

Pairwise sequence alignment

Jiankui He

SUSTC, Shenzhen, China



Image by Nicolas Bouvier; courtesy of Genevieve Almouzni,
Curie Institute, Paris, France



南方科技大学

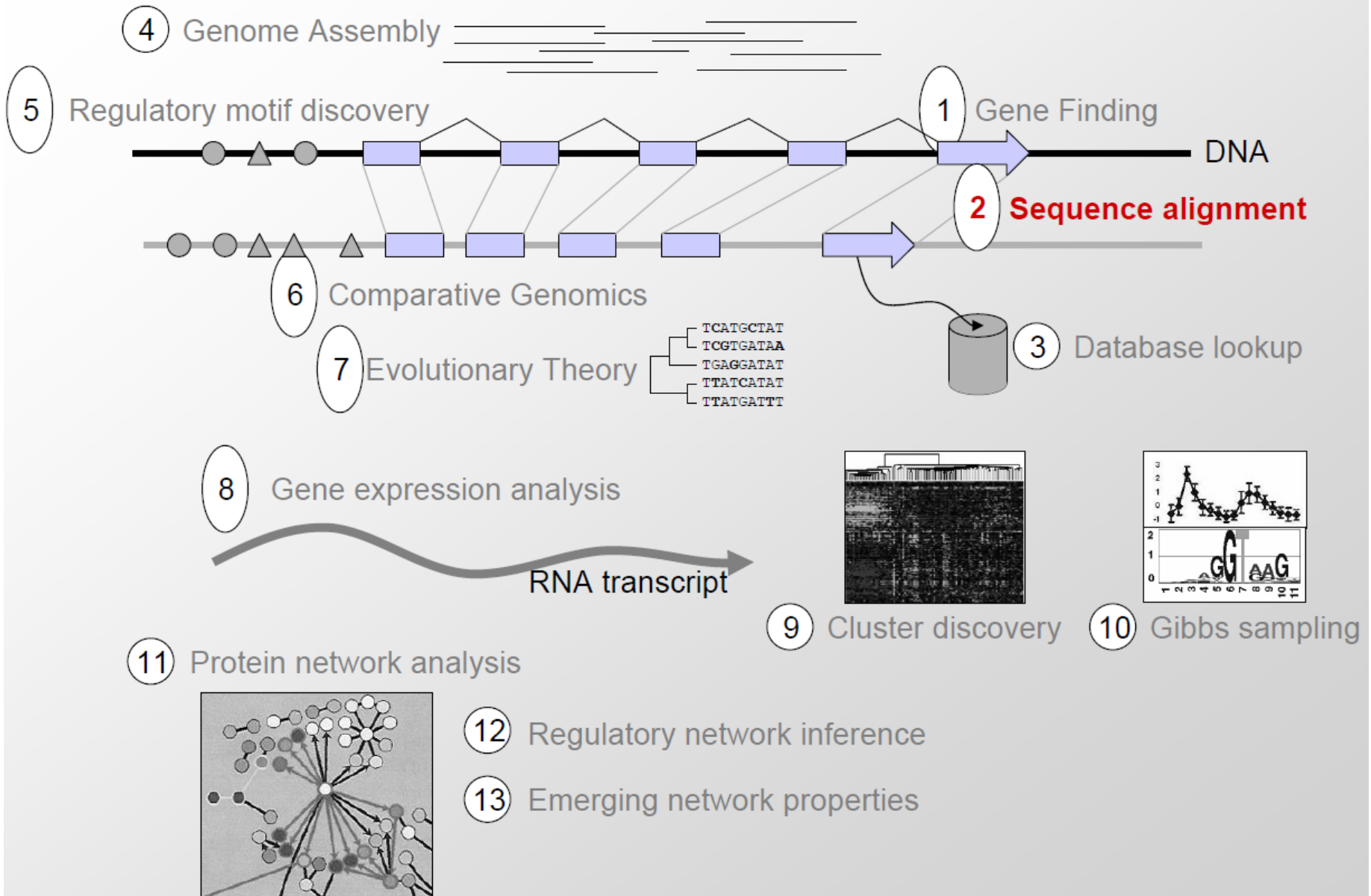
SOUTH UNIVERSITY OF SCIENCE AND TECHNOLOGY OF CHINA

PPT adopted from John Hopkins University

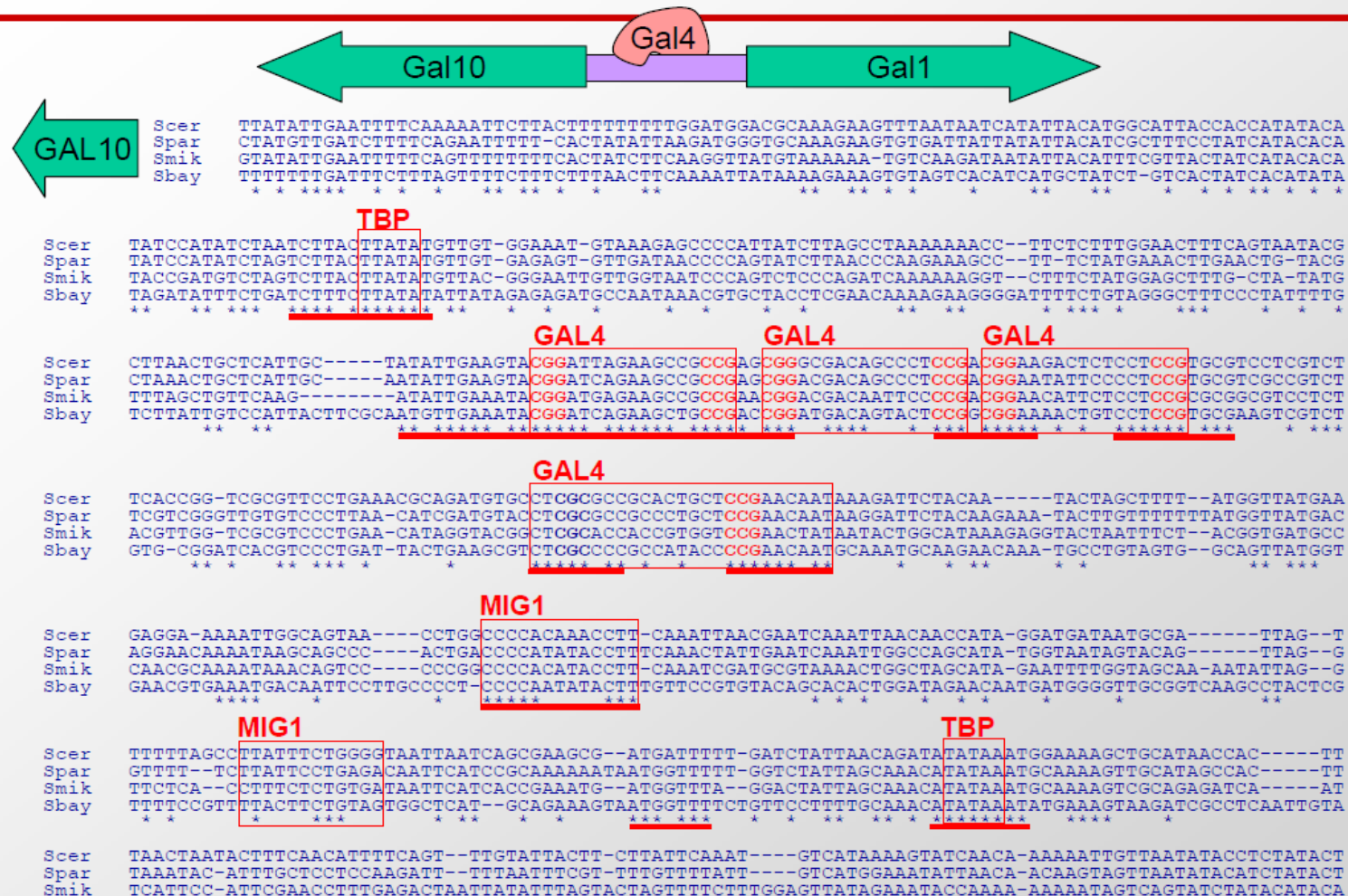
Schedule

- Class I: lecture
- Class II: homework, programming
- Class III: Demo

Challenges in Computational Biology



Evolution preserved functional elements!



