



ALGORITHMS OF BIOINFORMATICS

2 Hidden Messages

23 October 2025

Prof. Dr. Sebastian Wild

2.1 Biology Big Picture

Biology

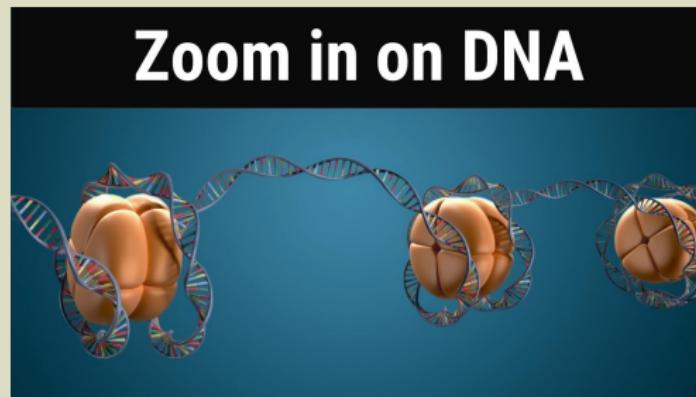
- ▶ *biology* = the scientific study of *living* things
 - ▶ originally *naturalists*: individual people manually **observing** plants and animals
e.g., *Darwin's finches*
 - ▶ gradually more scientific: controlled experiments, isolated mechanisms
e.g., *Mendel's inheritance experiments on peas*
 - ▶ gradually more focus on molecular/chemical mechanisms: microscopes, biochemistry

Biology

- ▶ *biology* = the scientific study of *living* things
 - ▶ originally *naturalists*: individual people manually **observing** plants and animals
e.g., *Darwin's finches*
 - ▶ gradually more scientific: controlled experiments, isolated mechanisms
e.g., *Mendel's inheritance experiments on peas*
 - ▶ gradually more focus on molecular/chemical mechanisms: microscopes, biochemistry
- ▶ now clear: fundamental mechanisms (and origins!) of life are microscopic
- ~~ fundamental mechanisms to be found in *molecular biology*

Bioinformatics

- ▶ 20th Century: discovery of DNA and genes
 - ▶ DNA stores information about biomolecules in **discrete form**
human genome: 3.055 billion letter string over alphabet {A, C, G, T} (!)
 - ~~ genetic information can **copied** precisely
mutations are errors in the copying
 - ▶ double strands (backup!) and “coiling up” into chromosomes protects data
 - ▶ production of chemicals in living cells (**proteins**) is determined by **genes** (parts of DNA)



▶ Zoom in on DNA
<https://youtu.be/wZoz0rFluiw>

Bioinformatics

- ▶ 20th Century: discovery of DNA and genes
 - ▶ DNA stores information about biomolecules in **discrete form**
human genome: 3.055 billion letter string over alphabet {A, C, G, T} (!)
 - ~~ genetic information can **copied** precisely
mutations are errors in the copying
 - ▶ double strands (backup!) and “coiling up” into chromosomes protects data
 - ▶ production of chemicals in living cells (**proteins**) is determined by **genes** (parts of DNA)
- ~~ *Life itself has inherently computational components!* 🤖



Zoom in on DNA
<https://youtu.be/wZozOrFluiw>

Bioinformatics

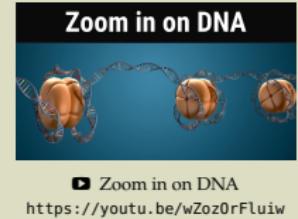
- ▶ 20th Century: discovery of DNA and genes
 - ▶ DNA stores information about biomolecules in **discrete form**
human genome: 3.055 billion letter string over alphabet {A, C, G, T} (!)
 - ~~ genetic information can **copied** precisely
mutations are errors in the copying
 - ▶ double strands (backup!) and “coiling up” into chromosomes protects data
 - ▶ production of chemicals in living cells (*proteins*) is determined by *genes* (parts of DNA)
- ~~ *Life itself has inherently computational components!* 🐾
- ~~ Computer science can contribute to the understanding these! ~~ *bioinformatics*



Zoom in on DNA
<https://youtu.be/wZozOrFluiw>

Bioinformatics

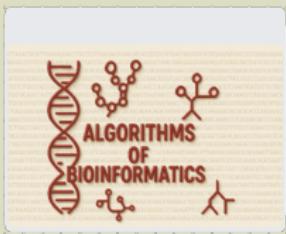
- ▶ 20th Century: discovery of DNA and genes
 - ▶ DNA stores information about biomolecules in **discrete form**
human genome: 3.055 billion letter string over alphabet {A, C, G, T} (!)
 - ~~ genetic information can **copied** precisely
mutations are errors in the copying
 - ▶ double strands (backup!) and “coiling up” into chromosomes protects data
 - ▶ production of chemicals in living cells (*proteins*) is determined by *genes* (parts of DNA)
 - ~~ *Life itself has inherently computational components!* 🐾
 - ~~ Computer science can contribute to the understanding these! ~~ *bioinformatics*
- ▶ But also: biology increasingly a data-centric field
 - ▶ much of knowledge discovery intrinsically reliant on computational analysis of collected data
 - ▶ e. g., reading the 3 billion letters of DNA is not possible with current lab techniques
 - ~~ use computers to puzzle it together (see *Sequencing Unit*)
 - ▶ “*in silico*” experiments



Collection of (more or less) Fun Sources

Collaborative Mindmap
on  infinity maps

- ▶ Share useful resources
- ▶ Structure knowledge hierarchically
- ▶ Link on Campuswire / ILIAS



Algorithms of Bioinformatics

BIOLOGY MINDMAP & SOURCES

Microbiology

The Origin of Life

Bioinformatics Lectures

Pop science

Cooperation

Microscopy to watch

A collage of various biology-related images and links, including a DNA helix, a tree of life, a book cover for 'A Brief History of Life on Earth', a YouTube channel thumbnail for 'Biology & Medicine', a course titled 'Microscopy', a book cover for 'SURVIVAL OF THE FITTEST', and a YouTube channel for 'microCOSMOS'.

*There's tons to learn,
new things discovered every day,
and it's about life itself!*

Molecular Biology 101

Molecular Biology (Britannica concise)

- ▶ concerned with chemical structures and processes of biological phenomena at the molecular level
- ▶ developed out of biochemistry, genetics, and biophysics
- ▶ particularly concerned with the study of **proteins**, nucleic acids, and enzymes

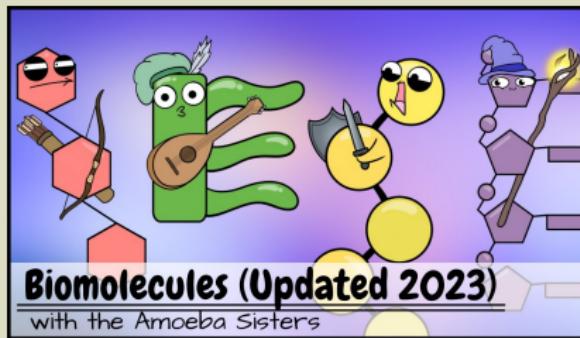
Molecular Biology 101

Molecular Biology (Britannica concise)

- ▶ concerned with chemical structures and processes of biological phenomena at the molecular level
- ▶ developed out of biochemistry, genetics, and biophysics
- ▶ particularly concerned with the study of **proteins**, nucleic acids, and enzymes

Biology = lots of terminology and names . . .

