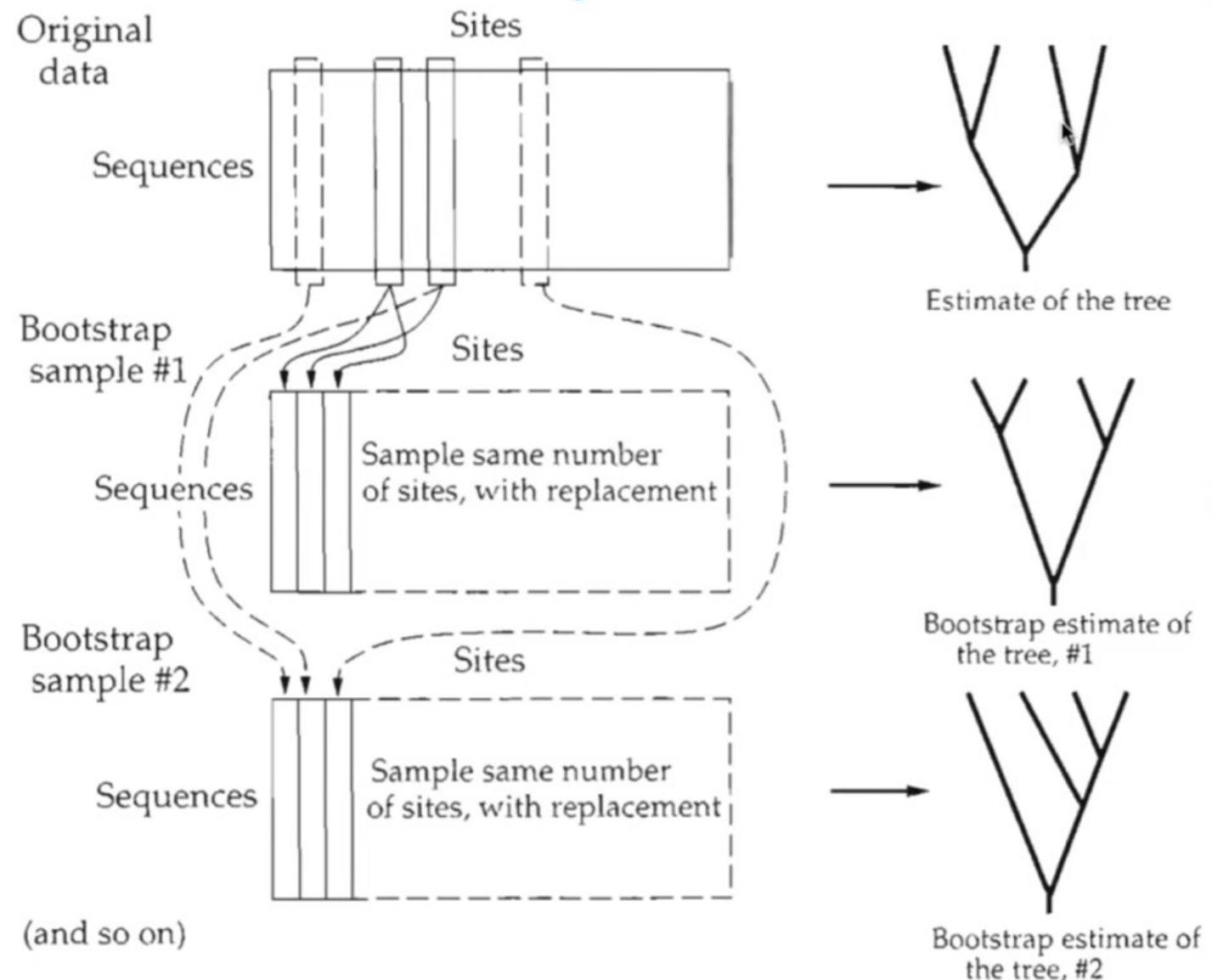
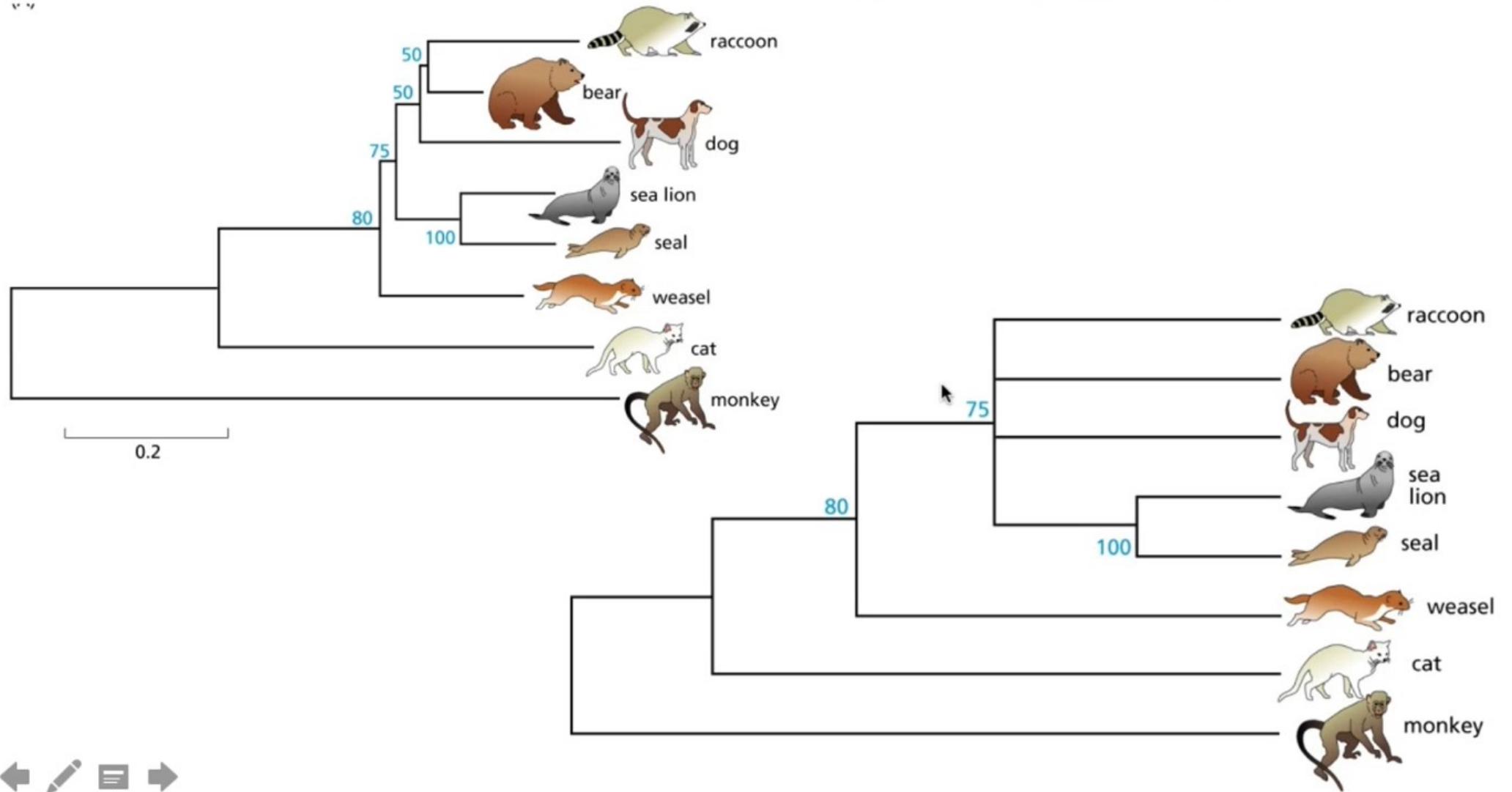
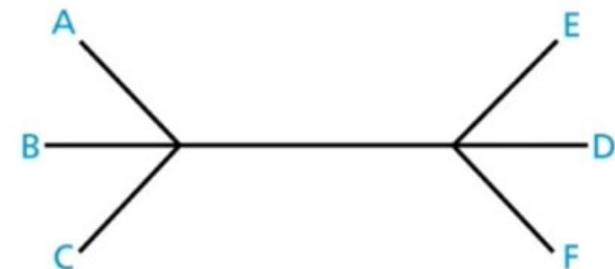
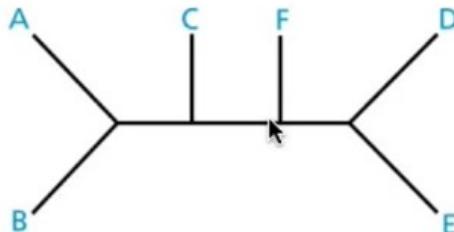
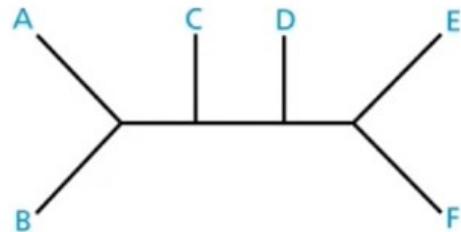


Image Credit: Felsenstein, 2004

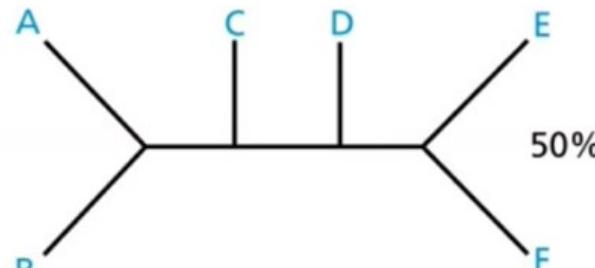
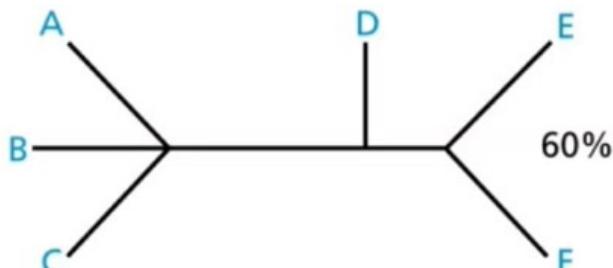
Condensed tree - internal branches not supported by bootstrap values removed



Consensus trees – show features that are consistent between trees



Strict Consensus



Majority Rule Consensus



Genome Analysis



Five approaches to genomics

Approach I: cataloguing genomic information

Genome size; number of chromosomes; GC content; isochores; number of genes; repetitive DNA; unique features of each genome

Approach II: cataloguing comparative genomic information

Orthologs and paralogs; COGs; lateral gene transfer

Approach III: function; biological principles; evolution

How genome size is regulated; polyploidization; birth and death of genes; neutral theory of evolution; positive and negative selection; speciation; epigenetics

Approach IV: Human disease relevance

Approach V: Bioinformatics aspects

Algorithms, databases, websites

Prominent web resources

Ensembl genomes

<http://www.ensembl.org>

<http://ensemblgenomes.org>

NCBI Genome

<https://www.ncbi.nlm.nih.gov/genome>

Dept. of Energy Joint Genome Institute (DOE JGI)

<http://jgi.doe.gov>

Genomes On Line Database (GOLD)

<https://gold.jgi.doe.gov>

UCSC

<https://genome.ucsc.edu>

<http://microbes.ucsc.edu>

Chronology of genome sequencing projects

1976: first viral genome

Fiers et al. sequence bacteriophage MS2 (3,569 base pairs, accession NC_001417). Bacteriophage are viruses that infect bacteria. →

1977:

Sanger et al. sequence bacteriophage φX174.
This virus is 5,386 base pairs (encoding 11 genes).
See accession J02482; NC_001422.

