



3000788 Intro to Comp Molec Biol

Lecture 5: Genome assembly and annotation

Fall 2025



Sira Sriswasdi, PhD

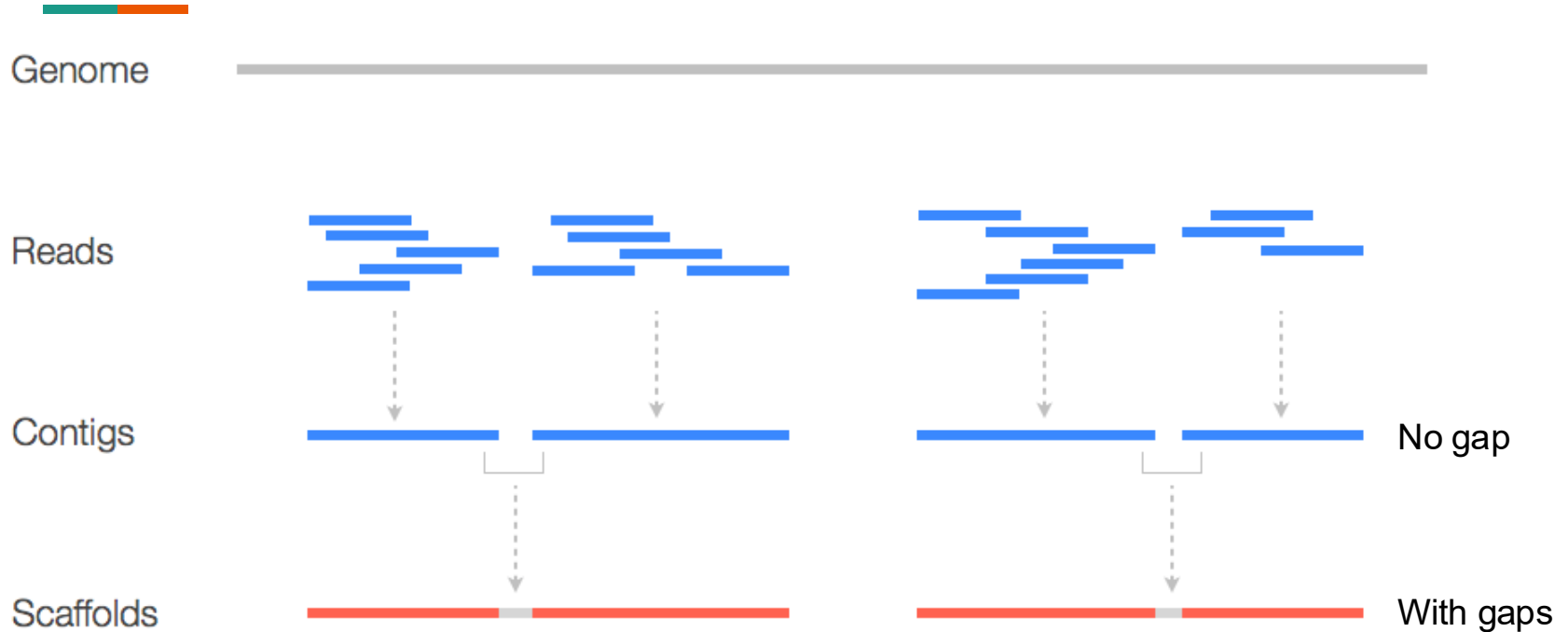
- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

Today's agenda



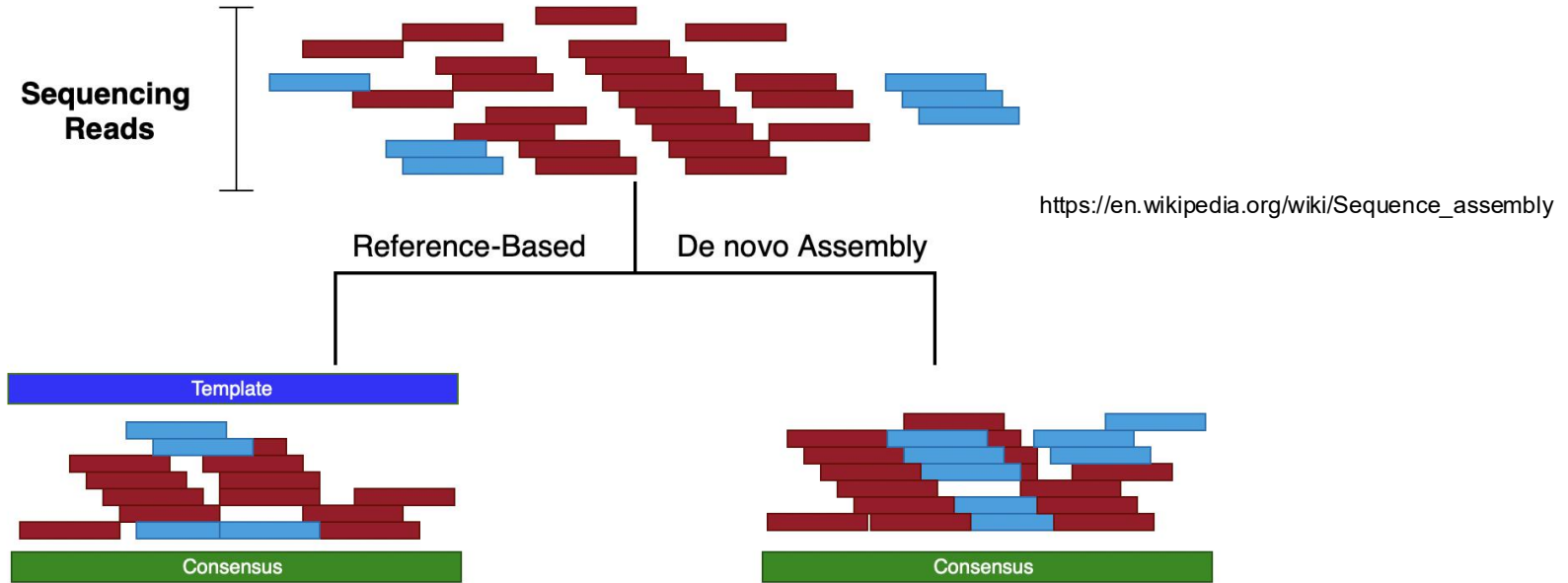
- Sequence assembly
- Structure of genome / genomic elements
- Annotation of genomic features