We will focus on mechanisms over terms, but a bit of context helps
let's make it at least whimsical (and maybe memorable)

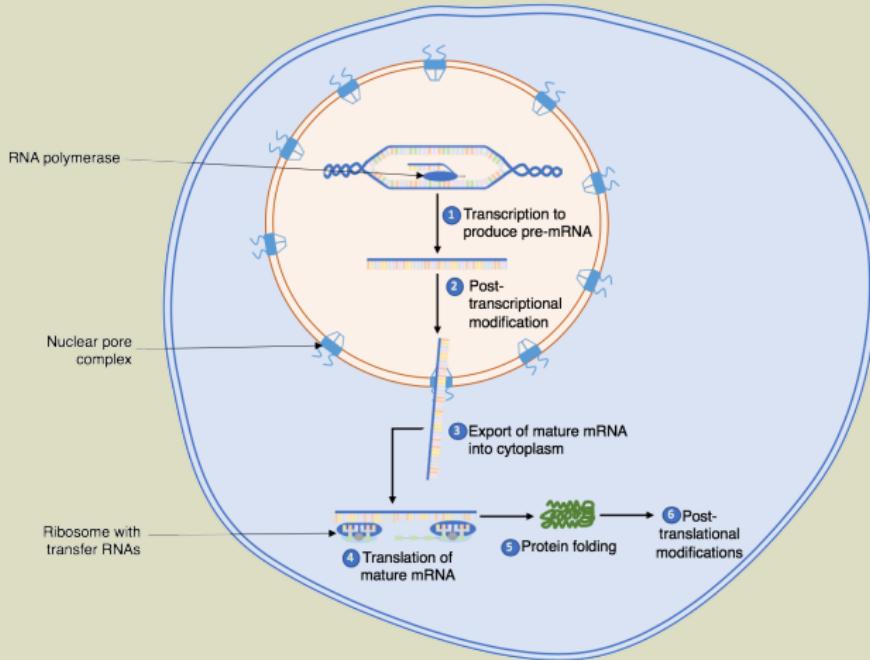


► Biomolecules (Updated 2023)
<https://youtu.be/1Dx7LDwINLU>

2.2 What are Genes?

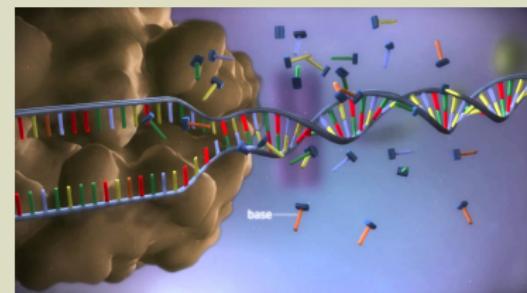
The Central Dogma of Molecular Biology

DNA makes RNA makes Protein



Protein Biosynthesis

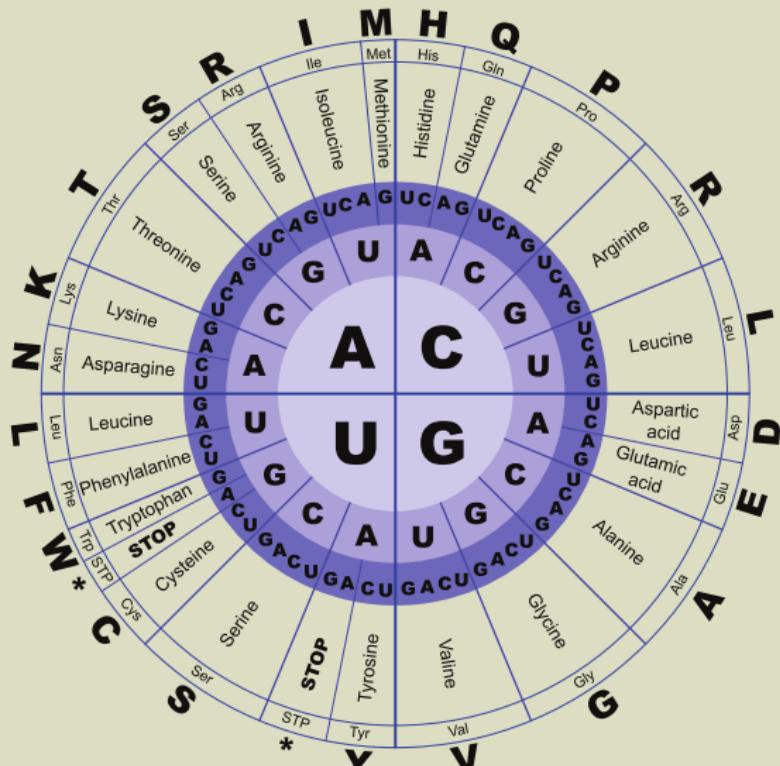
- mechanism to produce *protein* according to recipe stored in a *gene*



From DNA to protein - 3D
<https://youtu.be/gG7uCskU0rA>

https://commons.wikimedia.org/wiki/File:Summary_of_the_protein_biosynthesis_process.png

Genetic Code



Compeau & Pevzner, *Bioinformatics Algorithms*, Fig. 4.1
<https://cogniterra.org/lesson/29910/step/2?unit=22007>

- Within *ribosomes* (protein factories)

- ▶ translation
 - ▶ from RNA bases {A, C, G, U}
 - ▶ to amino acids (peptide)
{A, C, D, E, F, G, H, I, K, L,
M, N, P, Q, R, S, T, V, W, Y}
 - ▶ uses *transfer RNA*
"chemical finite state transducer"
 - ▶ *Genetic Code*:
3-base codons → amino acid

Inverse Codon Table

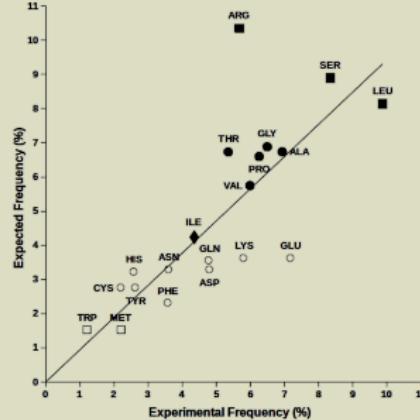
#Codons	Amino Acid (abbr.)	Codons
1	<i>Start</i>	> AUG
4	Ala	A A G G C U G C G A C G
2	Cys	C U G U U G C
2	Asp	D G A U G A C
2	Glu	E G A A G A G
2	Phe	F U U U U U C
4	Gly	G G Y G G C G G A G G G
2	His	H C A U C A C
3	Ile	I A U U A U C A U A
2	Lys	K A A A A A G
6	Leu	L C U U C U C U A C U G U U A U U G
1	Met	M A U G
2	Asn	N A A U A A C
4	Pro	P C C U C C C C C A C C G
2	Gln	Q C A A C A G
6	Arg	R C G U C G C C G A C G G A G A A G G
6	Ser	S U C U U C C U A C U G A G U A G C
4	Thr	T A C U A C C A C A A C G
4	Val	V G U U G U C G U A A G U
1	Trp	W U G G
2	Tyr	Y U A U U A C
3	<i>Stop</i>	< U A A U A G U G A
1	Sec	U (U G A)
1	Pyl	O (U A G)

Inverse Codon Table

#Codons		Amino Acid (abbr.)		Codons
1		Start	>	AUG
4		Ala	A	GCU GCC GCA GCG
2		Cys	C	UGU UGC
2		Asp	D	GAU GAC
2		Glu	E	GAA GAG
2		Phe	F	UUU UUC
4		Gly	G	GGU GGC GGA GGG
2		His	H	CAU CAC
3		Ile	I	AUU AUC AUA
2		Lys	K	AAA AAG
6		Leu	L	CUU CUC CUA CUG UUA UUG
1		Met	M	AUG
2		Asn	N	AAU AAC
4		Pro	P	CCU CCC CCA CCG
2		Gln	Q	CAA CAG
6		Arg	R	CGU CGC CGA CGG AGA AGG
6		Ser	S	UCU UCC UCA UCG AGU AGC
4		Thr	T	ACU ACC ACA ACG
4		Val	V	GUU GUC GUA GUG
1		Trp	W	UGG
2		Tyr	Y	UAU UAC
3		Stop	<	UAA UAG UGA
1		Sec	U	(UGA)
1		Pyl	O	(UAG)

Some amino acids have several codons
(most frequent amino acids receive strongest error protection!)