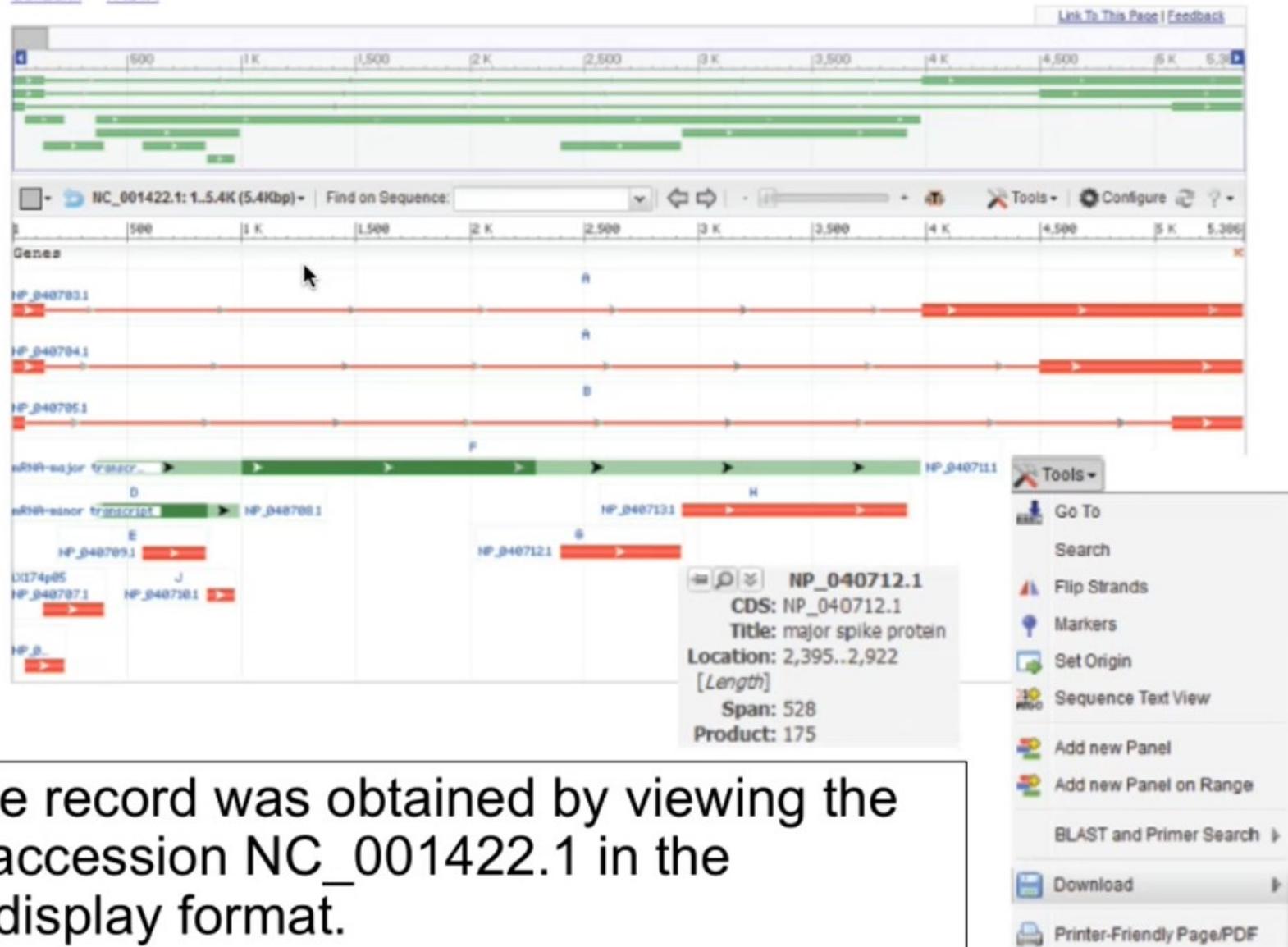
NCBI data for bacteriophage φX174

Enterobacteria phage phiX174 sensu lato, complete genome

NCBI Reference Sequence: NC_001422.1

GenBank FASTA

[Link To This Page](#) | [Feedback](#)



NCBI data for bacteriophage φX174

Sort the genomes list by taxonomy

Genome	Accession	Source information	Segm	Length	Protein Nbrs	Host	Created	Updated
Enterobacteria phage phiX174 sensu lato	NC_001422	-	-	5386 nt	11	77 bacteria	04/28/1993	04/17/2009

NCBI Genome record includes a summary of the accession number, length, number of proteins (11), sequence neighbors ($n = 77$), and host species.



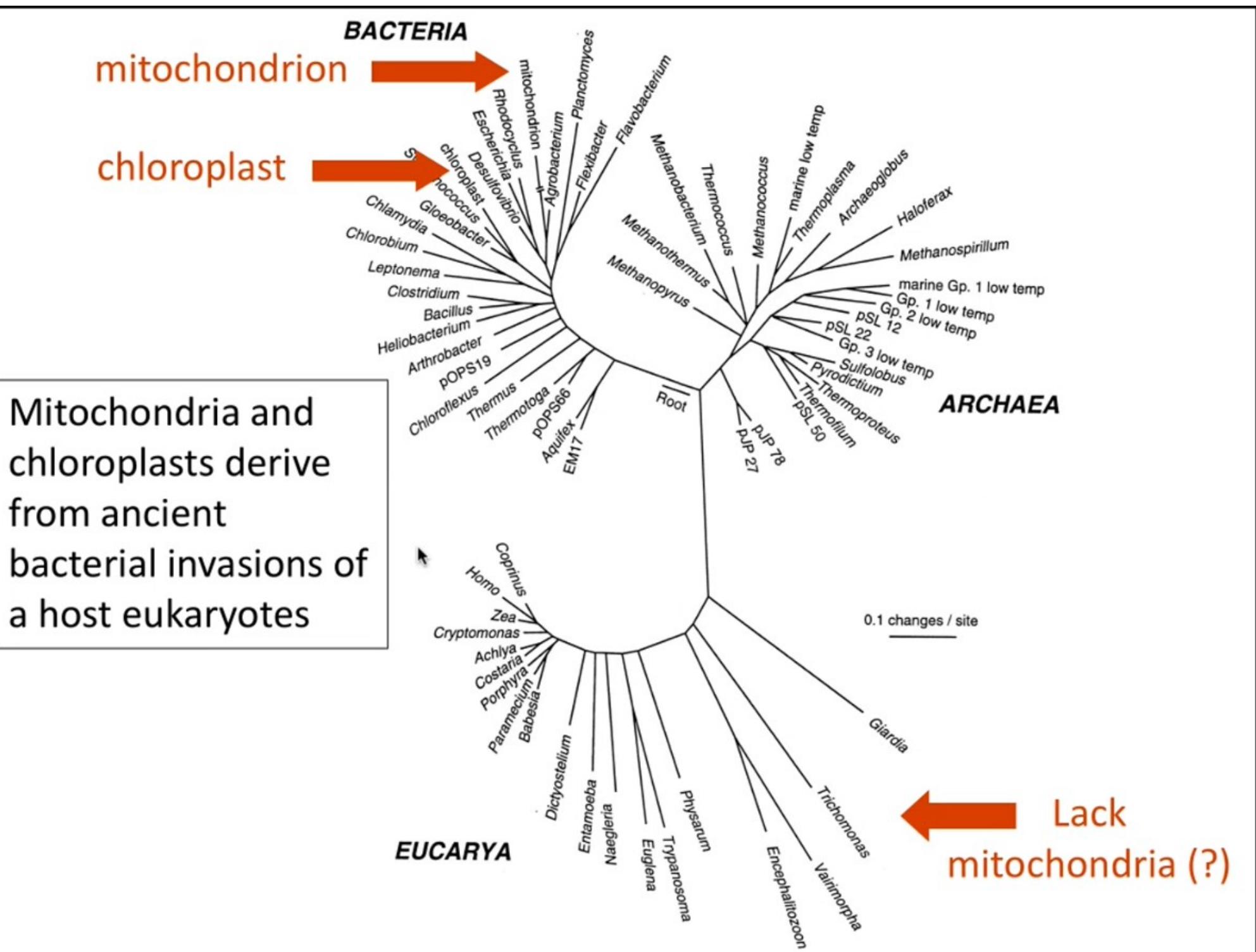
Chronology of genome sequencing projects

1981

Human mitochondrial genome: first eukaryotic organellar genome. 16,500 base pairs (encodes 13 proteins, 2 rRNA, 22 tRNA). Today (2021), over 12589 mitochondrial genomes sequenced.

1986

Chloroplast genome
156,000 base pairs (most are 120 kb to 200 kb)
>6000 chloroplast sequences (2021)



MitoMap: resource for organelle genomes

MITOMAP

A human mitochondrial genome database

A compendium of polymorphisms and mutations of the human mitochondrial DNA



Search MITOMAP for information on:

Perform search

Gene, disease, enzyme names may be abbreviated, truncated, etc. Examples of search words: ND1, ND4, NARP, LHON, 11778, 3243, Leu, Lys, etc.

