

RNAseq

Michael Schatz

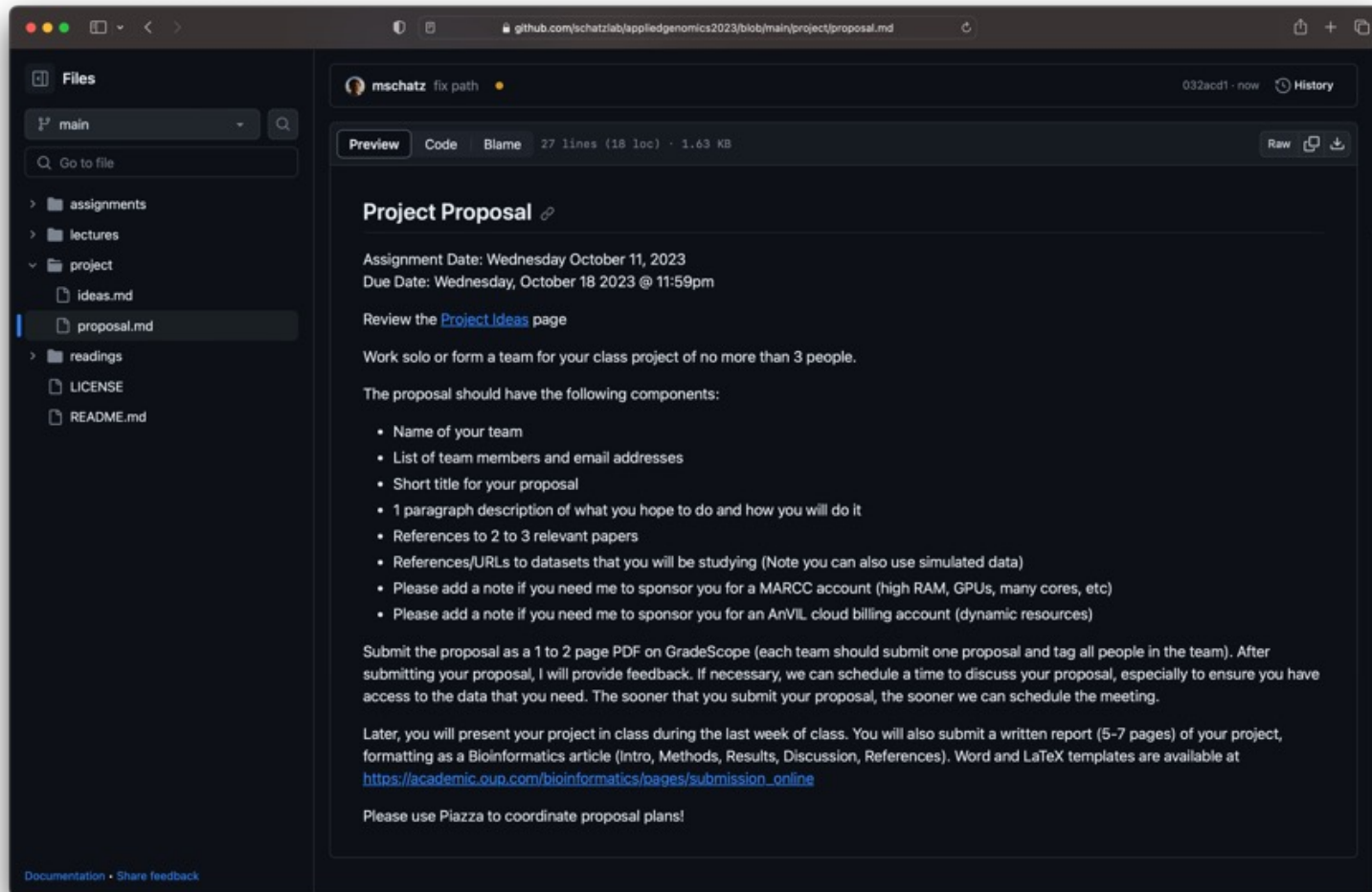
October 23, 2023

Lecture 16. Applied Comparative Genomics



Project Proposal

Due Wednesday Oct 18 by 11:59pm

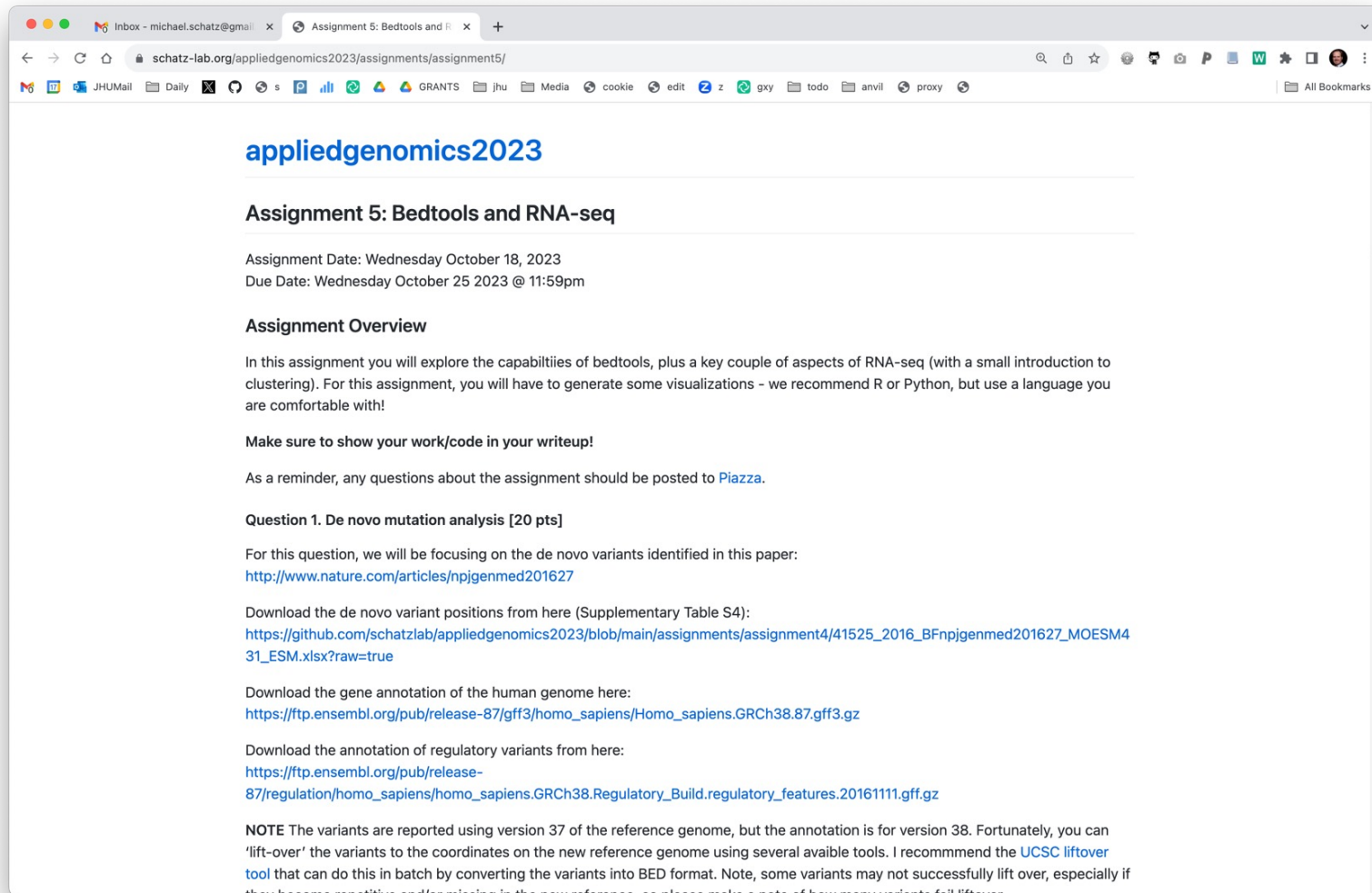


<https://github.com/schatzlab/appliedgenomics2023/blob/main/project/proposal.md>

Check Piazza for questions!

Assignment 5

Due: Wednesday Oct 25, 2023 by 11:59pm

A screenshot of a web browser displaying the 'Assignment 5: Bedtools and RNA-seq' page on the 'schatz-lab.org' website. The browser's address bar shows the URL 'schatz-lab.org/appliedgenomics2023/assignments/assignment5/'. The page content includes the title 'Assignment 5: Bedtools and RNA-seq', the assignment date (Wednesday October 18, 2023), and the due date (Wednesday October 25, 2023 @ 11:59pm). It also features an 'Assignment Overview' section with instructions on exploring bedtools and RNA-seq, a reminder to show work/code, and a link to Piazza for questions. The 'Question 1. De novo mutation analysis [20 pts]' section provides details on the de novo variants, including links to the source paper, the variant positions file, and the gene and regulatory variant annotations. A note at the bottom explains the reference genome version discrepancy and recommends the UCSC liftover tool.

appliedgenomics2023

Assignment 5: Bedtools and RNA-seq

Assignment Date: Wednesday October 18, 2023
Due Date: Wednesday October 25, 2023 @ 11:59pm

Assignment Overview

In this assignment you will explore the capabilities of bedtools, plus a key couple of aspects of RNA-seq (with a small introduction to clustering). For this assignment, you will have to generate some visualizations - we recommend R or Python, but use a language you are comfortable with!

Make sure to show your work/code in your writeup!

As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1. De novo mutation analysis [20 pts]

For this question, we will be focusing on the de novo variants identified in this paper:
<http://www.nature.com/articles/npjgenmed201627>

Download the de novo variant positions from here (Supplementary Table S4):
https://github.com/schatzlab/appliedgenomics2023/blob/main/assignments/assignment4/41525_2016_BFnpjgenmed201627_MOESM431_ESM.xlsx?raw=true

Download the gene annotation of the human genome here:
https://ftp.ensembl.org/pub/release-87/gff3/homo_sapiens/Homo_sapiens.GRCh38.87.gff3.gz

Download the annotation of regulatory variants from here:
https://ftp.ensembl.org/pub/release-87/regulation/homo_sapiens/homo_sapiens.GRCh38.Regulatory_Build.regulatory_features.20161111.gff.gz

NOTE The variants are reported using version 37 of the reference genome, but the annotation is for version 38. Fortunately, you can 'lift-over' the variants to the coordinates on the new reference genome using several available tools. I recommend the [UCSC liftover tool](#) that can do this in batch by converting the variants into BED format. Note, some variants may not successfully lift over, especially if they become repetitive and/or missing in the new reference, so please make a note of how many variants fail liftover.

<https://schatz-lab.org/appliedgenomics2023/assignments/assignment5/>

Check Piazza for questions!

Clustering Refresher

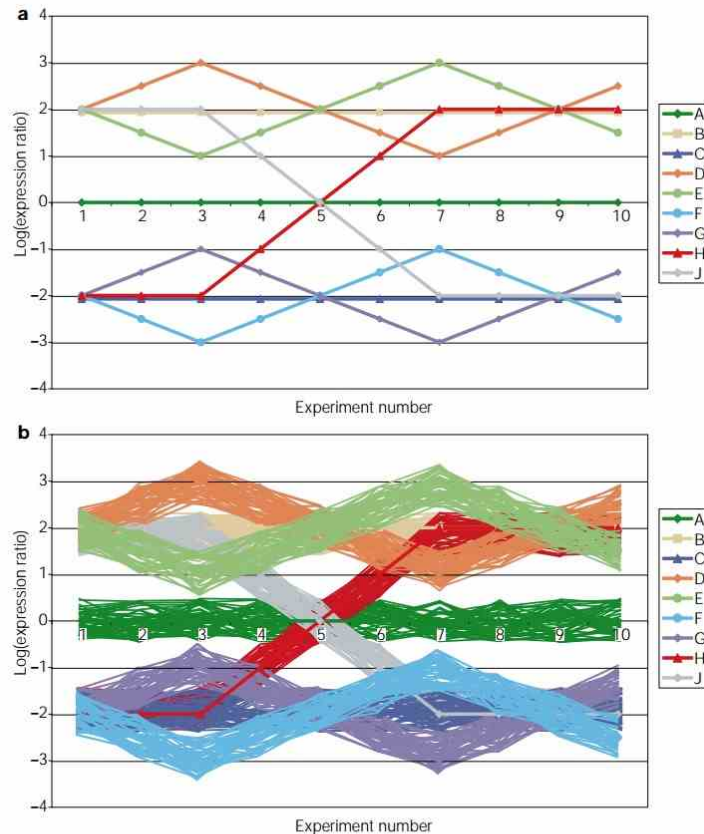
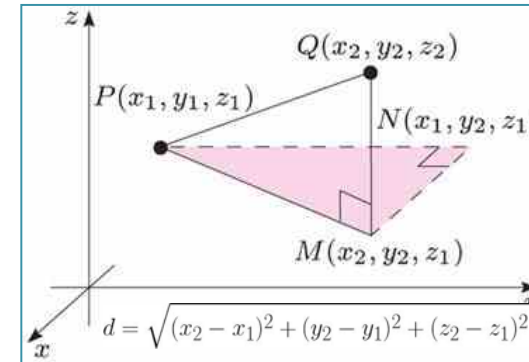
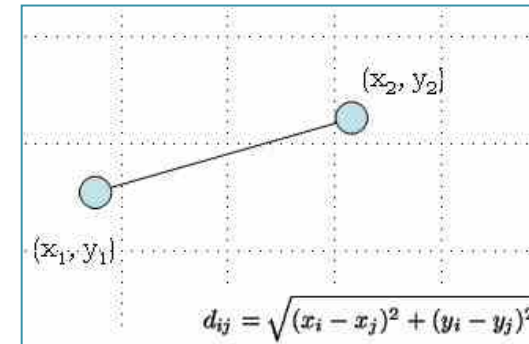


Figure 2 | **A synthetic gene-expression data set.** This data set provides an opportunity to evaluate how various clustering algorithms reveal different features of the data. **a** | Nine distinct gene-expression patterns were created with $\log_2(\text{ratio})$ expression measures defined for ten experiments. **b** | For each expression pattern, 50 additional genes were generated, representing variations on the basic patterns.

Euclidean Distance

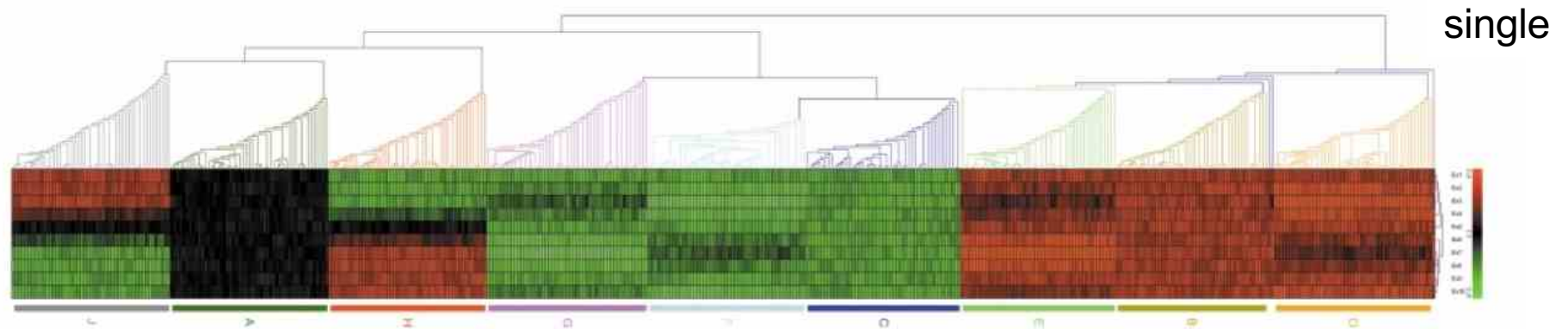
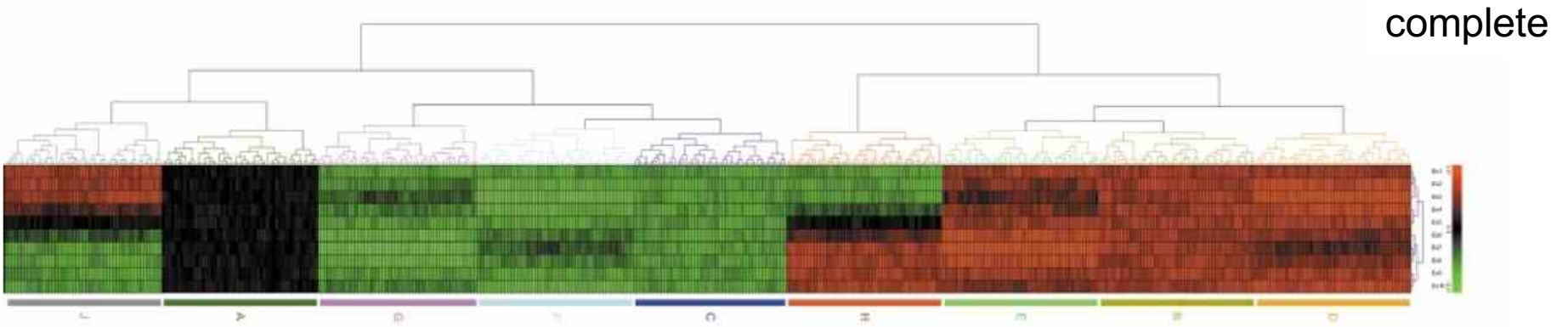
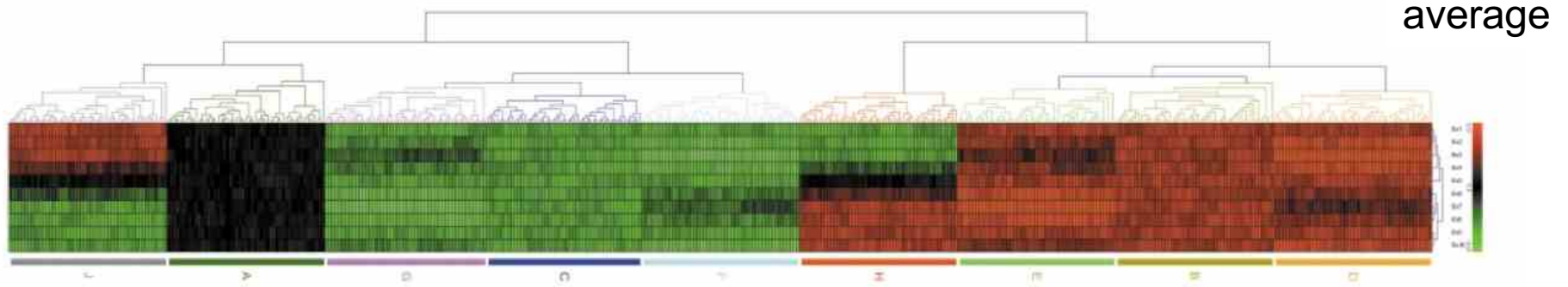