Sequence assembly



- Reconstruct the original genome from sequencing reads

Alignment versus assembly




- Alignment trusts the reference genome structure, assembly does not



Assembly of short reads

Assembling short reads is hard



GACCTACA
ACCTACAA
CCTACAAG
CTACAAGT
TACAAGTT
ACAAGTTA
CAAGTTAG
TACAAGTC
ACAAGTCC
CAAGTCCG



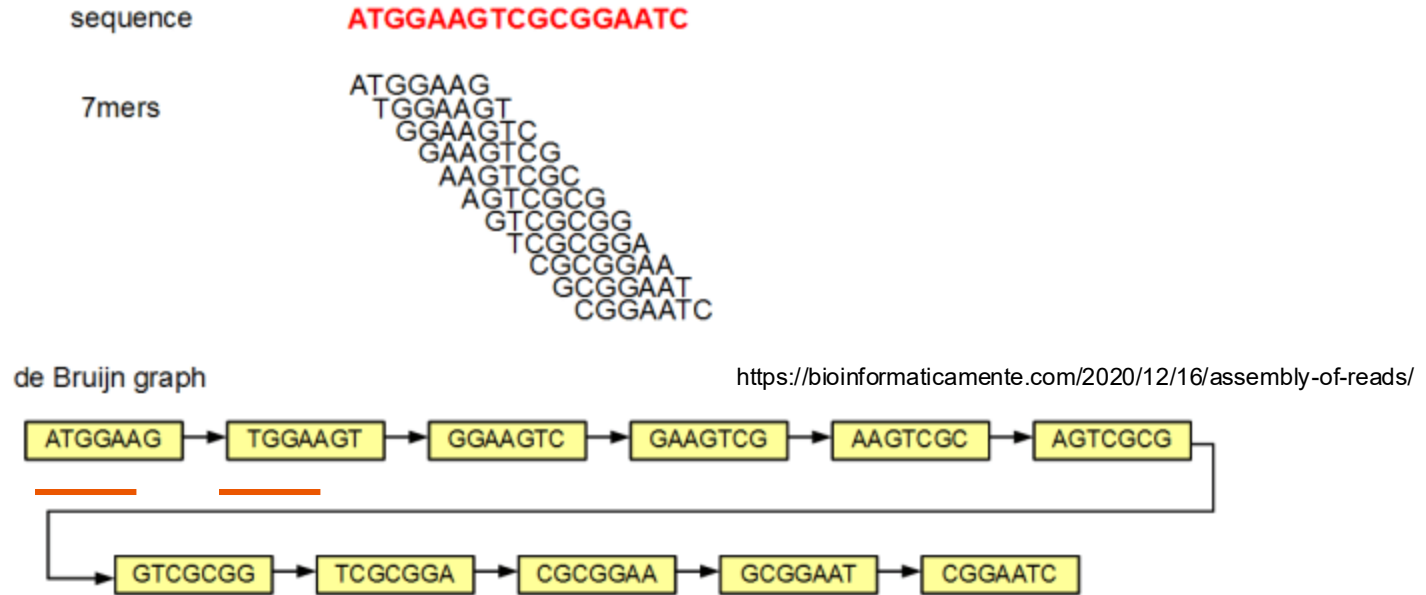
GACCTACAAGTTAG

Or

GACCTACAAGTCCG

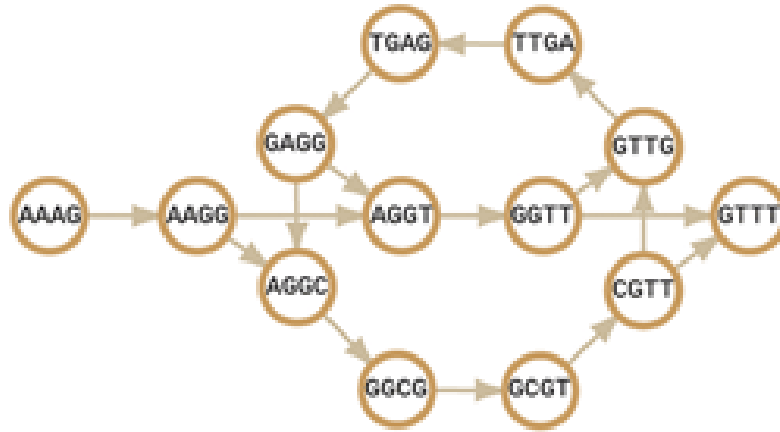
- Repetitive region? Polyploid genome?

De Bruijn graph



- Connect read whose **suffix** is identical to the next read's **prefix**
- An assembly is a walk along a **path** in this graph

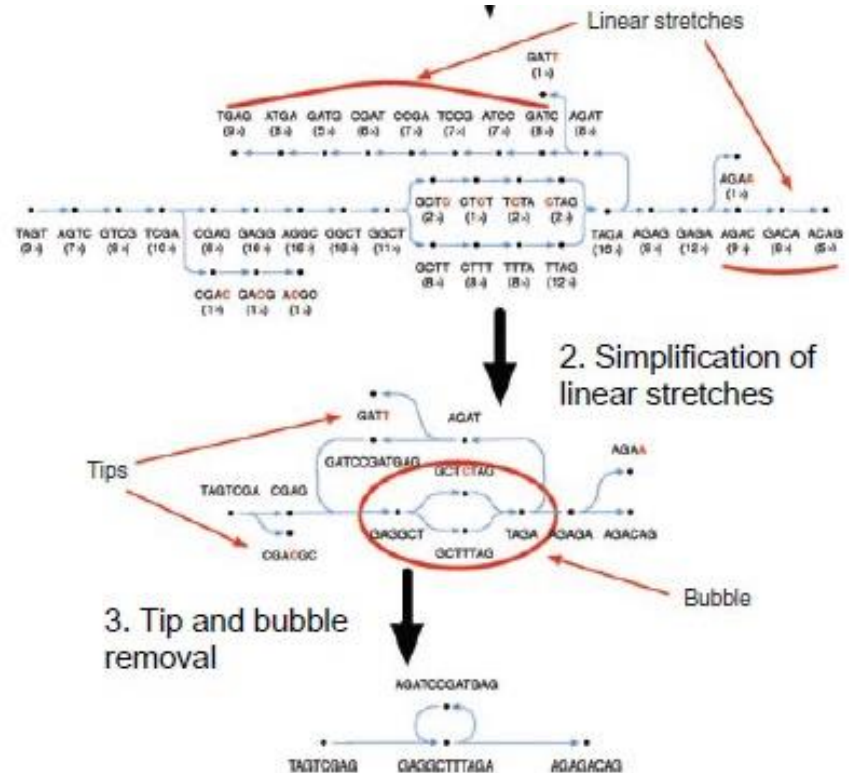
De Bruijn graph with branches and cycles



- Which path is best?
 - The longest one that visits as many nodes as possible
 - Cannot visit the same node again (cannot reuse read)
 - Related to the concept of **Hamiltonian Path**

Pruning less probable paths

- **Tips trimming:** remove the less probable branch (shorter sequence, lower read coverage)
- **Bubble removal:** use paired-end information to constrain the length of DNA between the forward read and the reverse read





Assembly of long reads

What changes with long read data?



- **Pros**

- Less number of reads, smaller de Bruijn graph
- Less ambiguity of overlapping sequences

- **Cons**

- More mismatches
- Multiple possibilities of overlapping sequences

AACACATACTCGACTACGACTACGACTAGCACT

Which one should be connected?

ACTACGACTAGCACTAGAC**CATCACGCATCA**

ACGACTAGCACTAGAC**TATAGCTACGACTACGACTACTA**


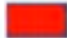

Overlap-Layout-Consensus algorithm



- **Overlap:** Connect reads with overlapping sequences (same as before)
- **Layout:** Summarize solvable paths into contigs
- **Consensus:** Call consensus sequence
- Why does OLC fit long read data?
 - Allow flexible overlapping sequences
 - Assume that the reads are not accurate