Mitomap Quick Reference

[The Human Mitochondrial Sequence](#) updated

[Amino Acid Translation Tables](#)

[Mitochondrial References \(A-Z\) \(>1 MB\)](#)

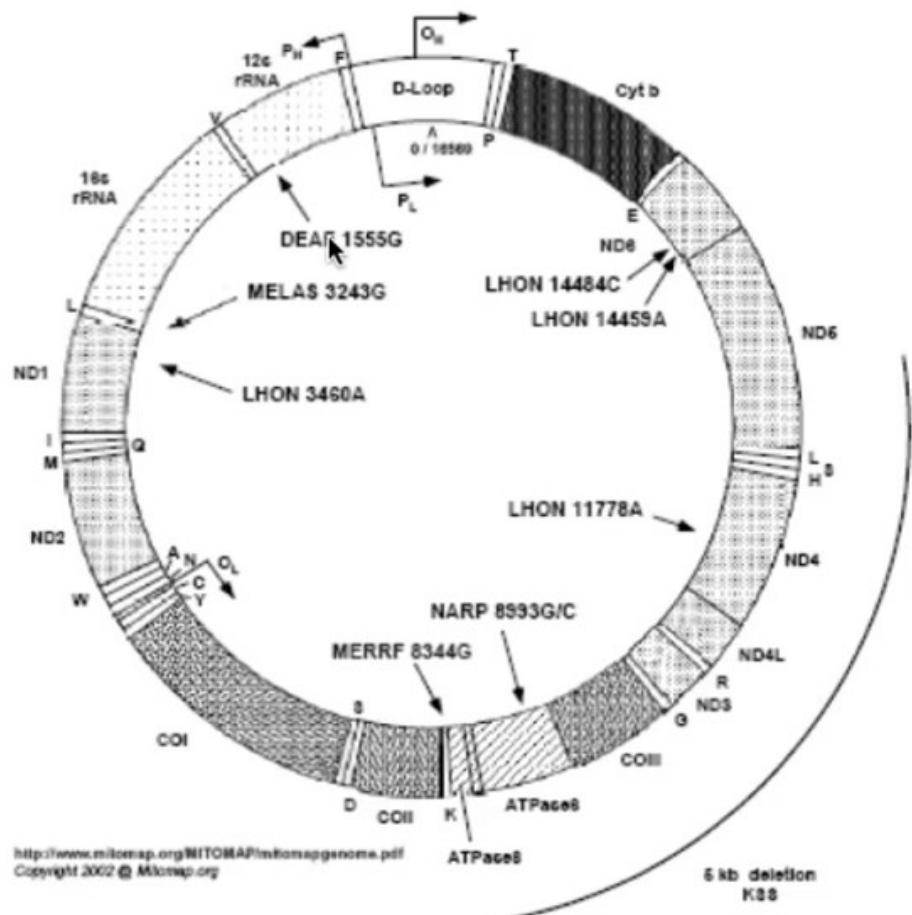
or view only [authors A-L](#) or [authors M-Z](#)

[Mitochondrial DNA Function Locations \(Gene Loci\)](#)

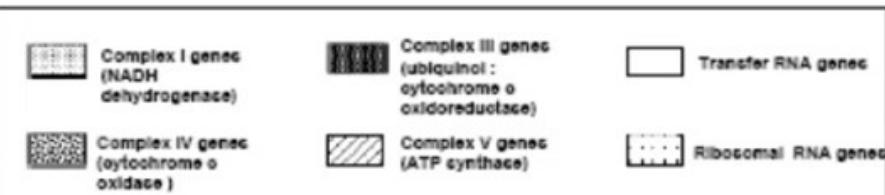
Illustrations

- Mitochondrial DNA Map
- Eleven pathological mutations in tRNA_{Leu(UUR)}
- Mitochondrial energetics
- Diabetes metabolism & the mitochondria
- World migrations
- mtDNA Tree

It is possible to map mutations in human mitochondrial DNA that are responsible for disease



Morbid map of the mitochondrial genome showing diseases associated with particular mutations.



Chronology of genome sequencing projects

1992: first eukaryotic chromosome.

Chromosome III of the budding yeast *S. cerevisiae* was sequenced.

1995: first genome of a free-living organism, the bacterium *Haemophilus influenzae*.

Accession: NC_000907.1

Size: 1,830,138 bp (1.8 Mb or megabase pairs)

Chronology of genome sequencing projects

1997:

More bacteria and archaea

Escherichia coli

4.6 megabases, 4200 proteins (38% of unknown function)

1998: first multicellular organism

Nematode *Caenorhabditis elegans*

97 Mb; 19,000 genes.

1999: first human chromosome

Chromosome 22 (49 Mb, 673 genes)

Chronology of genome sequencing projects

2000:

Fruitfly *Drosophila melanogaster* (13,000 genes)

Plant *Arabidopsis thaliana*

2001: draft sequence of the human genome
(public consortium and Celera Genomics).

Chronology of genome sequencing projects

- 2002: More completed genomes
- 2003: HapMap (catalog variation in human genome)
- 2004: chicken, rat; finished human chromosomes
- 2005: chimpanzee, dog
- 2006: sea urchin, honeybee
- 2007: rhesus macaque, first individual human
- 2008: platypus, first cancer genome
- 2009: bovine, methylome map
- 2010: 1000 Genomes pilot, Neanderthal
- 2011: genomics vision
- 2012: Denisovan, bonobo
- 2013: comb jelly, 700,000 year old horse
- 2014: primates, plants, ancient hominids
- 2015: diversity in Africa
- 2021: Complete Human Genome

Overview of genome analysis

- Selection of genomes for sequencing
- Sequence one individual genome, or several?
- How big are genomes?
- Genome sequencing centers
- Sequencing genomes: strategies
- When has a genome been fully sequenced?
- Repository for genome sequence data
- Genome annotation

Applications of genome sequencing

Purpose	Template	Example
De novo sequencing	Genome sequencing	Sequencing >1000 influenza genomes
	Ancient DNA	Extinct Neandertal genome
	Metagenomics	Human gut
Resequencing	Whole genomes	Individual humans
	Genomic regions	Assessment of genomic rearrangements or disease-associated regions
	Somatic mutations	Sequencing mutations in cancer
Transcriptome	Full-length transcripts	Defining regulated messenger RNA transcripts
	Noncoding RNAs	Identifying and quantifying microRNAs in samples
Epigenetics	Methylation changes	Measuring methylation changes in cancer

Genome sequencing has many applications and purposes.

Large-scale human genome sequencing projects

Recent large-scale projects include:

- Human-centered projects
- Human Microbiome Project
- For nonhuman species, large numbers of strains, inbred lines, or related organisms

Large-scale human genome sequencing projects

Project	URL	Goal
The 1000 Genomes Project	http://www.1000genomes.org/	Find human genetic variants having frequencies >1%
International Cancer Genome Consortium (ICGC)	http://www.icgc.org/	Catalog mutations in tumors from 50 cancer types
UK10K	http://www.sanger.ac.uk/about/press/2010/100624-uk10k.html	Sequence the genomes of 10,000 UK individuals
100,000 Genomes Project	http://www.genomicsengland.co.uk/	Sequence 100,000 individuals in the UK
Autism Genome 10K Project	http://autismgenome10k.org/	Sequence 10,000 autism-related genomes
Personal Genome Project	http://www.personalgenomes.org/	Effort to sequence 100,000 human genomes

Currently separate projects involving sequencing of hundreds of thousands of exomes/genomes are in progress.

Large-scale human genome sequencing projects

Project	URL	Goal
The 1000 Genomes Project	http://www.1000genomes.org/	Find human genetic variants having frequencies >1%
International Cancer Genome Consortium (ICGC)	http://www.icgc.org/	Catalog mutations in tumors from 50 cancer types
UK10K	http://www.sanger.ac.uk/about/press/2010/100624-uk10k.html	Sequence the genomes of 10,000 UK individuals
100,000 Genomes Project	http://www.genomicsengland.co.uk/	Sequence 100,000 individuals in the UK
Autism Genome 10K Project	http://autismgenome10k.org/	Sequence 10,000 autism-related genomes
Personal Genome Project	http://www.personalgenomes.org/	Effort to sequence 100,000 human genomes

Currently separate projects involving sequencing of hundreds of thousands of exomes/genomes are in progress.

Criteria for selecting genomes for sequencing

Criteria include:

- genome size (some plants are >>>human genome)
- cost
- relevance to human disease (or other disease)
- relevance to basic biological questions
- relevance to agriculture

Diversity of genome sizes

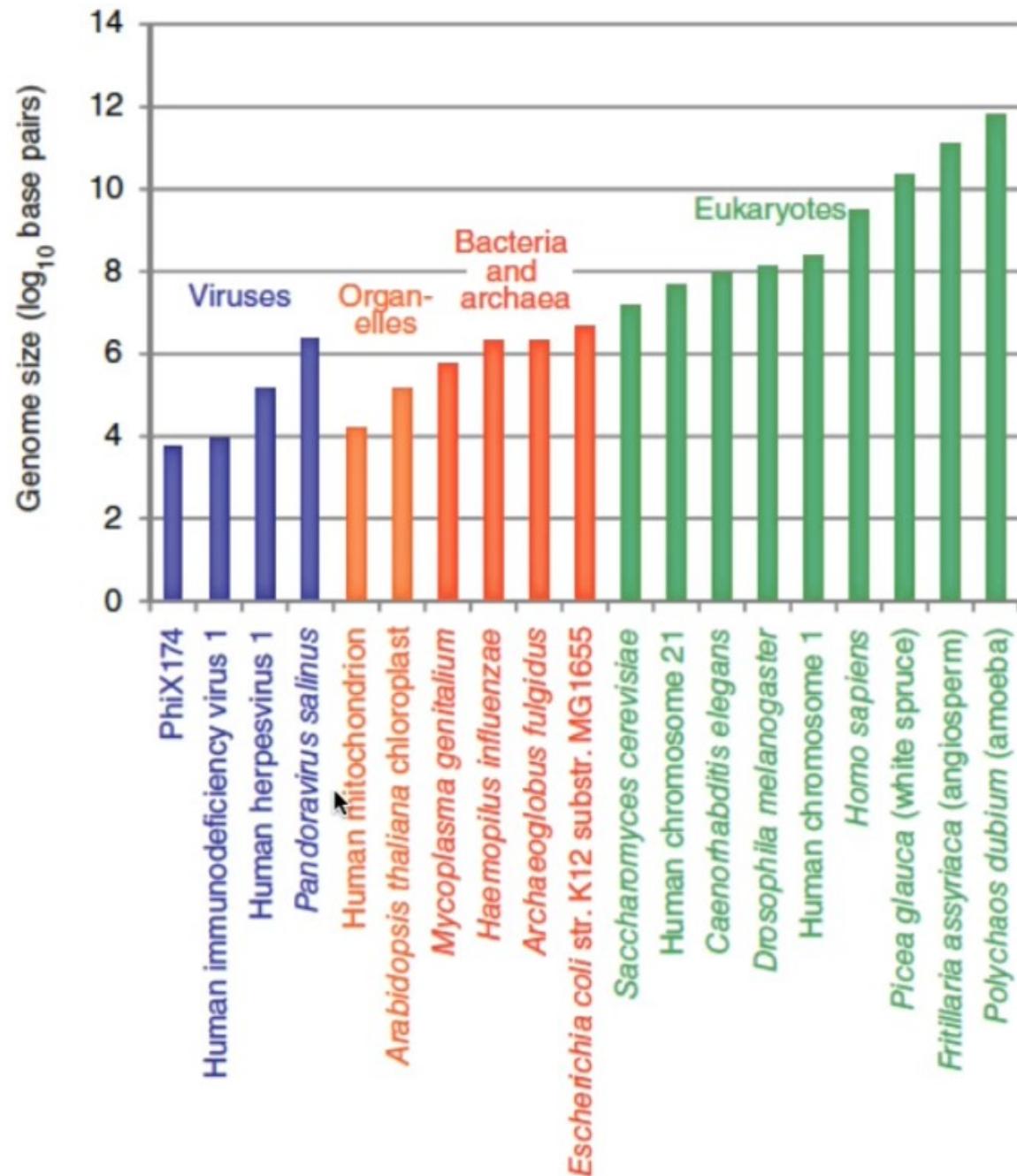
How big are genomes?