$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

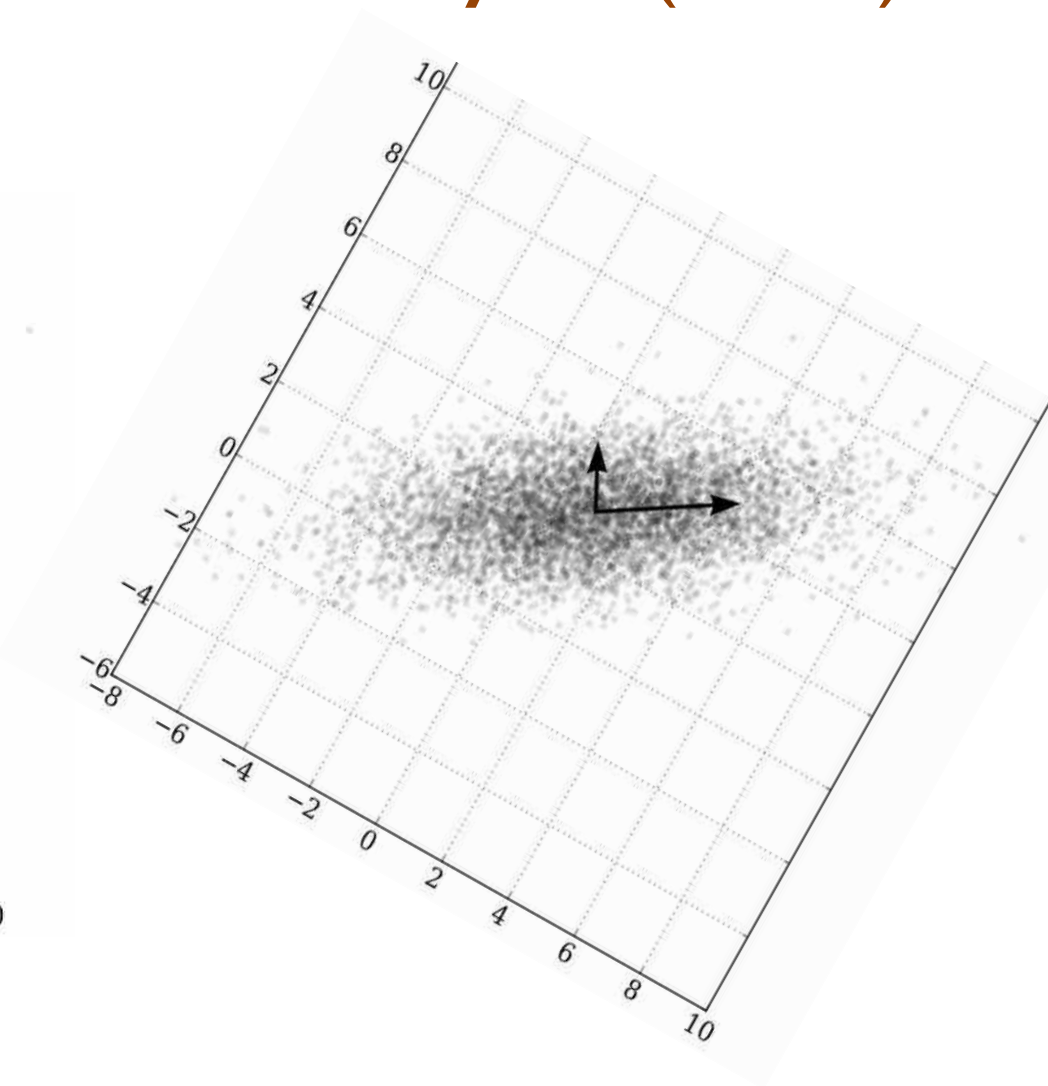
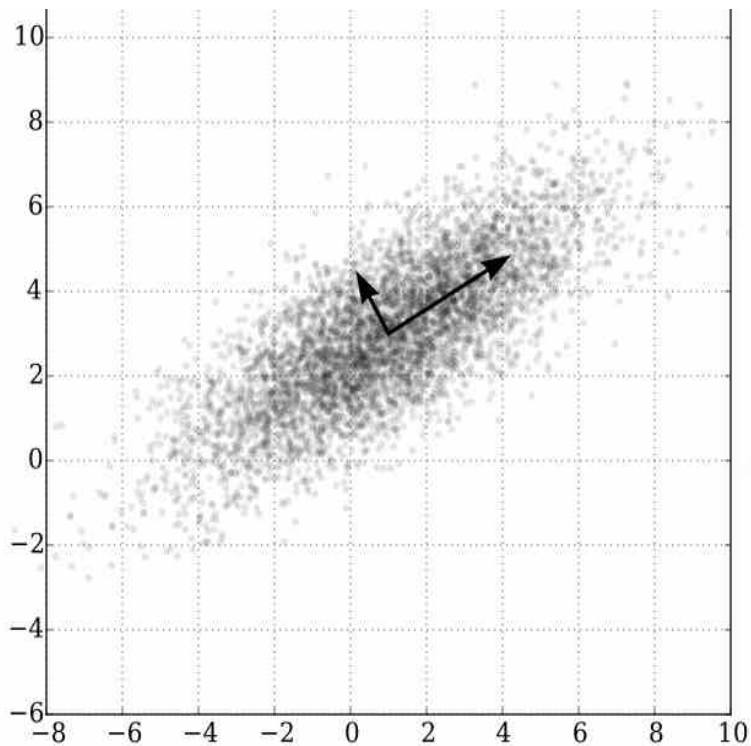
Computational genetics: Computational analysis of microarray data

Quackenbush (2001) *Nature Reviews Genetics*. doi:10.1038/35076576

Hierarchical Clustering



Principle Components Analysis (PCA)



PC1: "New X"- The dimension with the most variability
PC2: "New Y"- The dimension with the second most variability

Principle Components Analysis (PCA)

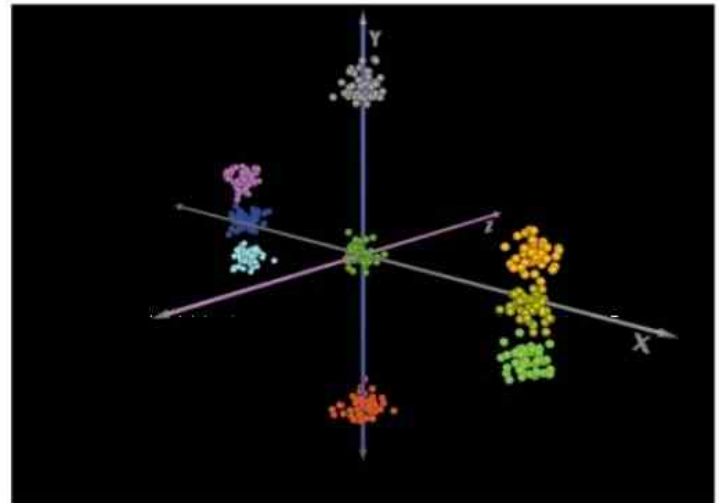
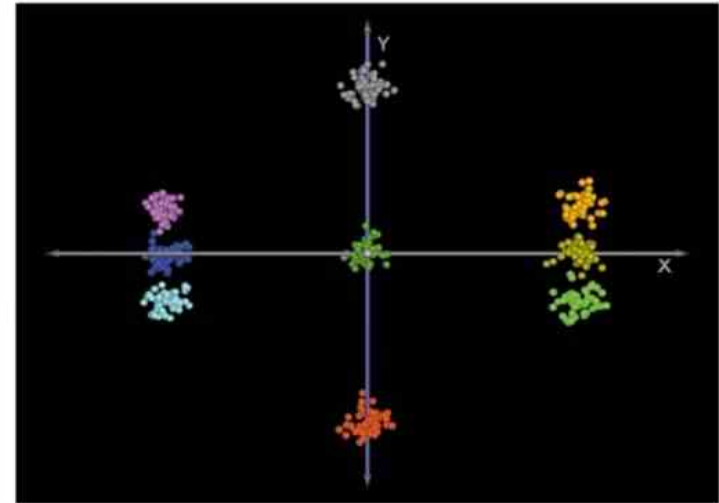
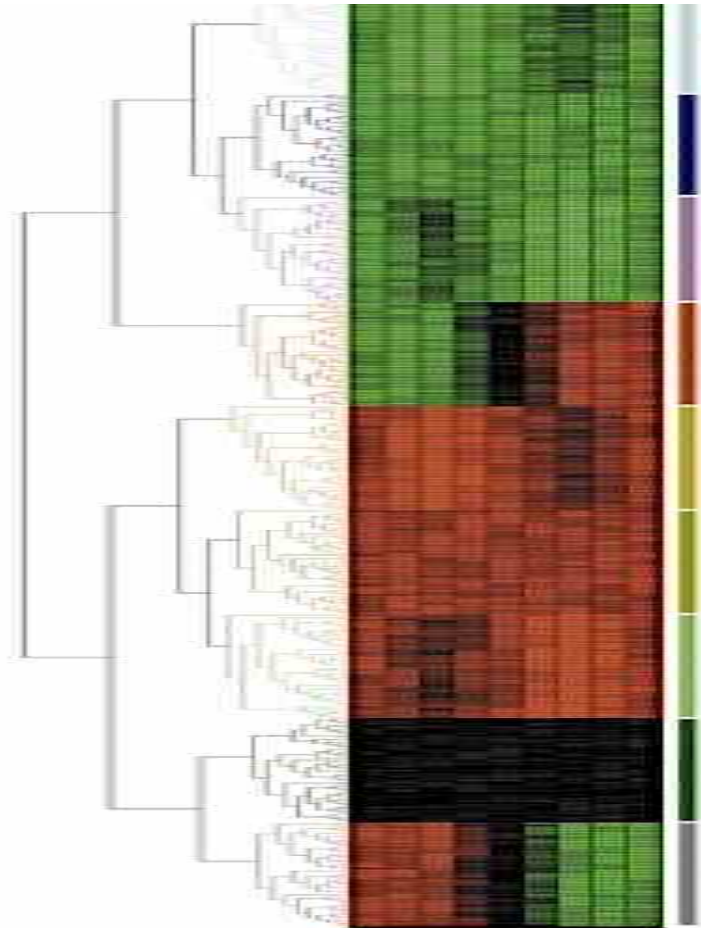
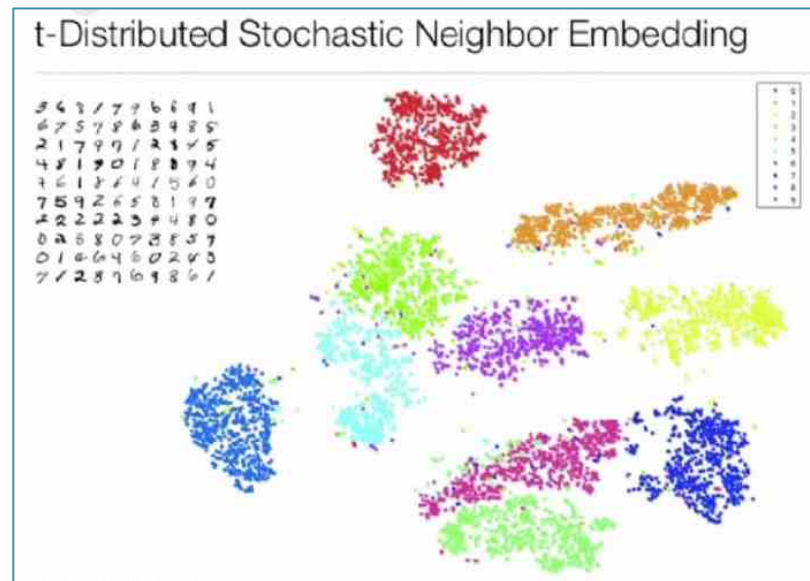
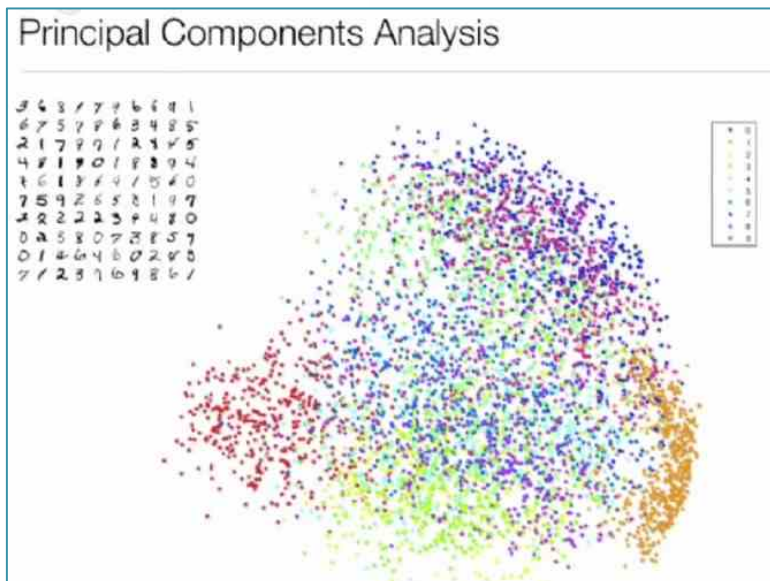
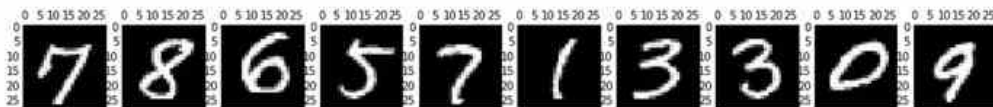


Figure 4 | **Principal component analysis.** The same demonstration data set was analysed using **a** | hierarchical (average-linkage) clustering and **b** | principal component analysis using Euclidean distance, to show how each treats the data, with genes colour coded on the basis of hierarchical clustering results for comparison.

PCA and t-SNE



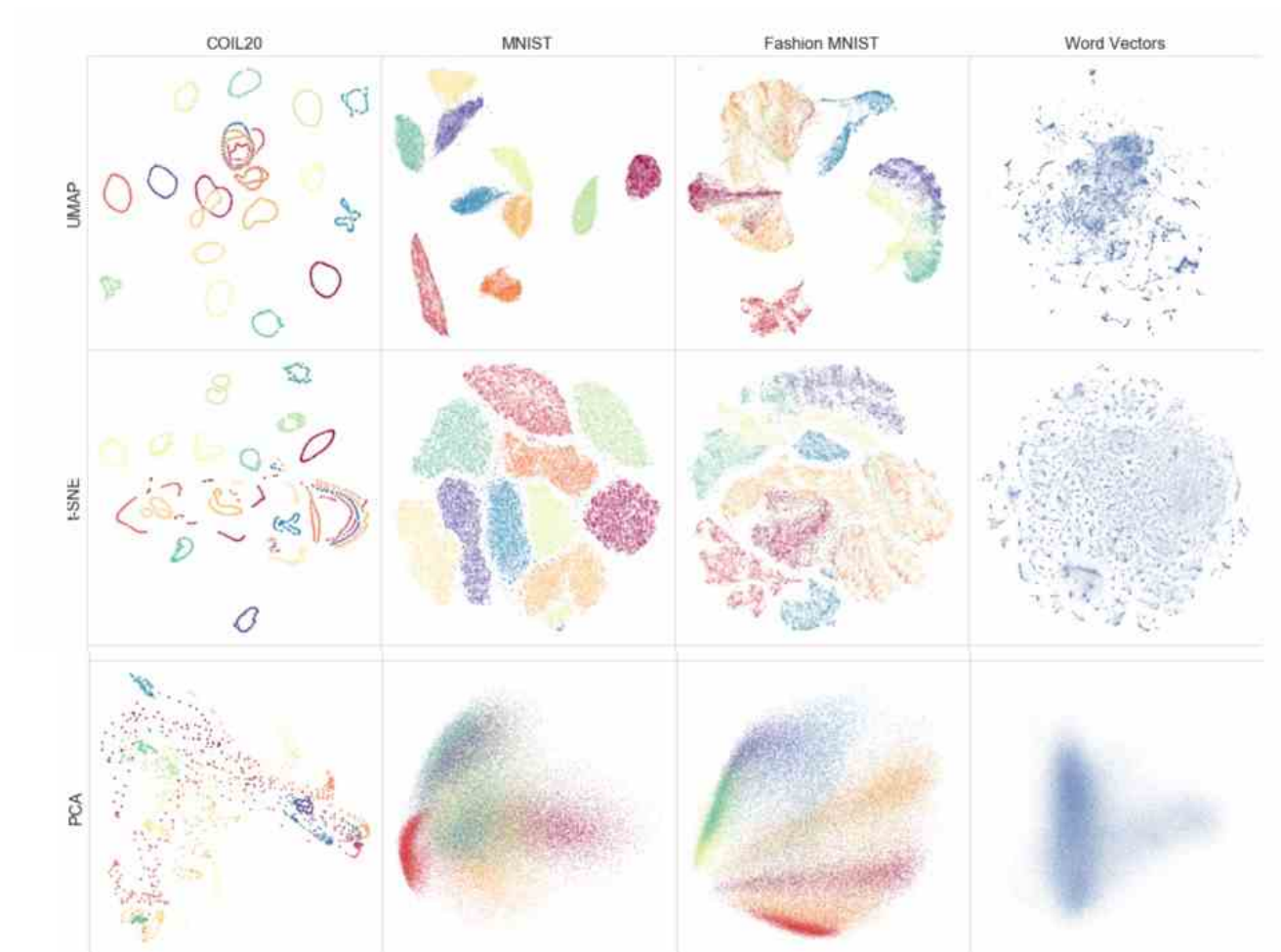
t-distributed Stochastic Neighborhood Embedding

