

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment  
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## Introduction to Alignments

Lec'06' slides

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## 1 More of the big picture

## 2 Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## 3 Summary

## 4 Reading for next class

# Classic example of paralogs and pseudogenes

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

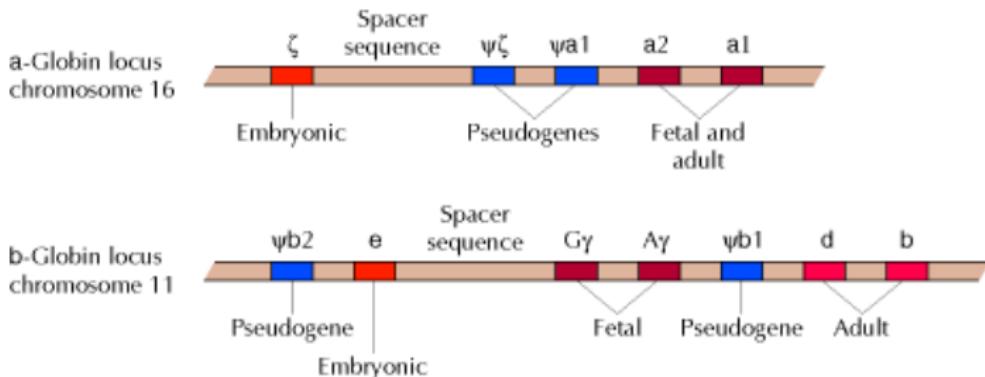
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class



**Figure 4.5. Globin gene families** Members of the human  $\alpha$ - and  $\beta$ -globin gene families are clustered on chromosomes 16 and 11, respectively. Each family contains genes that are specifically expressed in embryonic, fetal, and adult tissues, in addition to nonfunctional gene copies (pseudogenes).

## Review: Key idea from last lecture

## More of the big picture

- BLAST looks for **sequence similarity**
    - Example: AAACCG is similar to TAACCG
  - Homology + evolution creates **sequence similarity**

83% identity

Human catggccctgtggatgcgcctcctgccccctgctggcgctgctggccc  
||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||  
Mouse catggccctgtggatgcgcctcctgccccctgctggccctgctctcc

Therefore, we look for **sequence similarity** as evidence of homology.

Therefore, we use BLAST to find evidence of homology.

# Key idea of bioinformatics: Alignment → Similarity score → Homology

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

Line of reasoning:

- ① We have a DNA or protein sequence
- ② We **align** it against another DNA or protein sequence
- ③ We measure the similarity or **alignment score**
- ④ If the score is strong, we have **evidence for homology**.

But wait:

- ① What's the difference between sequence similarity and homology?
- ② If two sequences are similar, are the genes homologous?
- ③ How can we use this relationship? Why is this so important?

# What BLAST is really about

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

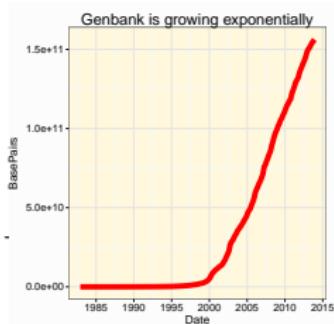
Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

- ① I have a DNA or protein sequence
- ② Because I'm a biologist,  
I want to find homologs
- ③ GenBank & Refseq & others have  
lots of sequence data → → → → →
- ④ Align (BLAST) my sequence against  
each sequence in database
- ⑤ High scoring alignments may be homologues.



More of the  
big picture

## Sequence Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

### ① More of the big picture

### ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

### ③ Summary

### ④ Reading for next class

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# Two DNA sequences optimally aligned

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

FASTA input:

```
>Rat
CAAGTA

>Cow
CAATTAA
```

align

Alignment output:

