

NTNU

Innovation and Creativity

TDT4287

Algorithms for bioinformatics

Autumn 2021

Pål Sætrom

Dept. of Computer Science

Dept. of Clinical and Molecular Medicine



Today

- (Practical matters)
- Bioinformatics?
- Genes and sequences
- Comparing sequences

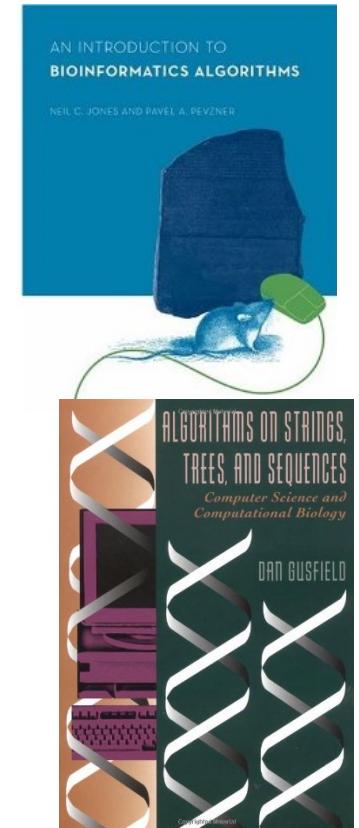
Goals:

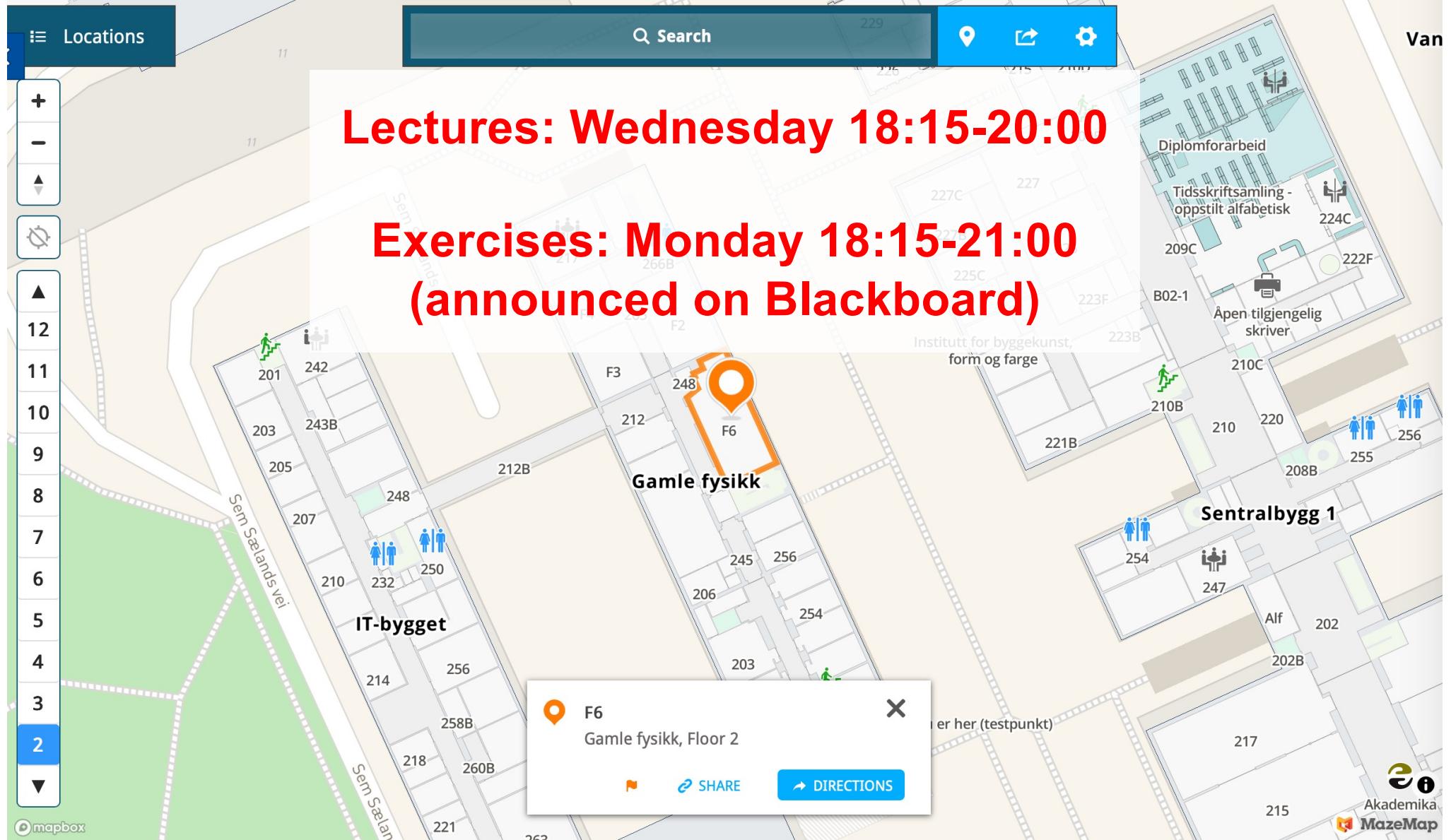
- Molecular biology for computer scientists
- Basic sequence comparison (and why)



Course material and contents

- Exam:
 - Written (or home)
 - Code D simple calculator only
 - Date: TBD
- Course material
 - An Introduction to Bioinformatics Algorithms
 - Algorithms on Strings, Trees and Sequences: Computer Science and Computational Biology
 - Handouts (blackboard)
- Mandatory project







Bioinformatics?

Webster:

bio·in·for·mat·ics (bī'ō-īn'fər-măt'iks)

Function: noun plural but singular in construction

: the collection, classification, storage, and analysis of biochemical and **biological information** using **computers** especially as applied in molecular genetics and genomics

Merriam-Webster Online Dictionary copyright © 2005 (<http://www.merriam-webster.com/>)

- Bioinformatics
 - Tools
 - Data
- Computational biology
 - Biological questions
 - Tools, models, and data are the means



Bioinformatics – this course

- Sequence (string) analysis
 - Alignment -> finding similar sequences
 - Suffix trees -> index-based searching
 - Motif finding -> similarities between sequences
 - Assembly -> reconstructing sequences from fragments
 - RNA structures -> similarities within a sequence