- Non-linear dimensionality reduction technique: distances are only locally meaningful
- Rather than Euclidean distances, for each point fits a Gaussian kernel to fit the nearest N neighbors (perplexity) that define the probabilities that two points should be close together
- Using an iterative spring embedding system to place high probability points nearby

Visualizing Data Using t-SNE

<https://www.youtube.com/watch?v=RJVL80Gg3IA>

UMAP



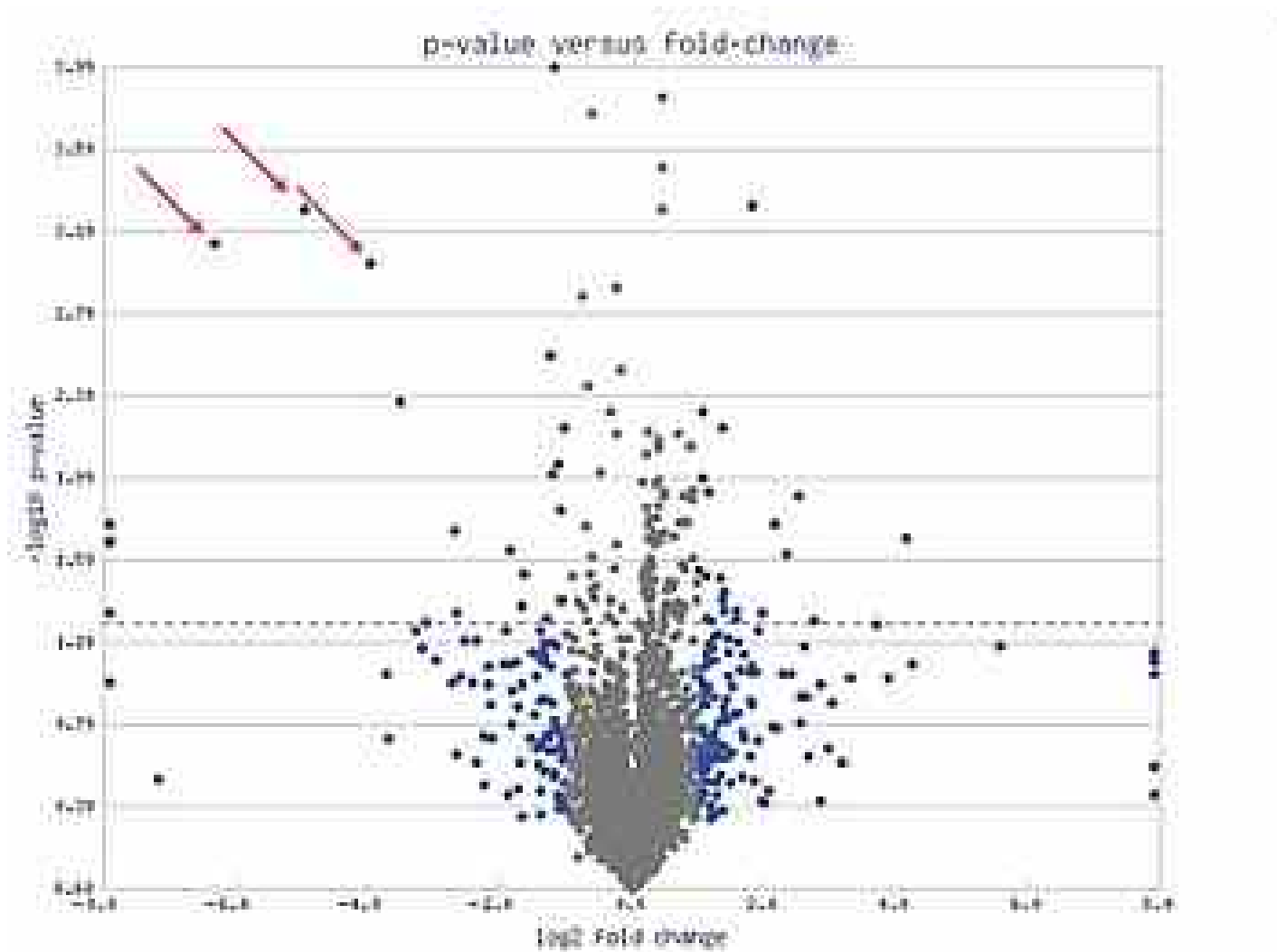
UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

McInnes et al (2018) arXiv. 1802.03426

<https://www.youtube.com/watch?v=nq6iPZVUxZU>

<https://towardsdatascience.com/how-exactly-umap-works-13e3040e1668>

Volcano Plot



https://en.wikipedia.org/wiki/Volcano_plot_%28statistics%29

JASPAR Database

The screenshot displays the JASPAR Database interface for matrix profile MA0002.1. The page is divided into several sections:

- Profile summary:** A table containing key information about the matrix profile.
- Sequence logo:** A visual representation of the sequence logo for the matrix profile, showing the relative frequency of nucleotides at each position.
- Frequency matrix:** A table showing the frequency of nucleotides (A, C, G, T) at each position (1-11).
- Binding sites information:** A section with buttons to download the HTML file or FASTA file.
- External links:** A section with links to PDB, UniProt, DRV, and TFBSshape.
- Version information:** A table showing the version history of the matrix profile.

Profile summary details:

Field	Value
Name	RUNX1
Matrix ID	MA0002.1
Class	Runt domain factors
Family	Runt-related factors
Collection	CORE
Taxon	Vertebrates
Species	Homo sapiens
Data Type	SELEX
Validation	8413232
Uniprot ID	Q01196
Source	
Comment	Matrix changed since last release: removal of primers and sites overlapping primers

Sequence logo: A sequence logo showing the relative frequency of nucleotides at each position. The logo is color-coded: A (green), C (blue), G (yellow), and T (red). The positions are numbered 1 to 11.

Frequency matrix:

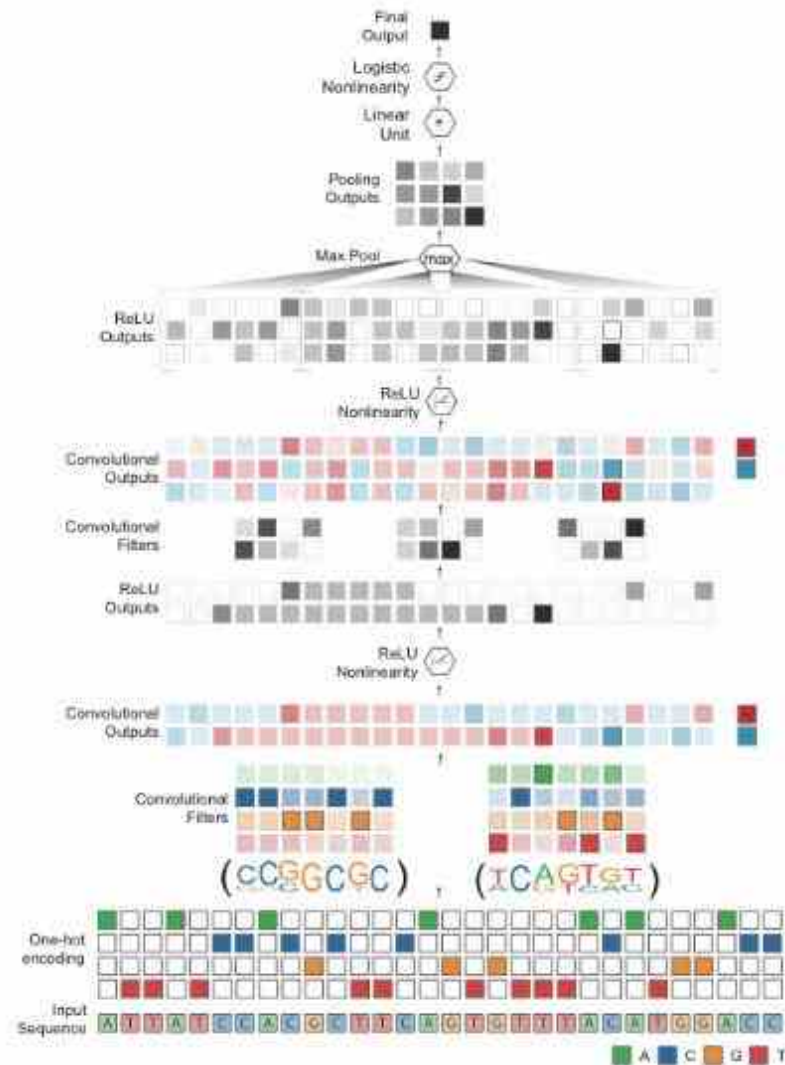
	1	2	3	4	5	6	7	8	9	10	11
A	10	12	4	1	2	2	0	0	0	8	13
C	2	2	7	1	0	8	0	0	1	2	2
G	3	1	1	0	23	0	26	26	0	0	4
T	11	11	14	24	1	16	0	0	25	16	7

Version information:

Matrix ID	Base ID	Version	Name	Species	Family	Class	Sequence logo
MA0002.2	MA0002	2	Runx1	Mus musculus	Runt-related factors	Runt domain factors	

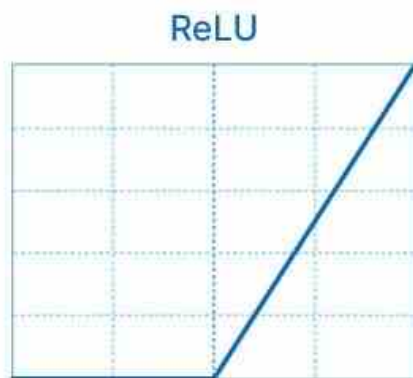
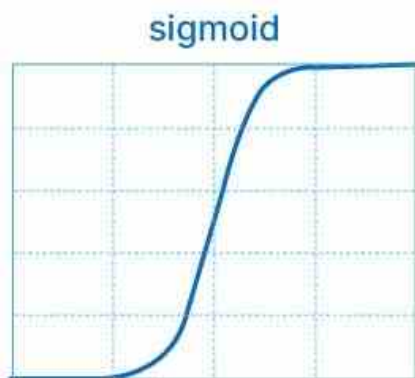
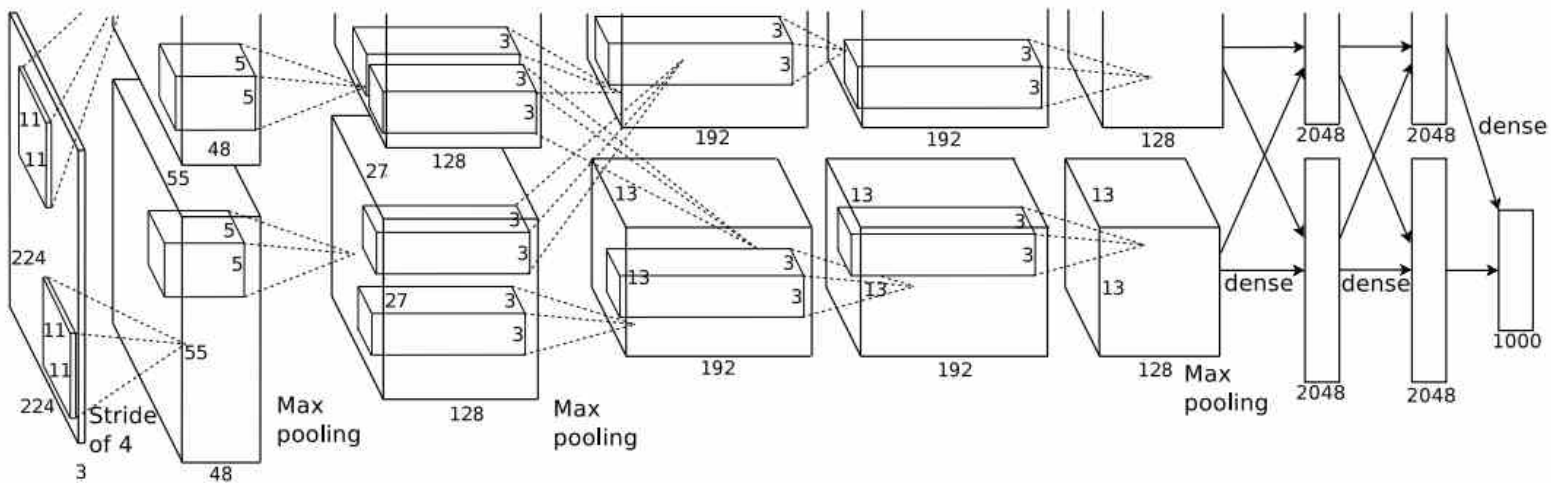
<https://jaspar.genereg.net/matrix/MA0002.1/>

ML with Strings

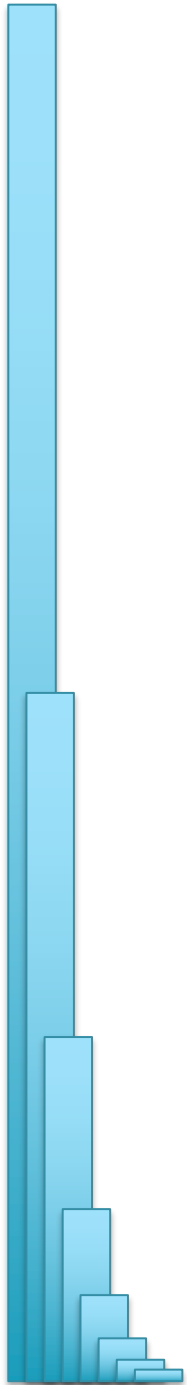


One hot encoding to sequence classification

<https://kundajelab.github.io/dragonntutorials.html>



ImageNet Classification with Deep Convolutional Neural Networks
 Krizhevsky et al. (2012) Advances in Neural Information Processing Systems 25

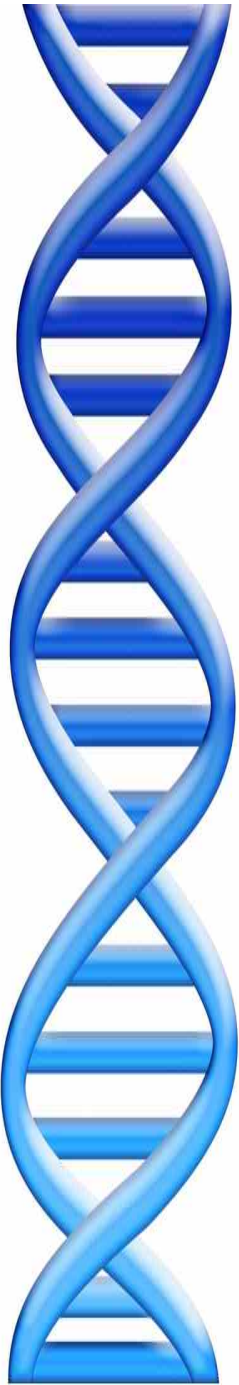


Annotation

Goal: Genome Annotations

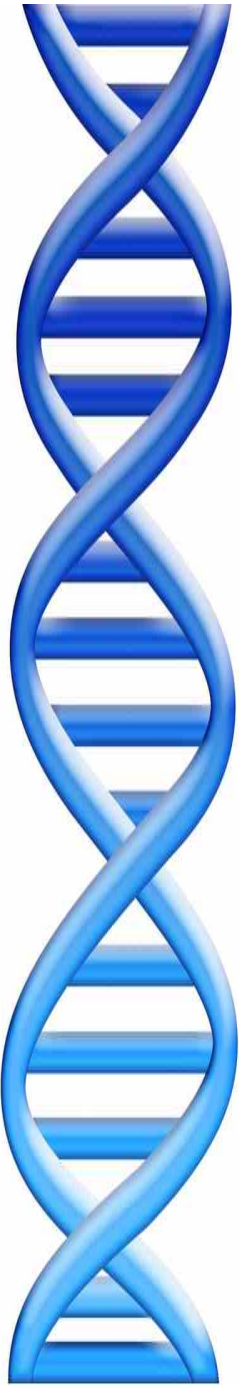
aatgcatgctggctatgctaagcatgctggctatgctaagctgggatccgatgacaatgcatgctggctatgctaag
gcatgctggctatgcaagctgggatccgatgactatgctaagctgggatccgatgacaatgcatgctggctatgct
aatgaatgggtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctggctatgctaagaa
tggtcttgggatttaccttgggaatgctaagcatgctggctatgctaagctgggatccgatgacaatgcatgctg
gctatgctaagcatgctggctatgcaagctgggatccgatgactatgctaagctgctggctatgctaagcatgctg
gctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatcc
gctggctatgctaagcatgctgggtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctg
atgctaagcatgctgggtcttgggatttaccttgggaatgctaagcatgctggctatgctaagcatgctg
gctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatccg
atgactatgctaagctgctggctatgctaagcatgctggctatgctaagcatgctggctatgctaagctgggaat
gcatgctggctatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctg
ggatccgatgactatgctaagctgctggctatgctaagcatgctggctatgctaagctgctggctatgctaagcatg
gtcttgggatttaccttgggaatgctaagctgggatccgatgacaatgcatgctggctatgctaagcatgctgg
gatttaccttgggaatgctaagcatgctggctatgctaagctgggaatgcatgctggctatgctaagctgggatc
cgatgacaatgcatgctggctatgctaagcatgctggctatgcaagctgggatccgatgactatgctaagctgctg
gctatgctaagcatgctggctatgctaagctcatgctg

Gene!



Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays



Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays

Very Similar Sequences

Query: HBA_HUMAN Hemoglobin alpha subunit

Sbjct: HBB_HUMAN Hemoglobin beta subunit

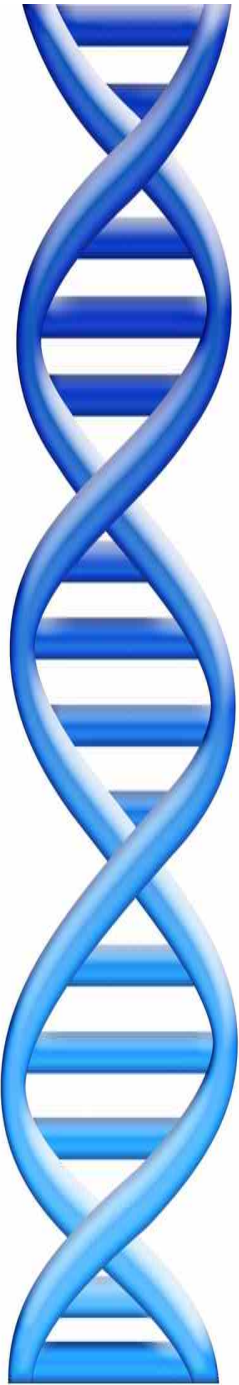
Score = 114 bits (285), Expect = 1e-26

Identities = 61/145 (42%), Positives = 86/145 (59%), Gaps = 8/145 (5%)

```
Query    2    LSPADKTNVKAANGKVGAGHAGEYGAELERMFLSFPTTKTYFPHF-----DLSHGSAQV 55
          L+P +K+ V A WGKV  +  E G EAL R+ + +P T+ +F  F      D    G+ +V
Sbjct    3    LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60

Query    56    KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA 115
          K HGKKV  A ++ +AH+D++      + LS+LH  KL VDP NF+LL + L+  LA H
Sbjct    61    KAHGKKVLGAFSDGLAHLNLRGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 120

Query    116   EFTPAVHASLDKFLASVSTVLTSKY 140
          EFTP V A+  K +A V+  L  KY
Sbjct    121   EFTPPVQAAYQKVVAGVANALAHKY 145
```

Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. Experimental & Functional Assays



Bacterial Gene Finding and Glimmer

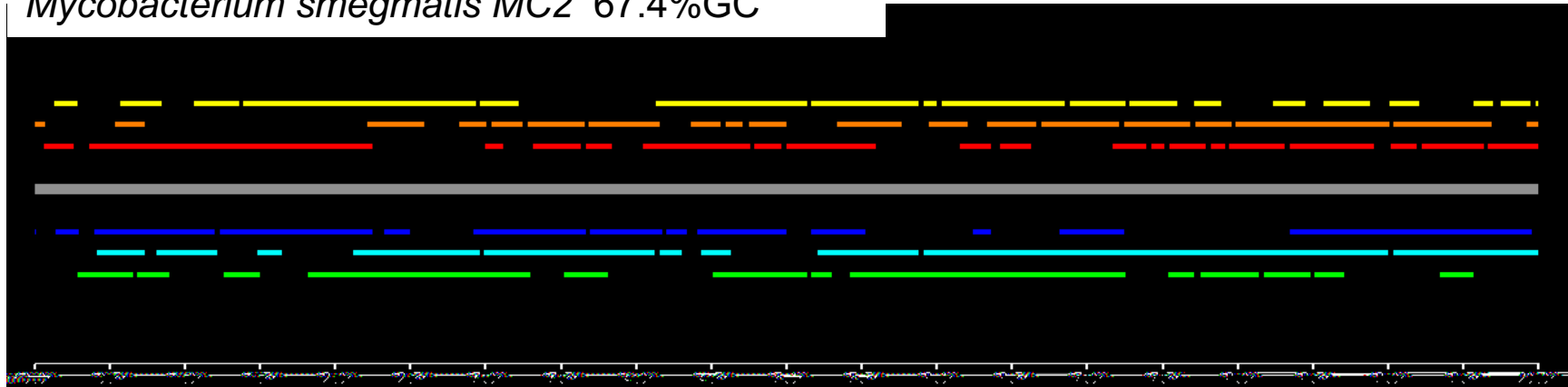
(also Archaeal and viral gene finding)

Arthur L. Delcher and Steven Salzberg

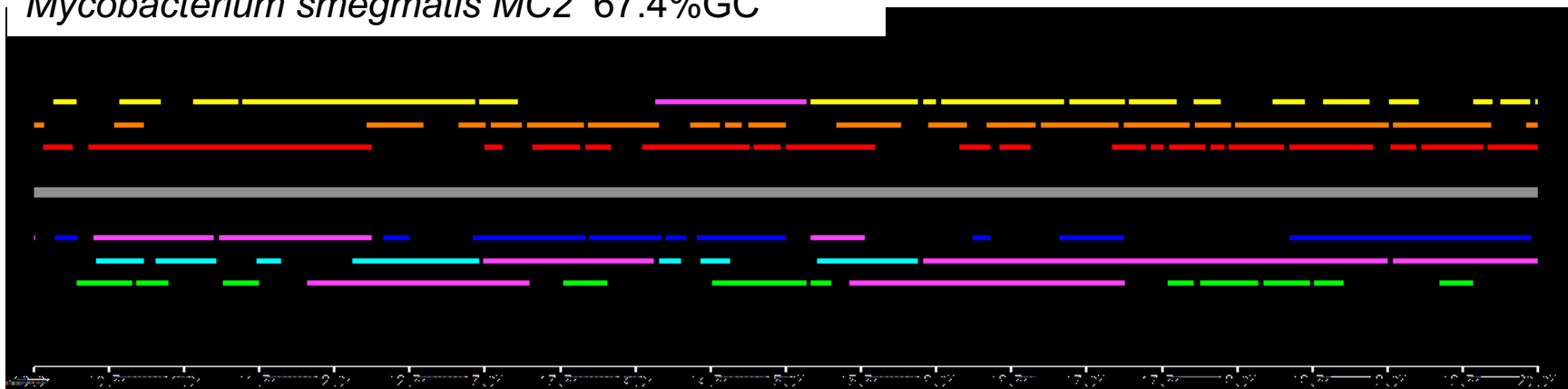
Center for Bioinformatics and Computational Biology

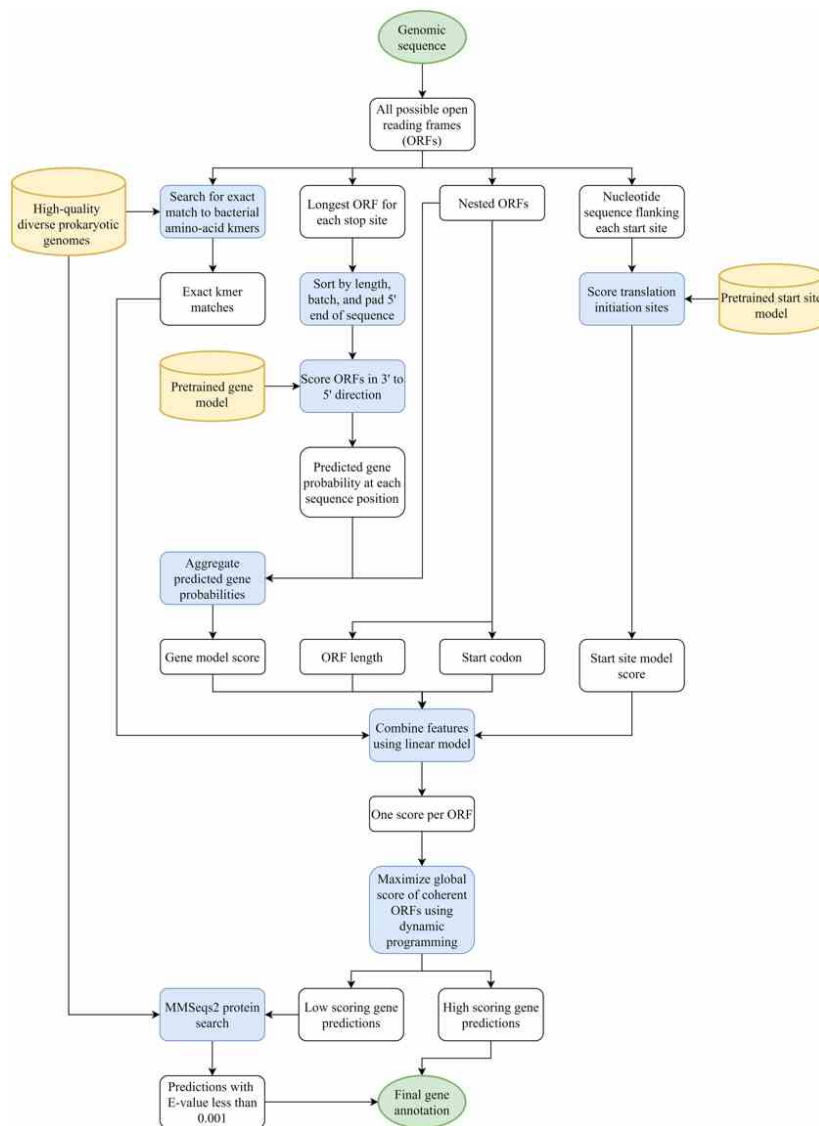
Johns Hopkins University

Mycobacterium smegmatis MC2 67.4%GC

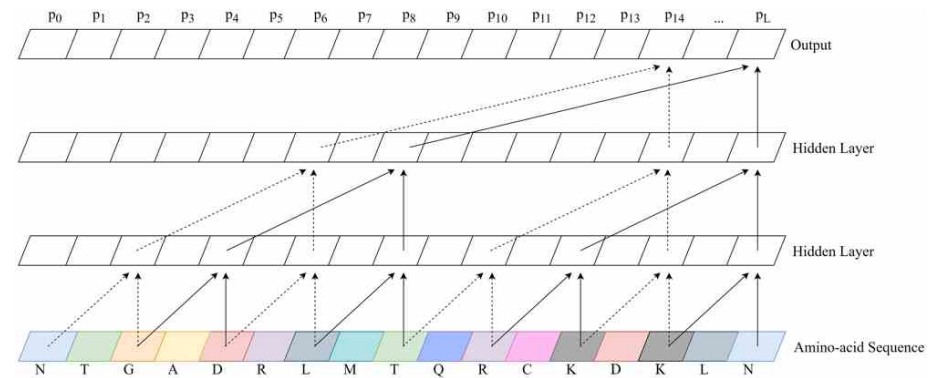


Mycobacterium smegmatis MC2 67.4%GC





Temporal Convolutional Network



Balrog: A universal protein model for prokaryotic gene prediction

Sommer, MJ, Salzberg, SL (2021) PLOS Comp. Bio. doi: 10.1371/journal.pcbi.1008727

Probabilistic Methods

- Create models that have a probability of generating any given sequence.
 - Evaluate gene/non-genome models against a sequence
- Train the models using examples of the types of sequences to generate.
 - Use RNA sequencing, homology, or “obvious” genes
- The “score” of an orf is the probability of the model generating it.
 - Most basic technique is to count how kmers occur in known genes versus intergenic sequences
 - More sophisticated methods consider variable length contexts, “wobble” bases, other statistical clues

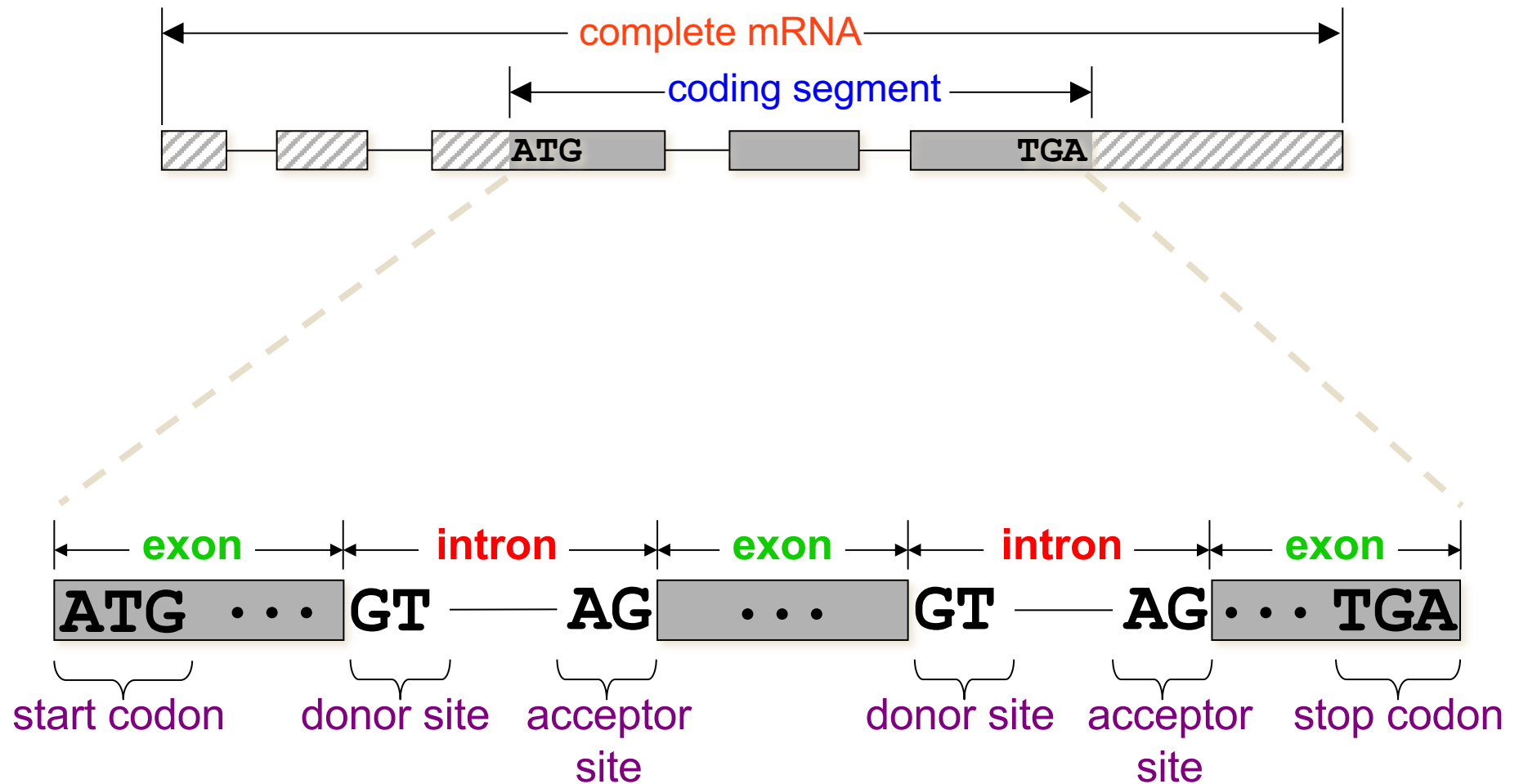


Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros

Eukaryotic Gene Syntax



Regions of the gene outside of the CDS are called **UTR**'s (*untranslated regions*), and are mostly ignored by gene finders, though they are important for regulatory functions.

What is an HMM?

- Dynamic Bayesian Network

- A set of states

- {Fair, Biased} for coin tossing
 - {Gene, Not Gene} for Bacterial Gene
 - {Intergenic, Exon, Intron} for Eukaryotic Gene
 - {Modern, Neanderthal} for Ancestry

- A set of emission characters

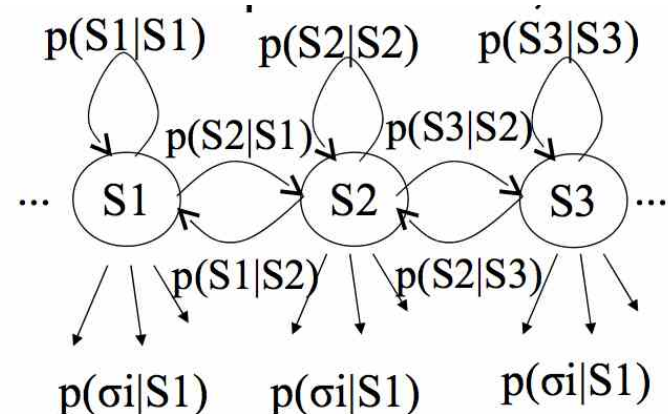
- $E=\{H,T\}$ for coin tossing
 - $E=\{1,2,3,4,5,6\}$ for dice tossing
 - $E=\{A,C,G,T\}$ for DNA

- State-specific emission probabilities

- $P(H \mid \text{Fair}) = .5, P(T \mid \text{Fair}) = .5, P(H \mid \text{Biased}) = .9, P(T \mid \text{Biased}) = .1$
 - $P(A \mid \text{Gene}) = .9, P(A \mid \text{Not Gene}) = .1 \dots$

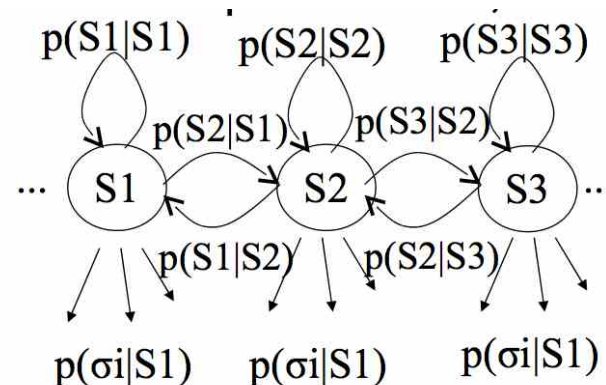
- A probability of taking a transition

- $P(s_i=\text{Fair} \mid s_{i-1}=\text{Fair}) = .9, P(s_i=\text{Bias} \mid s_{i-1} = \text{Fair}) .1$
 - $P(s_i=\text{Exon} \mid s_{i-1}=\text{Intergenic}), \dots$



Why Hidden?