Viral genomes: 1 kb to 350 kb (Mimivirus: 1181 kb)

Bacterial genomes: 0.5 Mb to 13 Mb

Eukaryotic genomes: 8 Mb to 686 Gb (human: ~3 Gb)

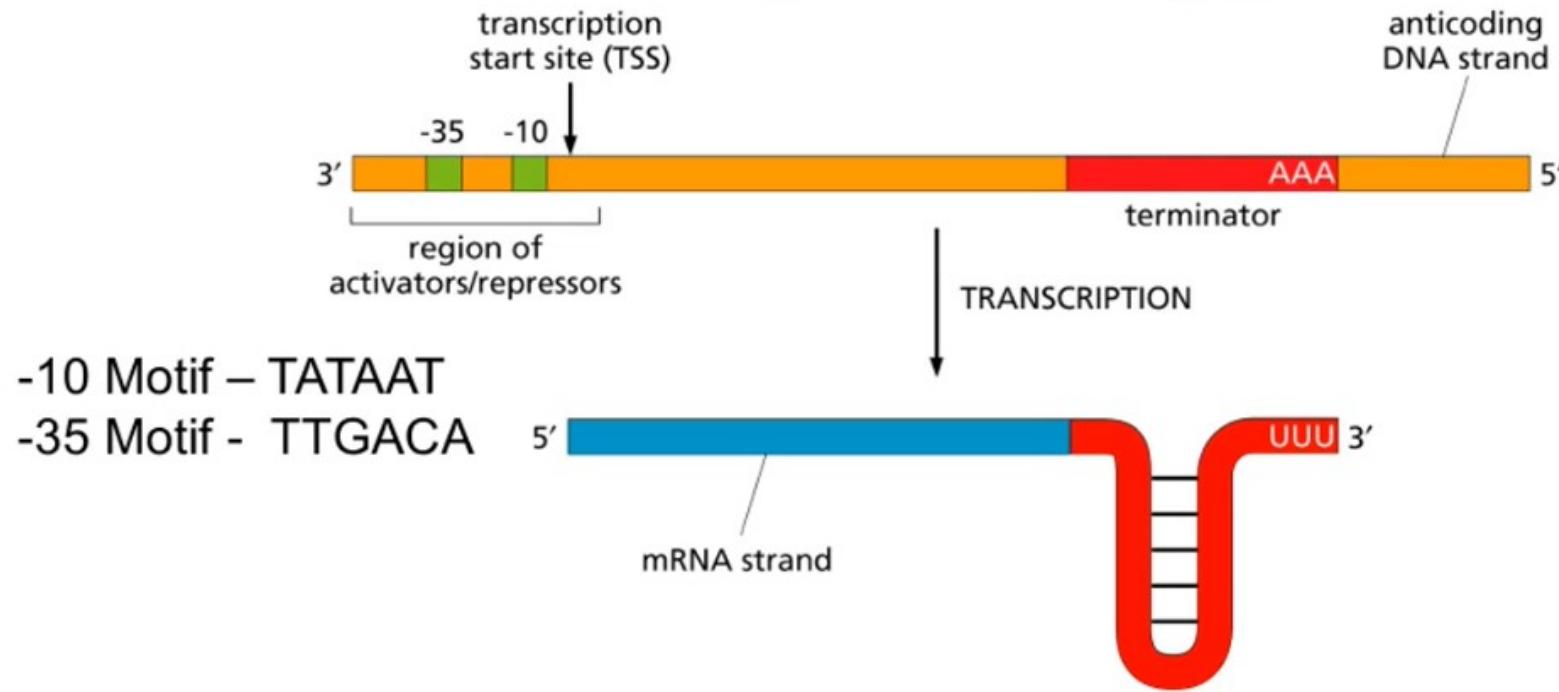
Genomes sizes from viruses to eukaryotes