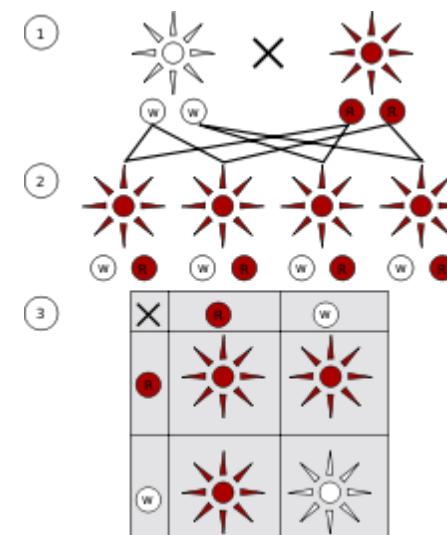
| Week | Topic | Curriculum |
|------|---------------------------------------|--|
| 34 | Introduction | |
| 35 | Alignment | Chps. 6.4-6.9 in Jones & Pevzner |
| 36 | Heuristic alignment | Chps. 9.3, 9.6-9.8 in Jones & Pevzner |
| 37 | BLAST/Statistical analysis | |
| 38 | Substring indexes | Chps. 9.4-9.5 in Jones & Pevzner; Handouts (available on itslearning). |
| 39 | Suffix tree applications | Chps. 9.4-9.5 in Jones & Pevzner; Handouts. |
| 40 | Sequencing and assembly; Project info | Chps. 8.1-8.9 in Jones & Pevzner |
| 41 | Markov Chains | Handouts; Ewens and Grant |
| 42 | Hidden Markov Models | Ewens and Grant |
| 43 | HMMs continued | |
| 44 | Alignment 2 | Chps. 5.5, 12.2-12.3; 6.10, 7.2-7.4 in Jones & Pevzner |
| 45 | Motif discovery | Chps. 4.4-4.9, 12.2-12.3 in Jones & Pevzner |
| 46 | Motif discovery - ctd. | |
| 47 | RNA secondary structure prediction | Chps. 5.1-5.4; 6.11-6.14 in Jones & Pevzner; Article 1 Article 2 |
| ... | Exam | |

- *Why sequence analysis?*



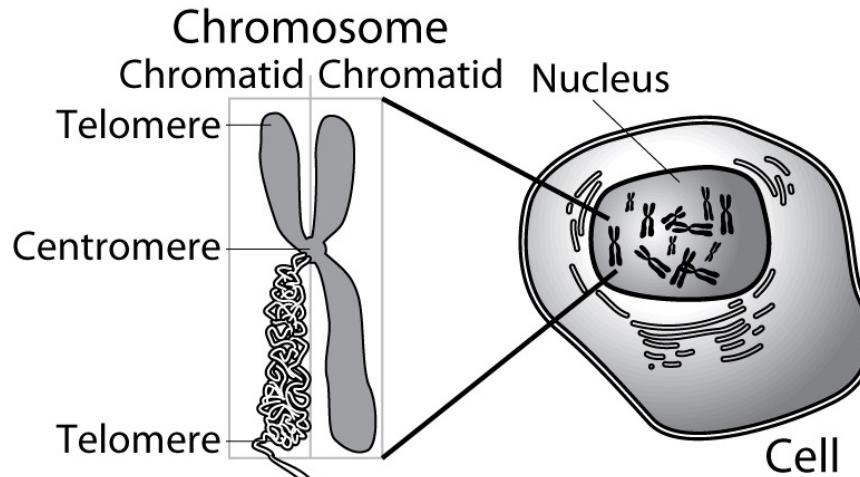
Genes control heritable characteristics

- Gene (Gregor J. Mendel - 1866)
 - Abstract unit of heredity
 - Alternative forms (alleles) -> different (physical) characteristics
 - Genotype + environment -> phenotype





Genes are in DNA



- Cells' cookbook
 - “Book of life”
 - 23 chromosomes
 - $3 * 10^9$ ordered base pairs
 - Self-copying
 - *Mechanism?*

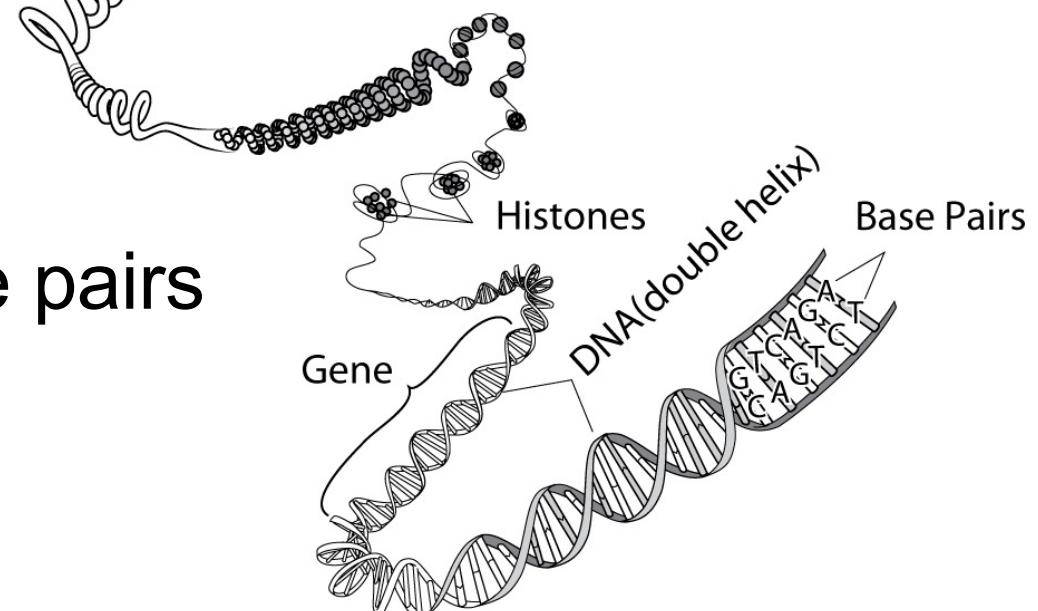
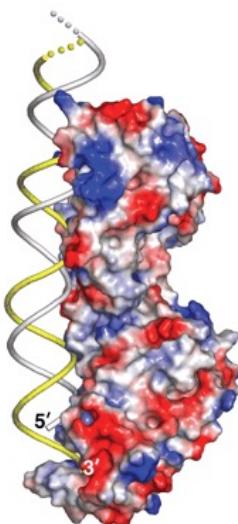
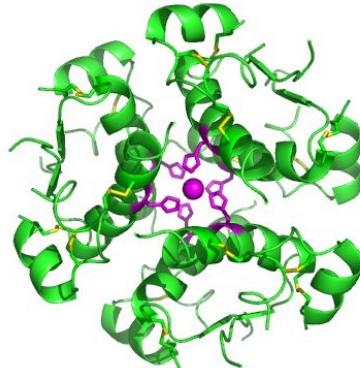
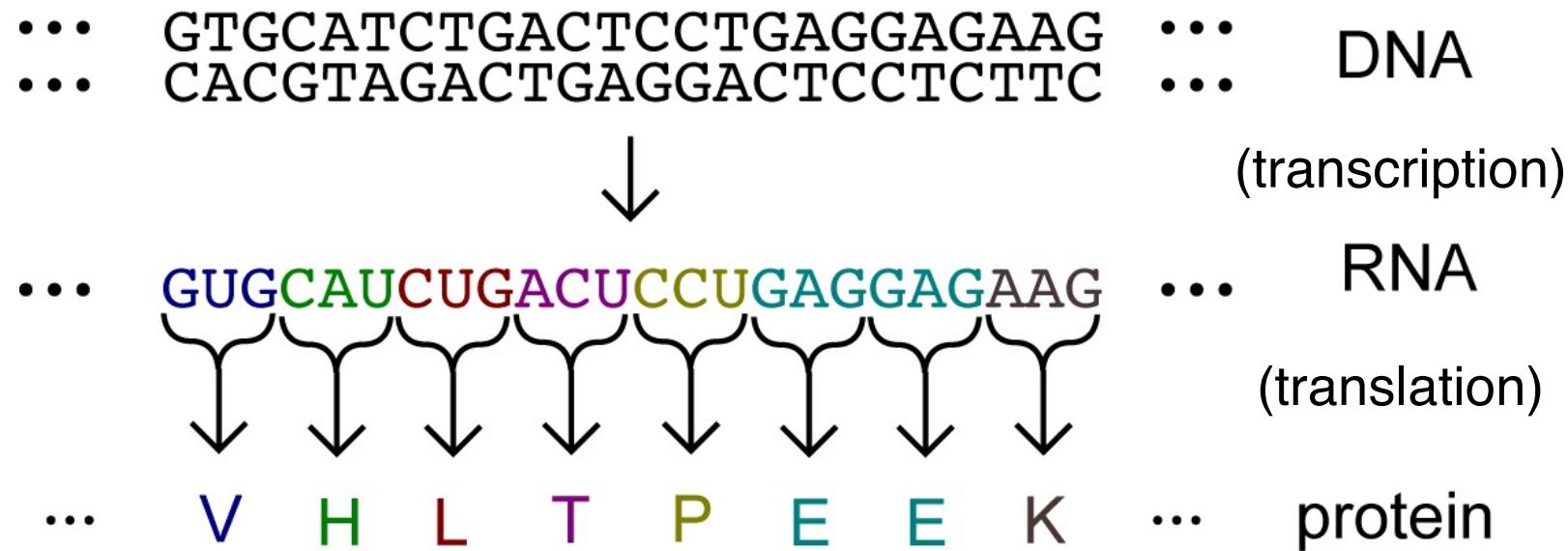