- Similar to Markov models used for prokaryotic gene finding, but system may transition between multiple models called states (gene/non-gene, intergenic/exon/intron)
- Observers can see the emitted symbols of an HMM (i.e., nucleotides) but have no ability to know which state the HMM is currently in.
 - But we can *infer* the most likely hidden states of an HMM based on the given sequence of emitted symbols.

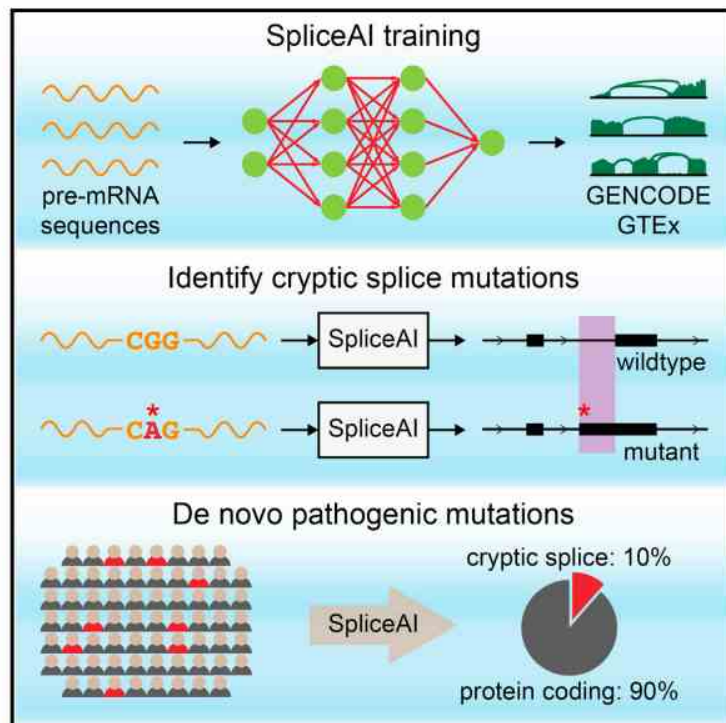


See
lecture
notes!

AAAGCATGCATTTAACGTGAGCACAAATAGATTACA

Predicting Splicing from Primary Sequence with Deep Learning

Graphical Abstract



Authors

Kishore Jaganathan,
Sofia Kyriazopoulou Panagiotopoulou,
Jeremy F. McRae, ..., Serafim Batzoglou,
Stephan J. Sanders, Kyle Kai-How Farh

Correspondence

kfarh@illumina.com

In Brief

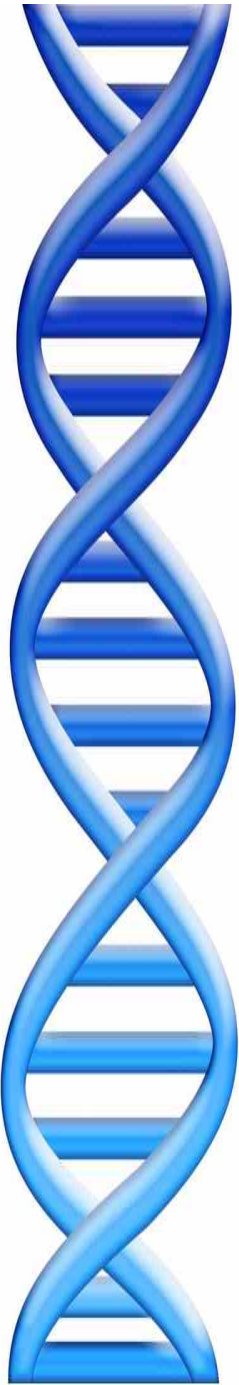
A deep neural network precisely models mRNA splicing from a genomic sequence and accurately predicts noncoding cryptic splice mutations in patients with rare genetic diseases.

Highlights

- SpliceAI, a 32-layer deep neural network, predicts splicing from a pre-mRNA sequence
- 75% of predicted cryptic splice variants validate on RNA-seq
- Cryptic splicing may yield ~10% of pathogenic variants in neurodevelopmental disorders
- Cryptic splice variants frequently give rise to alternative splicing

Gene Finding Overview

- Prokaryotic gene finding distinguishes real genes and random ORFs
 - Prokaryotic genes have simple structure and are largely homogenous, making it relatively easy to recognize their sequence composition
- Eukaryotic gene finding identifies the genome-wide most probable gene models (set of exons)
 - “Probabilistic Graphical Model” to enforce overall gene structure, separate models to score splicing/transcription signals
 - Accuracy depends to a large extent on the quality of the training data



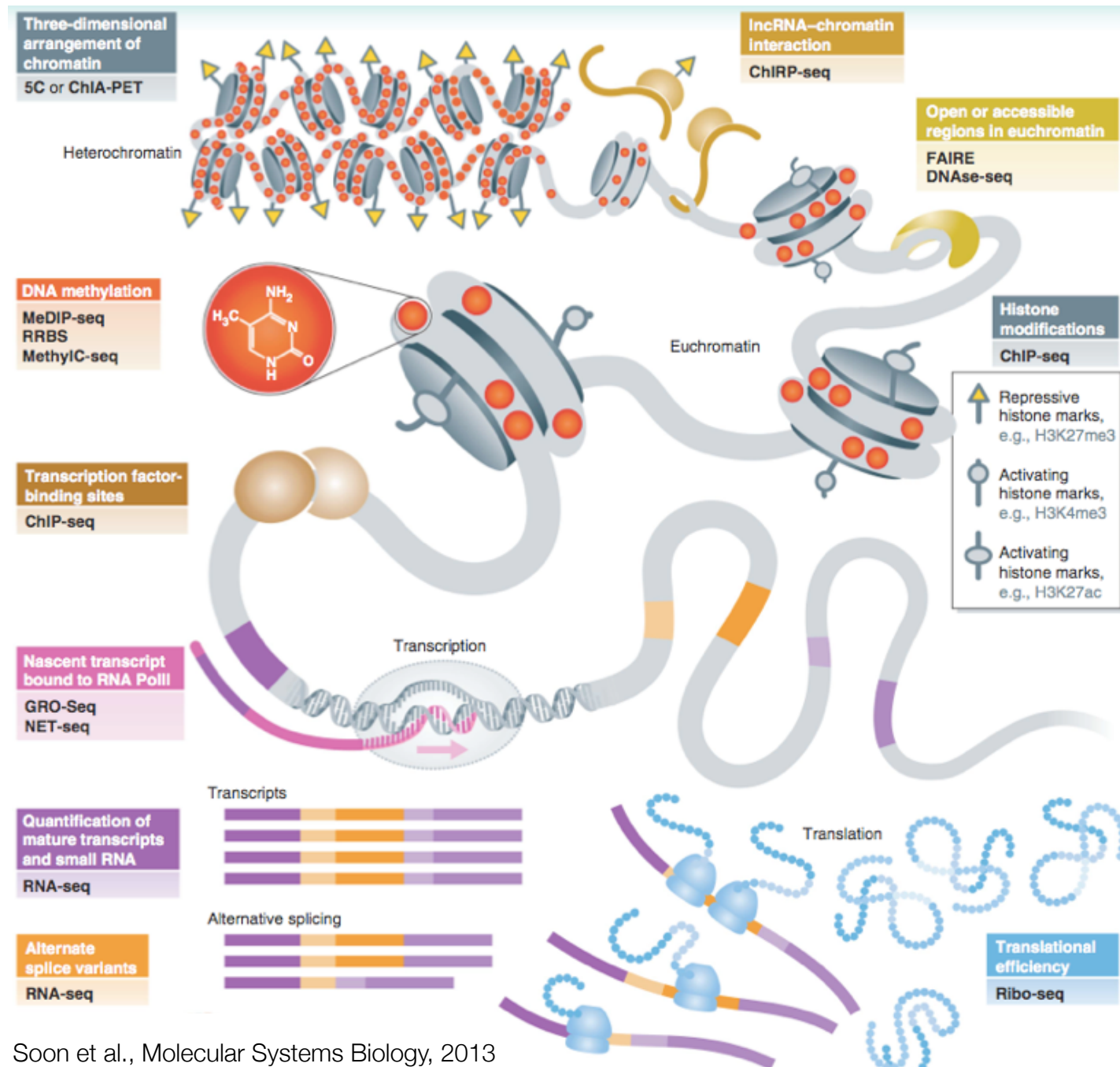
Outline

1. Alignment to other genomes
2. Prediction aka “Gene Finding”
3. **Experimental & Functional Assays**

Sequencing Assays

The *Seq List (in chronological order)

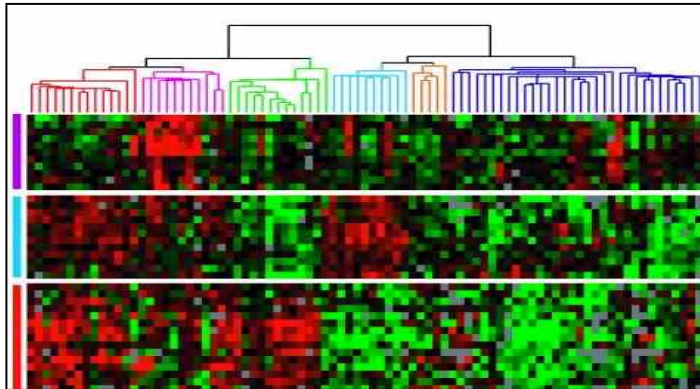
1. Gregory E. Crawford et al., "Genome-wide Mapping of DNase Hypersensitive Sites Using Massively Parallel Signature Sequencing (MPSS)," *Genome Research* 16, no. 1 (January 1, 2006): 123–131, doi:10.1101/gr.4074106.
2. David S. Johnson et al., "Genome-Wide Mapping of in Vivo Protein-DNA Interactions," *Science* 316, no. 5830 (June 8, 2007): 1497–1502, doi:10.1126/science.1141319.
3. Tarjei S. Mikkelsen et al., "Genome-wide Maps of Chromatin State in Pluripotent and Lineage-committed Cells," *Nature* 448, no. 7153 (August 2, 2007): 553–560, doi:10.1038/nature06008.
4. Thomas A. Down et al., "A Bayesian Deconvolution Strategy for Immunoprecipitation-based DNA Methylome Analysis," *Nature Biotechnology* 26, no. 7 (July 2008): 779–785, doi:10.1038/nbt1414.
5. Ali Mortazavi et al., "Mapping and Quantifying Mammalian Transcriptomes by RNA-Seq," *Nature Methods* 5, no. 7 (July 2008): 621–628, doi:10.1038/nmeth.1226.
6. Nathan A. Baird et al., "Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers," *PLoS ONE* 3, no. 10 (October 13, 2008): e3376, doi:10.1371/journal.pone.0003376.
7. Leighton J. Core, Joshua J. Waterfall, and John T. Lis, "Nascent RNA Sequencing Reveals Widespread Pausing and Divergent Initiation at Human Promoters," *Science* 322, no. 5909 (December 19, 2008): 1845–1848, doi:10.1126/science.1162228.
8. Chao Xie and Martti T. Tammi, "CNV-seq, a New Method to Detect Copy Number Variation Using High-throughput Sequencing," *BMC Bioinformatics* 10, no. 1 (March 6, 2009): 80, doi:10.1186/1471-2105-10-80.
9. Jay R. Hesselberth et al., "Global Mapping of protein-DNA Interactions in Vivo by Digital Genomic Footprinting," *Nature Methods* 6, no. 4 (April 2009): 283–289, doi:10.1038/nmeth.1313.
10. Nicholas T. Ingolia et al., "Genome-Wide Analysis in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling," *Science* 324, no. 5924 (April 10, 2009): 218–223, doi:10.1126/science.1168978.
11. Alayne L. Brunner et al., "Distinct DNA Methylation Patterns Characterize Differentiated Human Embryonic Stem Cells and Developing Human Fetal Liver," *Genome Research* 19, no. 6 (June 1, 2009): 1044–1056, doi:10.1101/gr.088773.108.
12. Mayumi Oda et al., "High-resolution Genome-wide Cytosine Methylation Profiling with Simultaneous Copy Number Analysis and Optimization for Limited Cell Numbers," *Nucleic Acids Research* 37, no. 12 (July 1, 2009): 3829–3839, doi:10.1093/nar/gkp260.
13. Zachary D. Smith et al., "High-throughput Bisulfite Sequencing in Mammalian Genomes," *Methods* 48, no. 3 (July 2009): 226–232, doi:10.1016/j.ymeth.2009.05.003.
14. Andrew M. Smith et al., "Quantitative Phenotyping via Deep Barcode Sequencing," *Genome Research* (July 21, 2009),



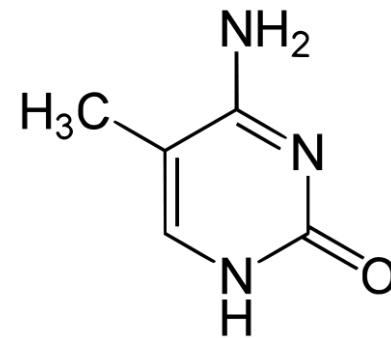
Soon et al., Molecular Systems Biology, 2013