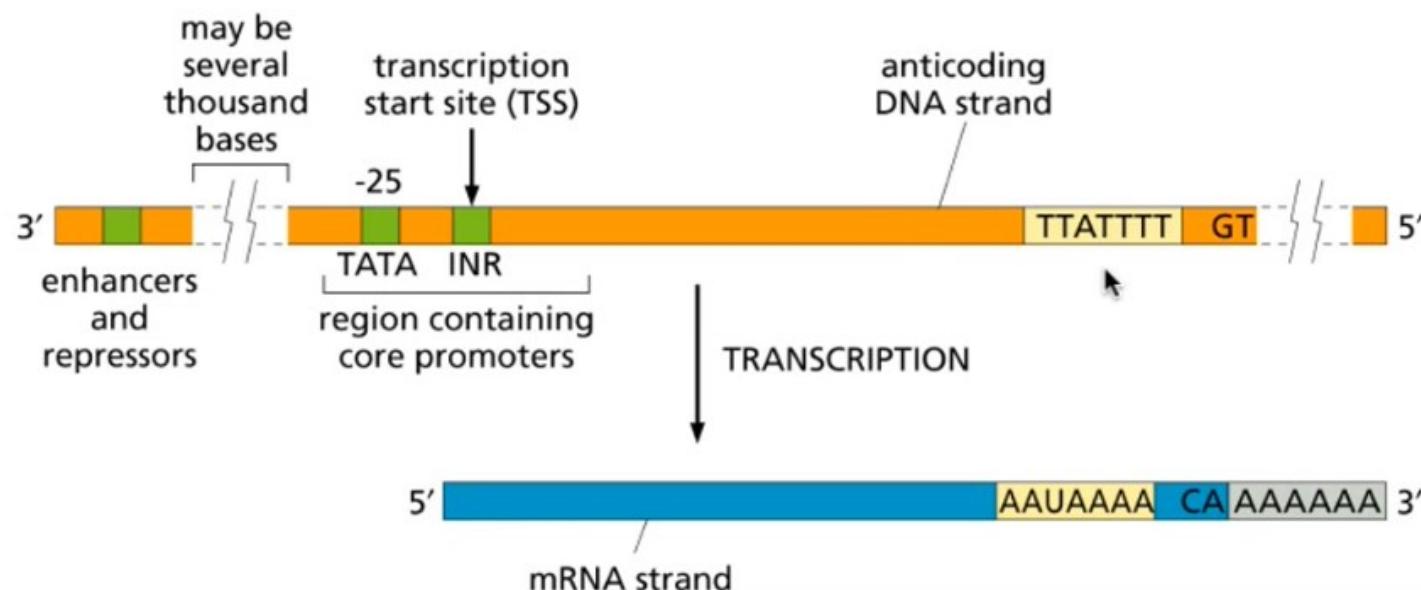
Genome Annotation

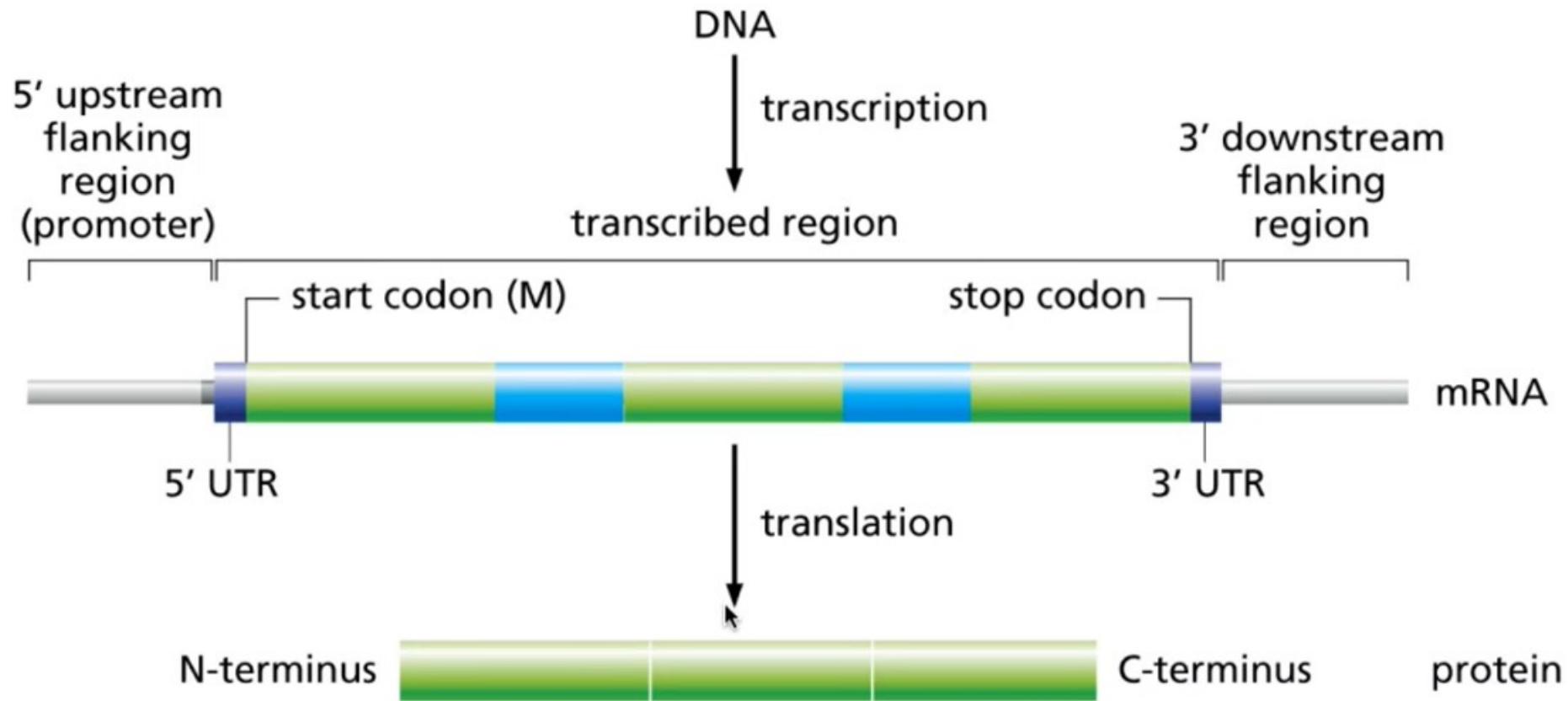


Prokaryotic Transcription



Eukaryotic Transcription



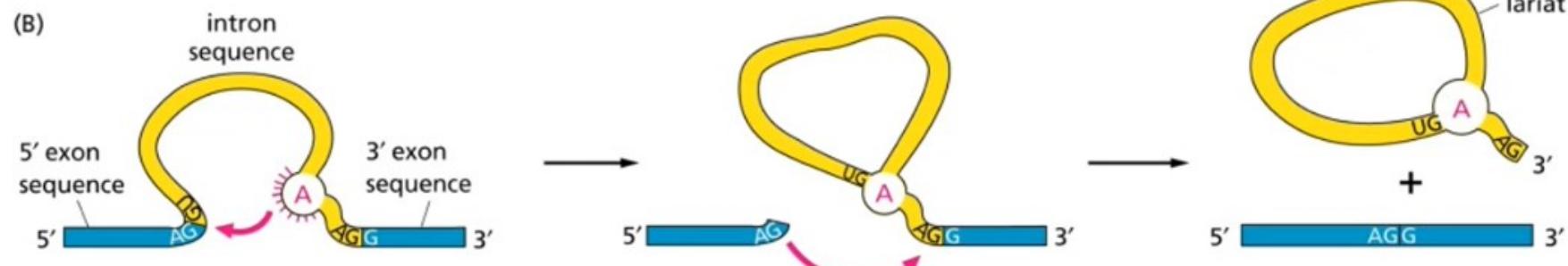


Modification

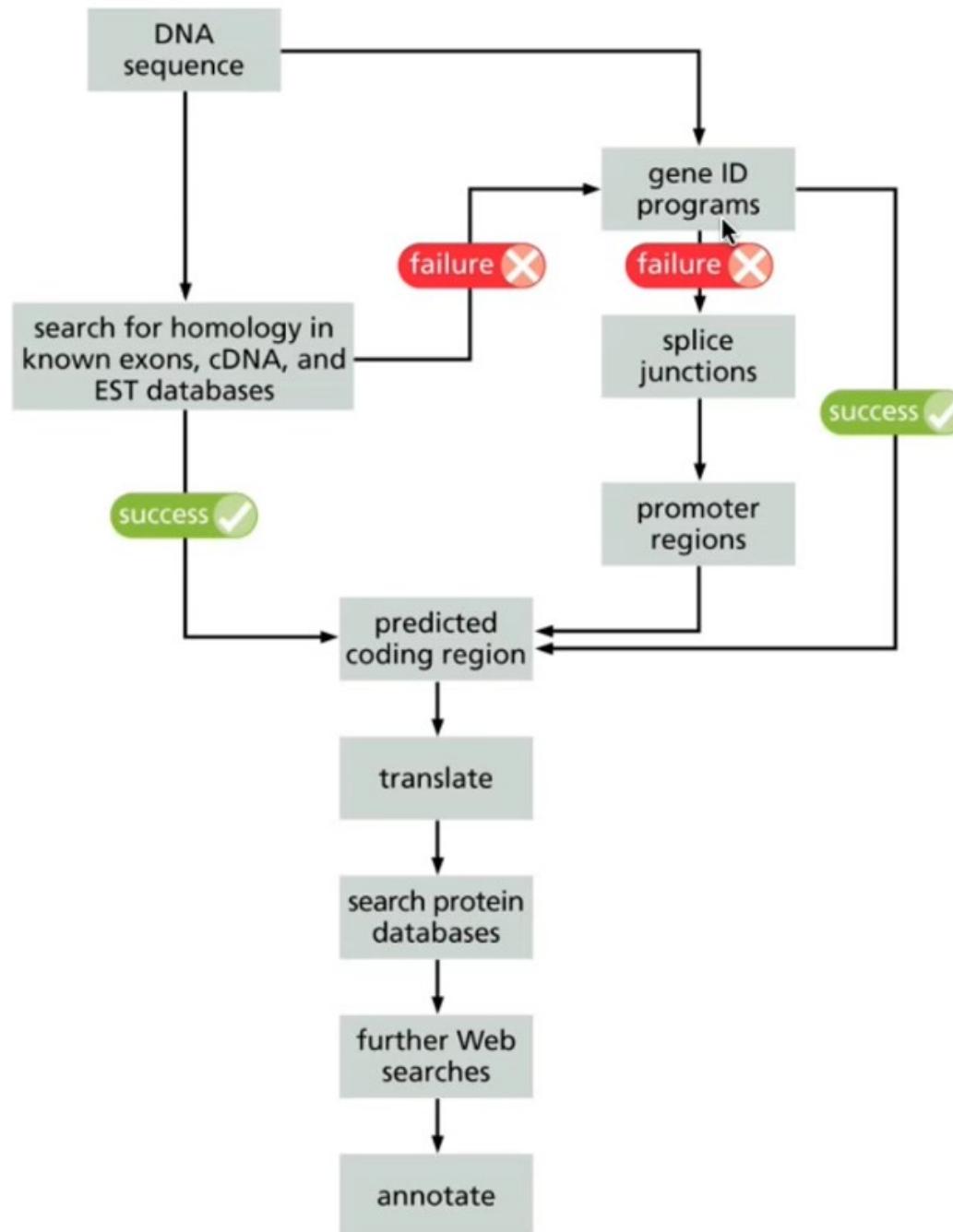
(A)



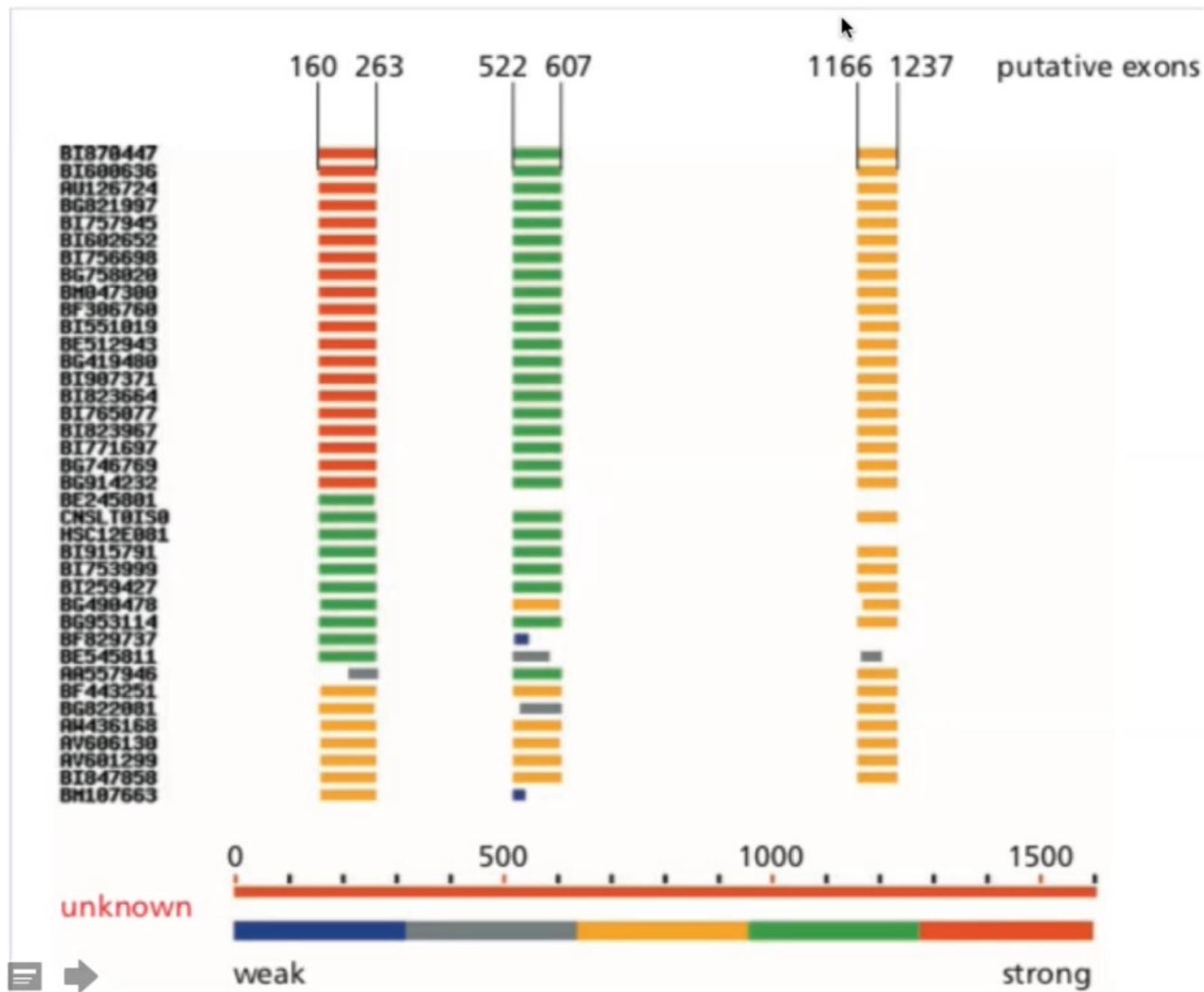
(B)



Steps Involved in the identification and annotation of gene sequences



Similarity search against EST database



Typical Series of Steps in Eukaryotic Gene Prediction

1. Submit DNA sequences to exon prediction programs
2. Take average or consensus exon prediction
3. Translate predicted exon into protein sequences in all frames and directions
4. Take translation with least or no stop codons
5. Search protein database with the translated segment
6. If hit found use protein homolog to delineate exon(s)
7. Repeat for all other exons
8. Annotate and splice exon to obtain putative protein sequence