Finding overlap with dynamic programming

		A	T	G	C	T
	0	0	0	0	0	0
A	0	1	0	0	0	0
G	0	0	0	1	0	0
C	0	0	0	0	2	0
T	0	0	0	0	0	3

Match : 1 
Mismatch : -1 
GAP : -2 

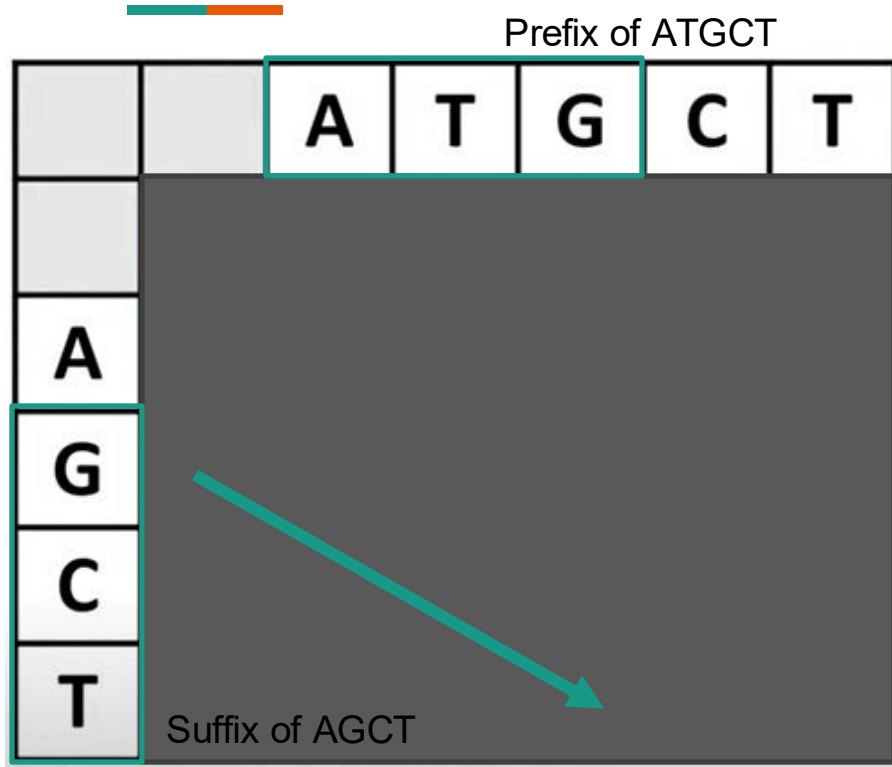
Seq1 : ~~A~~TGCT

|||

Seq2 : ~~A~~GCT


- For local alignment, we modify the algorithm to **reset negative scores to zeros**
- **How can we modify here?**

Where the path should start and end?



- Overlap = good match between the **suffix of the first sequence (AGCT)** and the **prefix of the second sequence (ATGCT)**
- The alignment path must **start from the first column** and **end on the last row**

How should the score be initialized?


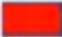



		A	T	G	C	T
	0	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
A	0					
G	0					
C	0					
T	0					

- The first column is all zeros because the alignment can start from anywhere in the first column
- The first row is minus infinity (minus large negative) because the alignment cannot start with gap on the second sequence

Putting everything together

		A	T	G	C	T
	0	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
A	0	1	-1	-3	-5	-7
G	0	-1	0	0	-2	-4
C	0	-1	-2	-1	1	-1
T	0	-1	0	-2	-2	2

Match : 1 
Mismatch : -1 
GAP : -2 


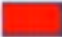

Seq1 : ATGCT

| | | |

Seq2 : A-GCT

Let's try one more example

		A	T	G	C	T
	0	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
A	0	1	-1	-3	-5	-7
T	0	-1	2	0	-2	-4
T	0	-1	0	1	-1	-1
G	0	-1	-2	1	0	-2

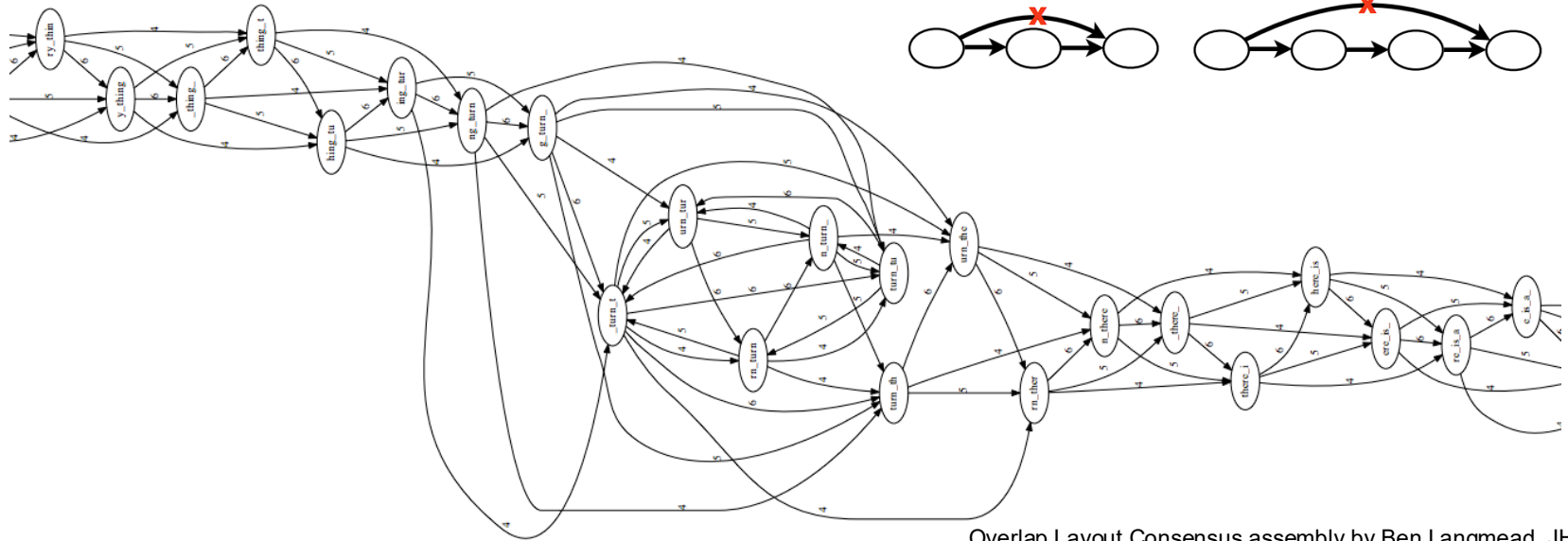
Match : 1 
Mismatch : -1 
GAP : -2 

Seq1: ATGCT

|||

Seq2: ATTG

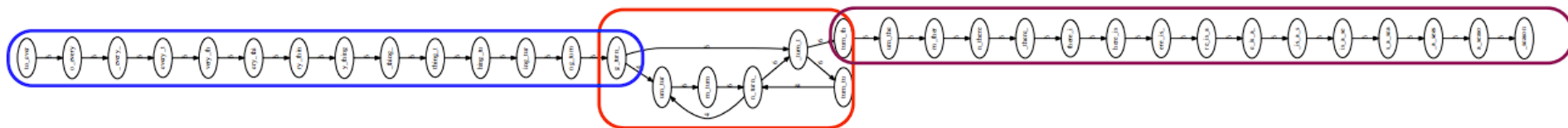
Laying out the contigs



Overlap Layout Consensus assembly by Ben Langmead, JHU

- Long reads will generate a combination of short and long overlaps
- Remove redundant edges to simplify the graph

Calling consensus sequence



TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAACTA
TAG TTACACAGATTATTGACTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