*-seq in 4 short vignettes

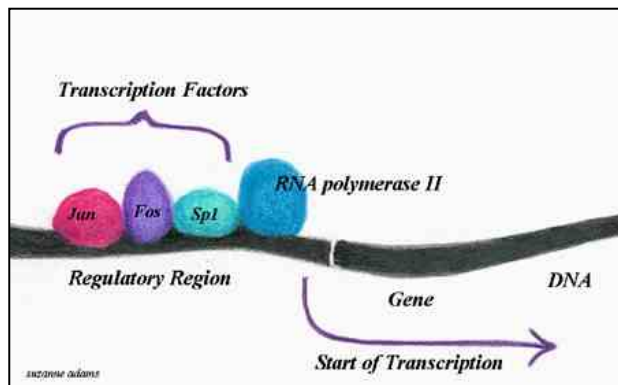
RNA-seq



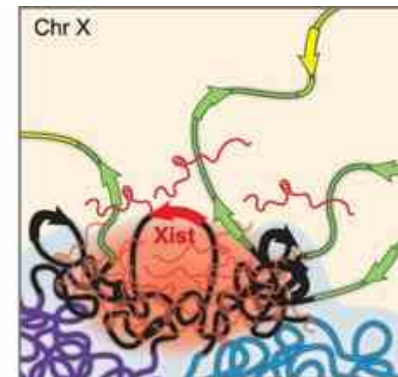
Methyl-seq



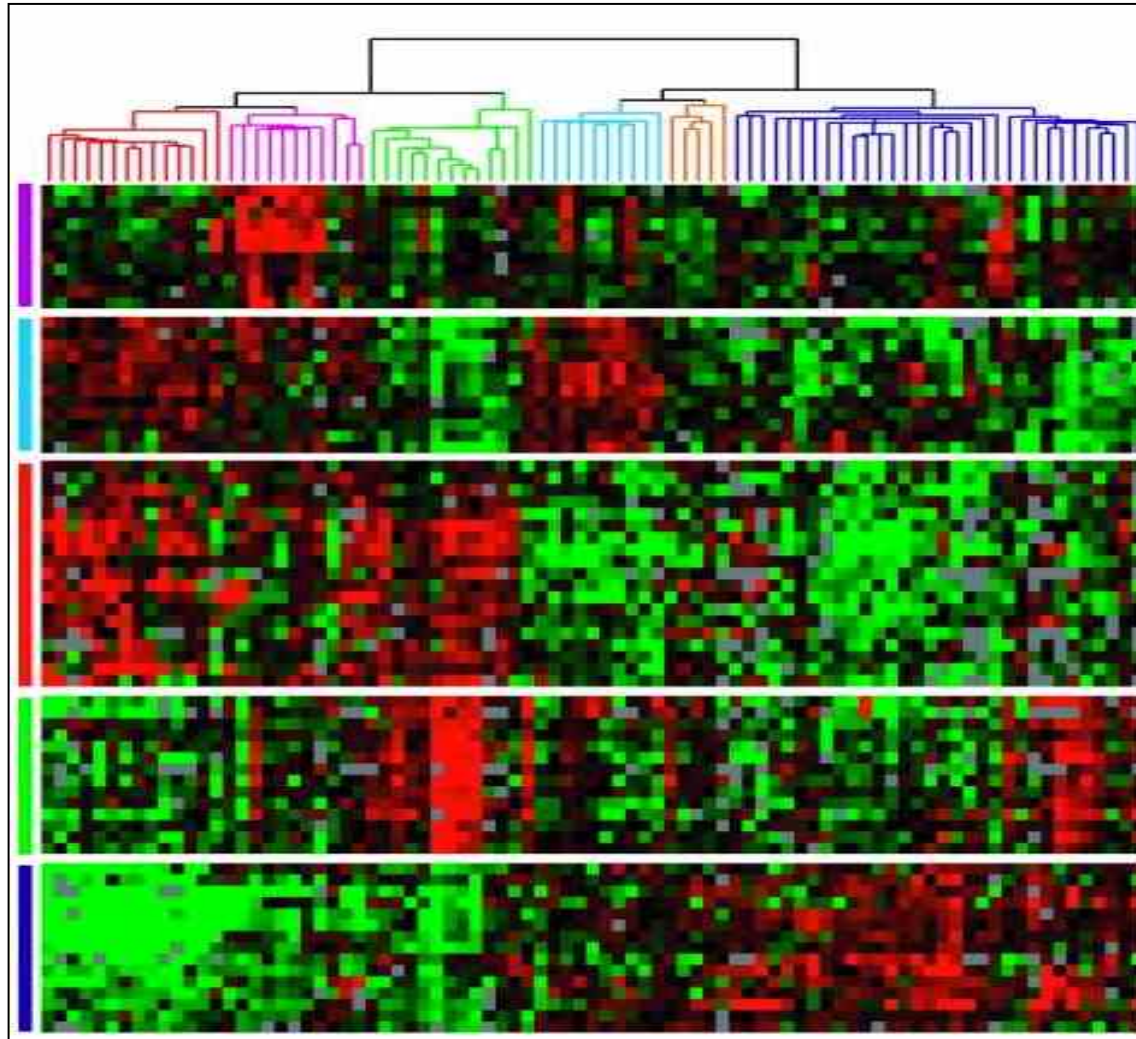
ChIP-seq



Hi-C

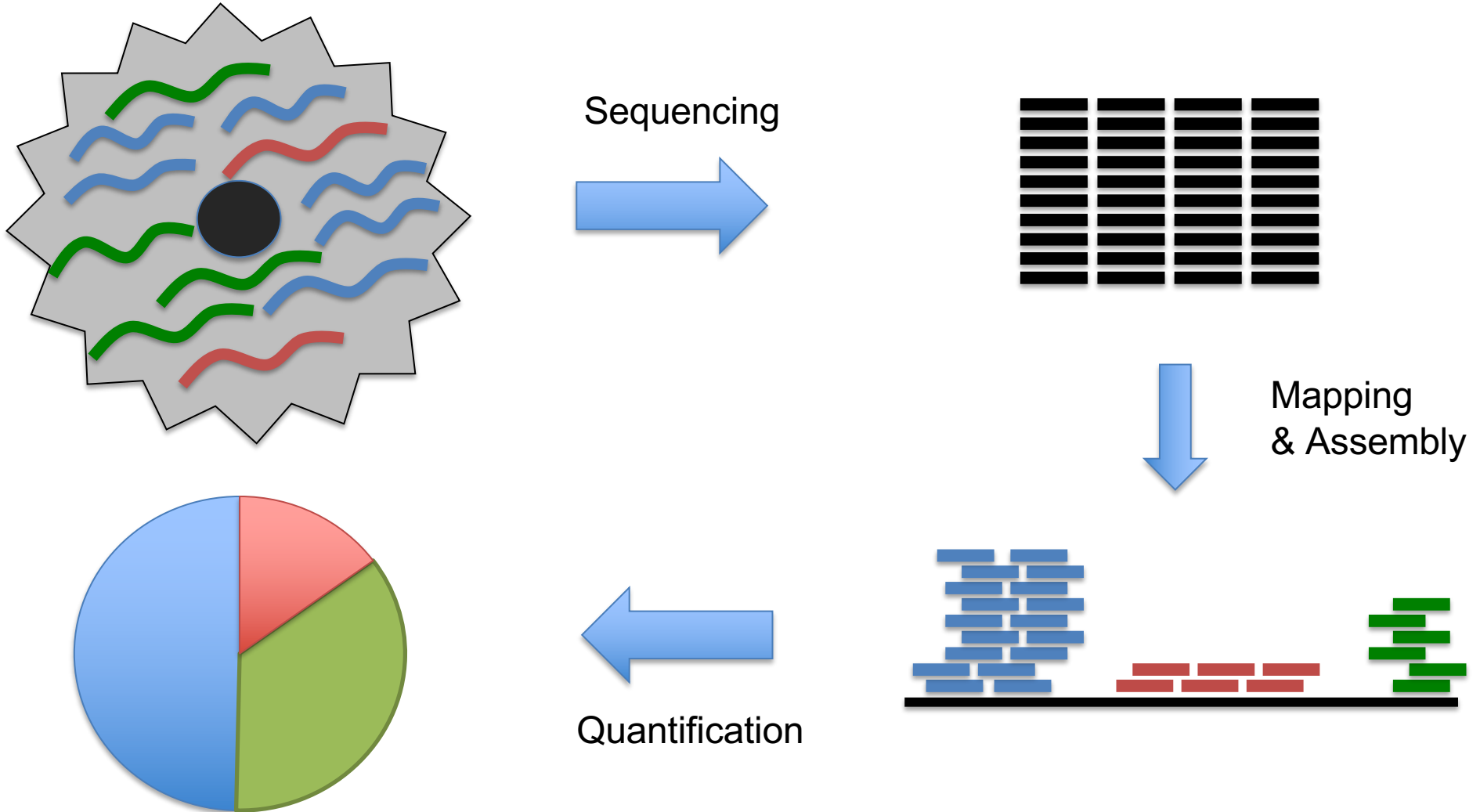


RNA-seq

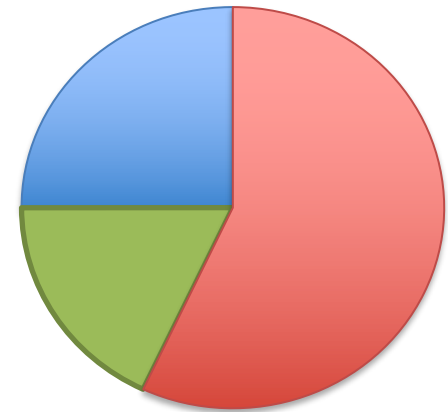
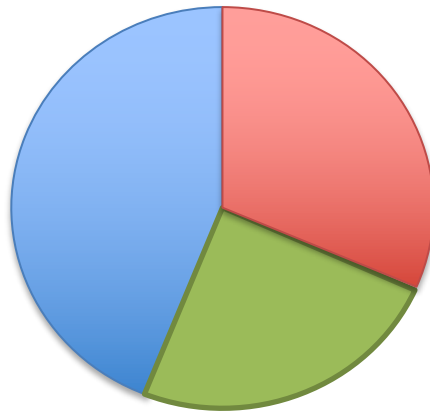
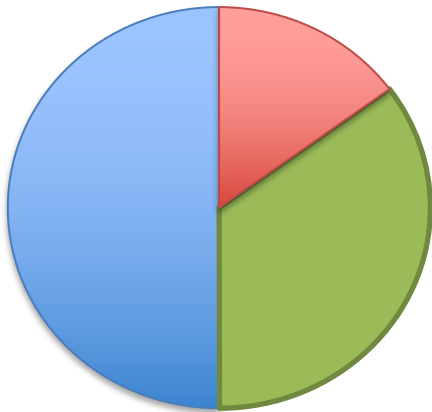
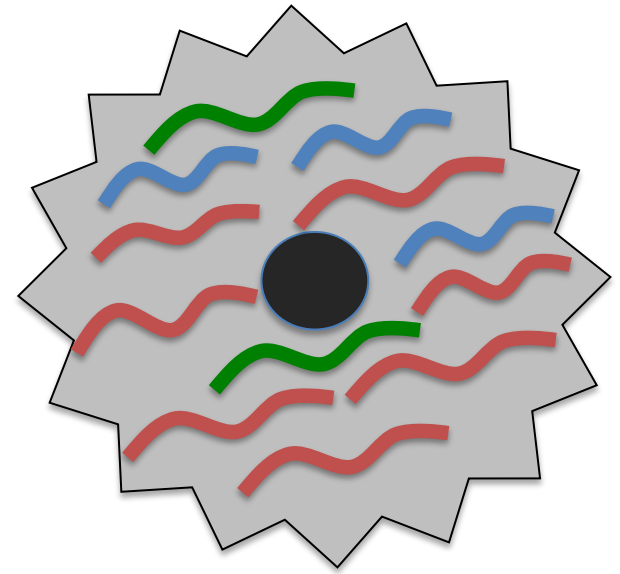
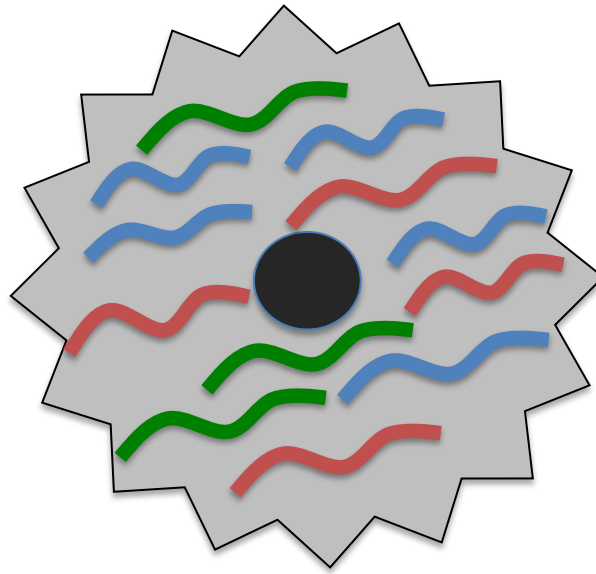
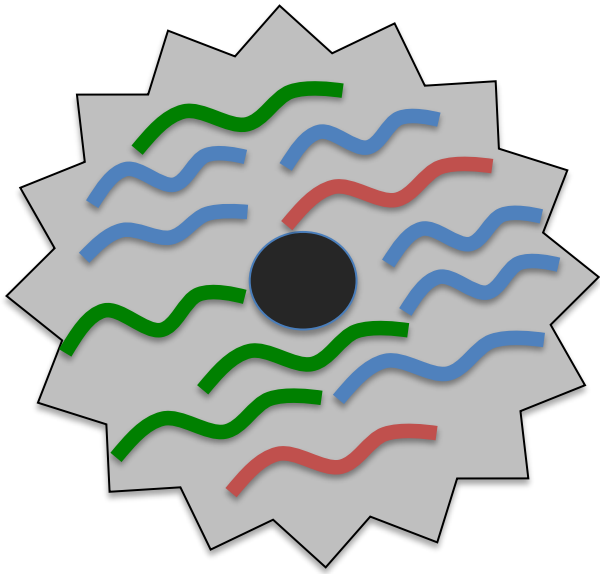


Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.
Sørli et al (2001) *PNAS*. 98(19):10869-74.

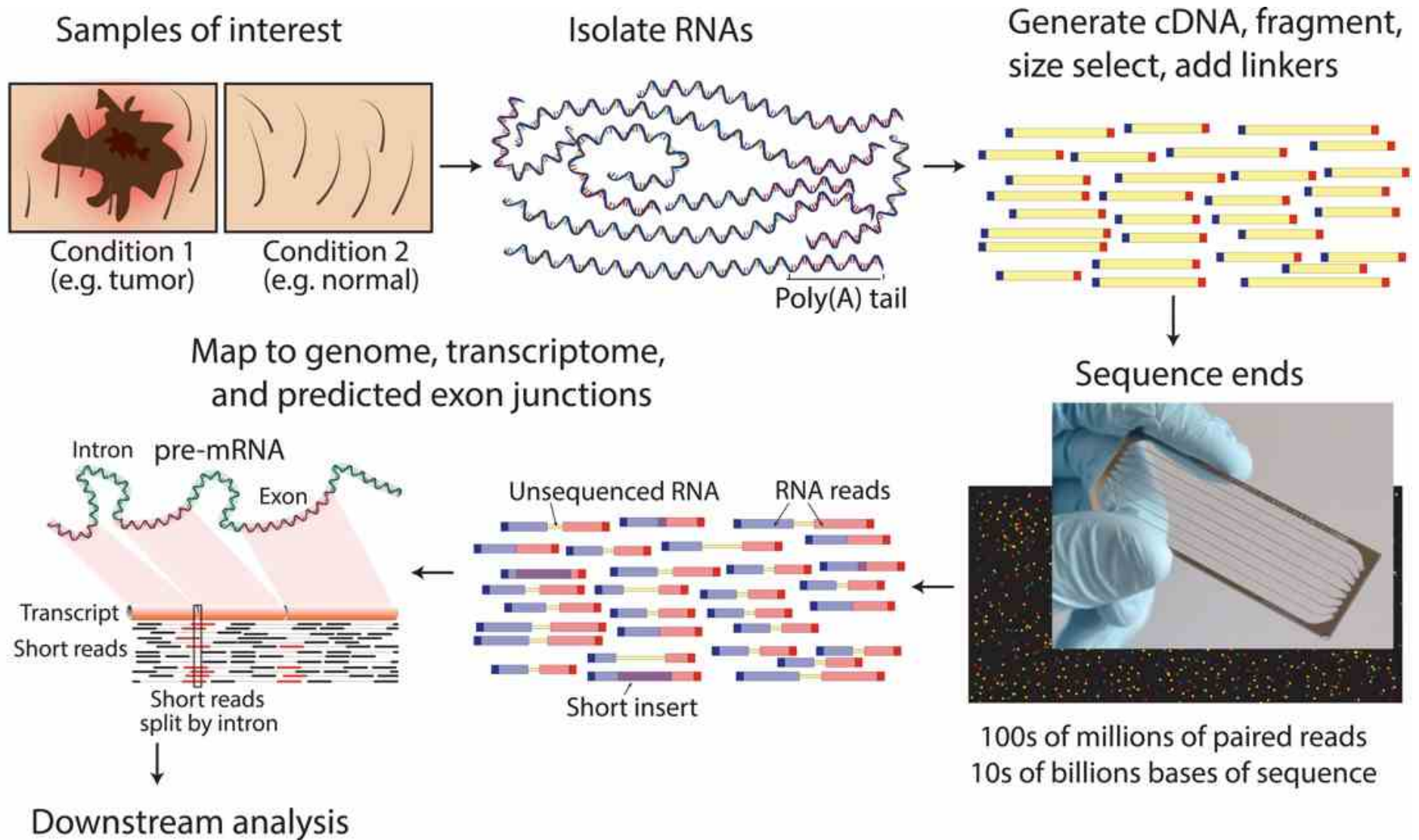
RNA-seq Overview



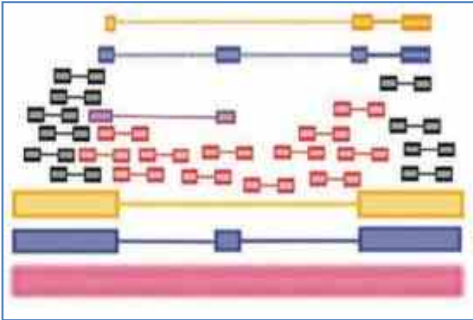
RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges



Challenge I: Eukaryotic genes are spliced

RNA-Seq Approaches

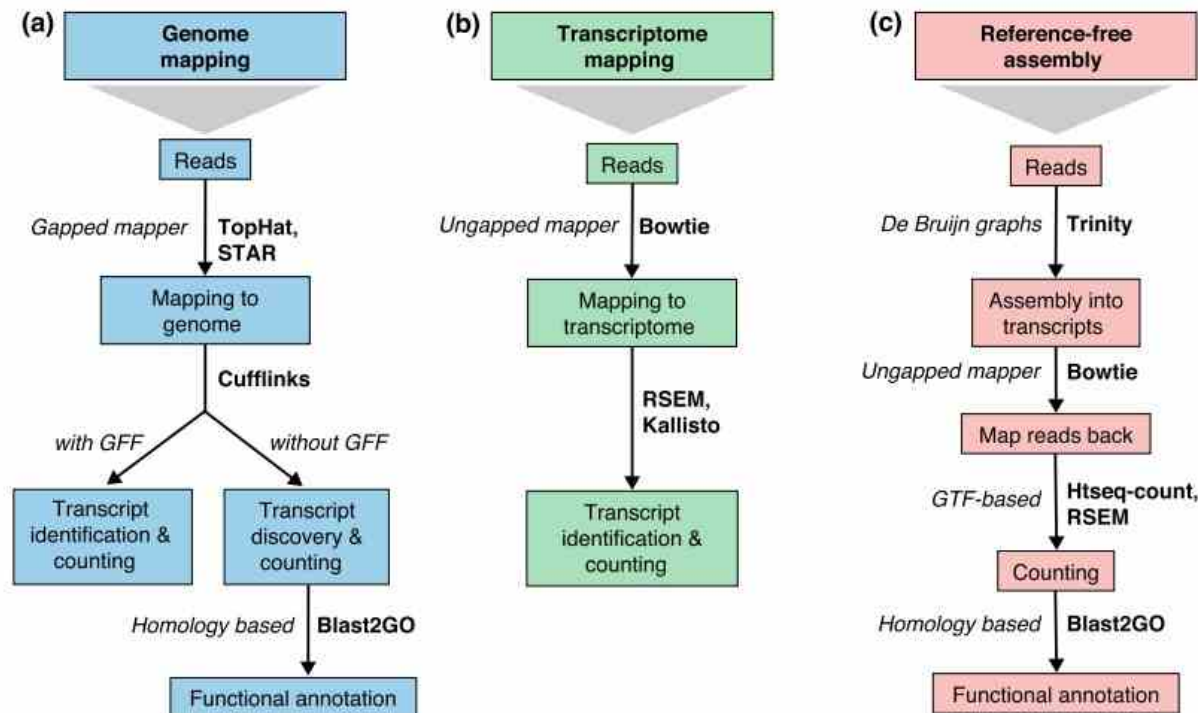


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation file. Novel transcripts are then functionally annotated. **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis proceeds as in **(b)** followed by the functional annotation of the novel transcripts as in **(a)**. Representative software that can be used at each analysis step are indicated in *bold text*. Abbreviations: *GFF* General Feature Format, *GTF* gene transfer format, *RSEM* RNA-Seq by Expectation Maximization