Gene Prediction



cgatccaagg	agcccgacgc	ctaggccgga	cccgcgggag	cgtctattga	gtaaccgttg
ttagggaga	cgaaagccgg	gaaggagctt	tcgcgcctgc	gccgcggggc	cgtcgcgtc
tgcgcctcg	cgcaagagag	gcgggg	at	qgcggagcca	gatctggagt
gacagccat	ccgtctgaag	tgtattcgta	aggagggctt	tcacgggt	cctccggAAC
acagggtgcg	ccccgggtcc	ccccggcagc	tctggccg	tcgcgtacgg	cactgccccgg
ctgggttcgg	ggggcctcg	gtcgcgtgc	cgcggcgg	cttccggcac	ggggggggaaac
gacagtccca	gagggtcccg	cgccgggggc	ggaaggccgg	gccccgggg	ctcaggggacc
ccgacagccg	gtccctggag	atctgagggg	ccggggccgg	ctcagggtatg	cttcgcggcc
gcggagagac	ggggccggga	cttgggaaa	gcaggctccg	ggatccagct	cgggcgctgc
tgggttcagc	gcccggagctg	ggctttgcag	gctgagccgc	gagccactt	tttgggggaa
gaatttacac	ccgcgacgag	ttcgagctt	aaggctcccg	ctggggcttgg	gctgcgtatgg
gcgggggtcca	gcccgtctgg	cggcttcac	caaacccctgg	gccccgtccct	tgggggggtgt
ccgggtccct	cctctggccg	cttcacgcg	accctgggc	tcccttgcgg	agggtcccggt
gcctccctct	cgacagacca	ggagagagcc	tctggggc	tggggcgtgt	tgcgggggtgt
cacggcttgg	gggggtggag	agccctgaac	tttggccgt	ggtttgtt	tttaactgtct
ggctgggtct	ctgagaggcc	aagccacgt	ttcagtaaga	atcattaaca	gatcgtggct
cgggcaggct	gctgcagctg	gaaacctga	gcttttagta	actccacagg	ggcagcagag
tcaggctcta	gcccatttt	tgcttccaac	tctgcctagc	tgtcaccac	ctgtggataa
tgaagttct	gtggcagata	ctggagggcc	cggagcacct	gcagcacat	gcaactgtacc
tgtcagagg	gtcagaggtt	ccagtcaagg	ttagcctca	ggggatgac	aacaggggcc
acacttact	cattccagacg	tgcagcaggc	gctaccgcg	ttgagagccg	tgtctggac
tctgggtgt	agtagtgagt	tggctttcc	ttttttttt	ttttttttag	acaggggtcta
gcgcctctgc	ccaggcttgg	atgcagtggc	gtgatcacgg	ctcagtcat	cgtcagcctc
ctgggtctcaa	gtgatcttcc	tgccctggcc	tcccaagtgc	tgaattaca	ggcgtgagcc
accggccccc	gcctggctg	gcattttttt	gagttcagg	agtgtgacaa	ggattttggac
acccagaaat	aagcgtgtcg	agaagagac	aagcagaggg	tgtgagaagg	.
.
.	.	.	gtttca	gctggggacga	tgcggggagtg
tgaaggagtt	tgagaagctg	aaccgcattt	gagagggtac	tcacggcatt	gtgtgtgatg
.
atcggggccc	gggacacccca	gacagatgag	attgtcgac	tgaagaaggt	gccccatggac
bagggagaagg	atg
.	cag
gcatccccat	cagcagctt	cgggagatca	cgcgtctgt	ccgcctgcgt	catccgaaca
tcgtggagct	gaaggaggtt	gttgtgggg	accacccctgg	gag	.
.
catcttcttg	gtgatgggtt	actgtgagca	ggacccctgg	agccctctgg	agaatatgcc
aacacccttc	tcggaggctc	ag.	.	.	.
.
gtcaagtgca	tcgtgtcgca	ggtgctccgg	ggcctccagt	atctgcacag	gaacttcatt
atccacacag
.
cagggacacttgaag	gtttccaaact	tgctcatgac	cgacaagggt	tgtgtgaaga	cag
.	cag
cggattttcg	cctggccccc	gcctatggt	tcccaagtaaa	gccaatgacc	ccccagggtgg
tcactctctg
.
gtaccgagcc	cctgaactgc	tgttgggaaac	caccacgcag	accaccagca	tcgacatgt
.
ggctgtgggc	tgcataactgg	ccgagctgt	ggcgacagg	ccttttctcc	ccggcacttc
cgagatccac	cagatcgact	tgatcgta	gctgtgggg	acgcccagt	agaacatctg
ggcg
.
ggcttttcca	agctgcccact	ggtcggccag	tacagccctc	ggaaggcagcc	ctacaacaac
ctgaagcaca	agtccccatg	gctgtcgag	gccgggctc	gcctgtcgca	cttccctgtt
atgtacgacc	ctaagaaaaag
.
ggcgacggcc	ggggactgtcc	tggagagctc	ctatttcaag	gagaagcccc	tac
.
c...g...g...cc	ggagctcatg	ccgacccctt	cccaccaccc	caacaaggcgg	ggcccccag
ccac...t...ga	ggggccagagc	aagcgtgt	aaccctgt	.	.
.
EXON 1					
Pred					
EXON 2					
EXON 3					
EXON 4					
EXON 5					
EXON 6					
EXON 7					
EXON 8					
EXON 9					
EXON 10					
EXON 11					
EXON 12					
EXON 13					
EXON 14					

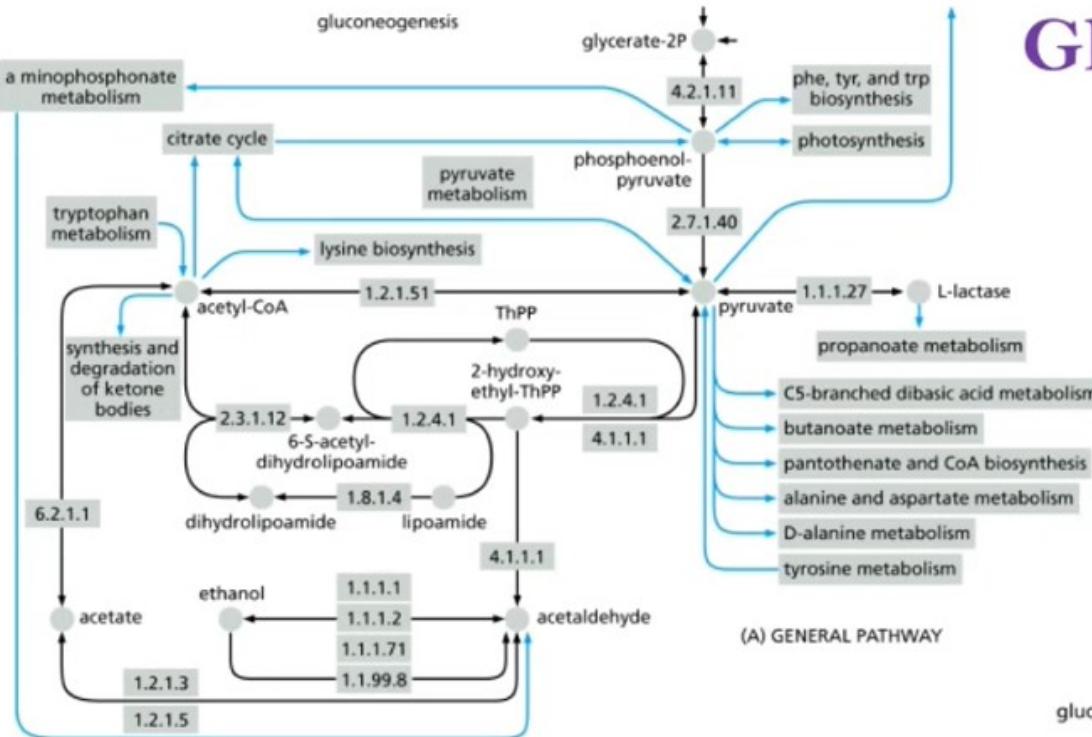
Prediction of start site is a nontrivial task



ttaaaaaggg	aggagcgggc	tggagggaa	agagggagaa	catggtcatt	actgaatcca	60	
cacattgcac	aatagaaaa	aggaacaggc	agggaaatag	tcaattatgt	attgcctcc	120	
tgctgttaa	atcagcactt	cagtaagata	aggtgaggac	agagcagcta	cctgtggga	180	
cattaacct	tttatctgt	gctatctgct	taggacata	gagaaaggca	gnttcttgc	240	
tgactcagct	ttttgcttaa	tttttcctt	ttggcatatg	aattgagctc	ccacngnnt	300	
ttggttggtt	ttgggcataa	gtggagagtt	caattggggc	cagggccccg	agtattttc	360	
tttcacaaa	tataccctt	agagttcaga	gaaaggctg	ggatttagagc	cttctttgag	420	
cactattcat	ggattacatg	aggagtngc	tgtaaaaga	ggcccccgg	ctaagccctg	480	
ggcacctgaa	caagtgagag	gtcaggaagg	ggagaaaaat	ccagcaaagg	aaccccagag	540	
gaccgggcc	gtagctaaaa	agaaaagtgc	ggaagagtgc	aggaaatgg	taaagaaaaag	600	
aaagctttcc	acaagacagg	accgatcagc	tgagccaact	gctgctgagt	acagtatgta	660	
gagctgaagt	ctctcgccc	aaaggcaata	aggccccgg	ggaaacctg	aggcagcgg	720	
gttctctctt	gggagcggga	acacggaacc	atggcagcgc	aggcaaatgc	agcttggga	780	
gaagctctga	gtaacggacg	tgcgaggatg	atgcttgc	aagaaggaaa	tgaccgctgt	840	
ccggctccg	cgggagacga	acccggctt	cccgcctca	gggactagct	ctccagggaa	900	
ctgggacggt	cagtgccgt	cgagggcagct	cctcgctgaa	ggaaggagca	ggggacaggg	960	
aacggaatgg	ggagcgctag	cccccaacta	cgtccatctg	gccatgtttg	aagagccana	1020	
aatggagga	gggaacccct	agagcgtgcc	agacggagac	tgctctccgt	ngcagtcggg	1080	
gcgcttccgg	cagggcgcnn	actcccagcc	gagccccc	cgcctgc	tccaggattc	1140	
ctcttcgccc	tttctggggc	cgcgggggg	cgctctcagn	aggatggcca	acaccttccc	1200	
tccatcccta	caccccgccg	ccccctgccc	gtggccgcgc	tcggctccc	cactgctcac	1260	
tccacccct	acatcccagc	ccgctgccag	agccggggag	aggcgggggg	ccgcgtggc	1320	
gagaccgtga	acagcggctg	tcacgtggc	cgcccaaggcc	aataggggtg	aggcttggg	1380	
tccagctcag	tcctcccccg	gcccctccga	ctggcagtgg	gactcagcgg	gcgtggaggt	1440	
cgcggctgag	cgagcgagcc	ctggggcaggt	gaattgtggc	tgtgggtga	cggtggagac	1500	
accccccga	gggaggcgga	gggaaaggag	gcgaggcctg	cacctgc	atg	1560	
cccactcccc	agcgcccccg	gaccgtgcag	ttctctgcag	gaccaggcca	TGGAGCTCGA	1620	
AGTCCGGCGG	GTCCGACAGG	CGTTCCTGTC	CGGCCGGTCG	CGACCTCTGC	GGTTTGGCT	1680	
GCAGCAGCTG	GAGGCCCTGC	GGAGGATGGT	GCAGGAGCGC	GAGAAGGATA	TCCTGACGGC	1740	
CATGCCGCC	GACCTGTGCA	AGgtancacg	cgtgcggcg	ggtgtggga	aactggccc	1800	
cggcngcac	tttgtggactg	gagtcttcgg	ctgggtttt	ttttgtctt	tacatttngg	1860	
attactccac	cactgggagt	atgatctcca	gcgatacaga	taaagccaaa	gttcccgcag	1920	
actttccagg	tcctcttagca	ctcagaaggg	catatgttac	ctagcttctg	tggttcctt	1980	
tctgtatatt	agagaattag	caagcccta	ccagggcgtg	aagggtgcaa	aaggagtctg	2040	
aatggcaaac	agctagtctg	ataatgccag	ttgttgtcac	tacaggtgt	cctggtnnn	2100	
gttctgacat	tnagggccaa	gtgtatcata	cttacnctgn	aagnntaact	gtgattctct	2160	
tataacagAG	TGAATTCAAT	GTGTACAGTC	AGGAAGTCAT	TACTGTCCTT	GGGGAAATTG	2220	
ATTTTATGCT	TGAGAATCTT	CCTGAATGGG	TTACTGCTAA	ACCAGTTAAG	AAGAACGTGC	2280	
TCACCATGCT	GGATGAGGCC	TATATTCA	CACAGCCTCT	GGGAGTGGTG	CTGATAATCG	2340	
GAGCTTGGAA	TTACCCCTTC	GTTCTCACCA	TTCA	GATAGGAGCC	ATCGCTGCAG	2400	
gtctgggtgc	caccttatgt	ctatatacct	ttttagggag	gcttattttc	tcatattaaat	2460	
tggnat	taag	gatagtggct	aattaaatac	attacttgg	tgatttgcc	ttgtttacac	2520
caccagtgt	ta	ctggaattca	tacatccata	cata			