Figure: NIH (modified)



Genes are recipes for molecules



Figures: Wikipedia, MacRae

MacRae, I.J. (2006). *Science* **311**: 195-198.



A cell as seen by a computer scientist



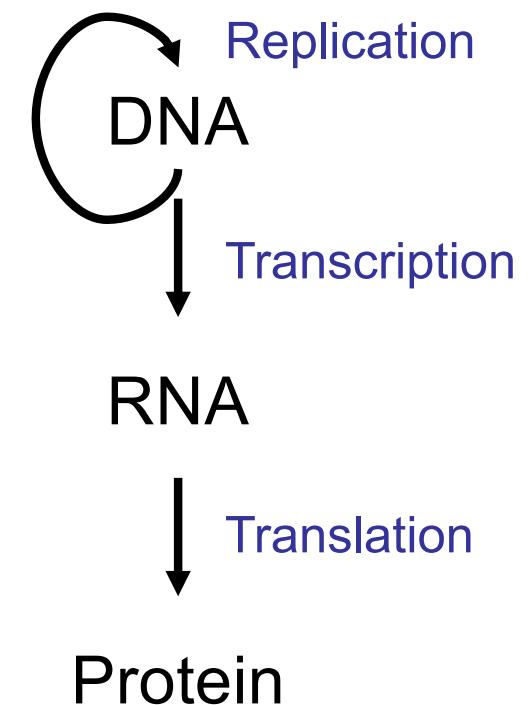
```
...CCAGCTGCTTCGGGCTGCCGAGGACCTCTGGGCCCA  
CATTAATGAGGCAGCCACCTGGCGAGTCTGACATGGCTGT  
CAGCGACGCGCTGCTCCATCTTCTCCACGTTCGCGTC...
```

f_{TR}

```
CAUUAUGAGGCAGCCACCUGGCGAGUCUGACAUGGCUGU  
CAGCGACGCGCUGCUCCCAUCUUUCUCCACGUUCGCGUC...
```

f_{TL}

```
MRQPPGESDMAVSDA...
```



If you know the DNA-sequence (the genome),
you will also know the genes!



...and explain traits and disease?



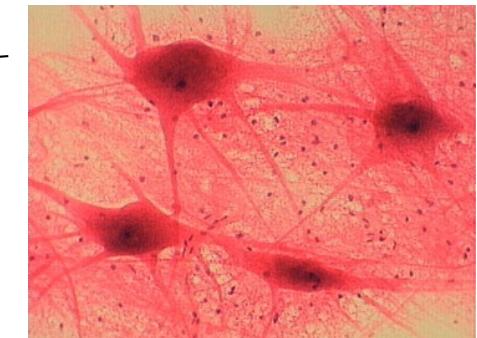
**Identical twins have the same genes, but
are still different**



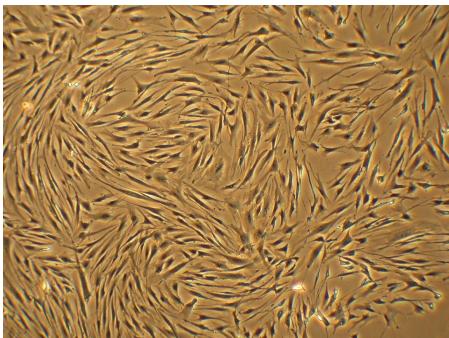


411 cell types, 10^{14} cells – same genes!