We can 'read' evolution to reveal functional elements

Today's goal:

How do we actually align two genes?

Genomes change over time

begin

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

mutation

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

deletion

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

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

insertion

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

end

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

Goal of alignment: Infer edit operations

begin

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

?



end

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

Importance of sequence comparison

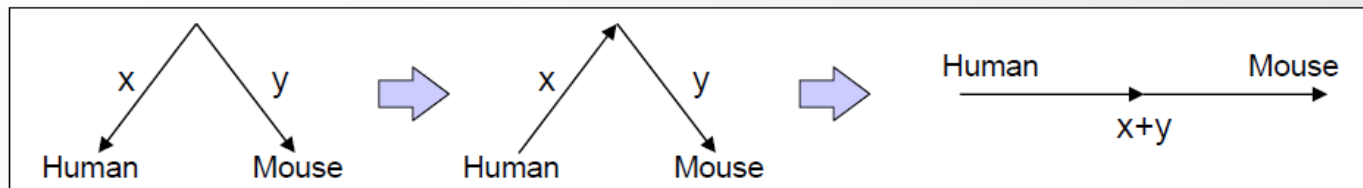
Discovering functional and evolutionary relationships in biological sequences:

- Similar sequences → evolutionary relationship
- evolutionary relationship → related function
- Orthologs → same (almost same) function in different organisms.

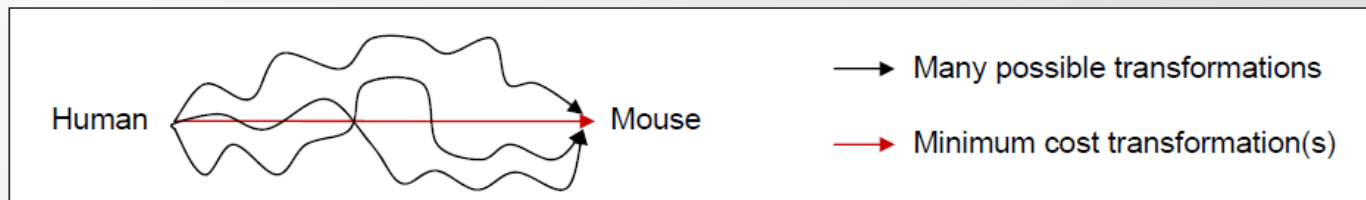
“→” should be read **usually** implies

From Bio to CS: Formalizing the problem

- Define set of evolutionary operations (insertion, deletion, mutation)
 - Symmetric operations allow time reversibility (part of design choice)

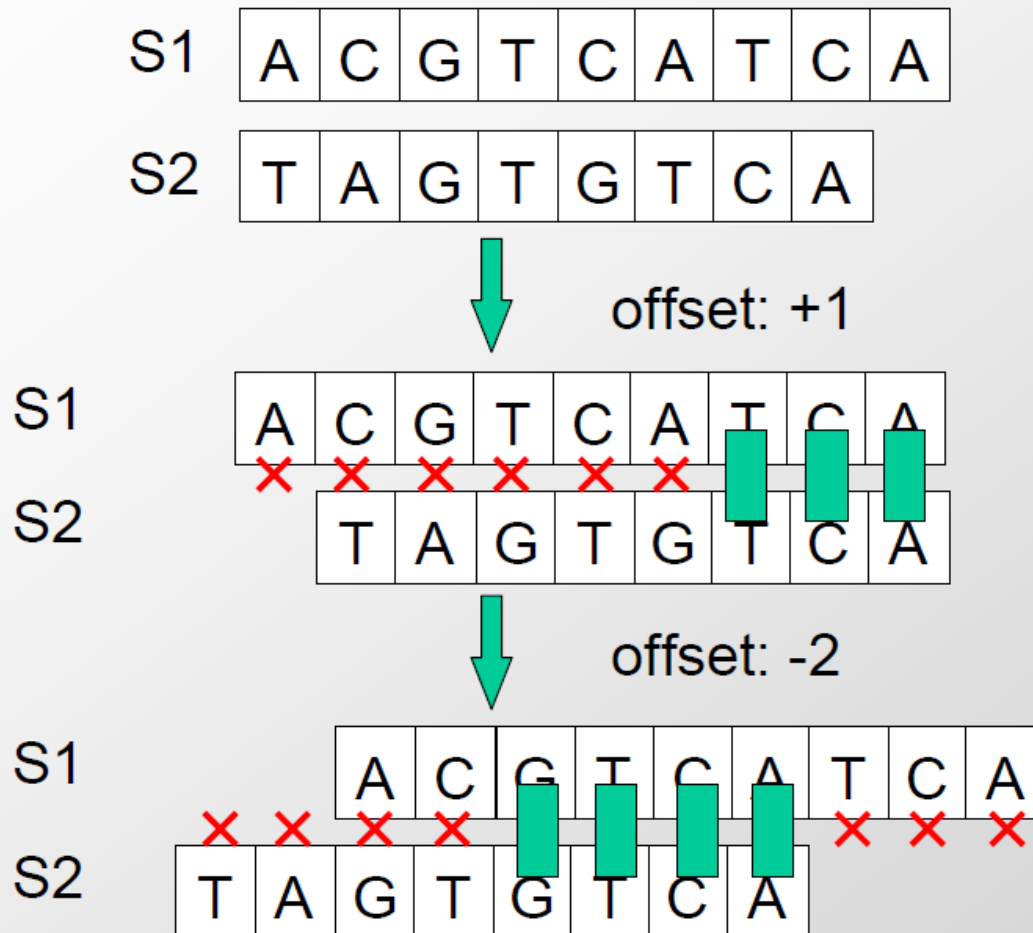


- Define optimality criterion (min number, min cost)
 - Impossible to infer exact series of operations (Occam's razor: find min)