Amino Acid Frequencies in Human Proteins

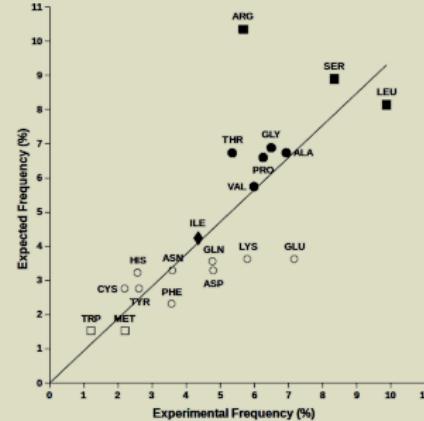


<https://doi.org/10.1371/journal.pone.0148174.g001>

Inverse Codon Table

#Codons		Amino Acid (abbr.)		Codons
1		Start	>	AUG
4		Ala	A	GCU GCC GCA GCG
2		Cys	C	UGU UGC
2		Asp	D	GAU GAC
2		Glu	E	GAA GAG
2		Phe	F	UUU UUC
4		Gly	G	GGU GGC GGA GGG
2		His	H	CAU CAC
3		Ile	I	AUU AUC AUA
2		Lys	K	AAA AAG
6		Leu	L	CUU CUC CUA CUG UUA UUG
1		Met	M	AUG
2		Asn	N	AAU AAC
4		Pro	P	CCU CCC CCA CCG
2		Gln	Q	CAA CAG
6		Arg	R	CGU CGC CGA CGG AGA AGG
6		Ser	S	UCU UCC UCA UCG AGU AGC
4		Thr	T	ACU ACC ACA ACG
4		Val	V	GUU GUC GUA GUG
1		Trp	W	UGG
2		Tyr	Y	UAU UAC
3		Stop	<	UAA UAG UGA
1		Sec	U	(UGA) ← Sometimes, stop codon UGA instead codes 21st amino acid Selenocystein...
1		Pyl	O	(UAG)

Amino Acid Frequencies in Human Proteins



<https://doi.org/10.1371/journal.pone.0148174.g001>

Some amino acids have several codons
(most frequent amino acids receive strongest error protection!)

But:

- ▶ non-ribosomal peptides (proteins not made according to central dogma)
- ▶ epigenetics (which genes are expressed)
- ▶ horizontal gene transfer (change genome during lifetime)
- ▶ retro viruses (inserts its own genes into host's genome!)
- ▶ proteins are also not the only active molecules (e. g., functional RNA)

Life finds a way . . . or a few dozen, just to be sure

2.3 Gene Detection

How can we find genes?

Recall: Gene = protein-coding region of DNA



Central options:

1. *ab initio* ("from the beginning"): just using the DNA
 - ▶ search for start (AUG) and stop codons (UAA, UAG, UGA) \rightsquigarrow open reading frame
 - ▶ search for promoter binding sites (docking station for transcription molecules)
 - ▶ bias of base frequencies in coding vs non-coding regions (*hidden Markov models*)
2. extrinsic methods: using additional (lab) data
 - ▶ e.g. sequencing messenger RNA from live cells (many more options)
 - ▶ comparison of genome to other species with known genes

Focus for today: Ab initio options

Why should there be any hope of finding hidden messages?

- ▶ Evolution!
 - ▶ Random mutations always at play
 - ▶ If functional part becomes dysfunctional, individual does not produce offspring
 - ▶ other parts might be subject to random modifications
- ~~ *signal*: property in a text that us unlikely to be present in random strings (noise)
- ~~ *noise / null model*: unused DNA is random

2.4 Waiting for Words

How big are genes?

UCSC Genome Browser on Human (GRCh38/hg38)

Move <<< << < > >> Zoom in 1.5x 3x 10x Base Zoom out 1.5x 3x 10x 100x

Multi-region chr16:89,919,211-89,919,373 163 bp | MC1R | Search [Examines](#)

chr16 (q24.3) 16p13.3 p13.2 13.13 13.12 13.11 16p12.3 16p12.2 16p12.1 16p11.2 16q11.2 16q12.1 16q12.2(q15) 16q21 16q22.1 16q22.2 16p23.1 q23.2 q23.3 16p24.1(q24.3)

Scale 50 bases → hg38

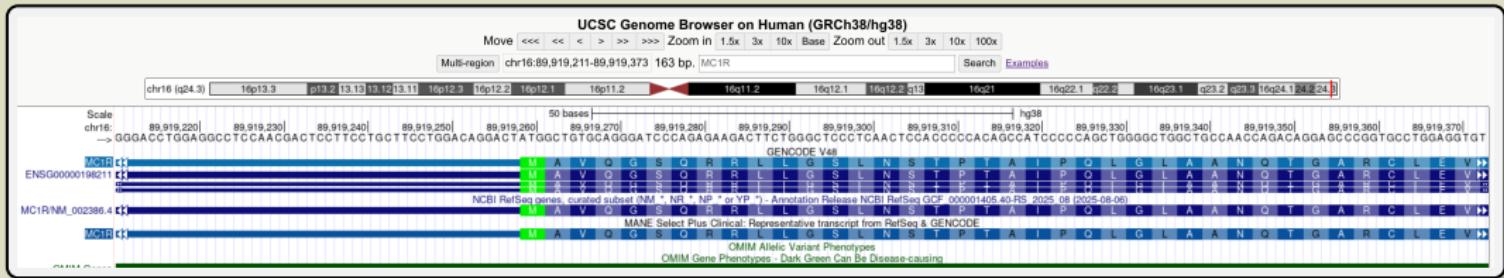
chr16:89,919,220 89,919,230 89,919,240 89,919,250 89,919,260 89,919,270 89,919,280 89,919,290 89,919,300 89,919,310 89,919,320 89,919,330 89,919,340 89,919,350 89,919,360 89,919,370

→ GGGACCTGGAAGCCTCCACGACTCTTCTGGCACAGGACTATGGCTGTGCAGGGATCCCAGAAAGACTCTGGCTCCCTGAACCTCCACCCCAACGCCATCCCCAGCTGGGGCTGGCTGCCAACCAAGACAGGCCCCGGTCCCTGGAGGTGT

ENSG00000198211 MC1R ENSG00000198211 MC1R NCBI RefSeq Genes, current build NCBI NM_1 NP_437777.1 Annotation Release NCBI RefSeq GCF_000001405.40 RS 2025 08 08 (2025-08-06)

MC1R_NM_002386.4 MC1R MC1R NM_002386.4 MC1R MANE Select Plus Clinical Representative transcript from RefSeq & GENCODE

OMIM Allelic Variant Phenotypes OMM Gene Phenotypes - Dark Green Can Be Disease-causing



How big are genes?