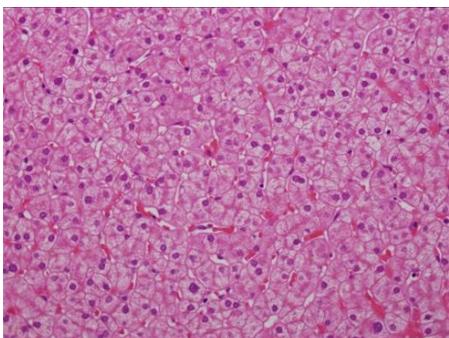
Brain



Connective tissue



Liver



Muscle

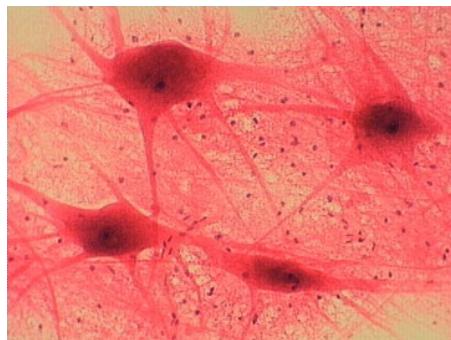


Pictures:
Photo Researchers, Inc., Iowa State Univ.,
Stephanie Saade, USA Today



Different genes turned on in different cell types

Brain



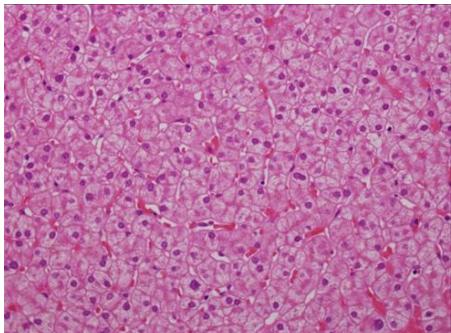
Doublecortin
Myosin
Albumin

Muscle



Doublecortin
Myosin
Albumin

Liver



Doublecortin
Myosin
Albumin



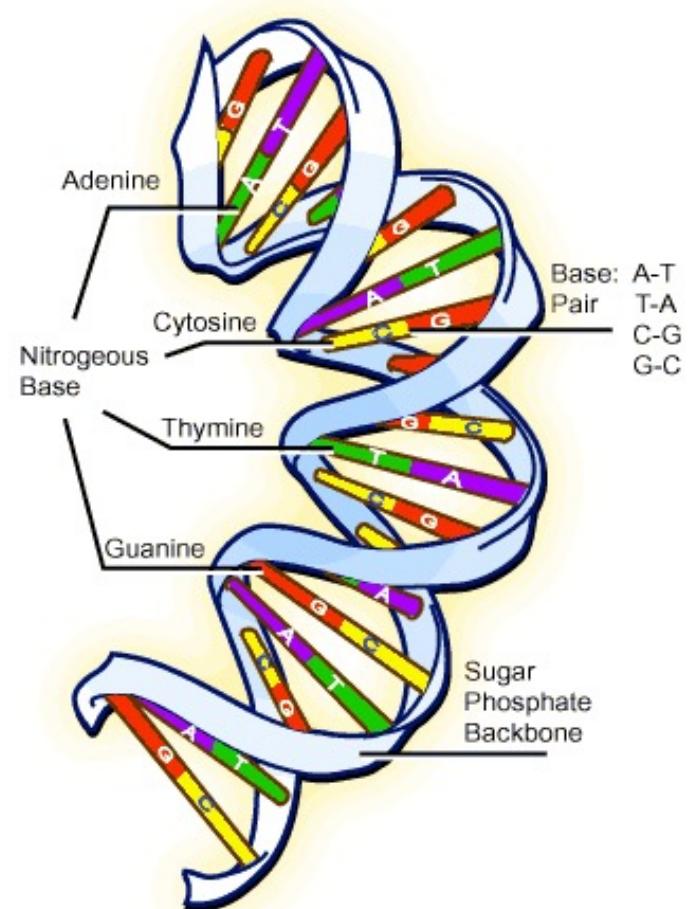
List of biological sequence problems

- How do genes encode proteins?
- What is the function of a protein?
- How do we recognize genes in DNA?
- How are genes controlled so that proteins are produced at the correct place and time?



Facts about genes (1960)

- Genes are stored in DNA
- Genes encode proteins (& RNAs)
- RNA is information carrier
- Proteins?





Protein structure

- Sequence (string) of amino acids
 - 20 different amino acids
 - Defined direction ($\text{H}_2\text{N}->\text{COOH}$ or N->C)
 - *Primary structure*
 - First sequence: Insulin (Sanger, 1955)
- Amino acid sequence folds into (unique), reproducible 3D structure
 - Tertiary structure
- Tertiary structure essential
 - Binding sites
 - Structure determines function
- How do genes encode proteins?



Genes (DNA) and proteins are sequences

- DNA (4 nucleotides)
- Proteins (20 amino acids)
- Genes (DNA) encode proteins:
 - Mapping $f : \text{DNA} \rightarrow \text{Protein}$
- $f = ?$



The genetic code

- ...must use at least 3 nucleotides to encode all 20 amino acids ($4^2 = 16$)
- ...has some redundancy ($4^3 = 64$)
- ...signal for start and stop



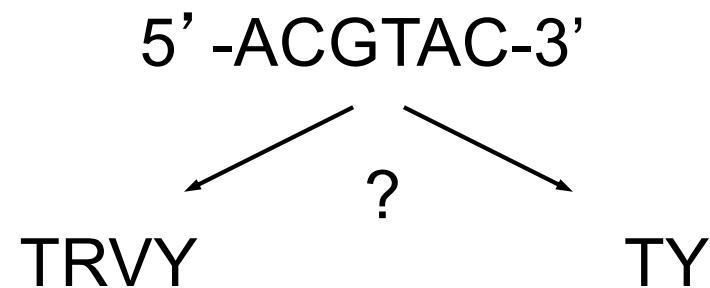
The genetic code (Nirenberg & Khorana, 1965)

| | | Second letter | | | | | |
|--------------|---|---------------------------------------|--------------------------------|--|---|------------------|--------------|
| | | U | C | A | G | | |
| First letter | U | UUU UUC UUA UUG } Phe | UCU UCC UCA UCG } Ser | UAU UAC UAA UAG } Tyr Stop Stop | UGU UGC UGA UGG } Cys Stop Trp | U C A G | Third letter |
| | C | CUU CUC CUA CUG } Leu | CCU CCC CCA CCG } Pro | CAU CAC CAA CAG } His Gln | CGU CGC CGA CGG } Arg | U C A G | |
| | A | AUU AUC AUA AUG } Ile Met | ACU ACC ACA ACG } Thr | AAU AAC AAA AAG } Asn Lys | AGU AGC AGA AGG } Ser Arg | U C A G | |
| | G | GUU GUC GUA GUG } Val | GCU GCC GCA GCG } Ala | GAU GAC GAA GAG } Asp Glu | GGU GGC GGA GGG } Gly | U C A G | |



Reading the genetic code?