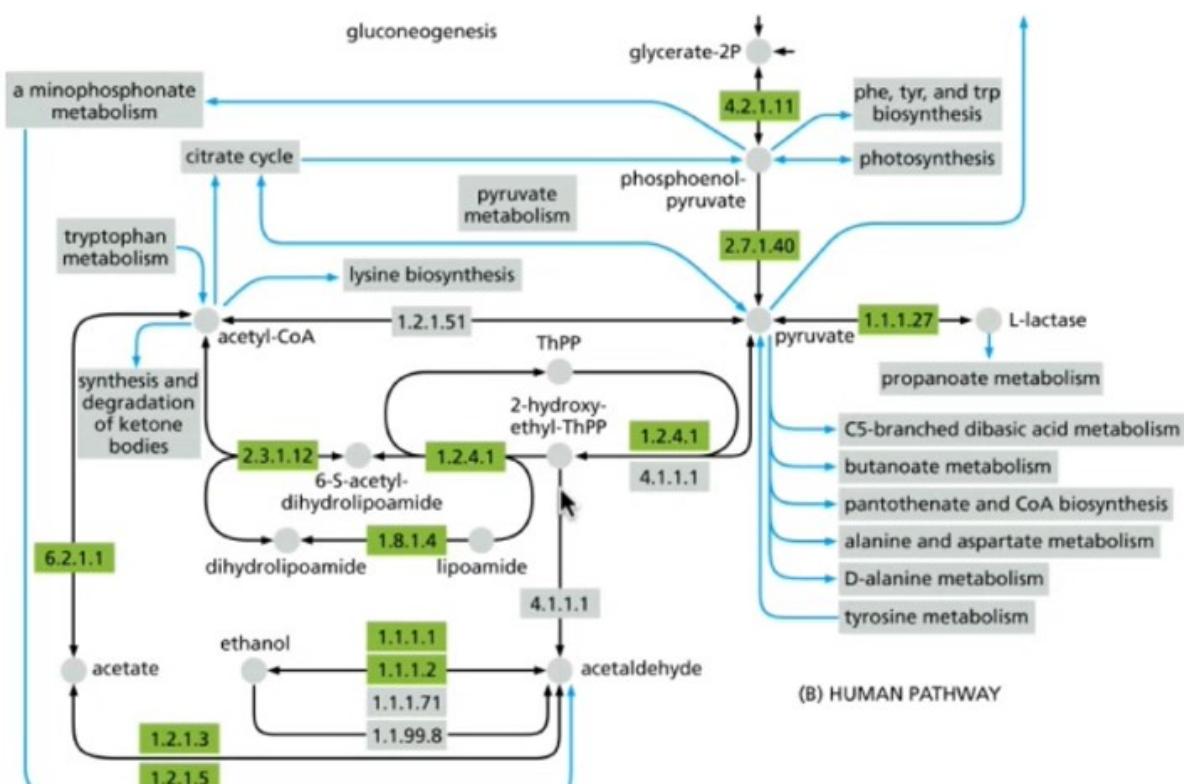


Gluconeogenesis Pathway



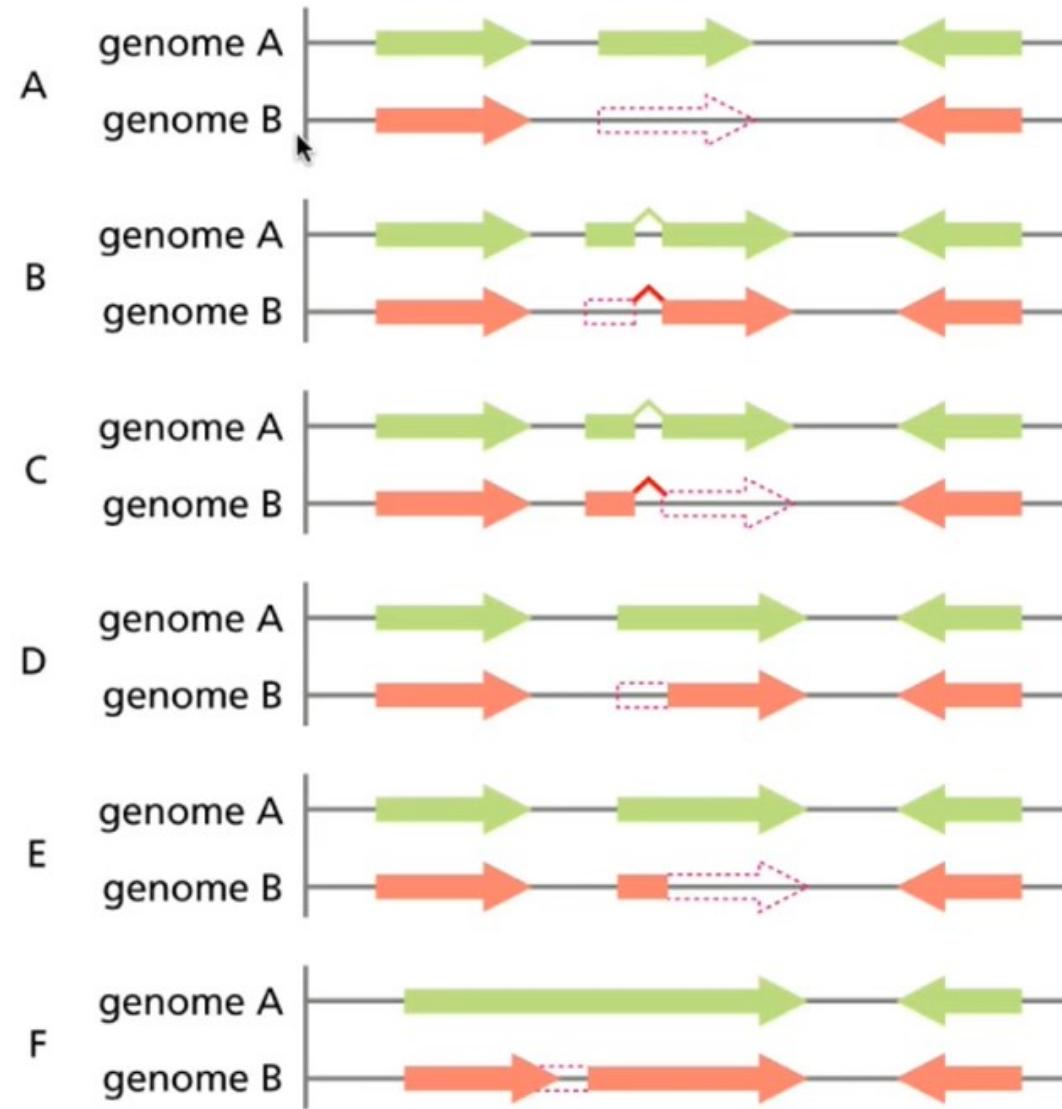
General Pathway

(A) GENERAL PATHWAY



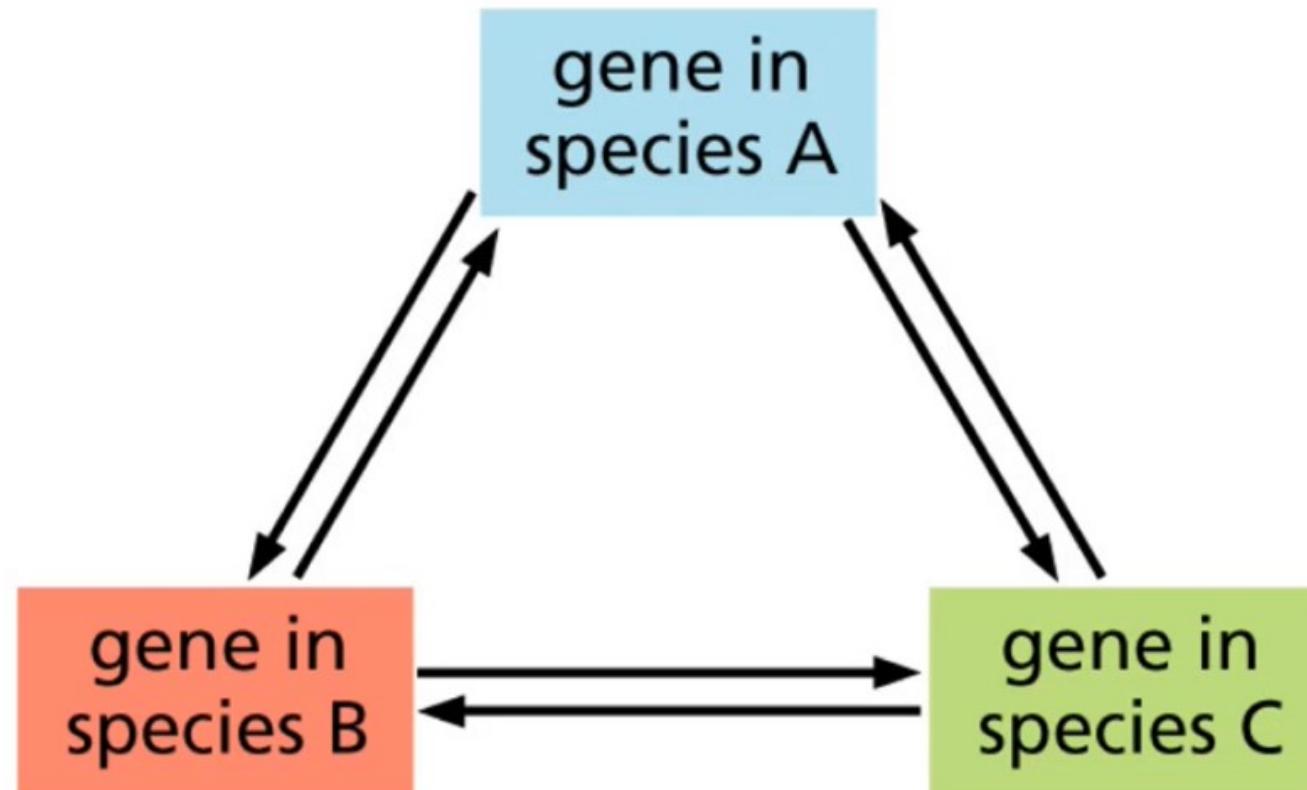
Human Pathway

Analysis of related genomic region helps to detect errors in annotation

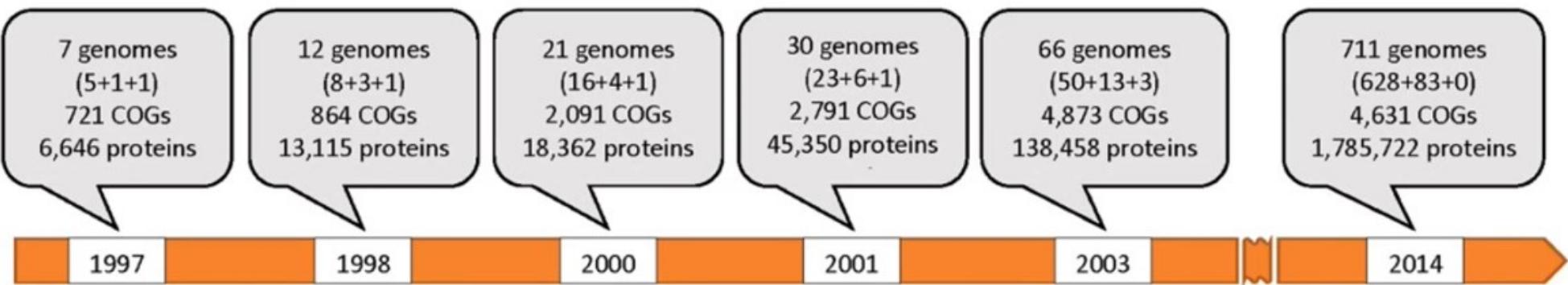


Comparative Genomics

Cluster of orthologous Groups (COG)