UCSC Genome Browser on Human (GRCh38/hg38)

Move <<< << < > >> Zoom in 1.5x 3x 10x Base Zoom out 1.5x 3x 10x 100x

Multi-region chr16:89,919,211-89,919,373 163 bp | MC1R | Search Examples

chr16 (q24.3) 16p13.3 p13.2 13.13 13.12 13.11 16p12.3 16p12.2 16p12.1 16p11.2 16q11.2 16q12.1 16q12.2(q15) 16q21 16p22.1 16p22.2 16p23.1 q23.2 q23.3 16p24.1(q24.3)

Scale 50 bases hg38 GENCODE V46

→ GGGACCTGAGGCCCTCACGACTCCTCTCGTTCCTGGACAGGACTATGGCTGTGCAGGGATCCAGAGAACCTCTGGGTCCTGACTCCA

ENSG00000198211 MC1R NCBI RefSeq Genes, current build NM_1 NP_0051405.4 RS 2025 08 08 (2025-08-06)

MC1R_NM_002386.4 MC1R MANE Select Plus Clinical Representative transcript from RefSeq & GENCODE

OMIM Allelic Variant Phenotypes OMM Gene Phenotypes - Dark Green Can Be Disease-causing

- ▶ only ~10% of human genome are genes
- ▶ length of (human) genes highly variable
 - ▶ shortest known gene (*U7 snRNA*) has only 63 bp
 - ▶ longest gene (*dystrophin*) over 2M bp
 - ▶ but: 99% of that are *introns* (cut out before translation)!
“split genes”
~~ transcription takes several hours (!)
 - ▶ more typical: ~20K bp

base pairs (DNA strands!)



How big are genes?

UCSC Genome Browser on Human (GRCh38/hg38)

Move <<< << < > >> Zoom in 1.5x 3x 10x Base Zoom out 1.5x 3x 10x 100x

Multi-region chr16:89,919,211-89,919,373 163 bp | MC1R | Search Examples

chr16 (q24.3) 16p13.3 p13.2 13.13 13.12 13.11 16p12.3 16p12.2 16p12.1 16p11.2 16q11.2 16q12.1 16q12.2(q15) 16q21 16p22.1 16p22.2 16p23.1 q23.2 q23.3 16p24.1(q24.3)

Scale 50 bases hg38 GENCODE V46

→ GGGACCTGAGGCCCTCACGACTCCTCTCGTTCCTGGACAGGACTATGGCTGTGCAGGGATCCAGAGAACCTCTGGCTCCCTGACTCCA

ENSG00000198211 MC1R NCBI RefSeq Genes, current build NM_1 NP_0051405.4 RS_2025_08_2025_08_06_01 MC1R_NM_002386.4 MC1R MANE Select Plus Clinical Representative transcript from RefSeq & GENCODE OMM Allelic Variant Phenotypes OMM Gene Phenotypes - Dark Green Can Be Disease-causing

NCBI RefSeq Genes, current build NM_1 NP_0051405.4 RS_2025_08_2025_08_06_01

MANE Select Plus Clinical Representative transcript from RefSeq & GENCODE

OMIM Allelic Variant Phenotypes

OMIM Gene Phenotypes - Dark Green Can Be Disease-causing

- ▶ only ~10% of human genome are genes
- ▶ length of (human) genes highly variable
 - ▶ shortest known gene (*U7 snRNA*) has only 63 bp
 - ▶ longest gene (*dystrophin*) over 2M bp
 - ▶ but: 99% of that are *introns* (cut out before translation)!
“split genes”
~~ transcription takes several hours (!)
 - ▶ more typical: ~20K bp
- ~~ base pairs (DNA strands!)



~~ Can we reliably distinguish genes from randomly occurring open reading frames?

Open Reading Frame Gene Detection

- ▶ **Random RNA Model:** String $D[0..N]$ generated i.i.d. uniformly
i.e., each $D[i] \stackrel{\mathcal{D}}{\sim} \text{Uniform}(\{\text{A, C, G, T}\})$
- ▶ **Random Open Reading Frame:** How many bp should we expect in **random RNA**
between occurrences of the start codon **ATG**
and **first** occurrence of any stop codon (**TAA, TAG, TGA**)?

(Recall: **U** in mRNA is **T** in DNA)

Clicker Question

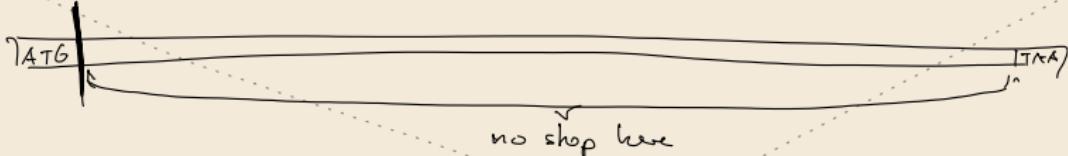


In a string of random (i.i.d. uniform) DNA, what is the expected length of an open reading frame?



→ *sli.do/cs594*

Back-of-the-envelope

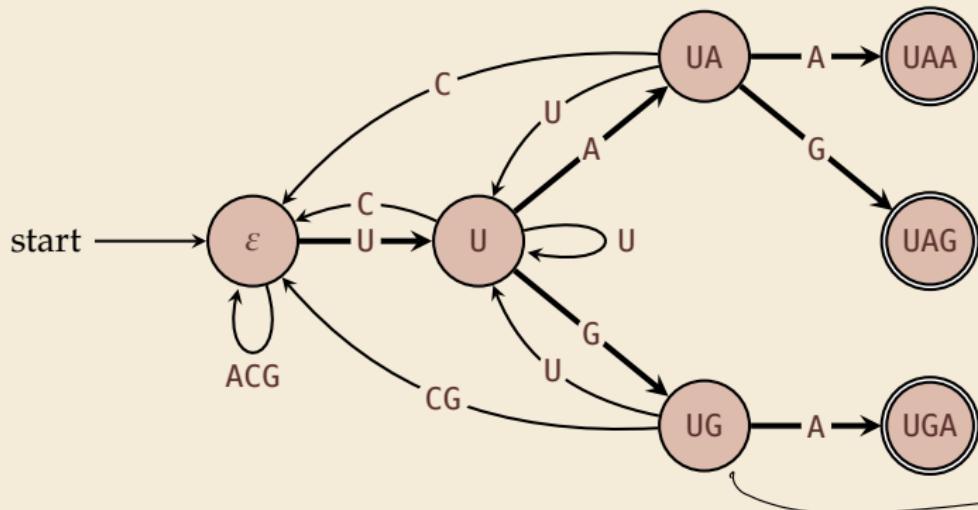


$$P \left[\boxed{?} \boxed{?} \boxed{?} \in \{ TAA, TGA, TAG \} \right] = \frac{3}{64} \approx 0.05$$

maybe ≈ 20 codons before stop

≈ 60 bp

Stop Codon automaton



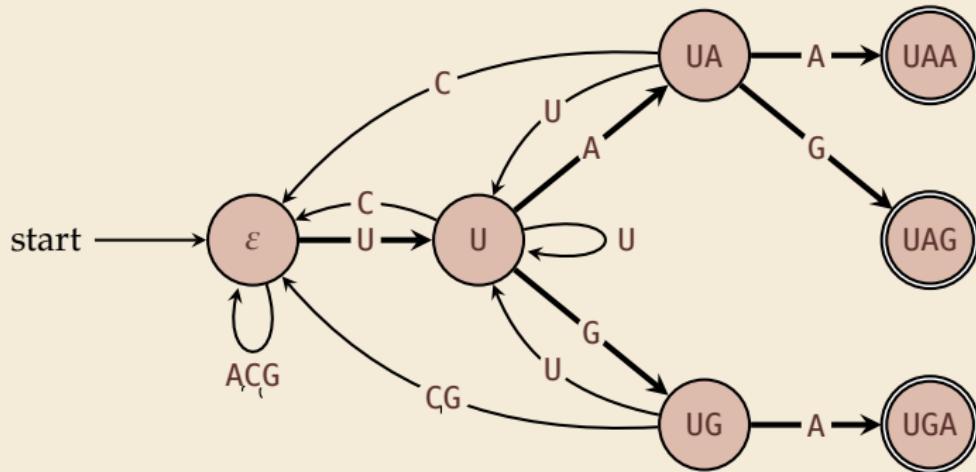
(Aho-Corasick
string-matching
automaton)



After seeing a start codon AUG, we accept the language of all strings that

- ▶ end with a stop codon **and**
- ▶ do not contain a stop codon earlier.