A survey of best practices for RNA-seq data analysis

Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-Seq Approaches

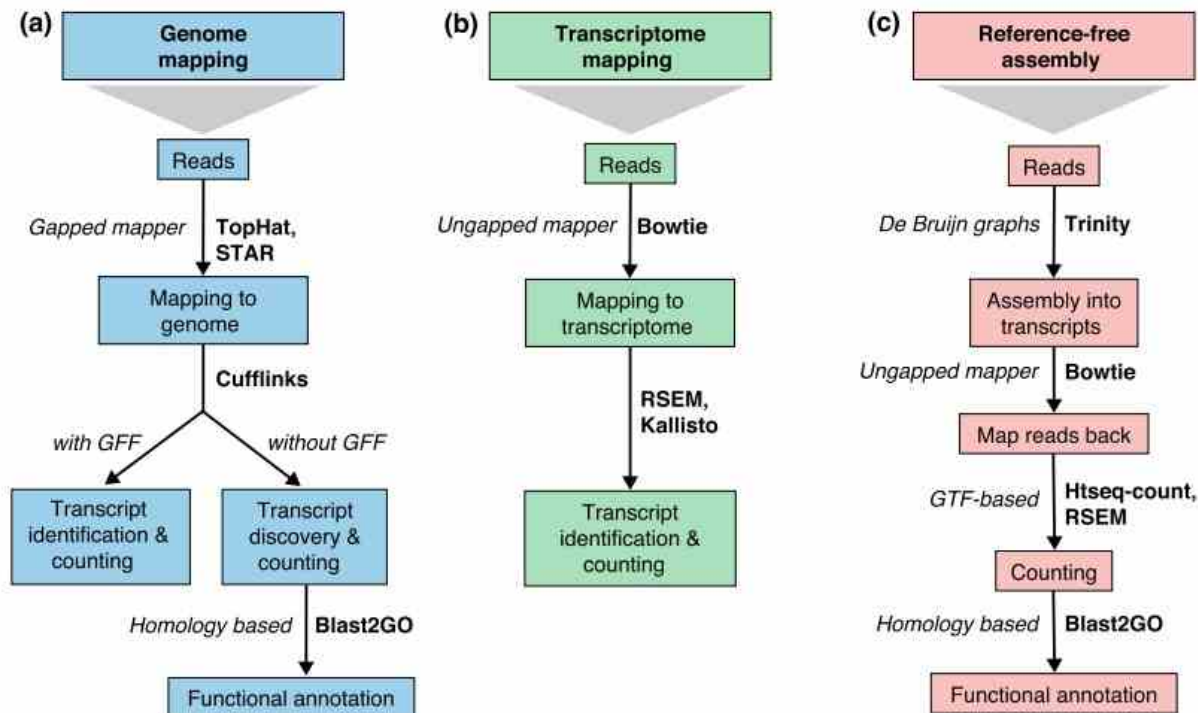


Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reference (or novel) transcript discovery and quantification can proceed with or without an annotation. **b** If a reference transcriptome is available and no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis is followed by the functional annotation of the novel transcripts as in **(a)**. Representative software that can be used at each analysis step are indicated in *bold text*. Abbreviations: *GFF* General Feature Format, *GTF* gene transfer format, *RSEM* RNA-Seq by Expectation Maximization

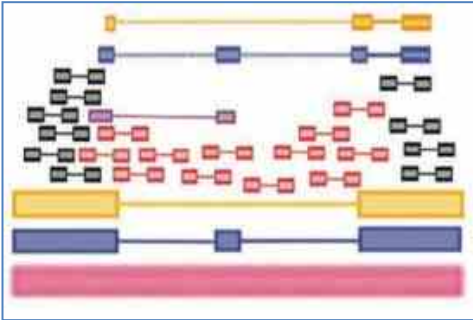
Which approach should we use?

It depends....

A survey of best practices for RNA-seq data analysis

Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-seq Challenges

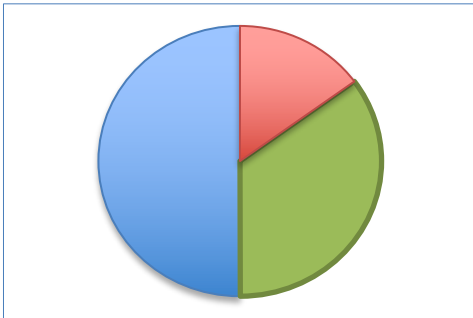


Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

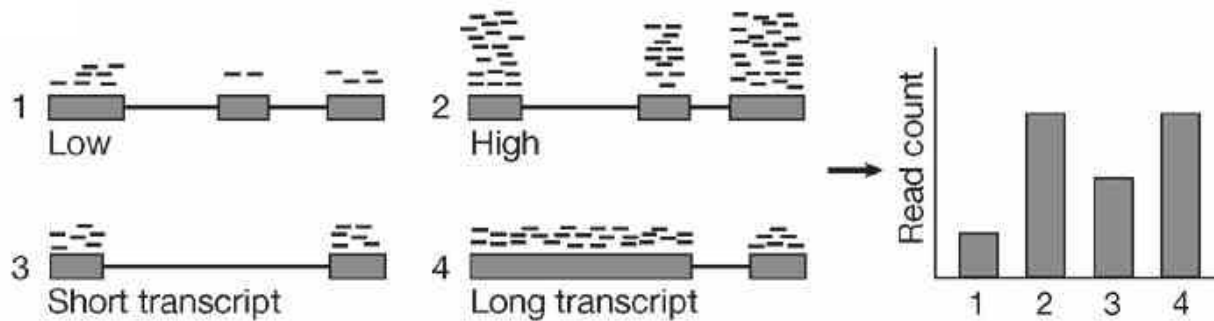
TopHat: discovering spliced junctions with RNA-Seq.

Trapnell et al (2009) *Bioinformatics*. 25:0 ||05-||||



Challenge 2: Read Count != Transcript abundance

RPKM, FPKM, TPM

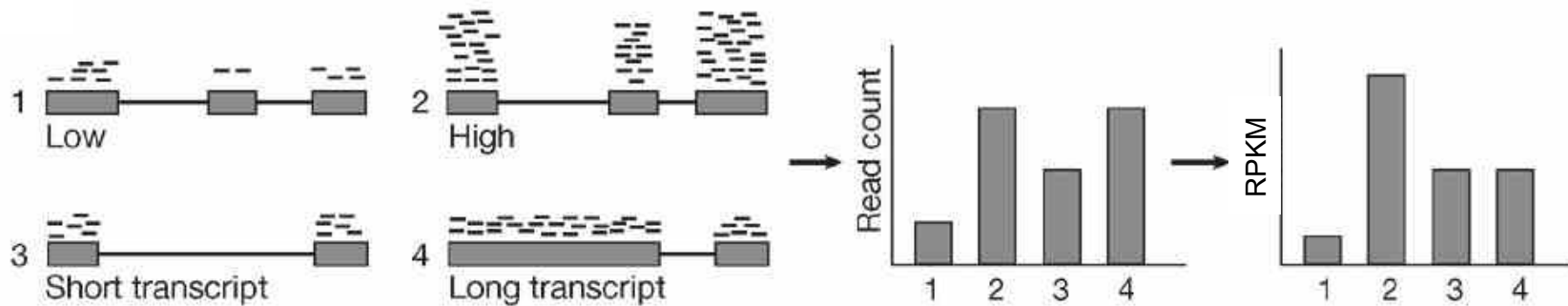


Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

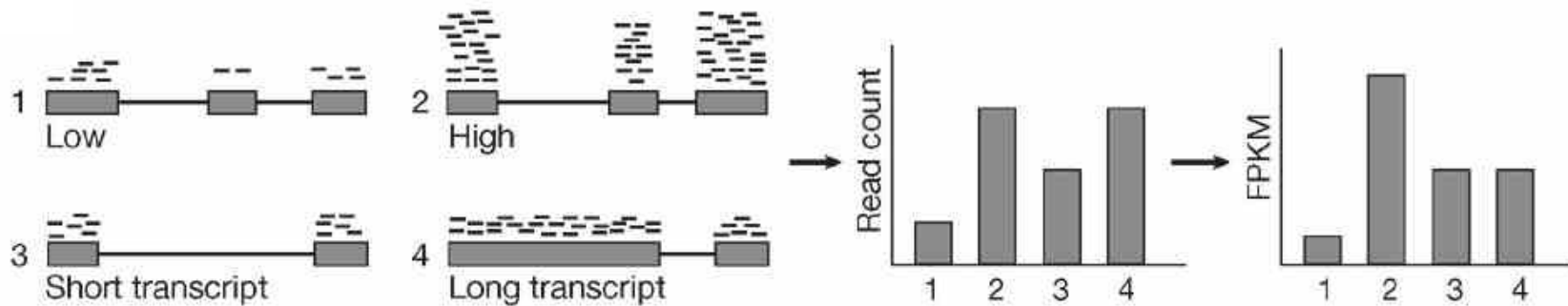
- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

(Count reads aligned to gene) / (length of gene in kilobases) / (# millions of read mapped)

=> Wait a second, reads in a pair arent independent!

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

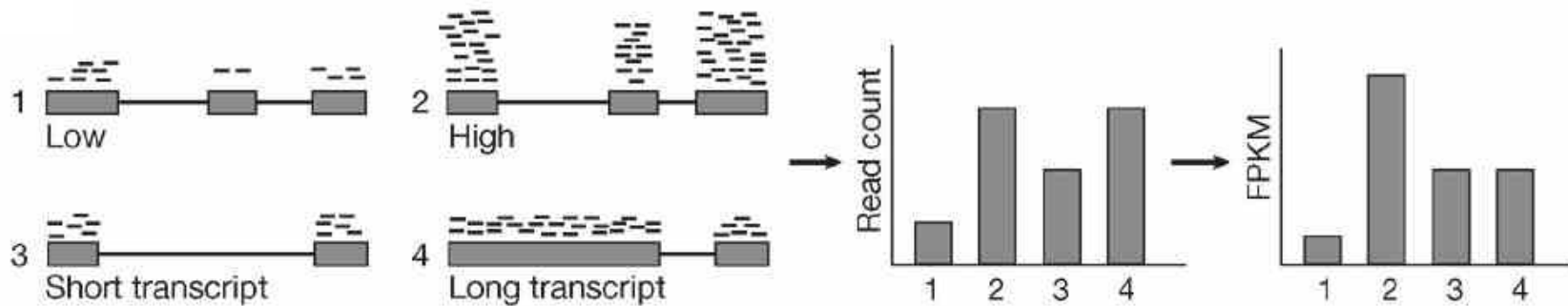
=> Wait a second, reads in a pair are not independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)

=> Does a much better job with short exons & short genes by boosting coverage

=> Wait a second, FPKM depends on the average transcript length!

RPKM, FPKM, TPM



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

=> Wait a second, reads in a pair are not independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)

=> Wait a second, FPKM depends on the average transcript length!

3. TPM: Transcripts Per Million (Li et al, 2011)

=> If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i , given the abundances of the other transcripts in your sample

=> Recommend you use TPM for all analysis, easy to compute given FPKM

$$TPM_i = \left(\frac{FPKM_i}{\sum_j FPKM_j} \right) \cdot 10^6$$

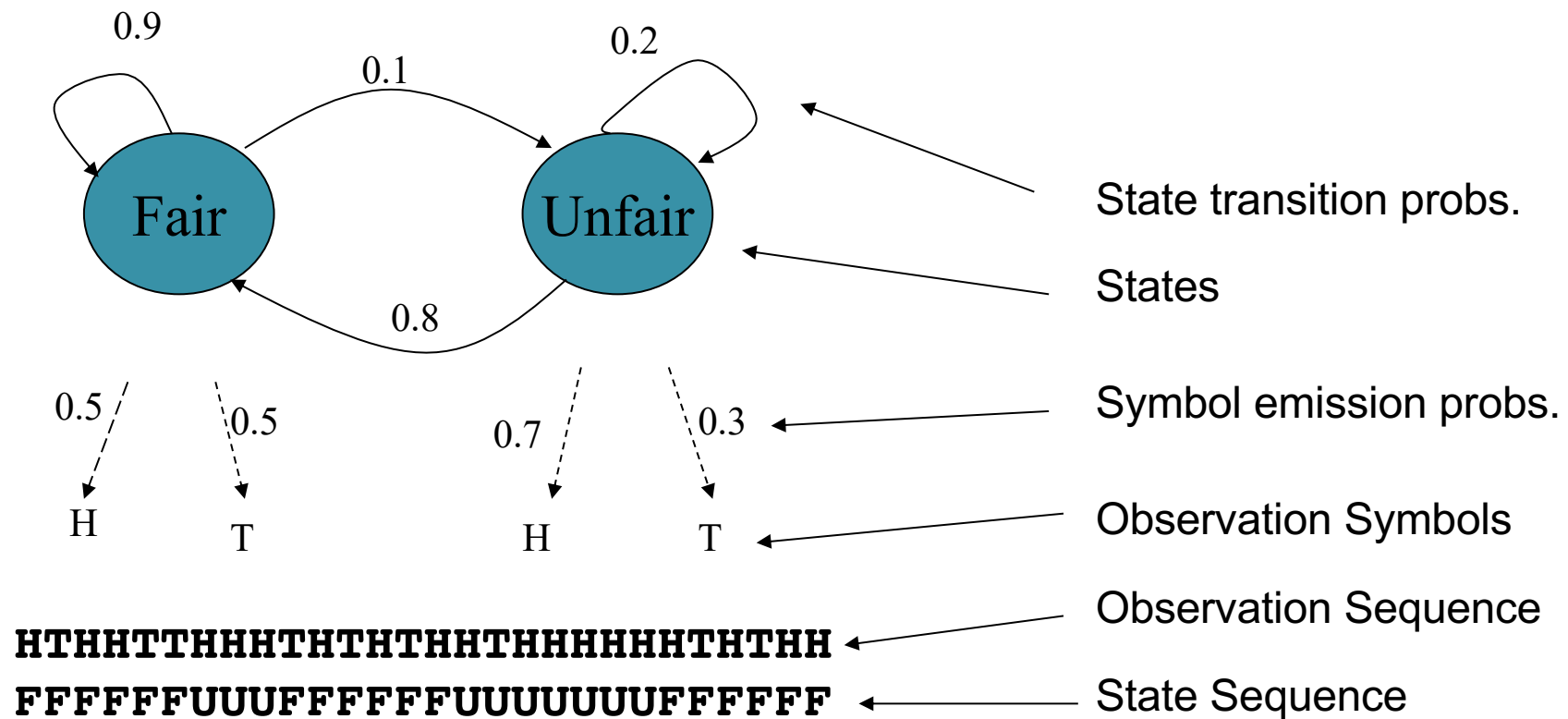


Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros

HMM Example - Casino Coin



Motivation: Given a sequence of H & Ts, can you tell at what times the casino cheated?

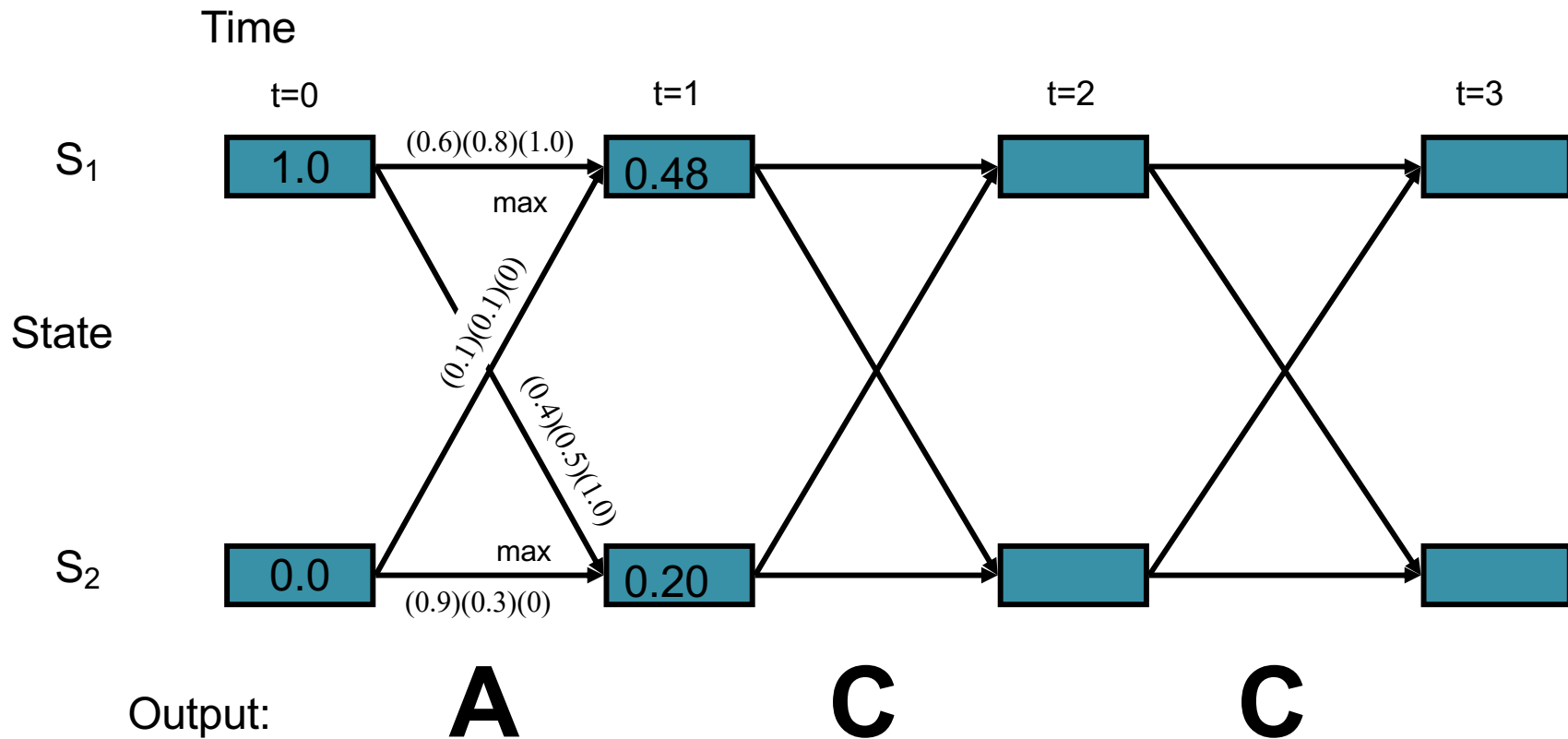
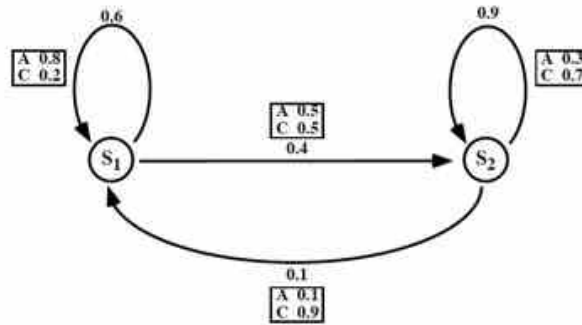
Solving the Decoding Problem: The Viterbi algorithm