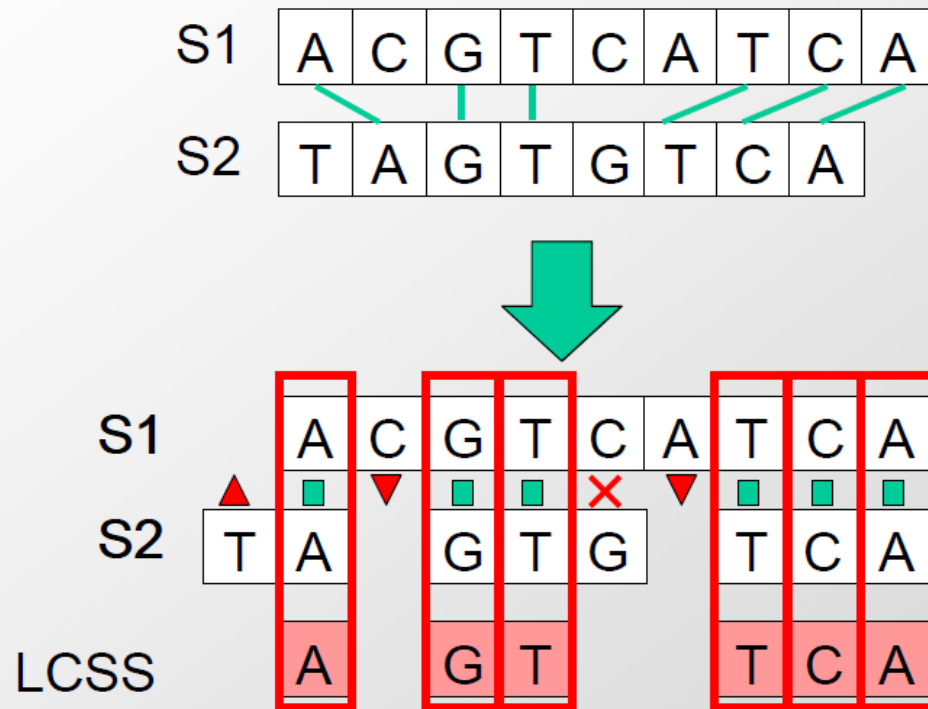
Formulation 1: Longest common substring

- Given two possibly related strings S1 and S2
 - What is the longest common substring? (no gaps)



Formulation 2: Longest common subsequence

- Given two possibly related strings S1 and S2
 - What is the longest common subsequence? (gaps allowed)



Edit distance:

- Number of changes needed for $S1 \rightarrow S2$
- Uniform scoring function

Formulation 3: Sequence alignment

- Allow gaps (fixed penalty)
 - Insertion & deletion operations
 - Unit cost for each character inserted or deleted
- Varying penalties for edit operations
 - Transitions (Pyrimidine \Leftrightarrow Pyrimidine, Purine \Leftrightarrow Purine)
 - Transversions (Purine \Leftrightarrow Pyrimidine changes)
 - Polymerase confuses Aw/G and Cw/T more often

Scoring function:

Match(x,x) = +1

Mismatch(A,G) = $-\frac{1}{2}$

Mismatch(C,T) = $-\frac{1}{2}$

Mismatch(x,y) = -1

	A	G	T	C
A	+1	$-\frac{1}{2}$	-1	-1
G	$-\frac{1}{2}$	+1	-1	-1
T	-1	-1	+1	$-\frac{1}{2}$
C	-1	-1	$-\frac{1}{2}$	+1

purine pyrimid.

Transitions:

A \Leftrightarrow G, C \Leftrightarrow T common
(lower penalty)

Transversions:

All other operations

Etc...
(e.g. varying gap penalties)

How can we compute best alignment

S1

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

S2

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

- Given additive scoring function:
 - Cost of mutation (AG, CT, other)
 - Cost of insertion / deletion
 - Reward of match
- Need algorithm for inferring best alignment
 - Enumeration?
 - How would you do it?
 - How many alignments are there?

Can we simply enumerate all possible alignments?

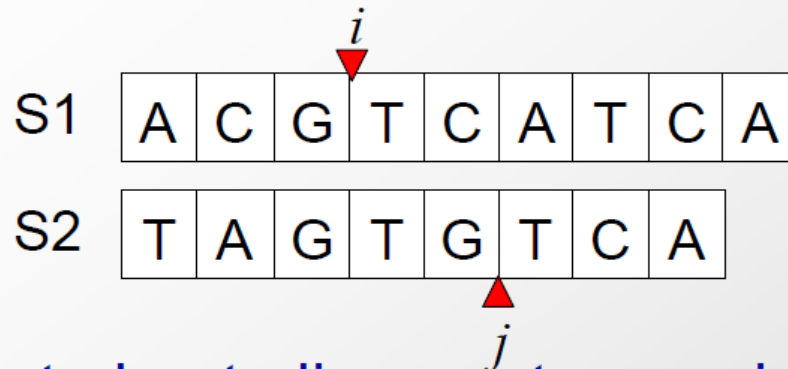
- Ways to align two sequences of length m, n

$$\binom{n+m}{m} = \frac{(m+n)!}{(m!)^2} \approx \frac{2^{m+n}}{\sqrt{\pi \cdot m}}$$