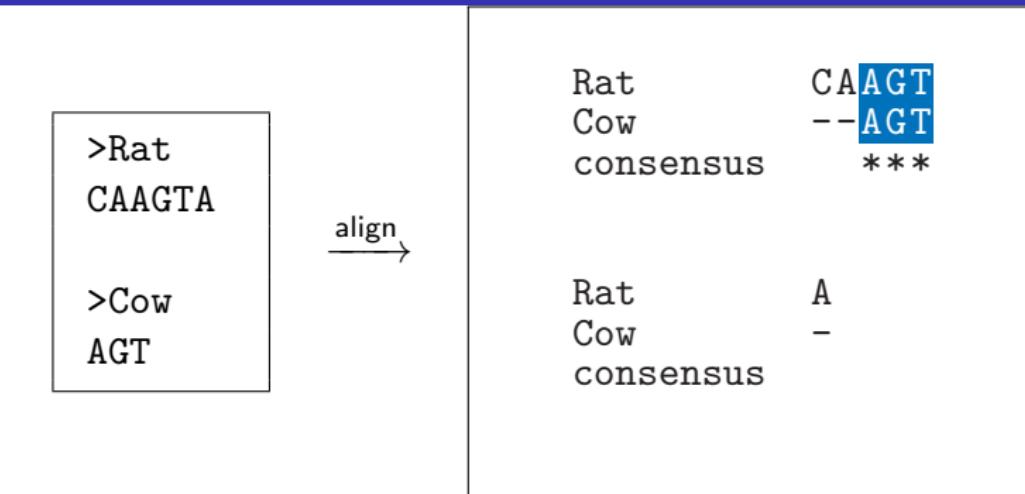
Rat	CAAGT
Cow	CAATT
consensus	**** *

Rat	A
Cow	A
consensus	*

How would you do this by hand?

# Slightly more complicated example

More of the  
big picture  
  
Sequence  
Alignment  
  
Principles of  
Sequence  
Alignment  
  
Scoring  
Alignments  
  
Substitution  
Matrices  
  
Inserting Gaps  
  
Types of  
Alignment  
  
Summary  
  
Reading for  
next class



How would you do this by hand?

How would you score this alignment's quality?

# Gaps hugely complicate finding optimal alignment

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

>Rat  
THISSEQUENCE

>Cow  
THISISASEQUENCE

align  
→

Rat	THIS	---	SEQUENCE	
Cow	THIS	IS	A	SEQUENCE
consensus	****	*****	*****	

How would you do this by hand?

How would you score this alignment's quality?

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# We can score alignment quality