- To solve the decoding problem (find the most likely sequence of states), we evaluate the Viterbi algorithm

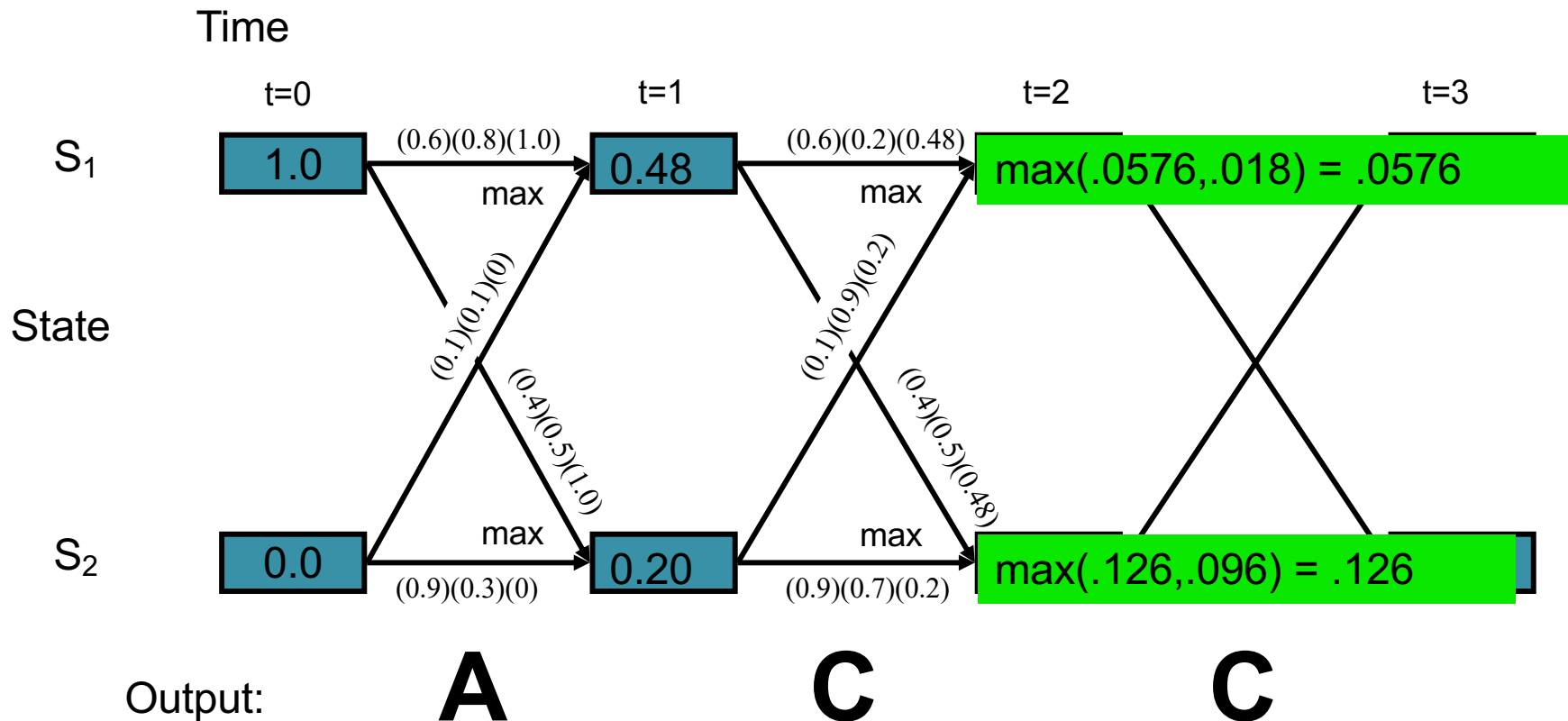
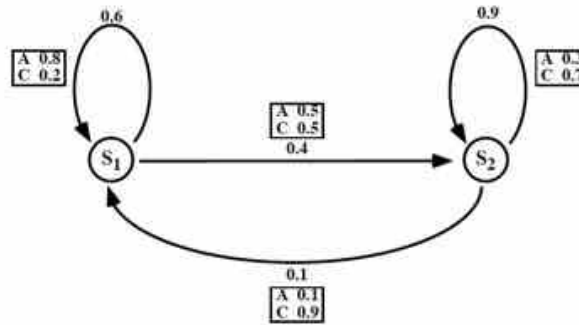
$$V_i(t) = \begin{cases} 0 & : t = 0 \wedge i \neq S_I \\ 1 & : t = 0 \wedge i = S_I \\ \max_j V_j(t-1) a_{ji} b_{ji}(y) & : t > 0 \end{cases}$$

Where $V_i(t)$ is the probability that the HMM is in state i after generating the sequence y_1, y_2, \dots, y_t , following the *most probable path* in the HMM

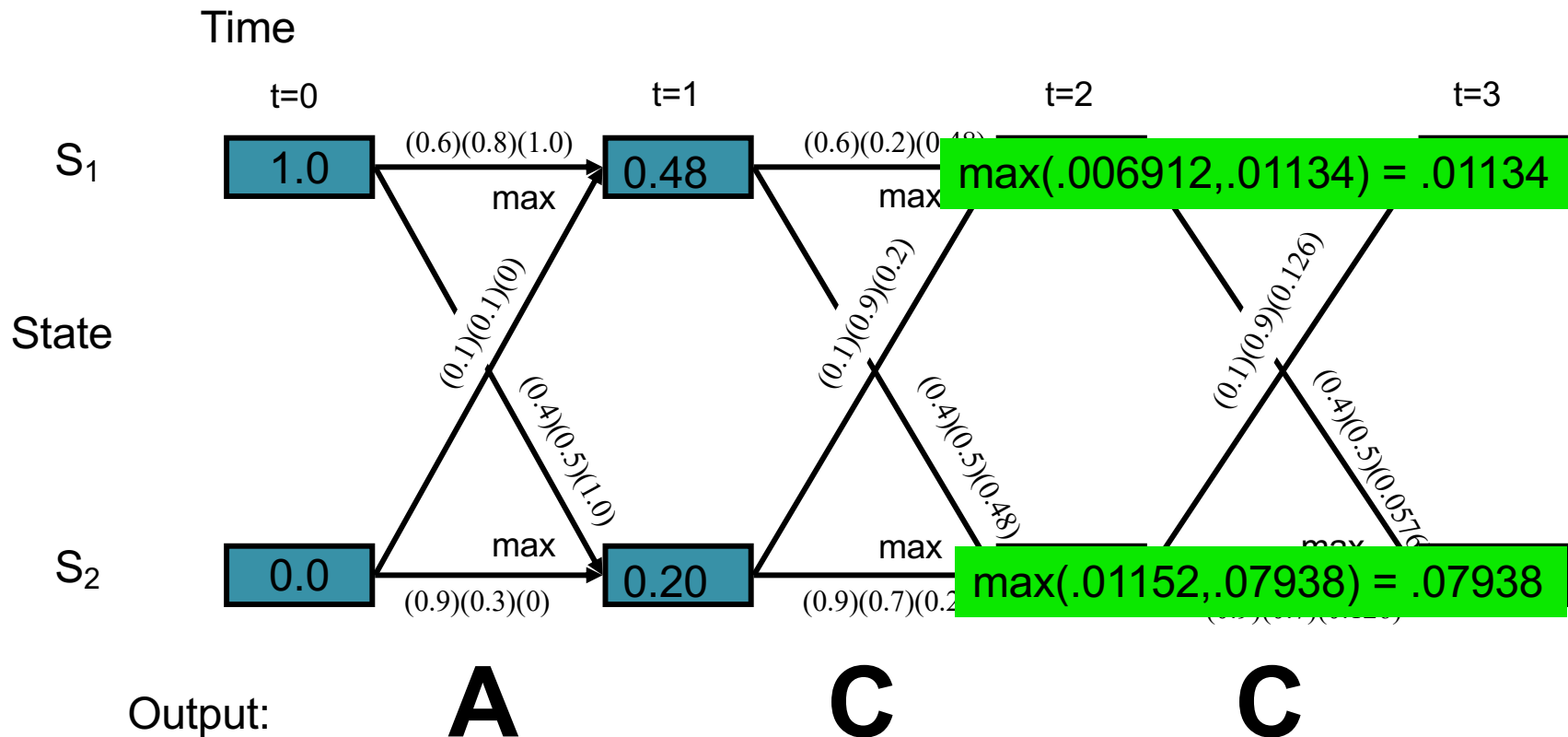
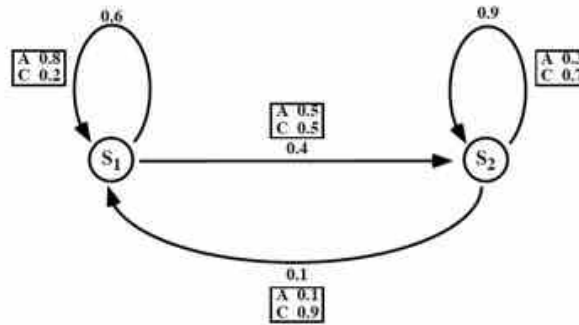
A trellis for the Viterbi Algorithm



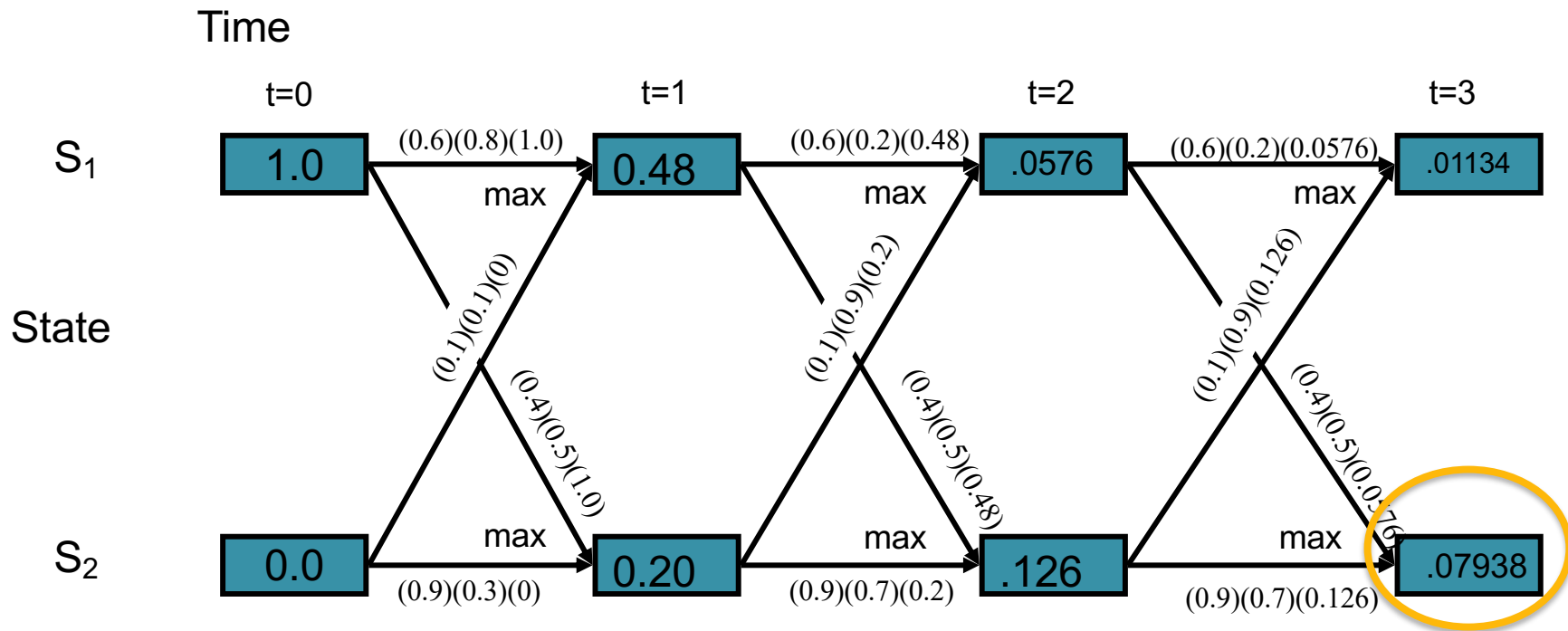
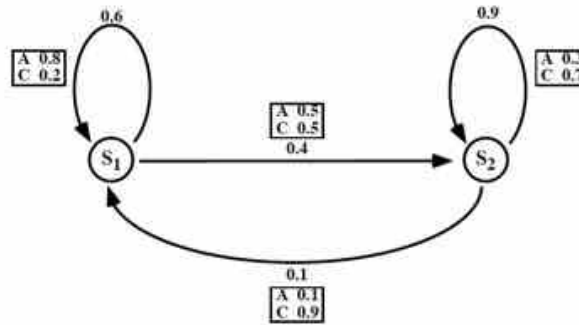
A trellis for the Viterbi Algorithm



A trellis for the Viterbi Algorithm

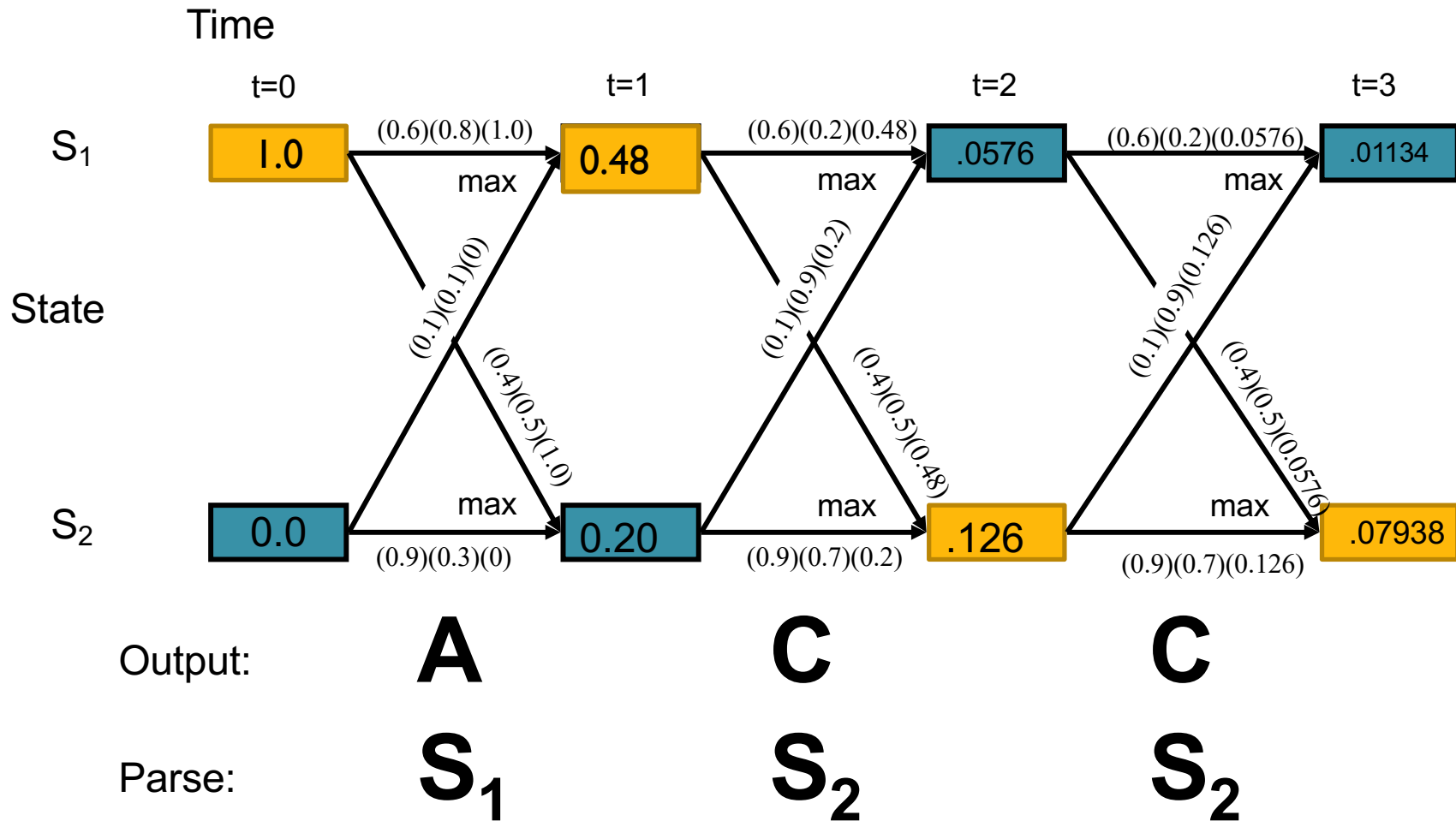


A trellis for the Viterbi Algorithm



S2 is final state → the most probable sequence of states has a 7.9% probability

A trellis for the Viterbi Algorithm



GlimmerHMM architecture