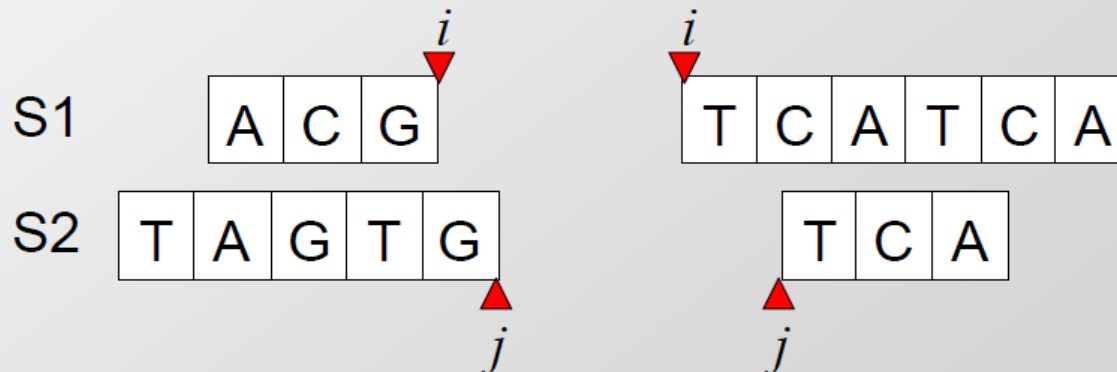
- For two sequences of length n

n	Enumeration	Today's lecture
10	184,756	100
20	1.40E+11	400
100	9.00E+58	10,000

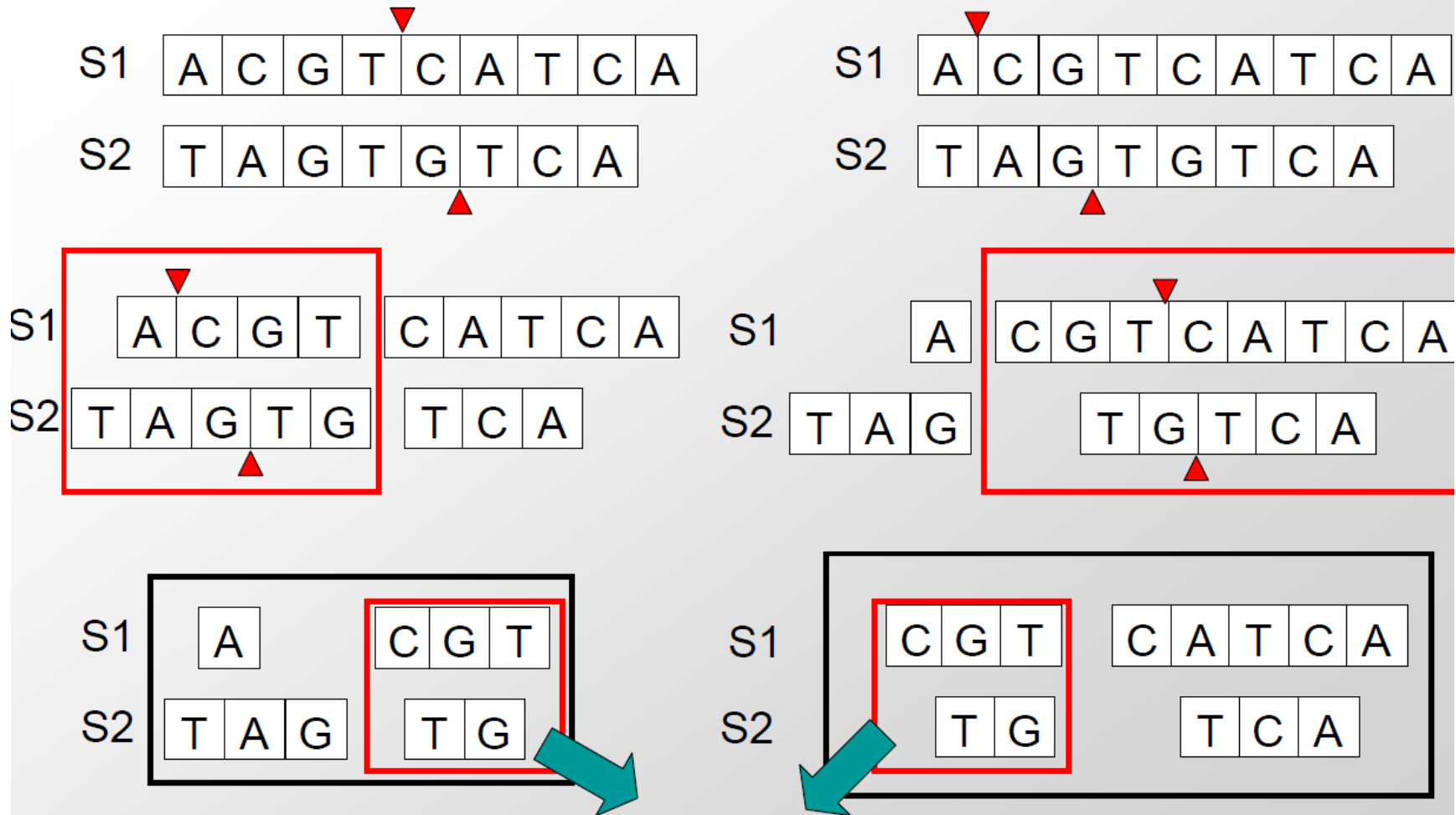
Key insight: score is additive!



- Compute best alignment recursively
 - For a given aligned pair (i, j) , the best alignment is:
 - Best alignment of S1[1.. i] and S2[1.. j]
 - + Best alignment of S1[i ..n] and S2[j ..m]
 - Proof: cut-and-paste argument (see 6.046)



Key insight: re-use computation



Identical sub-problems! We can reuse our work!

Solution #1 – Memoization

- Create a big dictionary, indexed by aligned seqs
 - When you encounter a new pair of sequences
 - If it is in the dictionary:
 - Look up the solution
 - If it is not in the dictionary
 - Compute the solution
 - Insert the solution in the dictionary
- Ensures that there is no duplicated work
 - Only need to compute each sub-alignment once!

Top down approach

Solution #2 – Dynamic programming

- Create a big table, indexed by (i,j)
 - Fill it in from the beginning all the way till the end
 - You know that you'll need every subpart
 - Guaranteed to explore entire search space
- Ensures that there is no duplicated work
 - Only need to compute each sub-alignment once!
- Very simple computationally!

Bottom up approach

Identity score

Let (x,y) be an aligned pair of elements of two sequences (at least one of x,y must not be a gap).

$$\text{id}(x, y) = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$$

Score of an alignment = sum of scores of aligned pairs

TGK - G

AGKVG

$$0+1+1+0+1 = 3$$

60 % identical

Identity score = number of matched elements / length of alignment

Gap penalties

Consider two pairs of alignments:

ATCG
ATTG

and

AT – C G
AT T - G

They have the same
identity score but
alignment on the left is
more likely to be correct

ATC - - T A
ATT T T TA

and

AT - C - T A
AT T T T TA

- The first problem is corrected by introducing “gap penalty”.
- Second problem is corrected by introducing additional penalty for opening a gap.

Example

Score the above alignment using identity score; gap penalty = 1
Gap opening penalty = 2

ATCG
1+1+0+1=3 ATTG

AT – C G
AT T - G

1+1-2-1-2-1+1=-3

ATC - - T A
ATT T T TA

AT - C - T A
AT T T T TA

1+1+0-2-1-1+1+1=0

1+1-2-1+0-2-1+1+1=-2

Problems with identity score

- In the two pairs of aligned sequence below there are mutations at the first and 6th position and insertion (or deletion) on the 4th position. However while V and A share significant biophysical similarity and we often see mutation between them, W and A do not often substitute one for the other.

VGK – GI... **W**GK – GI...

AGKVGL... **A**GKVGL

- What if I mutated to V and then back to I should this have the same score as when I was unchanged?

If we will like to use the score to estimate evolutionary distances it would be important to take into account possibility of such double mutation.

Second base

		Second base					
		U	C	A	G		
First base	U	UUU } Phenyl- alanine F UUC } UUA } Leucine L UUG }	UCU } UCC } Serine S UCA } UCG }	UAU } Tyrosine Y UAC } UAA } Stop codon UAG } Stop codon	UGU } Cysteine C UGC } UGA } Stop codon UGG } Tryptophan W	U	C
	C	CUU } CUC } Leucine L CUA } CUG }	CCU } CCC } Proline P CCA } CCG }	CAU } Histidine H CAC } CAA } Glutamine Q CAG }	CGU } CGC } Arginine R CGA } CGG }	U	C
	A	AUU } Isoleucine I AUC } AUA } AUG } Methionine start codon M	ACU } ACC } Threonine T ACA } ACG }	AAU } Asparagine N AAC } AAA } Lysine K AAG }	AGU } Serine S AGC } AGA } Arginine R AGG }	U	C
	G	GUU } GUC } Valine V GUA } GUG }	GCU } GCC } Alanine A GCA } GCG }	GAU } Aspartic acid D GAC } GAA } Glutamic acid E GAG }	GGU } GGC } Glycine G GGA } GGG }	U	C
						Third base	

Scoring Matrices

An amino-acid scoring matrix is a 20x20 table such that position indexed with amino-acids so that position X,Y in the table gives score of aligning amino-acid X with amino-acid Y

Identity matrix – Exact matches receive one score and non-exact matches a different score (1 on the diagonal 0 everywhere else)

Mutation data matrix – a scoring matrix compiled based on observation of protein mutation rates: some mutations are observed more often than other (PAM, BLOSUM).

Not used:

Physical properties matrix – amino acids with similar biophysical properties receive high score.

