



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
- ▶ 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)
- ~~ *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



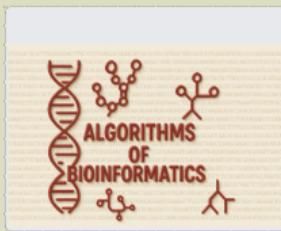
Zoom in on DNA
<https://youtu.be/wZoz0rFluiw>

Collection of (more or less) Fun Sources

Collaborative Mindmap
on  infinity maps

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

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



Algorithms of Bioinformatics

BIOLOGY MINDMAP & SOURCES

Microbiology



The Origin of Life



Bioinformatics Lectures



Pop science



Cooperation



Microscopy to watch



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)

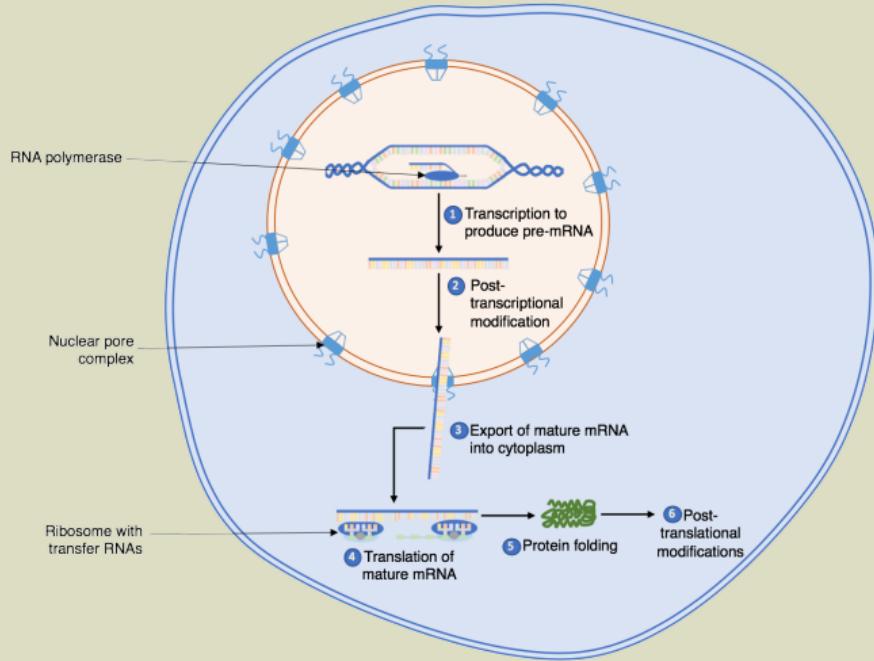


► Biomolecules (Updated 2023)
<https://youtu.be/1Dx7LDwINLU>

2.2 What are Genes?

The Central Dogma of Molecular Biology

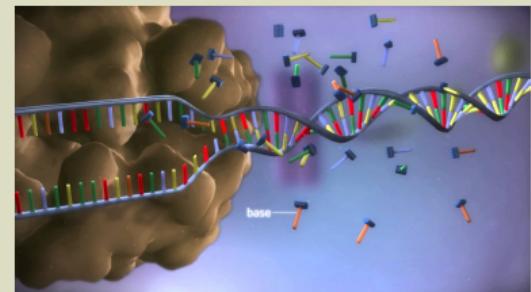
DNA makes RNA makes Protein



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

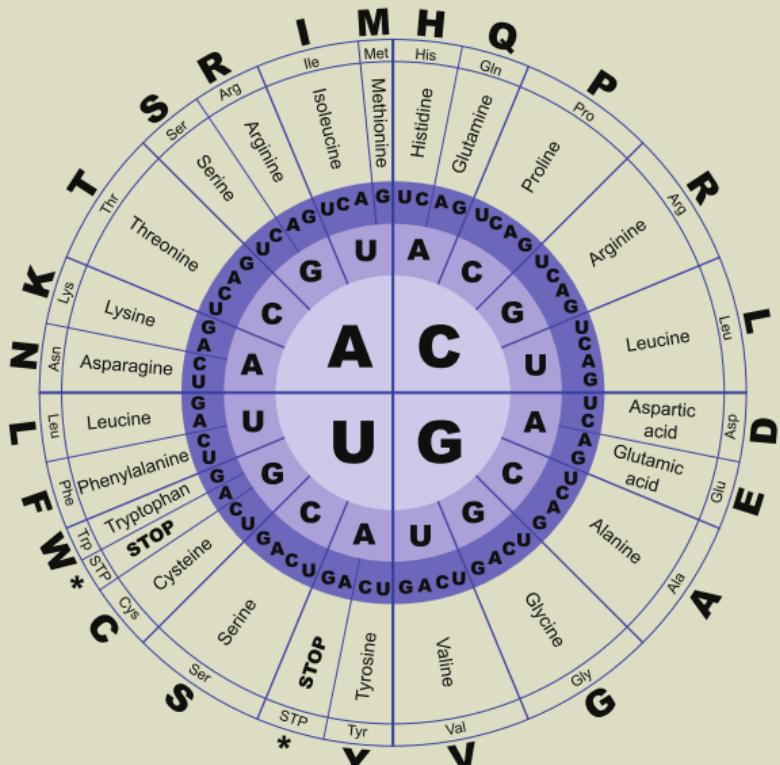
Protein Biosynthesis