More of the  
big picture  
  
Sequence  
Alignment  
  
Principles of  
Sequence  
Alignment  
  
Scoring  
Alignments  
  
Substitution  
Matrices  
  
Inserting Gaps  
  
Types of  
Alignment  
  
Summary  
  
Reading for  
next class

>Rat  
CAAGTA  
  
>Cow  
AGT

align →

Good alignment

$$\frac{3}{6} = 50\% \text{ identity}$$

Rat	CAAGT
Cow	--AGT
consensus	***

Rat	A
Cow	-
consensus	

Bad alignment

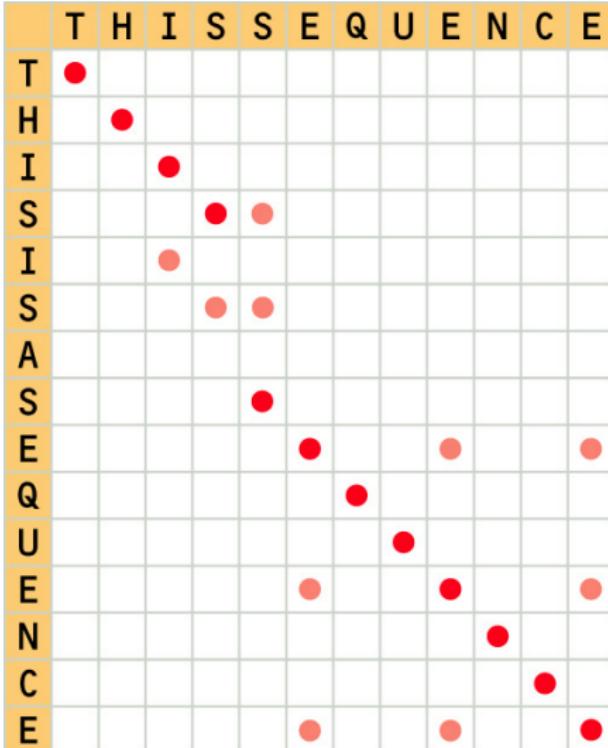
$$\frac{1}{6} = 17\% \text{ identity}$$

Rat	CAAGT
Cow	-AGT-
consensus	*

Rat	A
Cow	-
consensus	

# Dotplots align two sequences

More of the  
big picture  
  
Sequence  
Alignment  
  
Principles of  
Sequence  
Alignment  
  
Scoring  
Alignments  
  
Substitution  
Matrices  
  
Inserting Gaps  
  
Types of  
Alignment  
  
Summary  
  
Reading for  
next class



Red: Sequence match  
  
Pink: accidental match

# Remove noise with window and minimum identity

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment  
Scoring  
Alignments

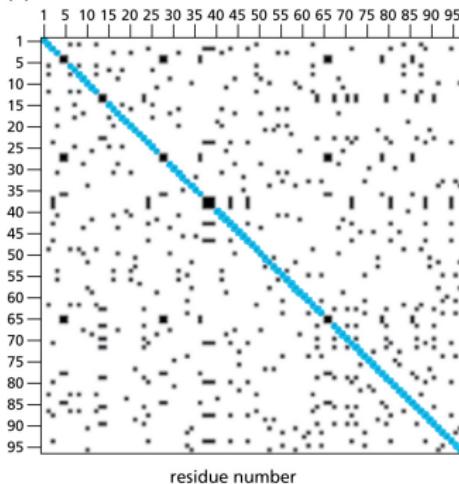
Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

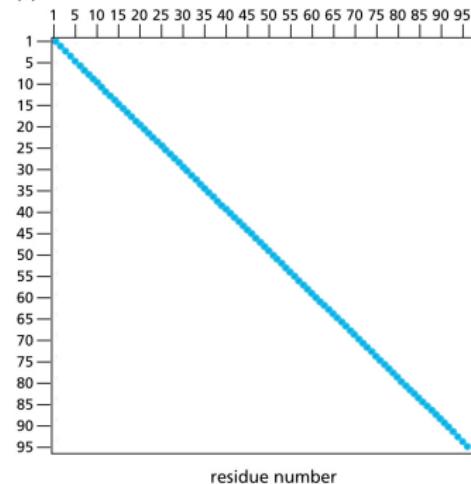
Summary

Reading for  
next class

(A)



(B)



# Chimp vs human dotplots, 2 different chromosomes

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment  
Scoring  
Alignments

Substitution  
Matrices

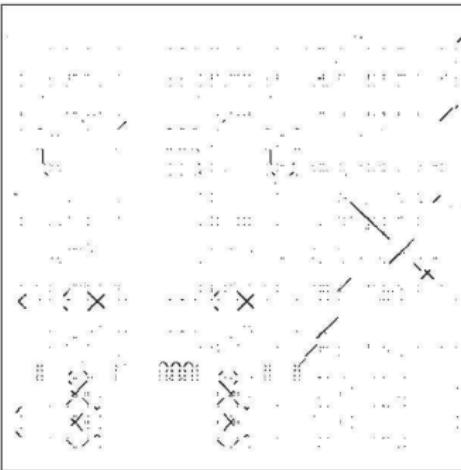
Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

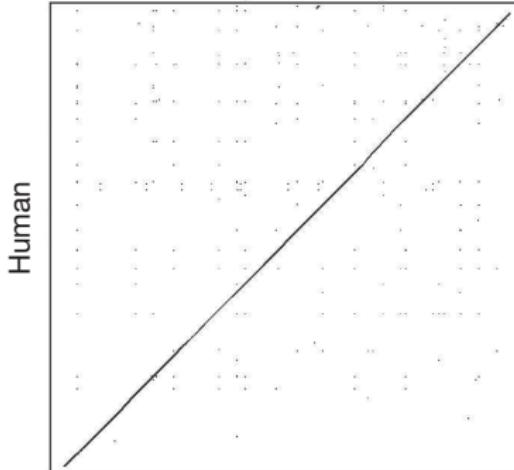
Human

Y chromosome



Chimpanzee

Chromosome 21

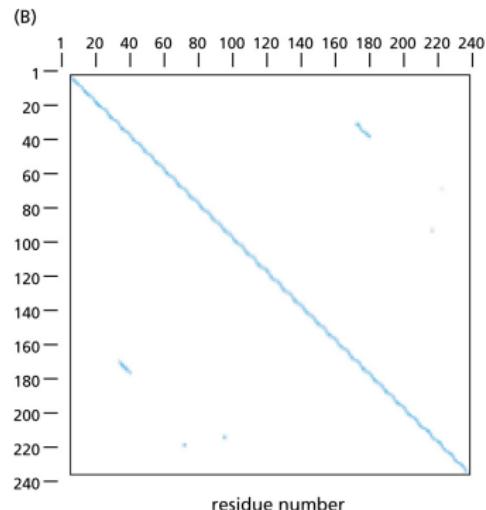
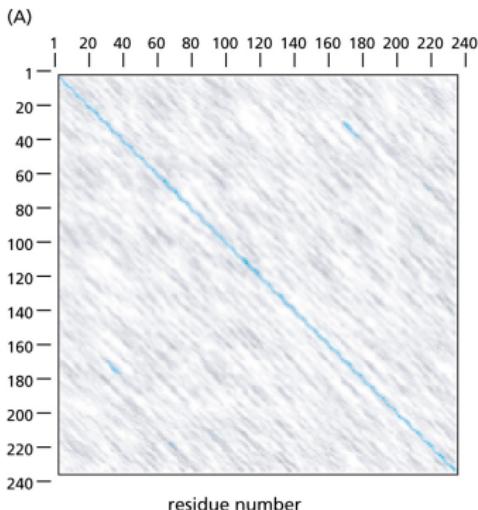