Genetic code matrix – amino acids are scored based on similarities in the coding triple.

Principles of Dynamic programming

- Need to figure out how to use solution to smaller problems for solving larger problem.
- We need to keep a reasonable bound on how many sub-problems we solve
- Make sure that each sub-problem is solved only once

Dynamic programming algorithm for computing the score of the best alignment

For a sequence $S = a_1, a_2, \dots, a_n$ let $S_j = a_1, a_2, \dots, a_j$

S, S' – two sequences

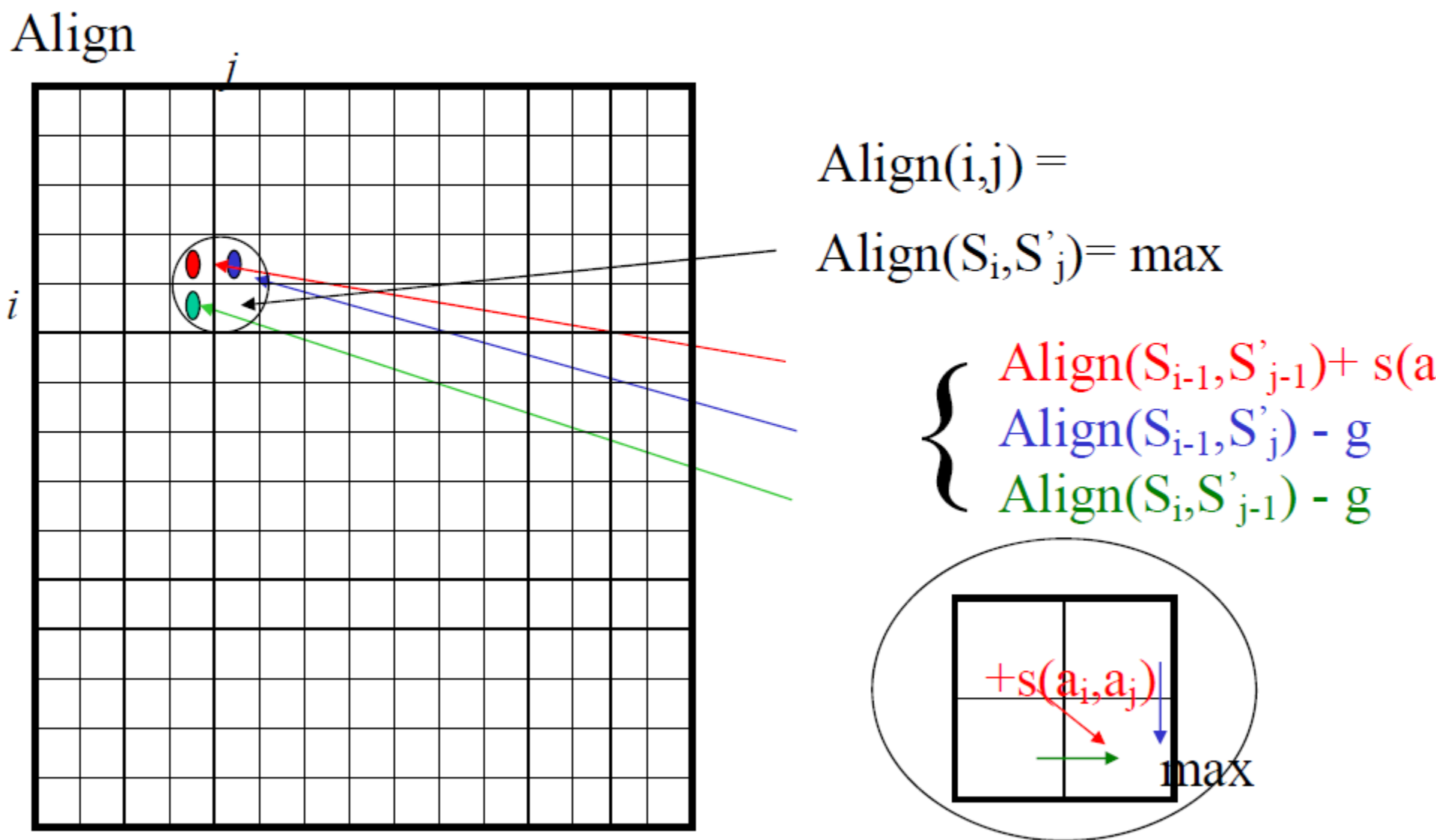
$\text{Align}(S_i, S'_j)$ = the score of the highest scoring alignment between S^1_i, S^2_j

$S(a_i, a'_j)$ = similarity score between amino acids a_i and a_j given by a scoring matrix like PAM, BLOSUM

g – gap penalty

$$\text{Align}(S_i, S'_j) = \max \begin{cases} \text{Align}(S_{i-1}, S'_{j-1}) + S(a_i, a'_j) \\ \text{Align}(S_i, S'_{j-1}) - g \\ \text{Align}(S_{i-1}, S'_j) - g \end{cases}$$

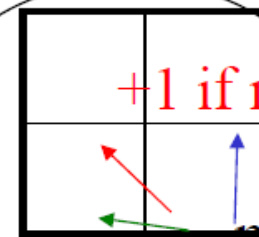
Organizing the computation – dynamic programming table



Example of DP computation with $g = 0$; match = 1; mismatch=0 Maximal Common Subsequence

initialization

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



+1 if match else 0

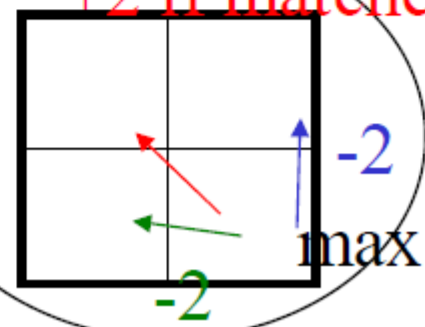
max

Example of DP computation with $g = 2$ match = 2; mismatch = -1

Initialization (penalty for starting with a gap)

	A	T	T	G	C	G	C	G	C	A	T	
A	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
T	-2	2	0	-2								
T	-4	0	4									
G	-6			6								
C	-8											
T	-10											
T	-12											
A	-14											
A	-16											
C	-18											
C	-20											
A	-22											

+2 if matched



The iterative algorithm

$m = |S|; n = |S'|$

for $i \leftarrow 0$ to m do $A[i,0] \leftarrow -i * g$

for $j \leftarrow 0$ to n do $A[0,j] \leftarrow -j * g$

for $i \leftarrow 1$ to m do

 for $j \leftarrow 1$ to n

$A[i,j] \leftarrow \max ($

$A[i-1,j] - g$

$A[i-1,j-1] + s(i,j)$

$A[i,j-1] - g$

$)$

return($A[m,n]$)

Complexity of the algorithm

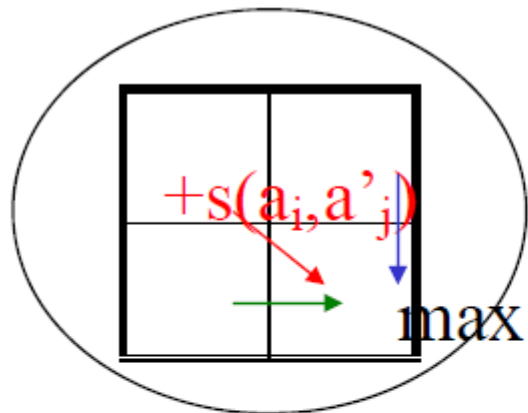
- Time $O(nm)$; Space $O(nm)$ where n , m the lengths of the two sequences.
- Space complexity can be reduced to $O(n)$ by not storing the entries of dynamic programming table that are no longer needed for the computation (keep current row and the previous row only).