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



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

Genetic Code



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

Within *ribosomes* (protein factories)

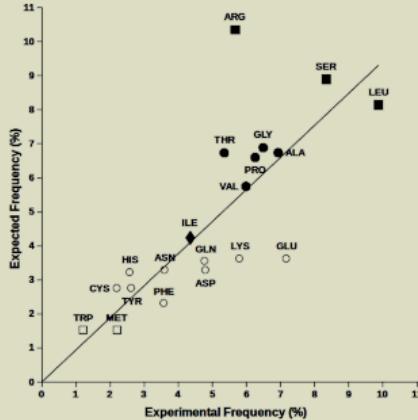
- ▶ translation
 - ▶ from RNA bases {A, C, G, U}
 - ▶ to amino acids (peptide)
{A, C, D, E, F, G, H, I, K, L,
M, N, P, Q, R, S, T, V, W, Y}
- ▶ uses *transfer RNA*
“chemical finite state transducer”
- ▶ *Genetic Code*:
3-base codons → amino acid

Inverse Codon Table

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

(UGA) ← Sometimes, stop codon UGA instead codes 21st amino acid Selenocystein...

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**) ↵ *open reading frame*
 - ▶ search for promoter binding sites (docking station for transcription molecules)
 - ▶ bias of base frequencies in coding vs non-coding regions
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

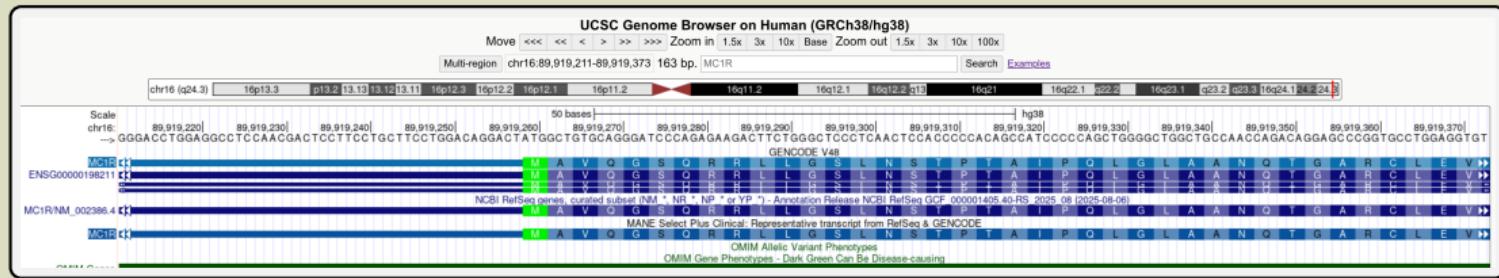
Focus for today: Ab initio options

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?