Stop Codon automaton

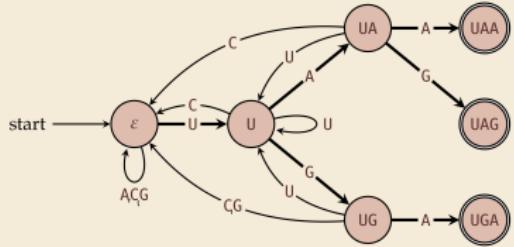


Expected number of
characters from state q
to some accepting state.
(We want q_ϵ)

$T_x :=$ expected # chars
from q_x to \textcircled{O}

After seeing a start codon **AUG**, we accept the language of all strings that

- ▶ end with a stop codon **and**
- ▶ do not contain a stop codon earlier.



$$T_{UAA} = T_{UAC} = T_{UCA} = \emptyset$$

$$T_{UA} = \frac{1}{4} \cdot T_{UAA} + \frac{1}{4} T_{UAG} + \frac{1}{4} T_U + \frac{1}{5} T_\varepsilon + 1$$

$$T_{UG} = \frac{1}{4} T_{UGA} + \frac{1}{4} T_U + \frac{1}{2} T_\varepsilon + 1$$

⋮

$$\bar{T}_\varepsilon = \frac{64}{3} = 21.\overline{3} \quad \ll 60$$

2.5 Probability Generating Functions

Probability Generating Functions

Expected values do not tell the full story . . . can we get at the distribution?

Probability Generating Functions

Expected values do not tell the full story . . . can we get at the distribution?

Definition 2.1 (pgf)

For $X \in \mathbb{N}_{\geq 0}$ a random variable, define its *probability generating function (pgf)* as

$\# \text{ chars}$ $G_X(z) = \sum_{k \geq 0} \mathbb{P}[X = k] \cdot z^k$



Probability Generating Functions

Expected values do not tell the full story . . . can we get at the *distribution*?

Definition 2.1 (pgf)

For $X \in \mathbb{N}_{\geq 0}$ a random variable, define its *probability generating function (pgf)* as

$$G_X(z) = \sum_{k \geq 0} \mathbb{P}[X = k] \cdot z^k$$

$$\mathbb{E}[X] = \sum_{k \geq 0} \mathbb{P}[X = k] \cdot k$$

$$\begin{aligned} G'_X(z) &= \sum_{k \geq 0} \underbrace{\frac{d}{dz} \mathbb{P}[X = k]}_{= \mathbb{P}[X = k] \cdot k} \cdot z^k \\ &= \mathbb{P}[X = k] \cdot k \cdot z^{k-1} \end{aligned}$$

Lemma 2.2 (Moments from pgf)

1. The expected value of X is $\mathbb{E}[X] = G'_X(1)$

2. The variance of X is $\text{Var}[X] = G''_X(1) + G'_X(1) - (G'_X(1))^2$

$$\frac{d^2}{dz^2} \mathbb{P}[X = k] \cdot z^k = \mathbb{P}[X = k] k \cdot (k-1) z^{k-2} \quad (k \geq 2)$$

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \sum_{k \geq 0} \mathbb{P}[X = k] (k - \mathbb{E}[X])^2 = \sum_{k \geq 0} \mathbb{P}[X = k] (k^2 - 2k \mathbb{E}[X] + \mathbb{E}[X]^2)$$

$$\begin{aligned}
 &= \sum_{k \geq 0} P\{X=k\} k^2 - 2 \underbrace{\mathbb{E}[X] \sum_{k \geq 0} P\{X=k\} k}_{\mathbb{E}[X]^2} + \underbrace{\mathbb{E}[X]^2 \sum_{k \geq 0} P\{X=k\}}_{=1} \\
 k^2 &= k(k-1) + k
 \end{aligned}$$

Example: Uniform Distribution

$U_n \stackrel{\mathcal{D}}{=} \text{Uniform}([0..n])$ (each value $u \in [0..n]$ with prob. $\frac{1}{n}$)

$$\rightsquigarrow \text{pgf: } U_n(z) = \frac{1}{n} (z^0 + z^1 + \cdots + z^{n-1})$$

Example: Uniform Distribution

$U_n \stackrel{\mathcal{D}}{=} \text{Uniform}([0..n])$ (each value $u \in [0..n]$ with prob. $\frac{1}{n}$)

$$\rightsquigarrow \text{pgf: } U_n(z) = \frac{1}{n} (z^0 + z^1 + \cdots + z^{n-1}) = \frac{1}{n} \frac{z^n - 1}{z - 1} \quad (n \geq 1)$$

Example: Uniform Distribution

$U_n \stackrel{\mathcal{D}}{=} \text{Uniform}([0..n])$ (each value $u \in [0..n]$ with prob. $\frac{1}{n}$)

$$\rightsquigarrow \text{pgf: } U_n(z) = \frac{1}{n} (z^0 + z^1 + \cdots + z^{n-1}) = \frac{1}{n} \frac{z^n - 1}{z - 1} \quad (n \geq 1)$$

► $\mathbb{E}[U_n] = U'_n(1) = \frac{n-1}{2}$

Example: Uniform Distribution

$U_n \stackrel{\mathcal{D}}{=} \text{Uniform}([0..n])$ (each value $u \in [0..n]$ with prob. $\frac{1}{n}$)

$$\rightsquigarrow \text{pgf: } U_n(z) = \frac{1}{n}(z^0 + z^1 + \cdots + z^{n-1}) = \frac{1}{n} \frac{z^n - 1}{z - 1} \quad (n \geq 1)$$

► $\mathbb{E}[U_n] = U'_n(1) = \frac{n-1}{2}$

► $\text{Var}[U_n] = U''_n(1) + U'_n(1) - U'_n(1)^2 = \frac{n^2 - 1}{12}$

Operations of pgfs

Lemma 2.3 (pgf of ind. r.v.)

Let $X, Y \in \mathbb{N}_{\geq 0}$ be *independent* random variables. Then $G_{X+Y}(z) = G_X(z) \cdot G_Y(z)$



$$\left(\sum_{k \geq 0} a_k z^k \right) \cdot \left(\sum_{k \geq 0} b_k z^k \right) = \sum_{k \geq 0} z^k \sum_{e=0}^k a_e b_{k-e}$$

convolution

Operations of pgfs

Lemma 2.3 (pgf of ind. r.v.)

Let $X, Y \in \mathbb{N}_{\geq 0}$ be *independent* random variables. Then $G_{X+Y}(z) = G_X(z) \cdot G_Y(z)$



- ▶ For $X_i \stackrel{\mathcal{D}}{=} B(p)$ *independent* (1 with prob. p , 0 otherwise)
we have $G_{X_i}(z) = pz + (1-p)z^0 = p(z-1) + 1$
- ▶ $Y = X_1 + \cdots + X_n \stackrel{\mathcal{D}}{=} \text{Bin}(n, p)$
we have $G_Y(z) = \prod_{i=1}^n G_{X_i}(z) = (pz + 1 - p)^n$

Operations of pgfs

Lemma 2.3 (pgf of ind. r.v.)