From computing the score to computing of the alignment

Desired output:

Sequence of substitutions/insertion/deletions leading to the optimal score.

ATTGCGTTATAT
AT- GCG- TATAT



Red direction = match

Blue direction = gap in horizontal sequence

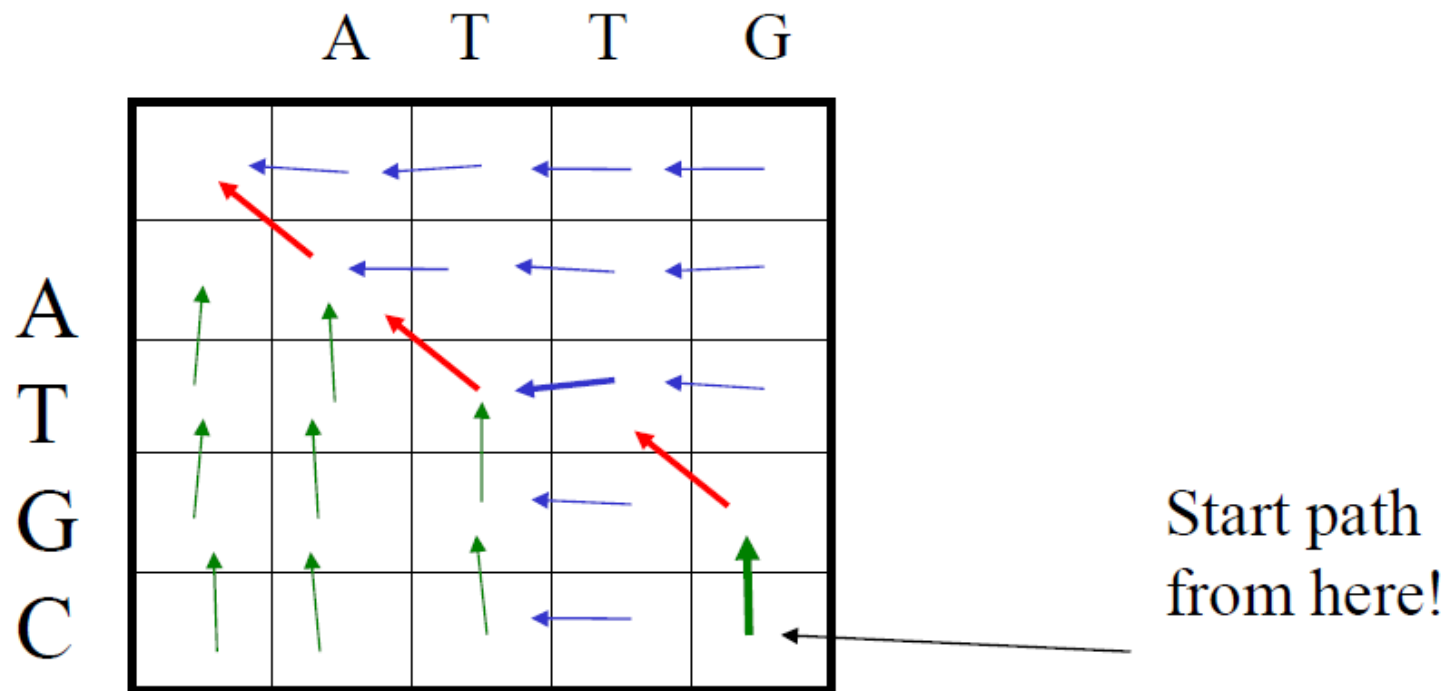
Green direction = gap in vertical sequence

a_1, a_2, \dots, a_i
 a'_1, a'_2, \dots, a'_j

$a_1, a_2, \dots, a_i -$
 a'_1, a'_2, \dots, a'_j

a_1, a_2, \dots, a_i
 $a'_1, a'_2, \dots, a'_j -$

Recovering the path



A T T G -

A T - G C

If at some position several choices lead to the same max value, the path need not be unique.

Local alignment

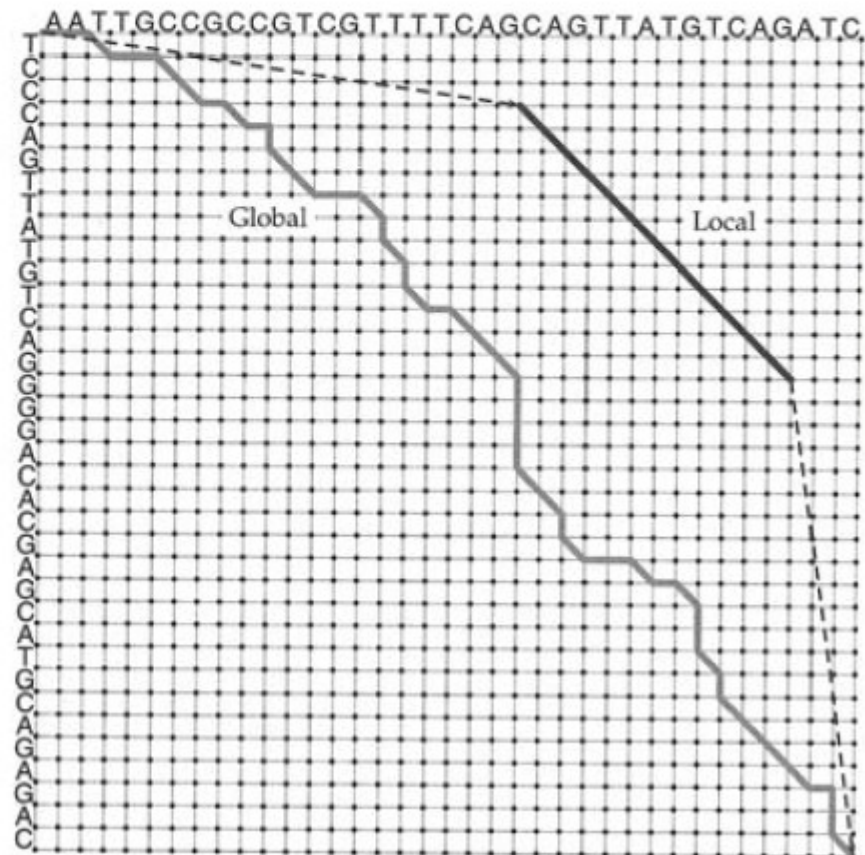
- The alignment techniques considered so far worked well for sequences which are similar over all their length
- This does not need to be the case: example gene from hox family have very short but highly conserved subsequence – the so called hox domain.
- Considered so far global alignment methods (that is algorithm that try to find the best alignment over whole length can miss this local similarity region

Global

```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-GAC
|  | | | | | | | | | | | | | | | |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG--T-CAGAT--C
```

Local

```
          tccCAGTTATGTCAGgggacacgagcatgcagaga
          | | | | | | | | | |
aattgccgccgctcggttttcagCAGTTATGTCAGatc
```



Local alignment (Smith, Waterman)

So far we have been dealing with **global alignment**.