- ▶ 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!)



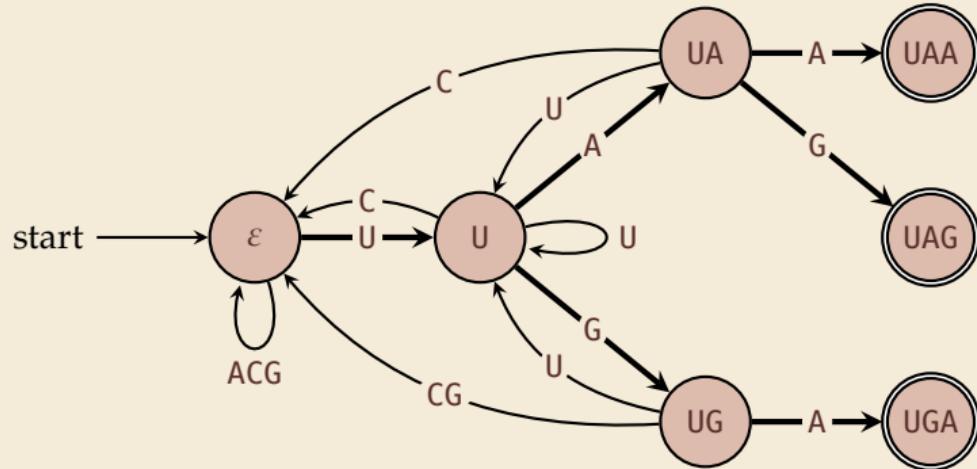
~~ Can we reliably distinguish genes from randomly occurring open reading frames?

Open Reading Frame Gene Detection

- ▶ **Random RNA Model:** String $D[0..N]$ generated i.i.d. uniformly
i.e., each $D[i] \stackrel{\mathcal{D}}{=} \text{Uniform}(\{\text{A, C, G, T}\})$
- ▶ **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)

Stop Codon 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.

2.5 Probability Generating Functions

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

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

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

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 $G_{X+Y}(z) = G_X(z) \cdot G_Y(z)$



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

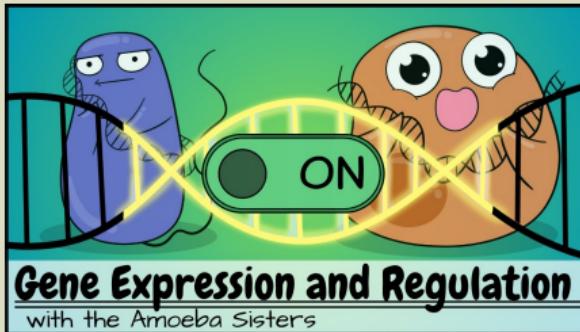
$$\rightsquigarrow \mathbb{E}[Y] = G'_Y(1) = p \cdot n(pz + 1 - p)^{n-1} \Big|_{z=1} = np$$

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

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*



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!*

↔ called a *motif*

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													
Motifs													
T	C	G	G	G	G	g	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		
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
c	C	C	G	G	t	G	A	c	T	T	a	C
a	C	C	G	G	G	G	A	T	T	T	t	C
<i>Motifs</i>	T	t	G	G	G	G	A	c	T	T	t	t
a	a	G	G	G	G	G	A	c	T	T	C	C
T	t	G	G	G	G	G	A	c	T	T	C	C
T	C	G	G	G	G	G	A	T	T	c	a	t
T	C	G	G	G	G	G	A	T	T	c	C	t
T	a	G	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

$$\text{PROFILE}(\text{Motifs})$$

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

$$\text{CONSENSUS}(\text{Motifs})$$

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 t (# strings) and constant σ

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 $O(\sigma^k t n)$ ($k \ll n, t$) \rightsquigarrow FPT wrt. parameter (σ, k)

Neither is 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.



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  $s_1^*, \dots, s_t^*$ 
```

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

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