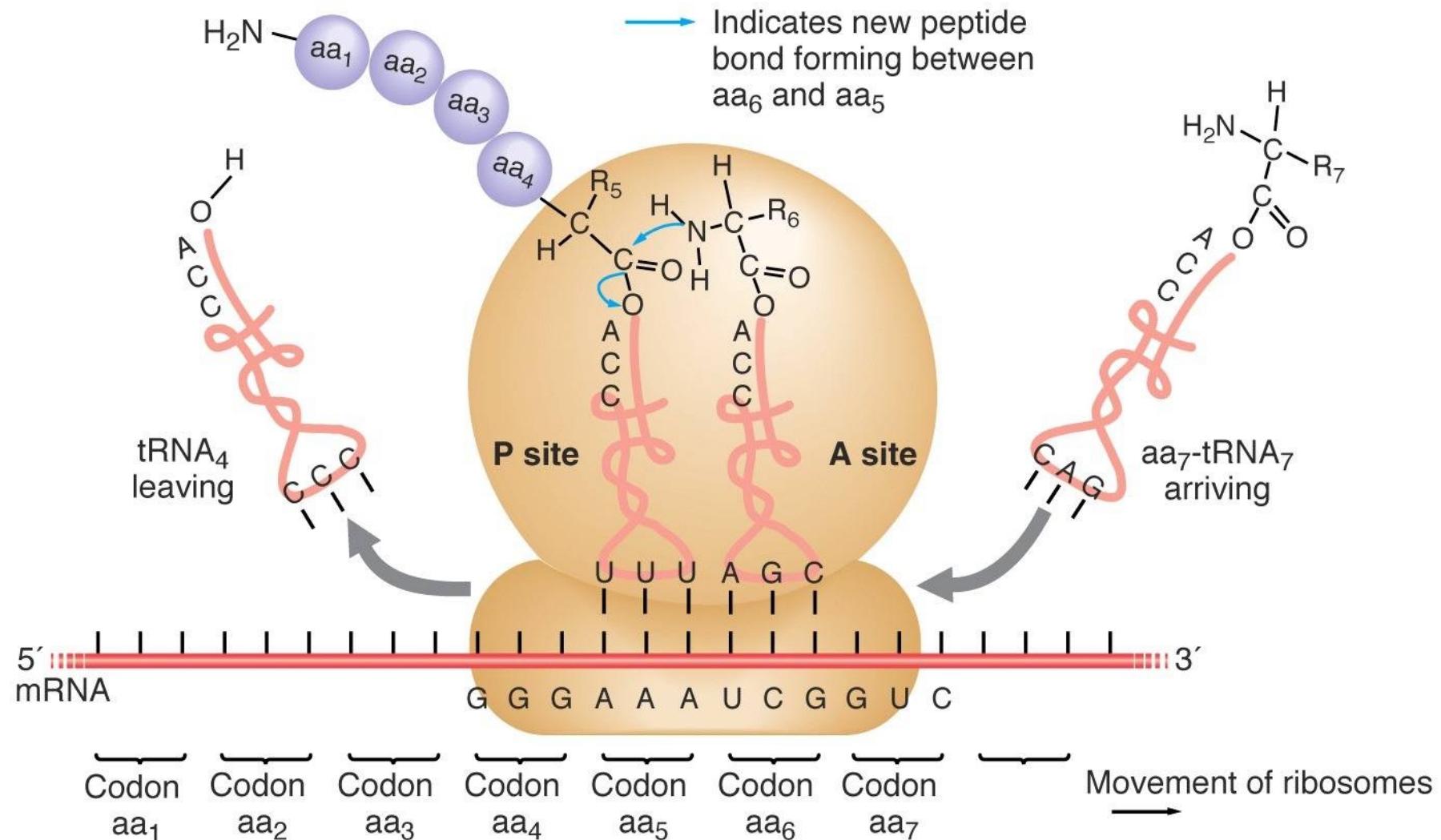
- Overlapping or exclusive words?



- Hints:
 - Adding or deleting 1 nucleotide destroys gene
 - Adding and deleting 1 nucleotide may be OK



Ribosomes translate mRNA to proteins





Changes in the recipe can affect function

... GTGCATCTGACTCCTGAGGGAGAAG ...

... GTGCATCTGAC**A**CCTGAGGGAGAAG ...

... GTGCATCTGAC-**C**CTGAGGGAGAAG ...

... GTGCATCTGAC**G**TCCTGAGGGAGAAG ...

Which of the above could have the strongest effect on function?



List of biological sequence problems

- ✓ How do genes encode proteins?
- What is the function of a protein?
- How do we recognize genes in DNA?
- How are genes controlled so that proteins are produced at the correct times?



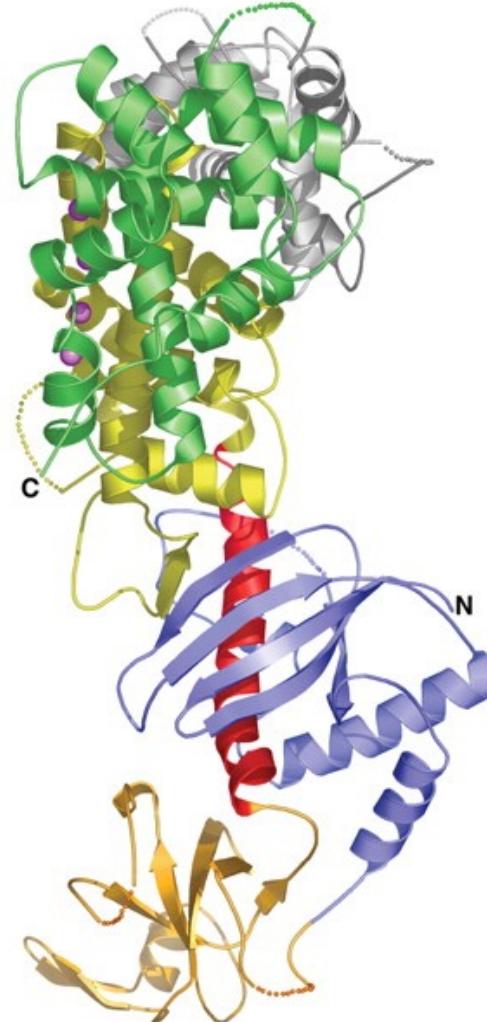
Protein sequence determines function

Sequence

```
>2FFL Giardia intestinalis Dicer
GAMHALGHCTVVTTRGPSHWLLLLDTHLGTLPGF
KVSAGRGLPAAEVYFEAGPRVLSRTDATIVAVYQ
SILFQLLGPTFPASWTEIGATMPHNEYTFPRFISN
PPQFATLAFPLLSPTSPLDLRLALMVTAQLMCDAK
RLSDEYTDYSTLSASLHGRMVATPEISWSLYVVLG
IDSTQTSLSYFTRANESITYMRYATAHNHLRAA
DLPLVAAVRLLDDLKDHQIPAPGSWDDLAPKLRFPLP
PELCLLLPEFDLIRVQALQFLPEIAKHICDIQNT
ICALDKSFDPDCRIGGERRYFAITAGLRLDQGRGRG
LAGWRTPFGPGFVGSSTDVFQRLELLGDAVLGFIVT
ARLLCLFPDASVGTIQLVELKMEILVRNEALNYLVQTL
GLPQLAEFSNNLVAKSKTWADMYEEIVGSIFTGPNN
GIYGCEEFLAKTLMSPHEHSKTVGSACPDAVTKASK
RVCMGGEAGAHEFRSLVDYACEQGISVFCSSRVSTM
FLERLRDIPAEDMLDWYRLGIQFSHRSGLSPGGV
VSVIDIMTHLARGLWLGS PGFYVEQQTDKNESACP
PTIPVLYIYHRSVQCPVLYGSLTETPTGPVASKVL
ALYEKILAYESSGGSKHIAAQTVSRSLAVPIPSGT
IPFLIRLIQIALTPHVYQKLELLGDAFIKCSLALH
LHALHPTILTEGALTRMRQSAETNSVLGRLTKRFPSS
VVSEVIIESHPKIQPDSKVYGDTFEAILAAILLAC
GEEAAGAFVREHVLQPQVADA
```