Local alignment – alignment between substrings.

Main idea: If alignment becomes too bad – drop it.

Set p and g so that **alignment of random strings gives negative score**

$$a[i,j] = \max \left\{ \begin{array}{l} a[i-1,j-1] + s(a_i, a_j) \\ a[i-1,j] + g \\ a[i,j-1] + g \\ 0 \end{array} \right.$$

Finding the alignment: find the highest scoring cell and trace it back

Example

- Sequence 1 = ACACACTA
- Sequence 2 = AGCACACA

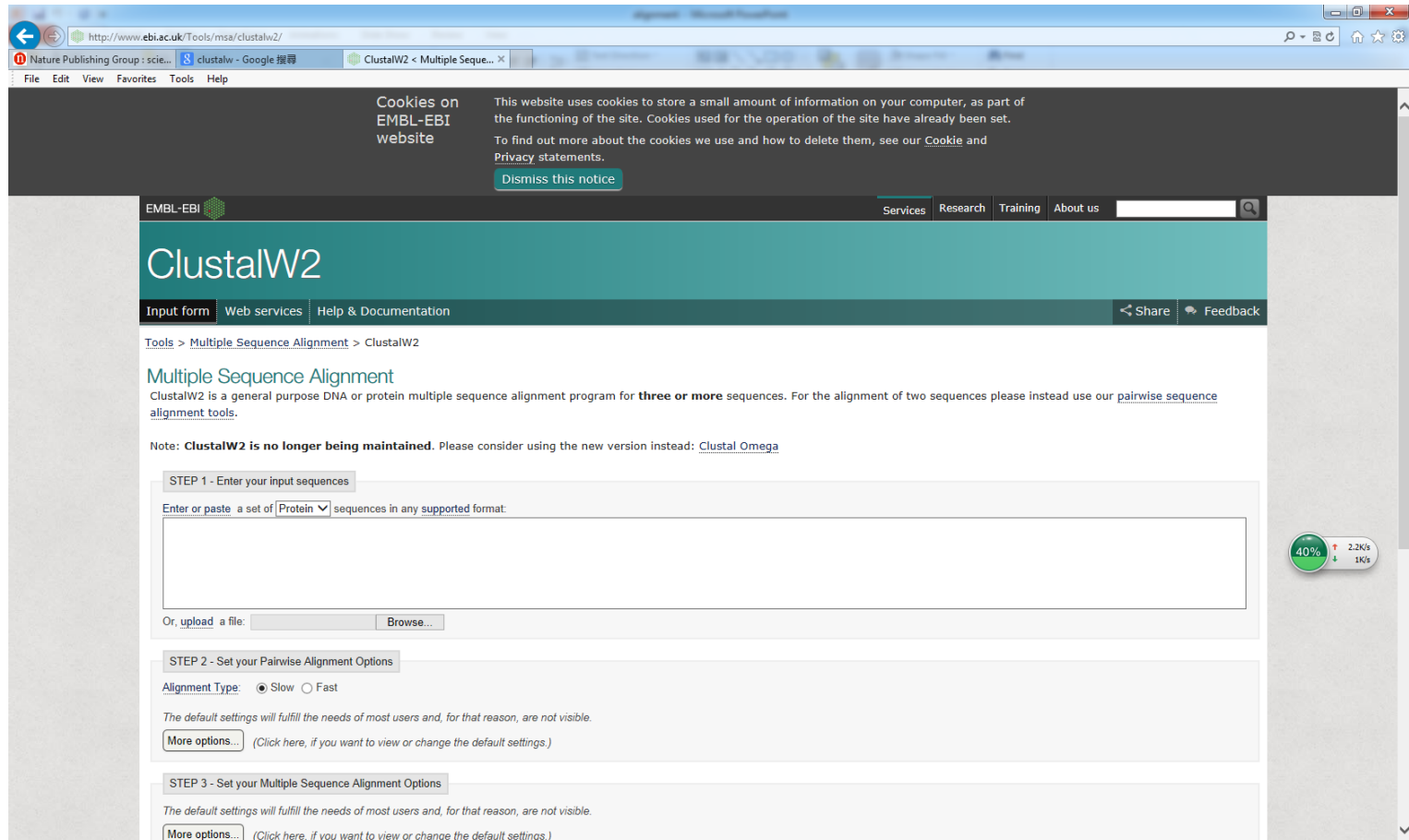
$s(a, b) = +2$ if $a = b$ (match), -1 if $a \neq b$ (mismatch)

$$H = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & 2 & 1 & 2 & 1 & 2 & 1 & 0 & 2 \\ G & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ C & 0 & 0 & 3 & 2 & 3 & 2 & 3 & 2 & 1 \\ A & 0 & 2 & 2 & 5 & 4 & 5 & 4 & 3 & 4 \\ C & 0 & 1 & 4 & 4 & 7 & 6 & 7 & 6 & 5 \\ A & 0 & 2 & 3 & 6 & 6 & 9 & 8 & 7 & 8 \\ C & 0 & 1 & 4 & 5 & 8 & 8 & 11 & 10 & 9 \\ A & 0 & 2 & 3 & 6 & 7 & 10 & 10 & 10 & 12 \end{pmatrix} \quad T = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow \\ G & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow \\ C & 0 & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow \end{pmatrix}$$

Sequence 1 = A--CACACTA

Sequence 2 = AGCACAC--A

ClustalW



The screenshot shows a web browser window with the URL <http://www.ebi.ac.uk/Tools/msa/clustalw2/>. The browser's address bar and tabs are visible at the top. A cookie notice banner is present, stating that the website uses cookies for site functionality and providing a link to the privacy statements. The main header of the website features the EMBL-EBI logo and navigation links for Services, Research, Training, and About us. The ClustalW2 title is prominently displayed in a teal banner. Below this, a navigation bar includes links for Input form, Web services, and Help & Documentation, along with Share and Feedback options. The breadcrumb trail indicates the path: Tools > Multiple Sequence Alignment > ClustalW2. The main content area is titled 'Multiple Sequence Alignment' and describes ClustalW2 as a general purpose DNA or protein multiple sequence alignment program for three or more sequences. A note informs users that ClustalW2 is no longer being maintained and suggests using Clustal Omega instead. The interface is divided into three steps: Step 1 - Enter your input sequences, Step 2 - Set your Pairwise Alignment Options, and Step 3 - Set your Multiple Sequence Alignment Options. Step 1 includes a text area for pasting sequences, a dropdown menu for selecting the sequence type (currently set to Protein), and a file upload option. Step 2 shows the Alignment Type set to Slow. Step 3 is partially visible. A system tray icon on the right side of the browser window indicates a 40% battery level.

http://www.ebi.ac.uk/Tools/msa/clustalw2/

Nature Publishing Group : scie... clustalw - Google 搜尋 ClustalW2 < Multiple Sequ...

File Edit View Favorites Tools Help

Cookies on EMBL-EBI website

This website uses cookies to store a small amount of information on your computer, as part of the functioning of the site. Cookies used for the operation of the site have already been set. To find out more about the cookies we use and how to delete them, see our [Cookie and Privacy statements](#).