Let $X, Y \in \mathbb{N}_{\geq 0}$ be *independent* random variables. Then $G_{X+Y}(z) = G_X(z) \cdot G_Y(z)$



- For $X_i \stackrel{\mathcal{D}}{=} \text{B}(p)$ (1 with prob. p , 0 otherwise)
we have $G_{X_i}(z) = pz + (1-p)z^0 = p(z-1) + 1$

- $Y = X_1 + \cdots + X_n \stackrel{\mathcal{D}}{=} \text{Bin}(n, p)$
we have $G_Y(z) = \prod_{i=1}^n G_{X_i}(z) = (pz + 1 - p)^n$

$$\rightsquigarrow \mathbb{E}[Y] = G'_Y(1) = p \cdot n(pz + 1 - p)^{n-1} \Big|_{z=1} = np$$

Operations of pgfs

Lemma 2.3 (pgf of ind. r.v.)

Let $X, Y \in \mathbb{N}_{\geq 0}$ be *independent* random variables. Then $G_{X+Y}(z) = G_X(z) \cdot G_Y(z)$



- For $X_i \stackrel{\mathcal{D}}{=} \text{B}(p)$ (1 with prob. p , 0 otherwise)

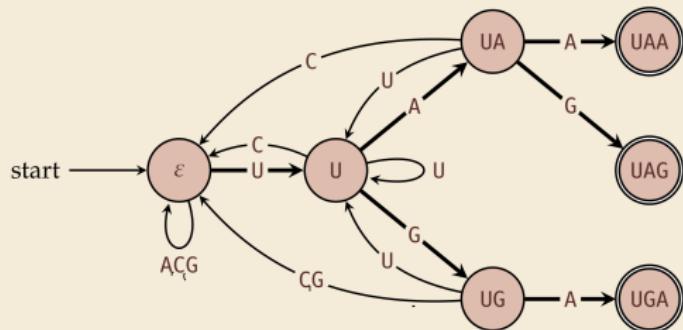
we have $G_{X_i}(z) = pz + (1-p)z^0 = p(z-1) + 1$

- $Y = X_1 + \dots + X_n \stackrel{\mathcal{D}}{=} \text{Bin}(n, p)$

we have $G_Y(z) = \prod_{i=1}^n G_{X_i}(z) = (pz + 1 - p)^n$

$$\rightsquigarrow \mathbb{E}[Y] = G'_Y(1) = p \cdot n(pz + 1 - p)^{n-1} \Big|_{z=1} = np$$

- $\text{Var}[Y] = G''_Y(1) + G'_Y(1) - G'_Y(1)^2 = p^2 n(n-1) + np - (np)^2 = np - np^2 = np(1-p)$



X_p = # steps from state p to $\textcircled{0}$

$G_p(z) = \text{pgf of } X_p$

$G_{\text{UA}}(z) = 1$

$$G_{\text{U}_A}(z) = \frac{1}{4} \cdot z \cdot G_{\text{UAA}}(z) + \frac{1}{4} z G_{\text{UAG}}(z)$$

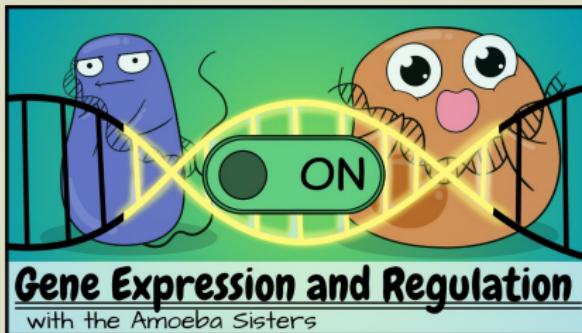
$$+ \frac{1}{4} z \cdot G_U(z)$$

$$+ \frac{1}{4} z \cdot G_\varepsilon(z)$$

2.6 Motif finding

Gene regulation

- ▶ For gene expression, *RNA polymerase* needs to bind to DNA at beginning of a gene (to start transcription of gene in DNA into messenger RNA)
- ▶ *promoter* molecule can **stop transcription** by binding to this start
 - ~~ RNA polymerase can't bind ~~ no mRNA created ~~ no protein made
 - ~~ *negative control*
- ▶ can also have promoters **enable** or **encourage** transcription ~~ *positive control*



► Gene Expression and Regulation
<https://youtu.be/ebIpkw3XapE>

Clock Genes in Plants

- ▶ most living beings have a *circadian rhythm*
Who controls that? "roughly 24h" (more details in Chapter 2 of *Bioinformatics Algorithms*)
- ▶ in plants, negative feedback loop of 3 proteins
 1. **TOC1 promotes** expression of **LHY** and **CCA1**
 2. sunlight triggers transcription **LHY** abd **CCA1**
 3. **LHY** abd **CCA1 repress** expression of **TOC1**
 4. without sunlight, **LHY** abd **CCA1** production diminishes
 5. **TOC1** no longer blocked, can accumulate at night
 6. **TOC1** triggers expression of **LHY** and **CCA1** (ahead of light!)
- ▶ **TOC1**, **CCA1**, and **LHY** turn other genes on or off (*promoters*)

Clock Genes in Plants

- ▶ most living beings have a *circadian rhythm*
Who controls that? "roughly 24h" (more details in Chapter 2 of *Bioinformatics Algorithms*)
- ▶ in plants, negative feedback loop of 3 proteins
 1. **TOC1 promotes** expression of **LHY** and **CCA1**
 2. sunlight triggers transcription **LHY** abd **CCA1**
 3. **LHY** abd **CCA1 repress** expression of **TOC1**
 4. without sunlight, **LHY** abd **CCA1** production diminishes
 5. **TOC1** no longer blocked, can accumulate at night
 6. **TOC1** triggers expression of **LHY** and **CCA1** (ahead of light!)
- ▶ **TOC1**, **CCA1**, and **LHY** turn other genes on or off (*promoters*)
 - ~~ genes with day/night rhythm should have **repeated binding sites** for **TOC1/CCA1/LHY!**
same substring
↓
 - ~~ called a **motif**

Motif Finding

Typical complication in bioinformatics: Nothing is exact . . .

Motif Finding

Typical complication in bioinformatics: Nothing is exact ...

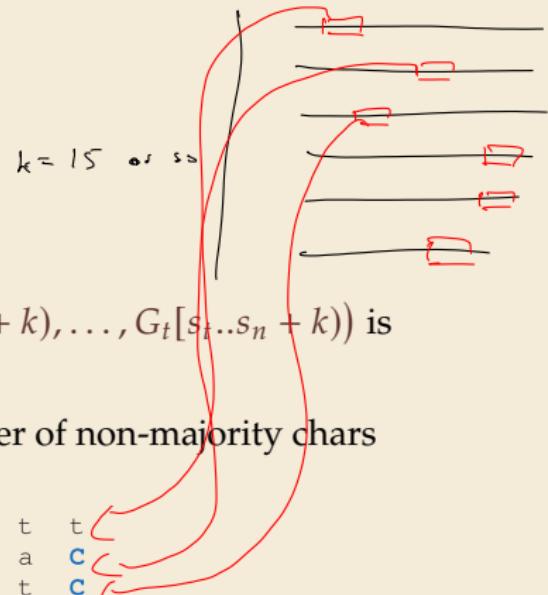
CONSENSUS PATTERN PROBLEM

- ▶ **Given:** Collection of strings $G_1, \dots, G_t \in \Sigma^n$, integer k
- ▶ **Goal:** Offsets $s_1, \dots, s_t \in [0..n - k]$ such that $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k])$ is minimized
- ▶ $d_H(T_1, \dots, T_n)$ = “total Hamming distance” = total number of non-majority chars

d _H motif score												
T	C	G	G	G	G	g	T	T	T	t	t	t
c	c	G	G	t	G	A	c	T	T	a	C	
a	C	G	G	G	G	A	T	T	T	t	C	
T	t	G	G	G	G	A	c	T	T	t	t	t
a	a	G	G	G	G	A	c	T	T	C	C	
T	t	G	G	G	G	A	c	T	T	C	C	
T	C	G	G	G	G	A	T	T	c	a	t	
T	C	G	G	G	G	A	T	T	c	C	t	
T	a	G	G	G	G	A	a	c	T	a	C	
T	C	G	G	G	t	A	T	a	a	C	C	