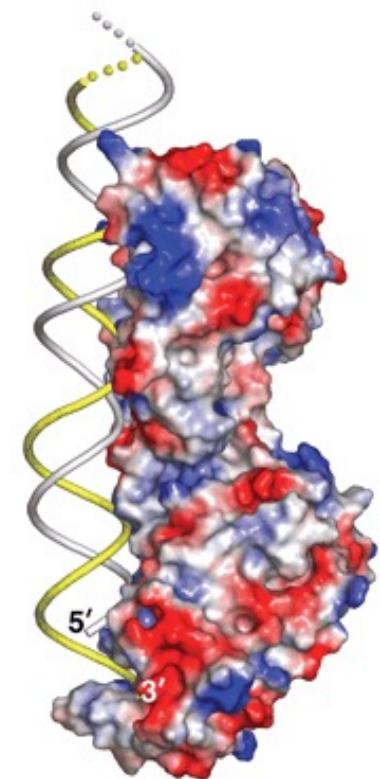
→

Structure



→

Function





Search to determine gene function

- Genes evolved from common ancestor – *homologs* – often have similar functions
- *Homologous genes*
 - *Orthologs* -> Genes in separate species
 - *Paralogs* -> Duplicated genes
- Homologous genes often have sequence similarity
 - Sequence similarity indicates homology
 - ...and similar function

Evaluating sequence similarity?



String terminology

- **String**

$S = s_1 \dots s_n, s_i \in A, |S| = n$

$S = \text{ACGTTAGCT}$

- **Substring of S**

$S_{i,m} = s_{1+i} \dots s_{m+i}, 0 \leq i, i + m \leq n$

$S_{1,4} = \text{CGTT}$

- **Prefix of S**

$S_i = s_1 \dots s_i, 0 \leq i \leq n$

$S_3 = \text{ACG}$

- **Suffix of S**

$S_j = s_{1+j} \dots s_n, 0 \leq j \leq n$

$S_3 = \text{TTAGCT}$



Subsequences

Subsequence: Ordered sequence of characters from string

AGCTTAGCTG

ACTGT, CTTA, GCG are subsequences

GACA, TCA, ATGA are not

Set of indexes:

ACTGT = {1,3,4,7,9}



Common subsequences

Common subsequence of two strings: subsequence in both strings

Formally:

$$S = s_1 \dots s_n, R = r_1 \dots r_m$$

$$1 \leq i_1 < i_2 < \dots < i_k \leq n$$

(subsequence in S)

$$1 \leq j_1 < j_2 < \dots < j_k \leq m$$

(subsequence in R)

$$s_{i_t} = r_{j_t}, \text{ for } 1 \leq t \leq k.$$

AG**C**T_A**GCTG**, T**CG**GAT**G**

CGTG is a common subsequence



Longest common subsequence (LCS)

Find the longest subsequence common to two strings

Input: Two strings, S and R

Output: Longest common subsequence of S and R



Recursive LCS

LCS(AGCTTAGCT**G**, TCGGAT**G**)?

Three possibilities:

1. Last letter from both strings part of LCS if identical; none of them otherwise
2. Last letter from left string is part of LCS
3. Last letter from right string is part of LCS



Recursive LCS (2)

Assume you know

$$\text{LCS}(\text{AGCTTAGCT}, \text{TCGGAT}) = \text{LCS}_{-1,-1}$$

$$\text{LCS}(\text{AGCTTAGCTG}, \text{TCGGAT}) = \text{LCS}_{0,-1}$$

$$\text{LCS}(\text{AGCTTAGCT}, \text{TCGGATG}) = \text{LCS}_{-1,0}$$

$$\text{LCS}(\text{AGCTTAGCTG}, \text{TCGGATG})$$

$$= \max\{\text{LCS}_{0,-1}, \text{LCS}_{-1,0}, \text{LCS}_{-1,-1} + f(\text{G}, \text{G})\}$$

$$f(r,s) = \begin{cases} 1 & r = s \\ 0 & r \neq s \end{cases}$$



Recursive LCS (3)

- Recursive LCS solves sub-problems multiple times
- Recursion table

$\text{LCS}(\text{AGCTTAGC}, \text{TCGGATG})$

$$= \max\{\text{LCS}_{0,-1}, \text{LCS}_{-1,0}, \text{LCS}_{-1,-1} + f(C, G)\}$$

Alternative solution?

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | ε | T | C | G | G | A | T | G |
| ε | | | | | | | | |
| A | | | | | | | | |
| G | | | | | | | | |
| C | | | | | | | | |
| T | | | | | | | | |
| T | | | | | | | | |
| A | | | | | | | | |
| G | | | | | | | | |
| C | | | | | | | | |

↑ $\text{LCS}_{-1,-1}$
↑ $\text{LCS}_{-1,0}$
↑ LCS
↑ $\text{LCS}_{0,-1}$



Dynamic programming (DP) LCS

- Solve LCS bottom up

LCS(S_i, R_j)