↓ ↓ ↓ ↓ ↓
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

Overlap Layout Consensus assembly by Ben Langmead, JHU

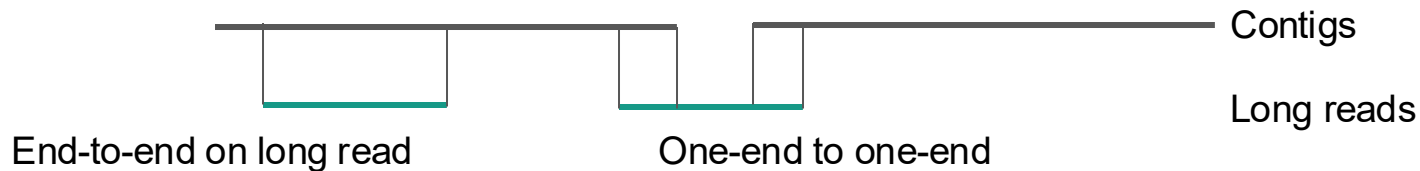


Hybrid assembly

Combining what we already know



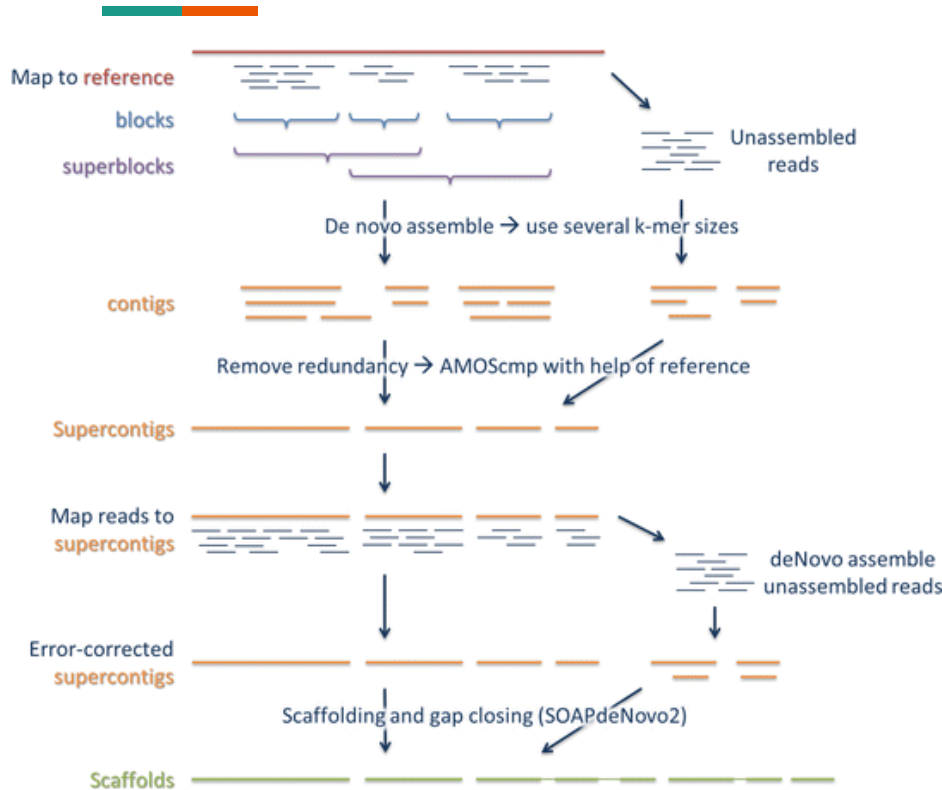
- Generate assembly from short reads
- Align long reads to the assembly
 - **Semi-global alignment**
 - Long read and assembly must reflect a consistent genome structure
 - Long read may span multiple short read contigs





Referenced-guided assembly

Guiding assembly by positions on reference genome



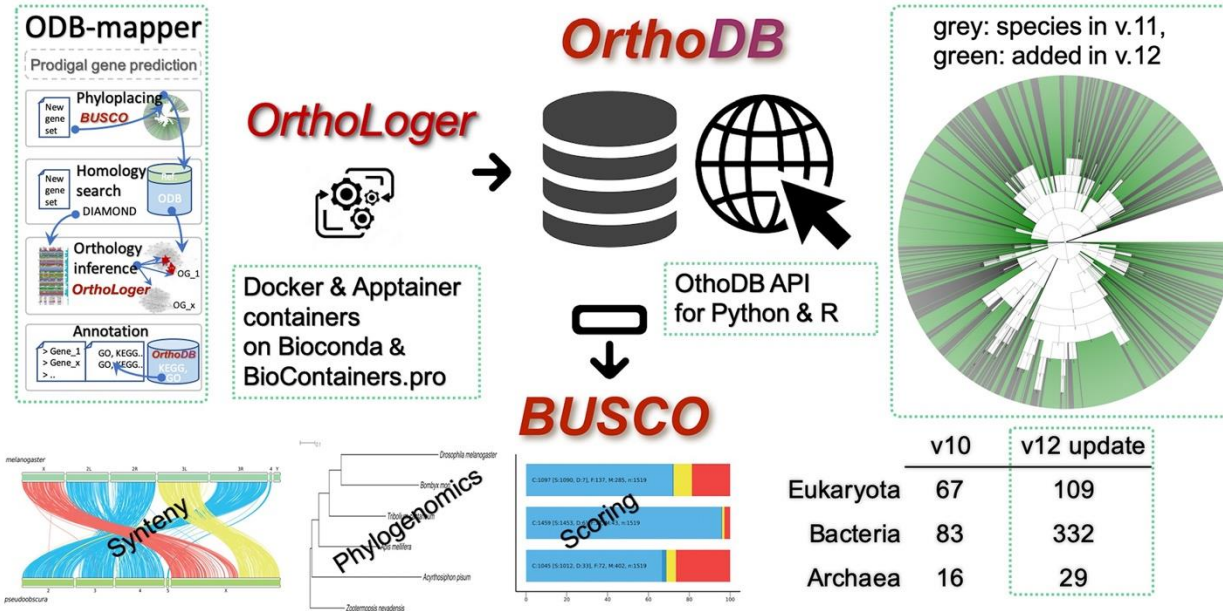
- Aligned reads can be grouped for local assembly
- Unaligned reads can indicate structural variations
- Useful when studying microorganisms, which evolve rapidly

How to check the quality of an assembly?



- Contig count
- **N50**: the length of the shortest contig that form 50% of total contig length
 - Contig lengths: 2, 3, 5, 7, 9
 - 50% of total lengths = $0.5 \times (2+3+5+7+9) = 13$
 - Contigs that form at least 50% of total length: 7 and 9
 - N50 is 7
- **U50**: N50 for a target genome, in the case where the assembly contains both host/pathogen, plasmid, or background DNA

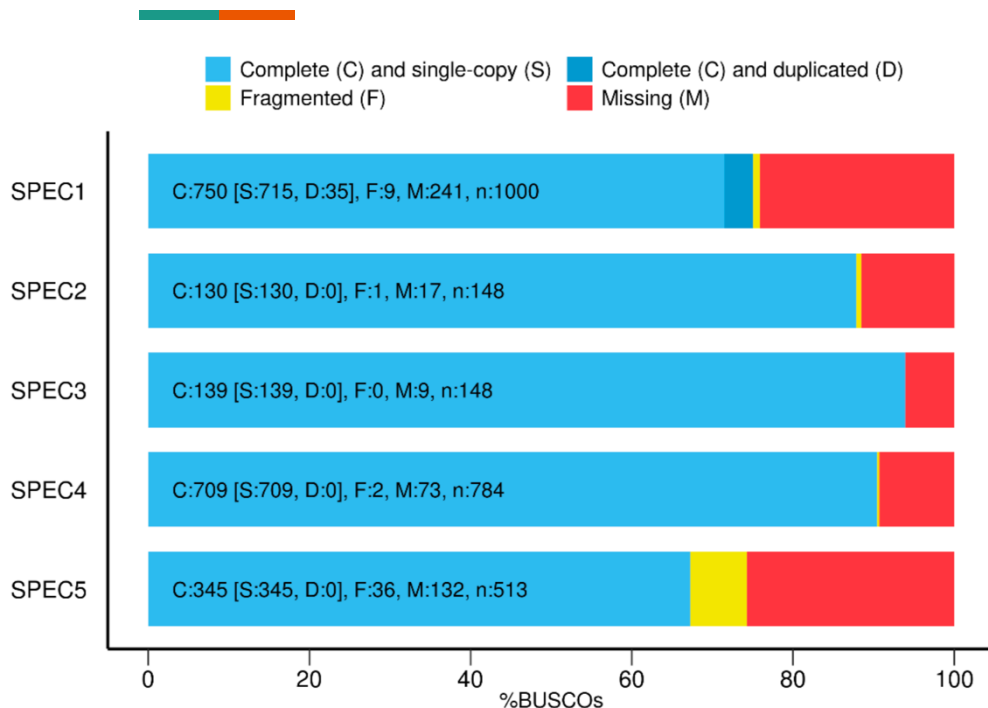
BUSCO: checking the presence of core genes



Simao, F.A. et al. Bioinformatics 31:3210-3212 (2015)