SCORE(Motifs) $3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30$

Compeau & Pevzner, Bioinformatics Algorithms, Fig. 2.2
<https://cogniterra.org/lesson/29868/step/2?unit=21966>



Median String

Motif with consensus and profile

	T	C	G	G	G	G	g	T	T	T	t	t	3
c	C	G	G	t	G	A	c	T	T	T	a	C	
a	C	G	G	G	G	A	T	T	T	T	t	C	
T	t	G	G	G	G	A	c	T	T	T	t	t	
a	a	G	G	G	G	A	c	T	T	T	C	C	
T	t	G	G	G	G	A	c	T	T	T	C	C	
T	C	G	G	G	G	A	T	T	T	c	a	t	
T	C	G	G	G	G	A	T	T	T	c	C	t	
T	a	G	G	G	G	A	a	c	T	a	C		
T	C	G	G	G	t	A	T	a	a	C	C	C	

► Equivalently:

$$d_H(T_1, \dots, T_t) = \sum_{i=1}^t d_H(\bar{T}, T_i)$$

for the *consensus string* \bar{T} :

for all $j \in [0..n]$:

$$\bar{T}[j] = \text{majority}(T_1[j], \dots, T_t[j])$$

$$\text{SCORE}(\text{Motifs}) = 3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30$$

A:	2	2	0	0	0	0	9	1	1	1	3	0
C:	1	6	0	0	0	0	0	4	1	2	4	6
G:	0	0	10	10	9	9	1	0	0	0	0	0
T:	7	2	0	0	1	1	0	5	8	7	3	4

A:	.2	.2	0	0	0	0	.9	.1	.1	.1	.3	0
C:	.1	.6	0	0	0	0	0	.4	.1	.2	.4	.6
G:	0	0	1	1	.9	.9	.1	0	0	0	0	0
T:	.7	.2	0	0	.1	.1	0	.5	.8	.7	.3	.4

$$\text{CONSENSUS}(\text{Motifs}) = \text{T C G G - G G A T T T C C}$$



Bad News

- ▶ CONSENSUS PATTERN PROBLEM is NP-hard (even for binary alphabet)



Elias: *Settling the Intractability of Multiple Alignment*, J. of Computational Biology 2006

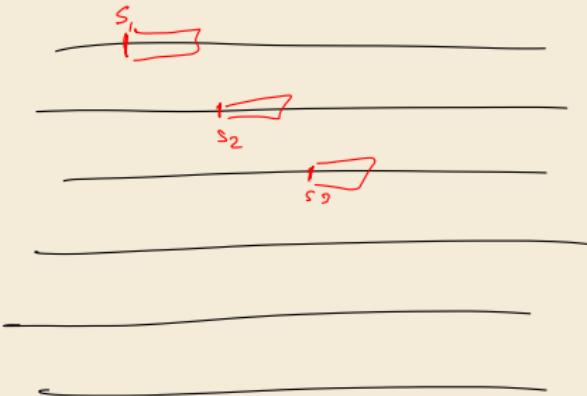
- ▶ even W[1]-hard for parameter ℓ (# strings) unbounded alphabet size

Brute-Force Options

There are two brute-force options:

1. Try all combinations of starting indices

- ▶ Each $s_i \in [0..n - k]$ \rightsquigarrow search space $(n - k + 1)^t$
 - ▶ Computing score for each is effort $O(t \cdot k)$
- \rightsquigarrow Total cost $O(n^t k)$ $(k \ll n, t)$

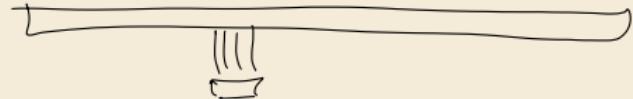


Brute-Force Options

There are two brute-force options:

1. Try all combinations of starting indices

- ▶ Each $s_i \in [0..n - k]$ \rightsquigarrow search space $(n - k + 1)^t$
 - ▶ Computing score for each is effort $O(t \cdot k)$
- \rightsquigarrow Total cost $O(n^t k)$ ($k \ll n, t$)



2. Try all consensus strings.

- ▶ Try all $\bar{T} \in \Sigma^k$ \rightsquigarrow search space σ^k (for $\sigma = |\Sigma|$)
- ▶ for each, \bar{T} and each string G_i , find best s_i $\rightsquigarrow t \cdot (n - k + 1)$ options

Note: Crucial that for given \bar{T} , scores from G_i are independent

- \rightsquigarrow Total cost $O(\sigma^k t n)$ ($k \ll n, t$)

Brute-Force Options

There are two brute-force options:

1. Try all combinations of starting indices

- ▶ Each $s_i \in [0..n - k]$ \rightsquigarrow search space $(n - k + 1)^t$
 - ▶ Computing score for each is effort $O(t \cdot k)$
- \rightsquigarrow Total cost $O(n^t t k)$ $(k \ll n, t)$

2. Try all consensus strings.

- ▶ Try all $\bar{T} \in \Sigma^k$ \rightsquigarrow search space σ^k (for $\sigma = |\Sigma|$)
- ▶ for each, \bar{T} and each string G_i , find best s_i $\rightsquigarrow t \cdot (n - k + 1)$ options

Note: Crucial that for given \bar{T} , scores from G_i are independent

- \rightsquigarrow Total cost $O(\sigma^k t n)$ $(k \ll n, t)$

Better, but still not feasible for $t \geq 15$ (or so) and $\sigma = 4 \dots$

2.7 Local search heuristics

Heuristic motif finding

*Exact solutions for CONSENSUS PATTERN PROBLEM
seem out of reach.*

↝ Give up optimality guarantee.

heuristic

+

approximation



Greedy Incremental

```

1 procedure greedyMotif( $G_1, \dots, G_t, k$ )
2    $s_1^*, \dots, s_t^* := 0$  // best so far
3   for  $s_1 := 0, \dots, n - k$  // try all  $s_1$ 
4     for  $i := 2, \dots, t$ 
5       Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k]$  for  $j \in [1..i)$ 
6        $s_i := \arg \max_s \mathbb{P}[G_i[s..s + k] | P]$  // most-likely next row
7       if  $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k]) < d_H(G_1[s_1^*..s_1^* + k], \dots, G_t[s_t^*..s_t^* + k])$ 
8          $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // better
9   return  $s_1^*, \dots, s_t^*$ 

```

$$\mathbb{P} \left[\text{AAC} \mid \begin{array}{|c|c|c|} \hline & A & C \\ \hline A & 0.5 & 0.2 & 0.1 \\ \hline C & 0.2 & 0.4 & 0.2 \\ \hline G & 0.2 & 0.1 & 0.2 \\ \hline T & 0.1 & 0.1 & 0.1 \\ \hline \end{array} \right] = 0.5 \cdot 0.2 \cdot 0.2$$