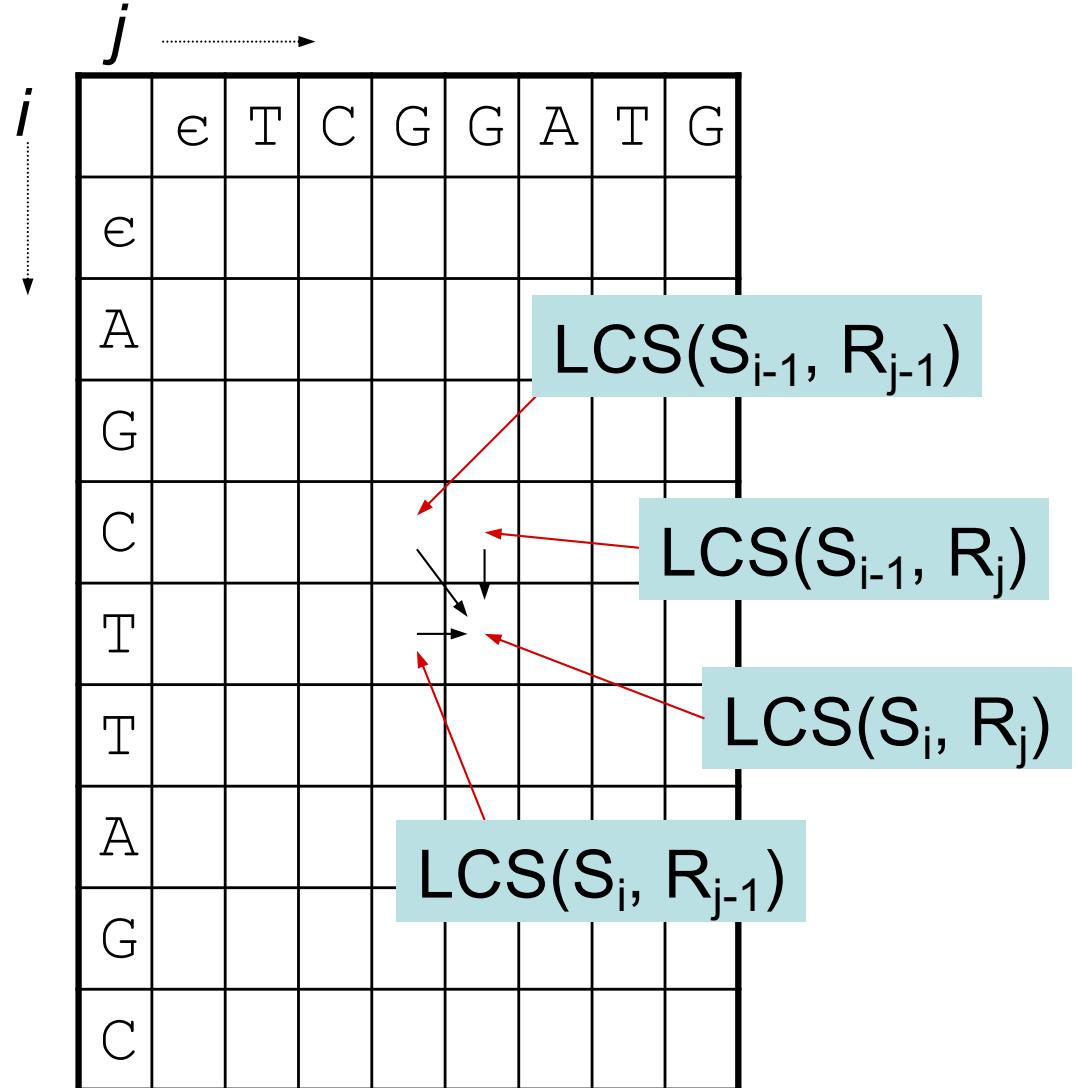
$$= \max\{$$

$$\text{LCS}(S_{i-1}, R_j),$$

$$\text{LCS}(S_i, R_{j-1}),$$

$$\text{LCS}(S_{i-1}, R_{j-1}) + f(s_i, r_j)$$

$$\}$$





Dynamic programming LCS (2)

- Solve LCS bottom up

$LCS(S_i, R_j)$

= max{

$LCS(S_{i-1}, R_j),$

$LCS(S_i, R_{j-1}),$

$LCS(S_{i-1}, R_{j-1}) + f(s_i, r_j)$

}

1. Initialize matrix
2. *Fill matrix using recurrence?*

| | ε | T | C | G | G | A | T | G |
|---|---|---|---|---|---|---|---|---|
| ε | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | | | | | | | |
| G | 0 | | | | | | | |
| C | 0 | | | | | | | |
| T | 0 | | | | | | | |
| T | 0 | | | | | | | |
| A | 0 | | | | | | | |
| G | 0 | | | | | | | |
| C | 0 | | | | | | | |



Dynamic programming LCS (3)

- Solve LCS bottom up

$LCS(S_i, R_j)$

= max{

$LCS(S_{i-1}, R_j),$

$LCS(S_i, R_{j-1}),$

$LCS(S_{i-1}, R_{j-1}) + f(s_i, r_j)$

}

1. Initialize matrix
2. Fill matrix using recurrence

| | ε | T | C | G | G | A | T | G |
|---|---|---|---|---|---|---|---|---|
| ε | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | | | | | | |
| G | 0 | 0 | | | | | | |
| C | 0 | 0 | | | | | | |
| T | 0 | 1 | | | | | | |
| T | 0 | | | | | | | |
| A | 0 | | | | | | | |
| G | 0 | | | | | | | |
| C | 0 | | | | | | | |



Dynamic programming LCS (3)

- Solve LCS bottom up

$LCS(S_i, R_j)$

= max{

$LCS(S_{i-1}, R_j),$

$LCS(S_i, R_{j-1}),$

$LCS(S_{i-1}, R_{j-1}) + f(s_i, r_j)$

}

1. Initialize matrix
2. Fill matrix using recurrence

| | ε | T | C | G | G | A | T | G |
|---|---|---|---|---|---|---|---|---|
| ε | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | | | | | |
| G | 0 | 0 | 0 | | | | | |
| C | 0 | 0 | 1 | | | | | |
| T | 0 | 1 | 1 | | | | | |
| T | 0 | 1 | 1 | | | | | |
| A | 0 | 1 | 1 | | | | | |
| G | 0 | 1 | 1 | | | | | |
| C | 0 | 1 | 2 | | | | | |



Dynamic programming LCS (3)

- Solve LCS bottom up

$LCS(S_i, R_j)$

= max{

$LCS(S_{i-1}, R_j),$

$LCS(S_i, R_{j-1}),$

$LCS(S_{i-1}, R_{j-1}) + f(s_i, r_j)$

}

1. Initialize matrix
2. Fill matrix using recurrence

| | ε | T | C | G | G | A | T | G |
|---|---|---|---|---|---|---|---|---|
| ε | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| G | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 |
| C | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 |
| A | 0 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |
| G | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 |
| C | 0 | 1 | 2 | 2 | 2 | 2 | 2 | 3 |



LCS DP algorithm

LCS(S, R):

1. $a[0,0] = 0$
2. for i in range(1, $|S|$):
 1. $a[i,0] = 0$
3. for j in range(1, $|R|$):
 1. $a[0,j] = 0$
4. for i in range(1, $|S|$):
 1. for j in range(1, $|R|$):
 1. $a[i,j] = \max\{a[i-1,j], a[i,j-1], a[i-1,j-1] + f(s_i, r_j)\}$



But what is the LCS?

$LCS(S_i, R_j)$
 $= \max\{$
 $LCS(S_{i-1}, R_j),$
 $LCS(S_i, R_{j-1}),$
 $LCS(S_{i-1}, R_{j-1}) + f(C, G)$
 $\}$

- LCS is path through matrix
- TTG
- “Remember where you’ve been”

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | ε | T | C | G | G | A | T | G |
| ε | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| G | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 |
| C | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 |
| A | 0 | 1 | 1 | 1 | 1 | 2 | 2 | 2 |
| G | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 |
| C | 0 | 1 | 2 | 2 | 2 | 2 | 2 | 3 |



LCS DP algorithm with backtrack matrix

LCS(S, R):

1. $a[0,0] = 0$
2. for i in range(1, $|S|$):
 1. $a[i,0] = 0$
3. for j in range(1, $|R|$):
 1. $a[0,j] = 0$
4. for i in range(1, $|S|$):
 1. for j in range(1, $|R|$):
 1. $a[i,j] = \max\{a[i-1,j], a[i,j-1], a[i-1,j-1] + f(s_i, r_j)\}$
 2. $b[i,j] = \{0 \text{ if } a[i,j] == a[i-1,j], 1 \text{ if } a[i,j] == a[i,j-1], 2 \text{ if } a[i,j] == a[i-1,j-1] + 1\}$

| | |
|------|--|
| 0 == | |
| 1 == | |
| 2 == | |



Does the path fit the algorithm?