- A good assembly should contain all core genes of the related taxonomy

BUSCO assessment result

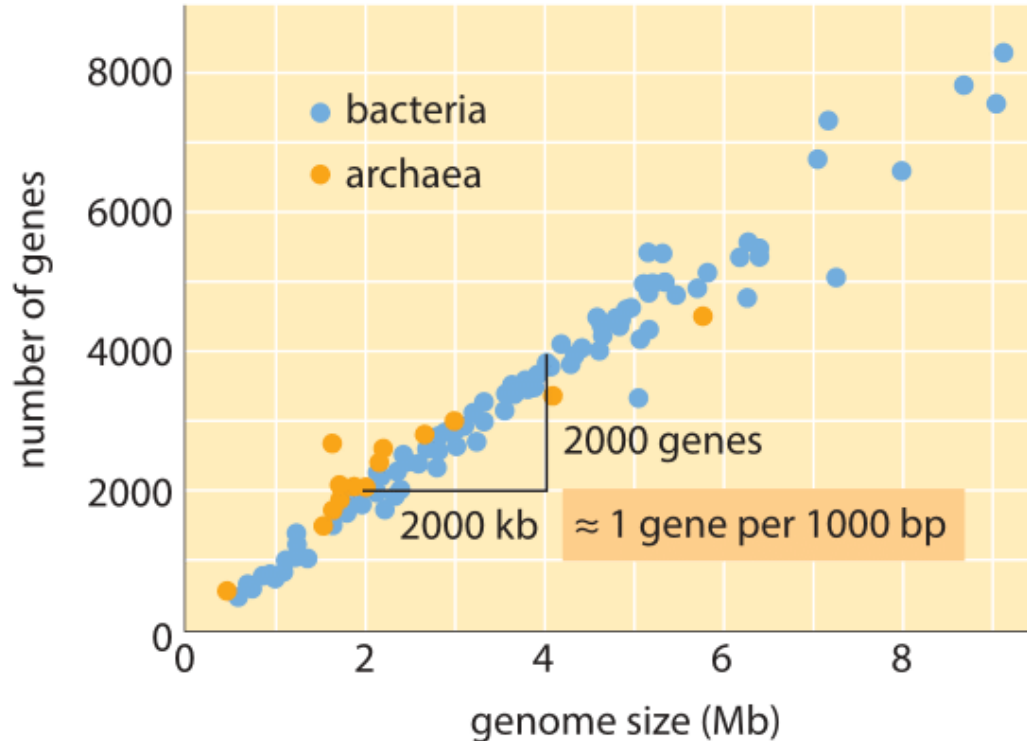


- **C:** Found the complete length of core genes
 - **S/D:** Whether the genes were found as single-copy or duplicated
 - Duplicated can indicate error
- **F:** Found incomplete genes (partial assembly)
- **M:** Did not find the genes



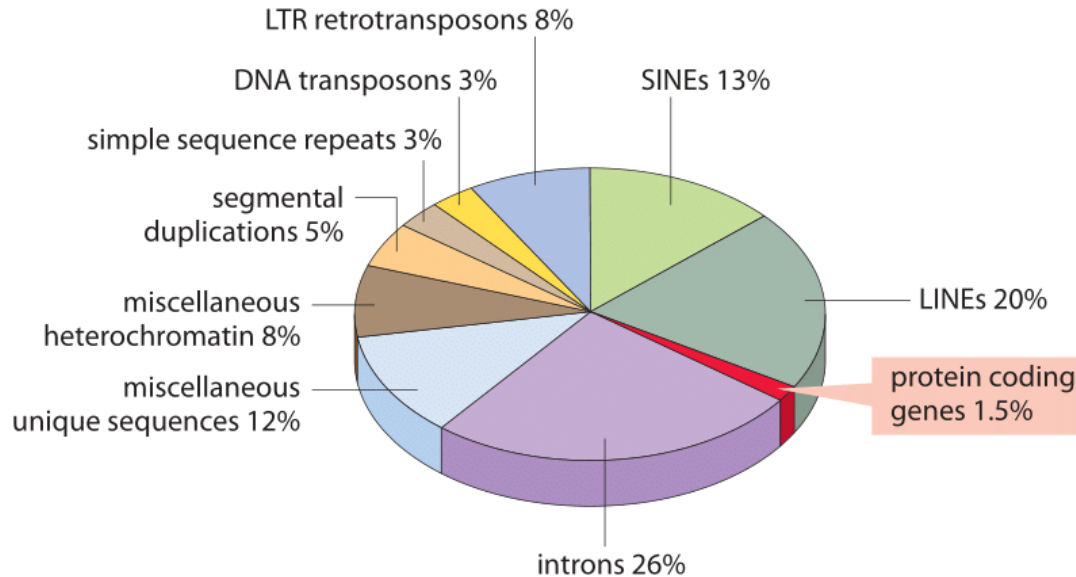
Genomic elements

In microbes, most genomic DNA are genes



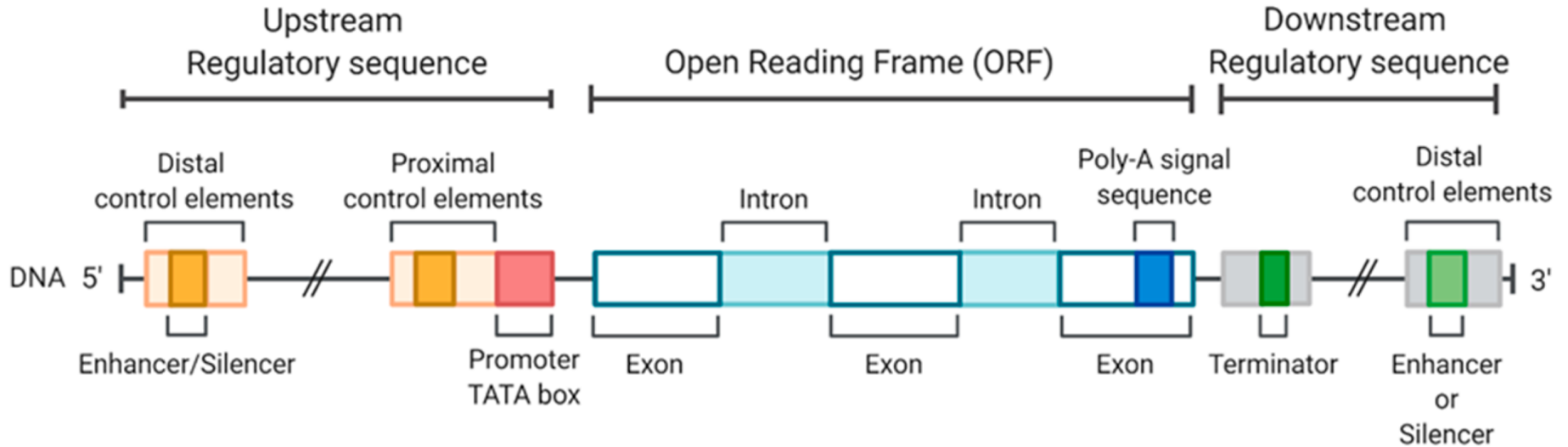
Eukaryotes have complex genomes

main components of the human genome



	Organism	# of protein-coding genes	# of genes naïve estimate: (genome size /1000)
viruses	HIV 1	9	10
	Influenza A virus	10-11	14
	Bacteriophage λ	66	49
	Epstein Barr virus	80	170
prokaryotes	<i>Buchnera sp.</i>	610	640
	<i>T. maritima</i>	1,900	1,900
	<i>S. aureus</i>	2,700	2,900
	<i>V. cholerae</i>	3,900	4,000
	<i>B. subtilis</i>	4,400	4,200
	<i>E. coli</i>	4,300	4,600
	<i>S. cerevisiae</i>	6,600	12,000
eukaryotes	<i>C. elegans</i>	20,000	100,000
	<i>A. thaliana</i>	27,000	140,000
	<i>D. melanogaster</i>	14,000	140,000
	<i>F. rubripes</i>	19,000	400,000
	<i>Z. mays</i>	33,000	2,300,000
	<i>M. musculus</i>	20,000	2,800,000
	<i>H. sapiens</i>	21,000	3,200,000
	<i>T. aestivum</i> (hexaploid)	95,000	16,800,000