La Place

start counting at 1 instead of 0
removes 0s from profile

Motif with consensus and profile

Motifs	T	C	G	G	G	g	T	T	T	t	t
c	C	C	G	G	t	G	A	c	T	T	a
a	C	G	G	G	G	A	T	T	T	t	C
T	t	G	G	G	G	A	c	T	T	t	t
a	a	G	G	G	G	A	c	T	T	C	C
T	t	G	G	G	G	A	c	T	T	C	C
T	C	G	G	G	G	A	T	T	c	a	t
T	C	G	G	G	G	A	T	T	c	C	t
T	a	G	G	G	G	A	a	c	T	a	C
T	C	G	G	G	t	A	T	a	a	C	C

SCORE(Motifs)

$$3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30$$

A: 2 2 0 0 0 0 9 1 1 1 3 0

COUNT(Motifs) C: 1 6 0 0 0 0 0 4 1 2 4 6

G: 0 0 10 10 9 9 1 0 0 0 0 0

T: 7 2 0 0 1 1 0 5 8 7 3 4

PROFILE(Motifs) A: .2 .2 (0) (0) (0) (5) .9 .1 .1 .1 .3 (0)

C: .1 .6 (0) (0) (0) (0) 0 .4 .1 .2 .4 .6

G: (0) (0) 1 1 .9 .9 .1 (0) (0) (0) (0) 0

T: .7 .2 (0) (0) 1 .1 (0) 5 .8 .7 3 .4

CONSENSUS(Motifs)

T C G G - G G A T T T C C



www.cs.tufts.edu/comp/bioinfo/logo.html

Greedy Incremental

```
1 procedure greedyMotif( $G_1, \dots, G_t, k$ )
2    $s_1^*, \dots, s_t^* := 0$  // best so far
3   for  $s_1 := 0, \dots, n - k$  // try all  $s_1$ 
4     for  $i := 2, \dots, t$ 
5       Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k]$  for  $j \in [1..i)$ 
6        $s_i := \arg \max_s \mathbb{P}[G_i[s..s + k] \mid P]$ 
7       if  $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k]) < d_H(G_1[s_1^*..s_1^* + k], \dots, G_t[s_t^*..s_t^* + k])$ 
8          $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // better
9   return  $s_1^*, \dots, s_t^*$ 
```

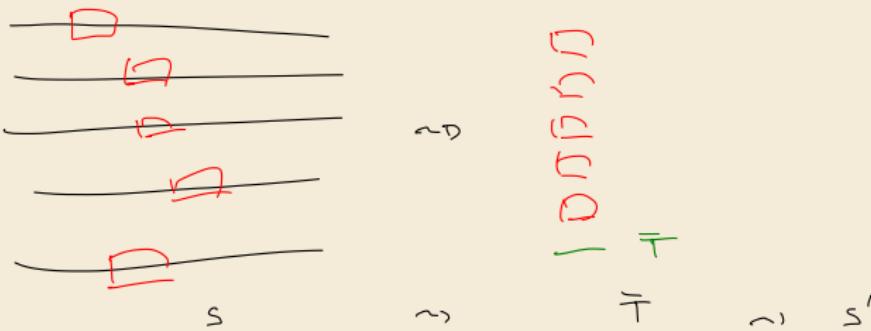
👍 deterministic

👎 highly sensitive to order of genomes ...

👎 easy to get stuck in local optimum (wrt to order)

Hill Climbing

```
1 procedure randomLocalSearch( $G_1, \dots, G_t, k$ )
2   Randomly choose  $s_1, \dots, s_t \in [0..n - k]$ 
3    $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // Remember best so far
4   repeat forever
5     Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k)$  for  $j \in [1..i)$ 
6     for  $i := 1, \dots, t$ :
7        $s_i := \arg \max_s \mathbb{P}[G_i[s..s + k) \mid P]$ 
8       if  $d_H(G_1[s_1..s_1 + k), \dots, G_t[s_t..s_t + k)) < d_H(G_1[s_1^*..s_1^* + k), \dots, G_t[s_t^*..s_t^* + k))$ 
9          $s_1^*, \dots, s_t^* := s_1, \dots, s_t$ 
10    else
11      return  $s_1^*, \dots, s_t^*$ 
```



Hill Climbing

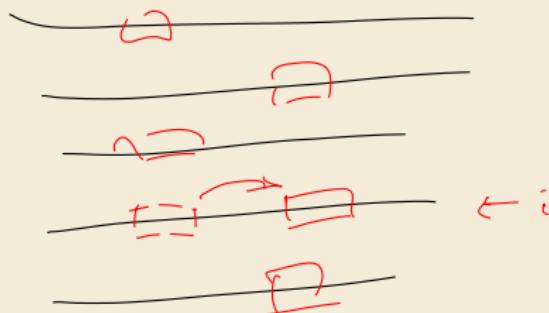
```
1 procedure randomLocalSearch( $G_1, \dots, G_t, k$ )
2   Randomly choose  $s_1, \dots, s_t \in [0..n - k]$ 
3    $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // Remember best so far
4   repeat forever
5     Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k]$  for  $j \in [1..i)$ 
6     for  $i := 1, \dots, t$ :
7        $s_i := \arg \max_s \mathbb{P}[G_i[s..s + k] \mid P]$ 
8       if  $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k]) < d_H(G_1[s_1^*..s_1^* + k], \dots, G_t[s_t^*..s_t^* + k])$ 
9          $s_1^*, \dots, s_t^* := s_1, \dots, s_t$ 
10    else
11      return  $s_1^*, \dots, s_t^*$ 
```

- 👍 deterministic for a given starting point
- 👍 always terminates in local optimum
- 👎 must be repeated many times to not be stuck in bad local optimum

\Rightarrow repeat
and maybe
stop hill climb
early

Gibbs Sampler

```
1 procedure gibbsSampler( $G_1, \dots, G_t, k, R$ )
2   Randomly choose  $s_1, \dots, s_t \in [0..n - k]$ 
3    $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // Remember best so far
4   repeat  $R$  times:
5      $i := \text{random } [1..t]$ 
6     Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k]$  for  $j \in [t] \setminus \{i\}$ 
7      $s_i := \text{random in } [0..n - k] \text{ w/p } \propto \mathbb{P}[G_i[s..s + k] | P]$ 
8     if  $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k]) < d_H(G_1[s_1^*..s_1^* + k], \dots, G_t[s_t^*..s_t^* + k])$ 
9        $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // better
10  return  $s_1^*, \dots, s_t^*$ 
```



Gibbs Sampler

```
1 procedure gibbsSampler( $G_1, \dots, G_t, k, R$ )
2     Randomly choose  $s_1, \dots, s_t \in [0..n - k]$ 
3      $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // Remember best so far
4     repeat  $R$  times:
5          $i := \text{random } [1..t]$ 
6         Compute profile  $P[0..k]$  from  $G_j[s_j..s_j + k]$  for  $j \in [t] \setminus \{i\}$ 
7          $s_i := \text{random in } [0..n - k] \text{ w/p } \propto \mathbb{P}[G_i[s..s + k] | P]$ 
8         if  $d_H(G_1[s_1..s_1 + k], \dots, G_t[s_t..s_t + k]) < d_H(G_1[s_1^*..s_1^* + k], \dots, G_t[s_t^*..s_t^* + k])$ 
9              $s_1^*, \dots, s_t^* := s_1, \dots, s_t$  // better
10    return  $s_1^*, \dots, s_t^*$ 
```

👍 Less prone to get stuck in local optima

👎 still no performance guarantee