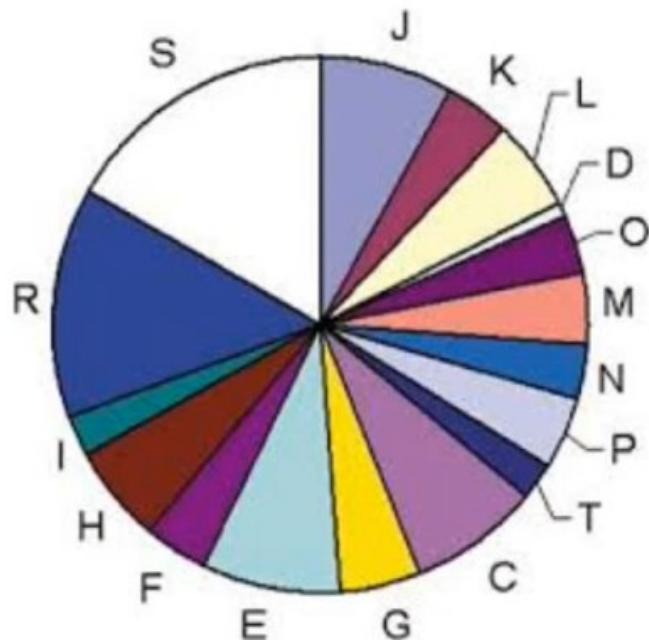


Orthology is established by bidirectional Best-scoring BLAST Hit (BeT).



2020 = 1300 genomes + 4877 COGs





One-letter abbreviations for the functional categories:

J, translation, including ribosome structure and biogenesis;

L, replication, recombination and repair;

K, transcription;

O, molecular chaperones and related functions;

M, cell wall structure and biogenesis and outer membrane;

N, secretion, motility and chemotaxis;

T, signal transduction;

P, inorganic ion transport and metabolism;

C, energy production and conversion;

G, carbohydrate metabolism and transport;

E, amino acid metabolism and transport;

F, nucleotide metabolism and transport;

H, coenzyme metabolism;

I, lipid metabolism;

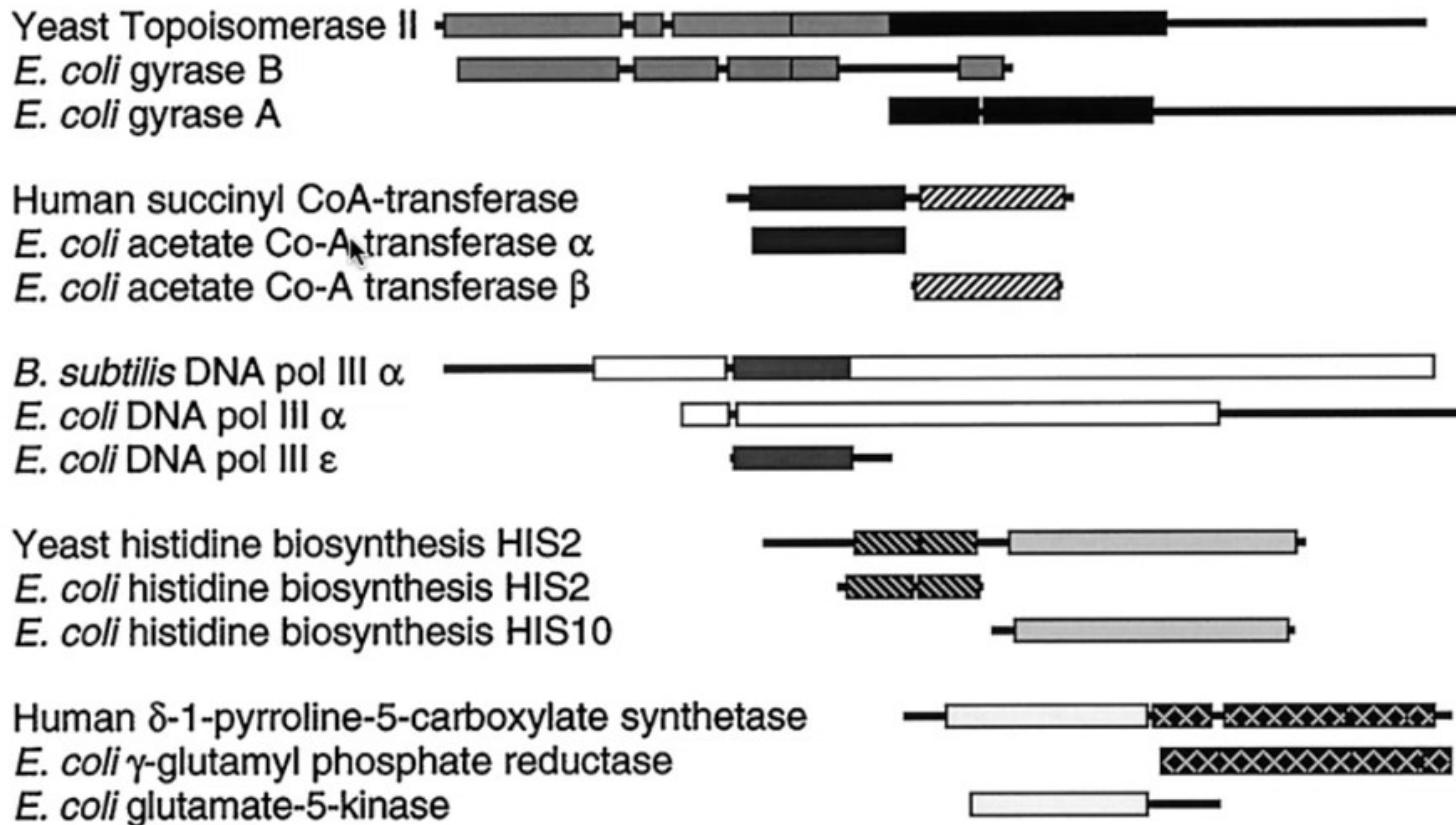
D, cell division and chromosome partitioning;

R, general functional prediction only;

S, no functional prediction

Genomic Context Based Methods

Gene Fusion/Fission



Inference: The observation that two independent genes in one organism are fused into a single gene in another organism suggests a direct physical interaction or a functional link between the genes.

Marcotte et al, 1999

Co-occurrence of Genes

Species 1



Species 2



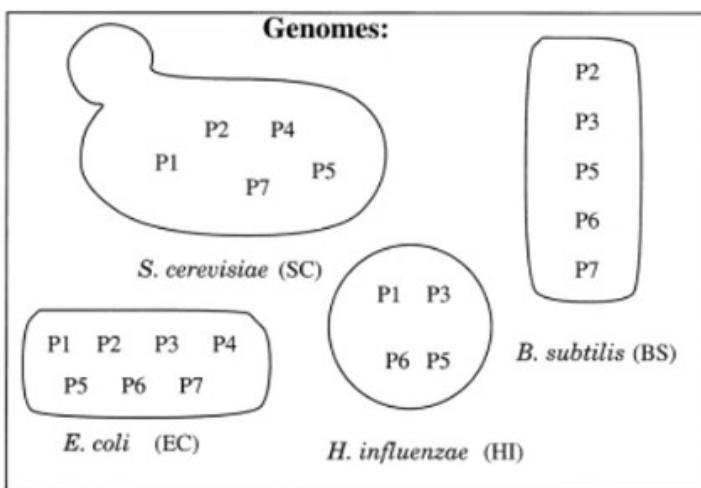
Species 3



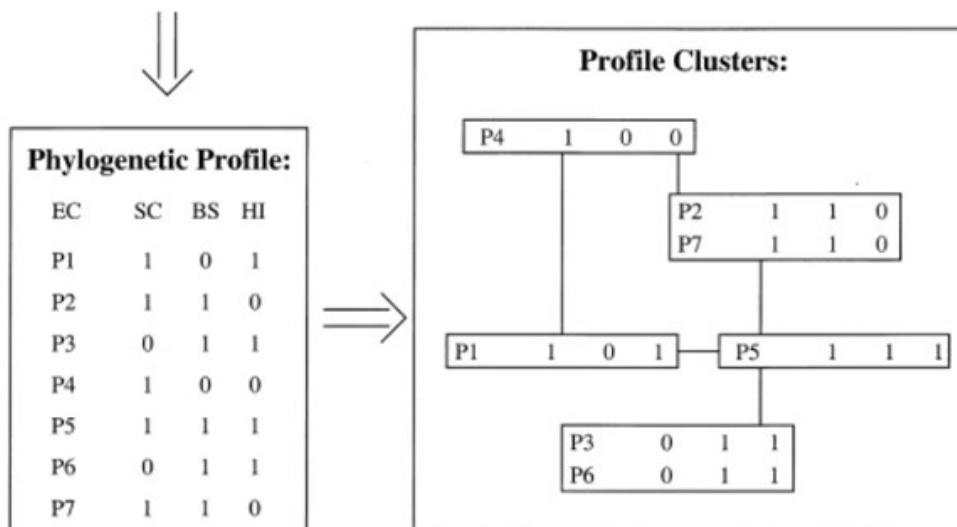
Inference: The observation that Gene B and Gene C have conserved their gene order in more than one organism suggests a selective pressure to regulate them together. Hence one can infer that the two genes must be functionally linked.

Overbeek et al, 1999

Phylogenetic Profiles



Inference: Since Protein 2 and Protein 7 cluster together, suggesting coevolution, it may be inferred that they are functionally linked.



Conclusion: P2 and P7 are functionally linked.
P3 and P6 are functionally linked

Pellegrini et al., 1999