An example of gene structure

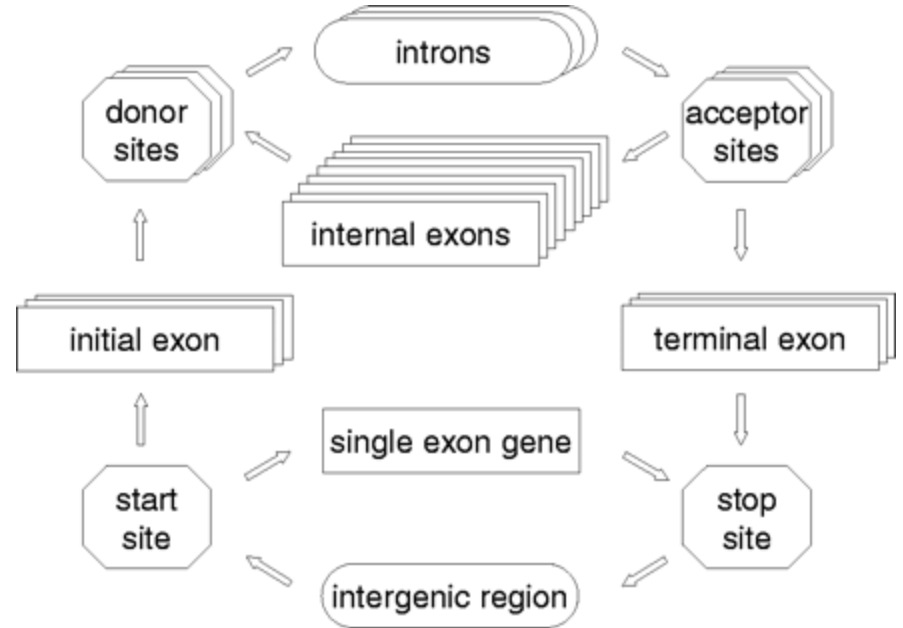
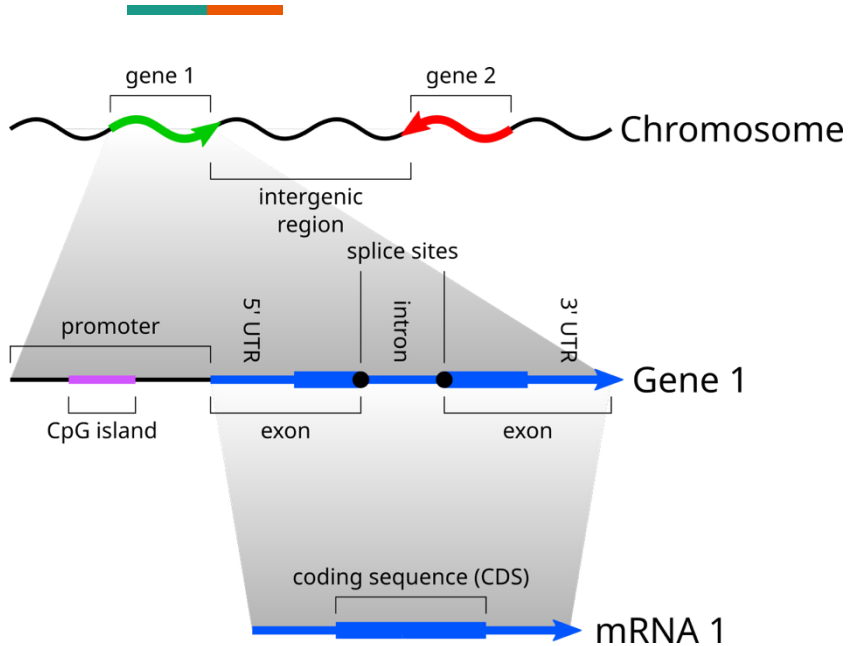


Genome annotation



- Protein coding genes
 - Exon/intron/splice sites
- Regulatory elements: Promoter/enhancer/silencer
- Repeats
- Non-coding RNA
- **Best if guided by data from other assays, such as RNA-seq and ChIP-seq**

Protein-coding gene annotation

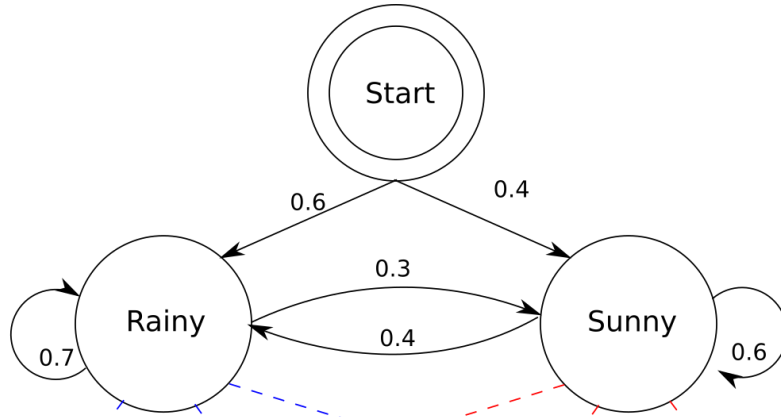


<https://cs.rice.edu/~ogilvie/comp571/hidden-markov-models/>

Lomsadze, A. et al. NAR 33:6494-6506 (2005)

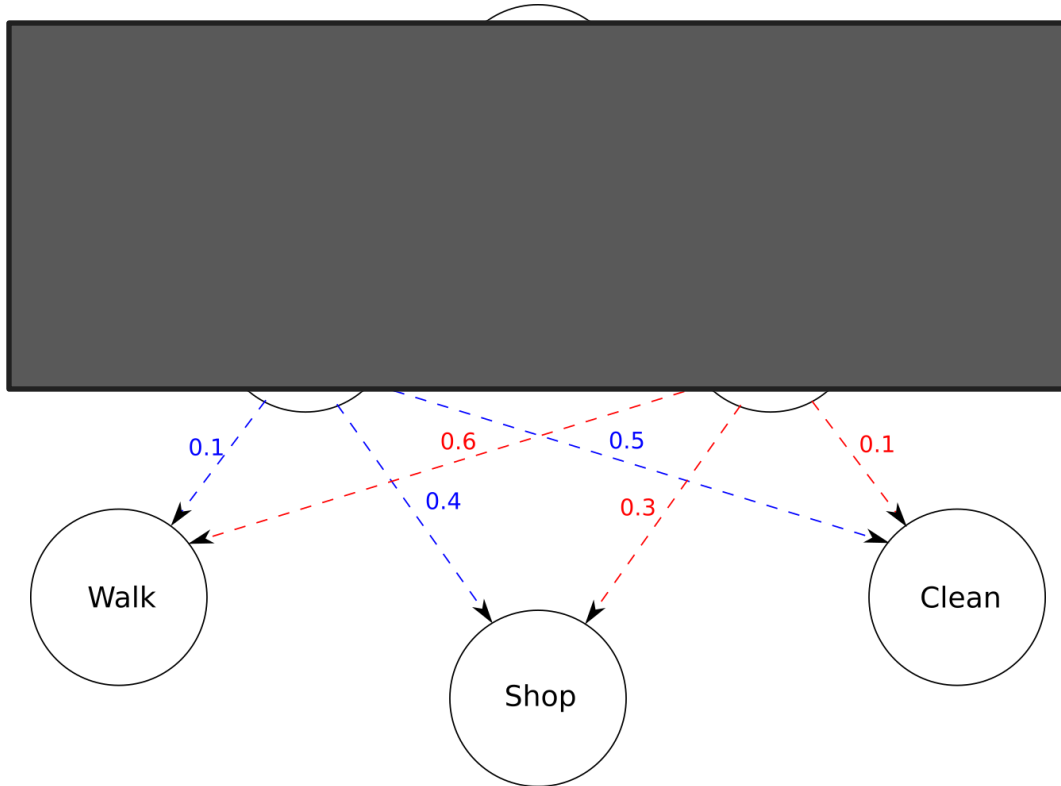
- Formulate the general structure of protein-coding genes