4.1.2 $b[i,j] = \{$

0 if $a[i,j] == a[i-1,j]$,

1 if $a[i,j] == a[i,j-1]$,

2 if $a[i,j] == a[i-1,j-1]+1$

}

| | |
|------|---|
| 0 == | ↑ |
| 1 == | ← |
| 2 == | ↖ |

| | $j \rightarrow$ | | | | | | | | |
|----------------|-----------------|---|---|---|---|---|---|---|--|
| $i \downarrow$ | ϵ | T | C | G | G | A | T | G | |
| ϵ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| A | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | |
| G | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 2 | |
| C | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 2 | |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | |
| T | 0 | 1 | 1 | 1 | 1 | 1 | 2 | 2 | |
| A | 0 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | |
| G | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 | |
| C | 0 | 1 | 2 | 2 | 2 | 2 | 2 | 3 | |



Printing the LCS

PrintLCS(b, S, i, j):

1. if $i == 0$ or $j == 0$:
 1. return
2. if $b[i,j] == 2$:
 1. PrintLCS(b, S, i-1, j-1)
 2. print s_i
3. else if $b[i,j] == 0$:
 1. PrintLCS(b, S, i-1, j)
4. else:
 1. PrintLCS(b, S, i, j-1)

0 == ↑
1 == ←
2 == ↙



Homework

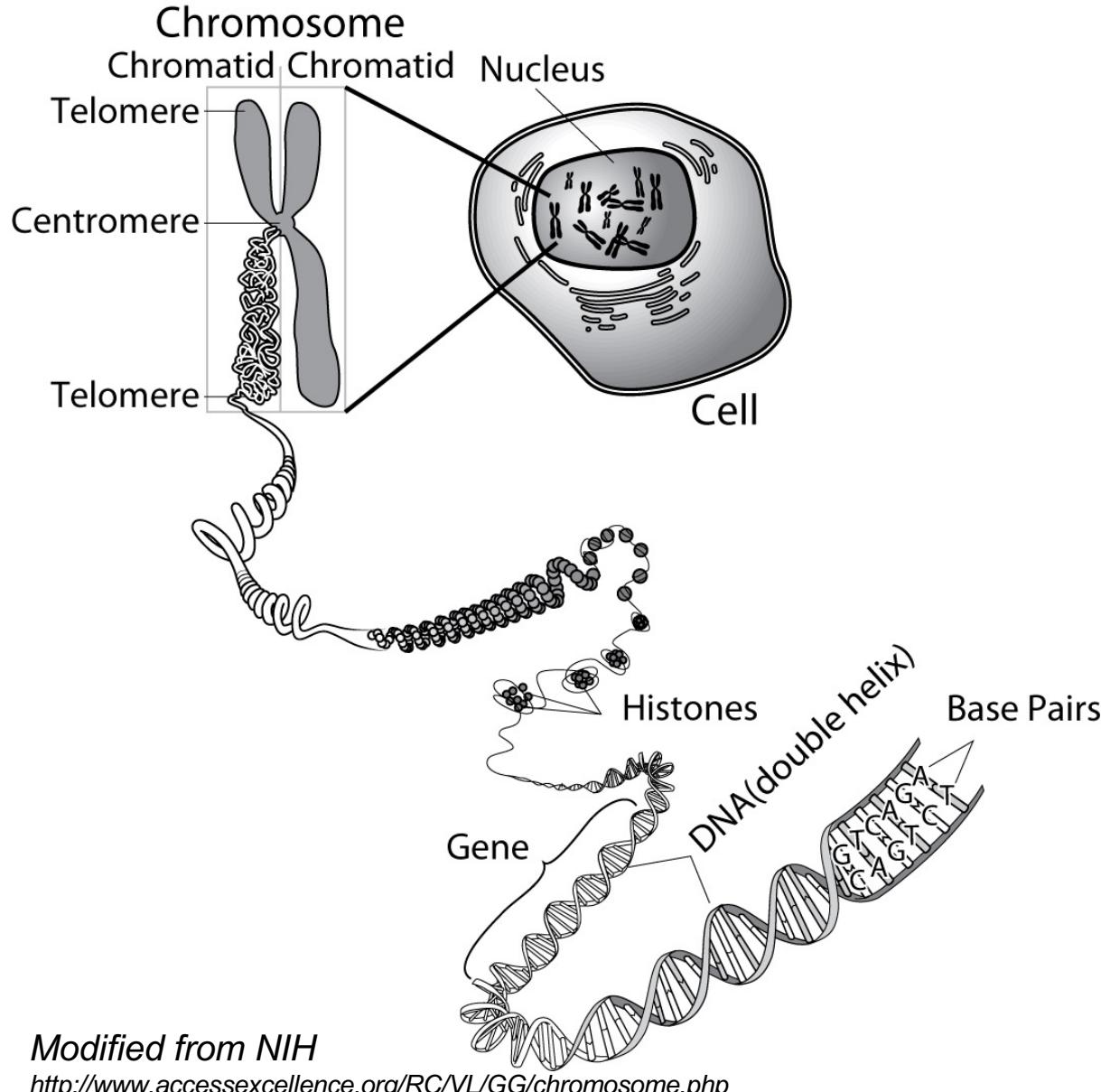
1. Compute the LCS DP matrix for $\text{LCS}(\text{ATTCGGGTTA}, \text{TAGTGATG})$.
2. Find the LCS without using a backtrack matrix.
3. $\text{LCS}(\text{AGCTTAGCTG}, \text{TCGGATG})$ has multiple solutions (for example, TTG or ATG). Find an algorithm that returns all the LCSs. (Hint: a stack could be useful)

Summary

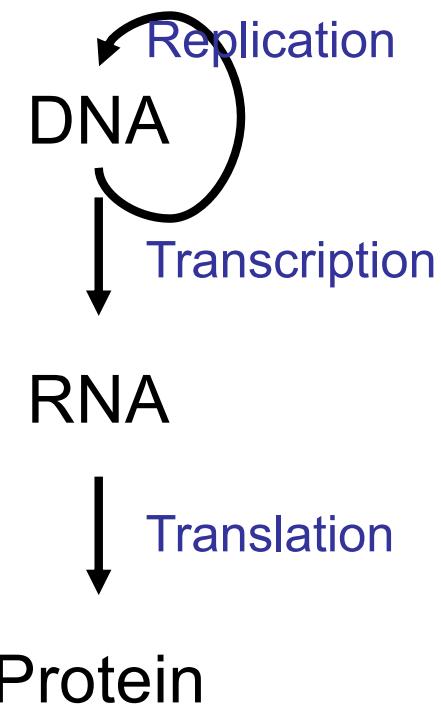
47



Sequences – basic data structures in cells



Sequence data



Modified from NIH

<http://www.accessexcellence.org/RC/VL/GG/chromosome.php>