[Dismiss this notice](#)

EMBL-EBI

Services Research Training About us

ClustalW2

Input form Web services Help & Documentation

Share Feedback

Tools > Multiple Sequence Alignment > ClustalW2

Multiple Sequence Alignment

ClustalW2 is a general purpose DNA or protein multiple sequence alignment program for **three or more** sequences. For the alignment of two sequences please instead use our [pairwise sequence alignment tools](#).

Note: **ClustalW2 is no longer being maintained**. Please consider using the new version instead: [Clustal Omega](#)

STEP 1 - Enter your input sequences

Enter or paste a set of Protein sequences in any supported format:

Or, upload a file: [Browse...](#)

STEP 2 - Set your Pairwise Alignment Options

Alignment Type: ☒ Slow ☐ Fast

The default settings will fulfill the needs of most users and, for that reason, are not visible.

[More options...](#) (Click here, if you want to view or change the default settings.)

STEP 3 - Set your Multiple Sequence Alignment Options

The default settings will fulfill the needs of most users and, for that reason, are not visible.

[More options...](#) (Click here, if you want to view or change the default settings.)

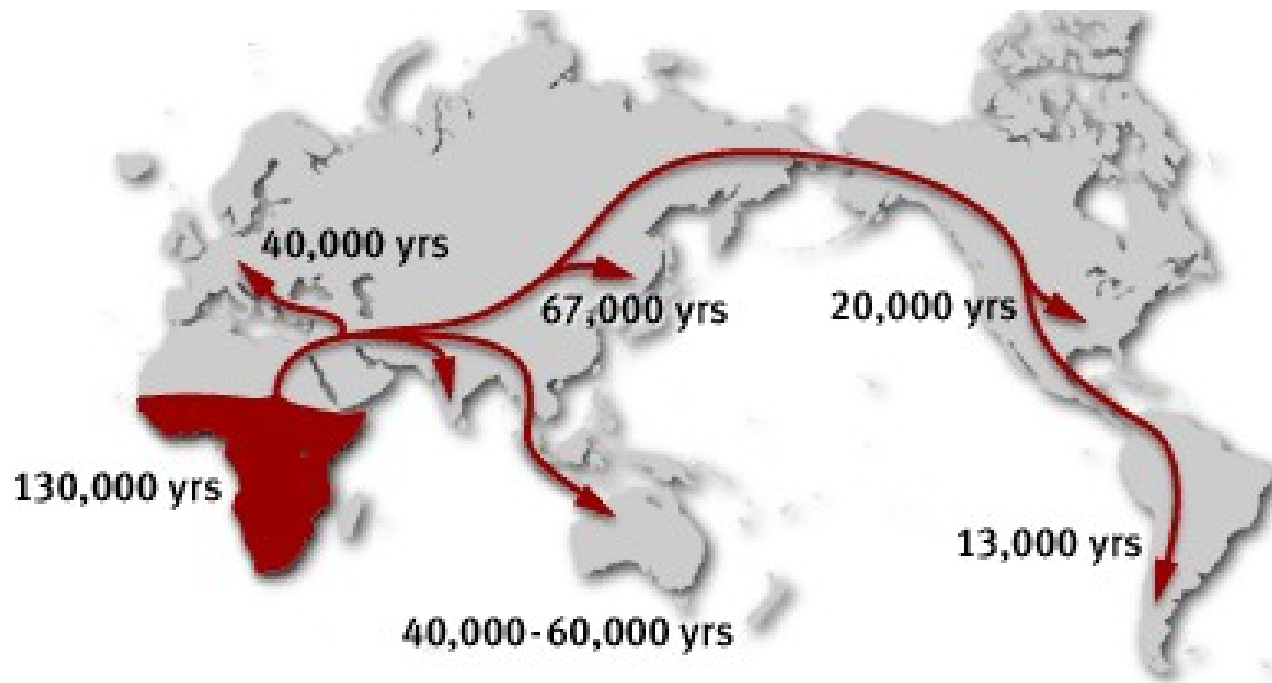
40% 2.2K/s 1K/s

Human Migration

Legend:
 □ Fossil or artifact site
 40,000 years ago
 Migration date
 → Generalized route

Key locations and dates:
 1. Omo Kibish, Ethiopia (200,000 years ago)
 2. Khadek River Mouth, China (100,000 years ago)
 3. Yana River, Siberia (30,000 years ago)
 4. Kostenki, Russia (40,000 years ago)
 5. Zhenkoudian (Shandong), China (15,000 years ago)
 6. Misotogawa, Japan (10,000 years ago)
 7. Lake Mungo, Australia (40,000 years ago)
 8. Monte Verde, South America (15,000 years ago)

Other locations marked: Spirit Cave (8,000-10,000 years ago), Kennewick (8,000 years ago), Clovis (13,000 years ago), Meadowcroft (14,000-15,000 years ago), Yana River (30,000 years ago), Kostenki (40,000 years ago), Zhenkoudian (Shandong) (15,000 years ago), Misotogawa (10,000 years ago), Omo Kibish (200,000 years ago), Khadek River Mouth (100,000 years ago), Yana River (30,000 years ago), Kostenki (40,000 years ago), Zhenkoudian (Shandong) (15,000 years ago), Misotogawa (10,000 years ago), Lake Mungo (40,000 years ago), Monte Verde (15,000 years ago).



- Evidence
 - Fossils
 - DNA sequence analysis

Homework

- **Insulin** is a peptide hormone, produced by beta cells of the pancreas, and is central to regulating carbohydrate and fat metabolism in the body.

>gi|386828|gb|AAA59172.1| insulin [Homo sapiens]

MAWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

>gi|82749718|gb|ABB89743.1| preproinsulin 1 [Rattus losea]

MALWMRFLPLLALLVWVEPKPAQAFVKQHLCGPHLVEALYLVCGERGFFYTPKSRREVEDPQVPQLELGGSPGAGDLQTLALEVARQKRGIVDQCCTSICSLYQLENYCN

>gi|402892452|ref|XP_003909428.1| PREDICTED: insulin isoform 4 [Papio anubis]

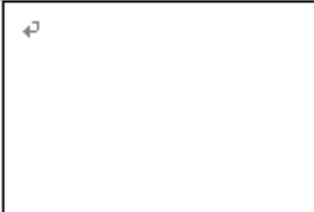
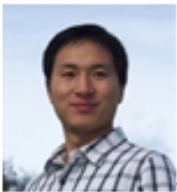







MALWMRLLPLLALLALWGPDSVPAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDPQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

>gi|344307503|ref|XP_003422420.1| PREDICTED: insulin-like isoform 1 [Loxodonta africana]

MALWTRLLPLLALLAVGAPPPARAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREVEDTQVGEVELGTGLQPFPAEAPKQKRGIVEQCCTGVCSLYQLENYCN

Identity score of insulin



Demo

- Write your own program, and demo in class
- A .fasta file will be provided to test your program
- Your program must be able to read fasta files

Score	Criteria
100	1, wrote a alignment program 2, demon in class 3, correctly calculate the identity score
80	1, wrote a alignment program 2, demon in class
0	1, no programming or copy other's program

Extra

- Write a similar program to estimate the similarity of two articles.
- Estimate the similarity of your program with others' program.