Markov model



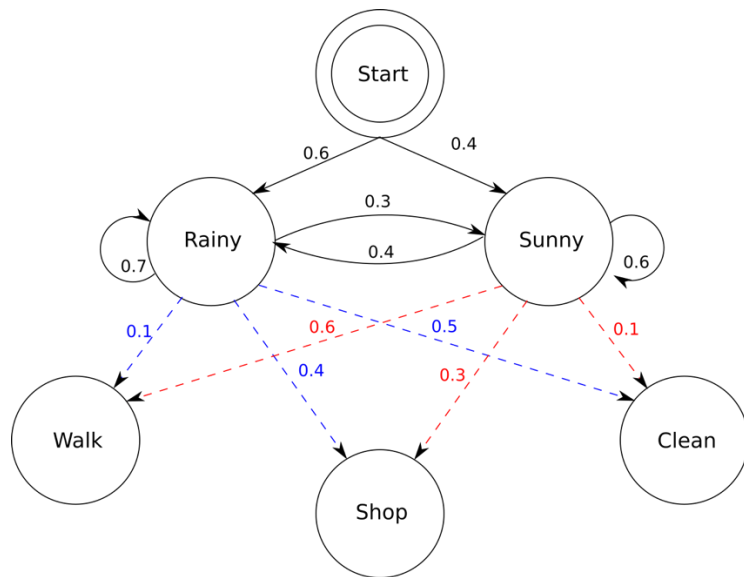
- State changes with specified probability
- $P(\text{Sunny} \mid \text{Rainy}) = 0.3$
- $P(\text{Rainy} \mid \text{Rainy}) = 0.7$
- **Data:** S,R,R,S,S,R

Hidden Markov model



- **Observation** depends on the **hidden** states
- **Data:** W,W,S,C,S,C
- $P(\text{Walk} \mid \text{Rainy}) = 0.1$
- $P(\text{Walk} \mid \text{Sunny}) = 0.6$
- Hidden state changes with specified probability

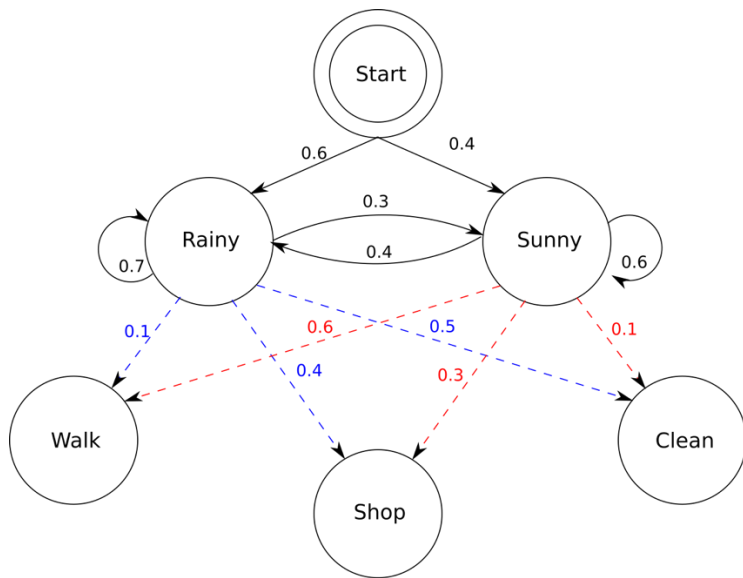
Decoding a Hidden Markov model



- Given the observation, what is the most likely sequence of hidden states (HS's)?
- Let's consider only the first two observations: Walk-Walk
- We have four possible HS's combinations: RR, RS, SR, and SS

Walk	Walk	Shop	Clean	Shop	Clean
HS1	HS2	HS3	HS4	HS5	HS6

Decoding a Hidden Markov model



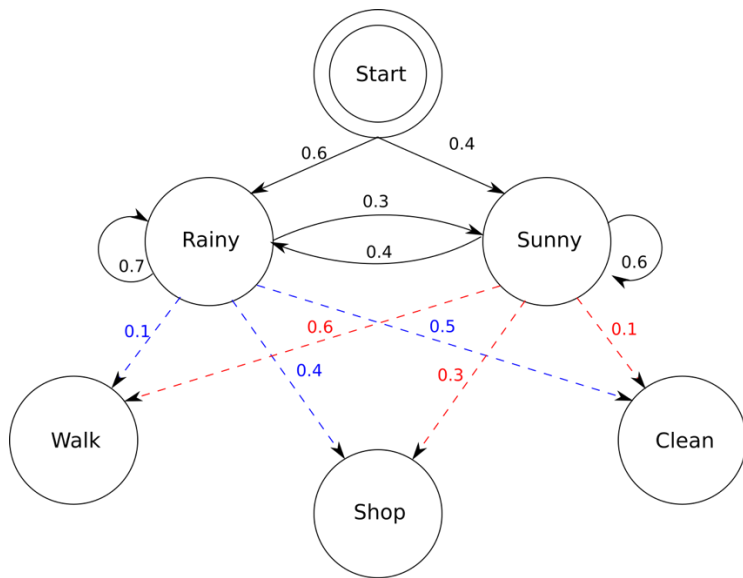
Conditional probability:

$$P(A | B) = P(A, B) / P(B)$$

- $P(RR | WW) = P(RR, WW) / P(WW)$
 - $P(RS | WW) = P(RS, WW) / P(WW)$
 - $P(SR | WW) = P(SR, WW) / P(WW)$
 - $P(SS | WW) = P(SS, WW) / P(WW)$
- It suffice to compare joint probabilities

Walk	Walk	Shop	Clean	Shop	Clean
HS1	HS2	HS3	HS4	HS5	HS6

Decoding a Hidden Markov model



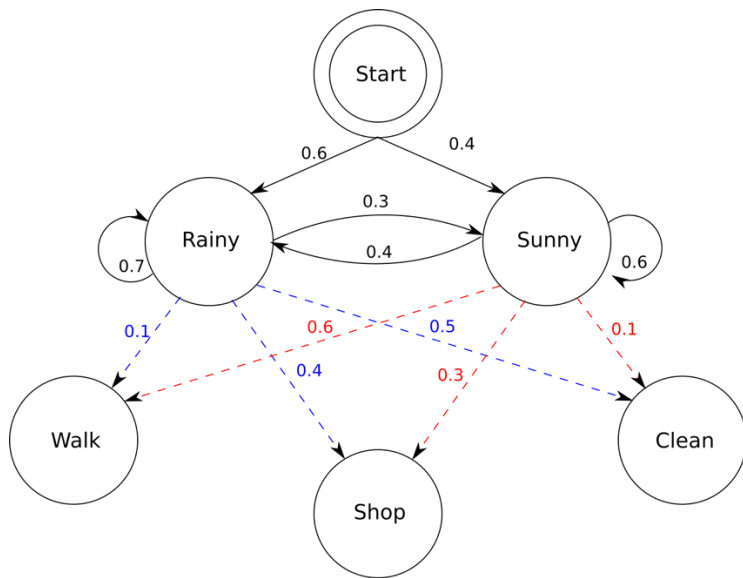
Breakdown the process sequentially:

1. HS1
2. Obs1 given HS1
3. HS2 given HS1
4. Obs2 given HS2

$$\begin{aligned} P(RS, WW) &= P(\text{Rainy}) \times P(W | \text{Rainy}) \times \\ &\quad P(\text{Sunny} | \text{Rainy}) \times \\ &\quad P(W | \text{Sunny}) \\ &= 0.6 \times 0.1 \times 0.3 \times 0.6 \end{aligned}$$

Walk	Walk	Shop	Clean	Shop	Clean
HS1	HS2	HS3	HS4	HS5	HS6

Decoding a Hidden Markov model



$$P(RR, WW) = 0.6 \times 0.1 \times 0.7 \times 0.1$$

$$P(RS, WW) = 0.6 \times 0.1 \times 0.3 \times 0.6$$

$$P(SR, WW) = 0.4 \times 0.6 \times 0.4 \times 0.1$$

$$P(SS, WW) = 0.4 \times 0.6 \times 0.6 \times 0.6$$

Intuition: $P(W \mid \text{Sunny}) \gg P(W \mid \text{Rainy})$

“Even though it’s less likely to observe a Sunny day, it’s even less likely to see someone take a walk on Rainy days”

Walk	Walk	Shop	Clean	Shop	Clean
HS1	HS2	HS3	HS4	HS5	HS6

Viterbi algorithm

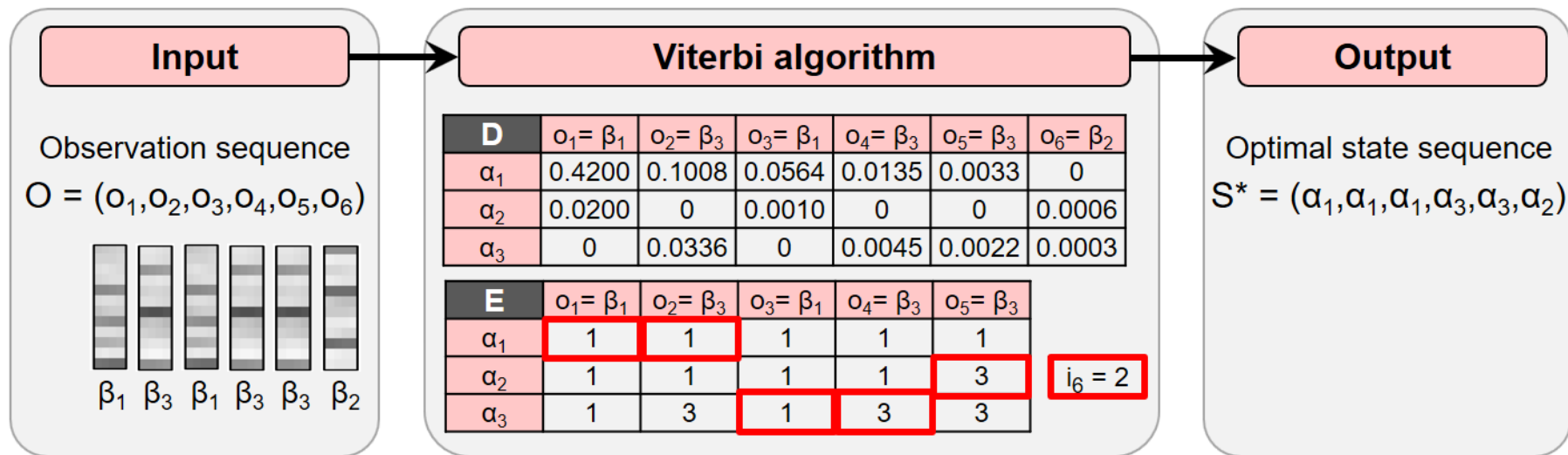
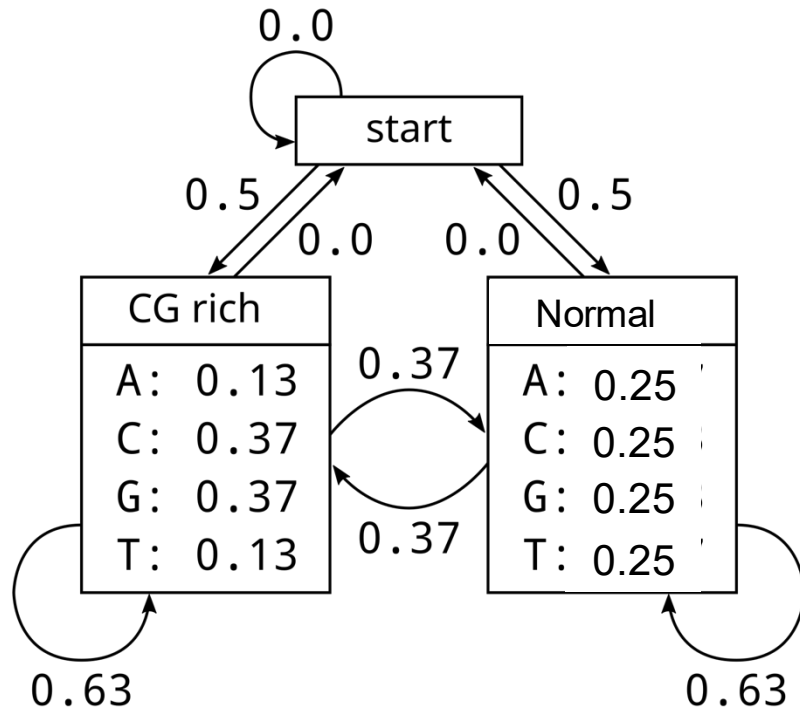


Figure 5.28b from [Müller, FMP, Springer 2015]

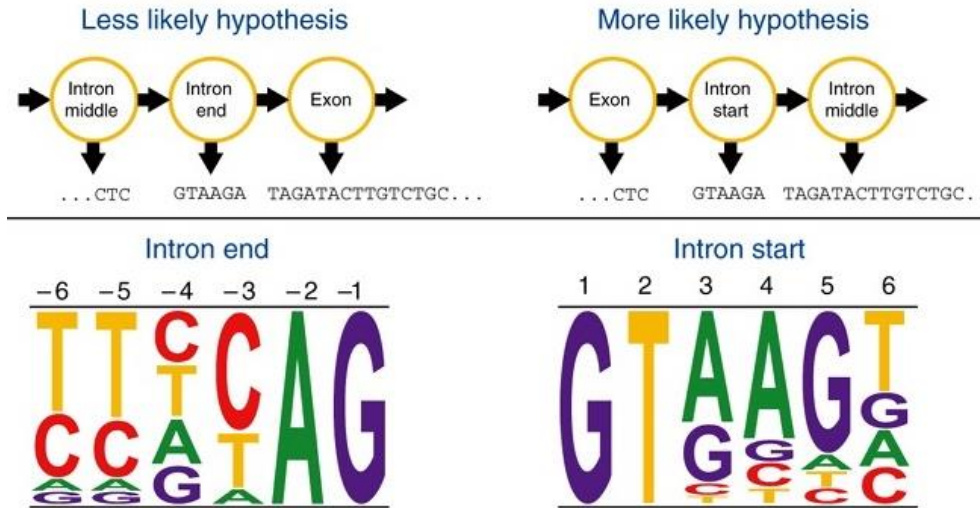
- Dynamic programming approach
- An optimal solution for n observations is based on solutions for $n-1$ obs.

HMM for CpG island annotation



- Different genomic elements emit different nucleotide profiles
- CpG islands would consist of more C/G than A/T
- Nucleotide frequencies can be adjusted based on the organism

HMM for protein-coding gene prediction



- Consider longer nucleotide segments, not just single bases
- What are the nucleotide frequencies/patterns of each state?
- Must be trained using known genes from related organisms

Any question?



- See you next time