Chimpanzee

What's going on here?

# Dotplots reveal repeats

More of the  
big picture  
  
Sequence  
Alignment  
  
Principles of  
Sequence  
Alignment  
  
Scoring  
Alignments  
  
Substitution  
Matrices  
  
Inserting Gaps  
  
Types of  
Alignment  
  
Summary  
  
Reading for  
next class



What's going on here biologically?

# What's going on here?

More of the  
big picture

Sequence  
Alignment

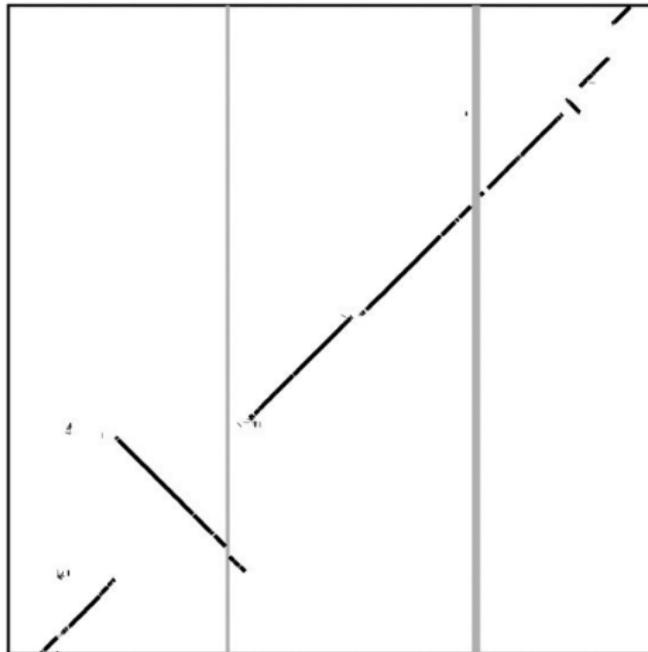
Principles of  
Sequence  
Alignment  
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class



More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## Traditional dot plotters:

- <http://myhits.isb-sib.ch/cgi-bin/dotlet>
- <http://www.vivo.colostate.edu/molkit/dnadot/>
- <http://bioinfo.lifl.fr/yass/yass.php>

Try BLAST 2 pairwise sequences at NCBI:

- <http://blast.ncbi.nlm.nih.gov/Blast.cgi>  
Click: **Align two or more sequences**

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

**Substitution Matrices**

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# Score can include amino acid similarity

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

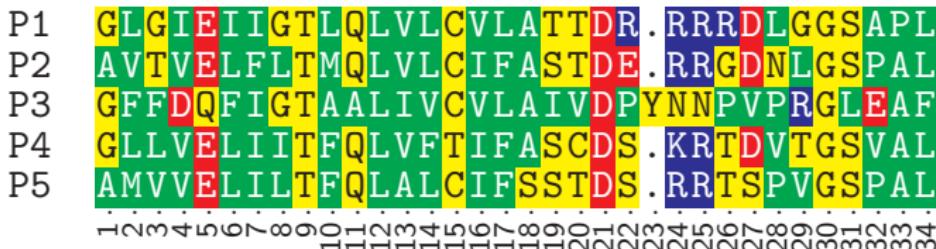
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class



- X acidic (-)
- X basic (+)
- X polar uncharged
- X hydrophobic nonpolar

What biological process happened here?

Hint: Look down column 2.

# PAM120 substitution matrix

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
9																			
-1	3																		
-3	2	4																	
-3	1	-1	6																
-3	1	1	1	3															
-5	1	-1	-2	1	5														
-5	1	0	-2	0	0	4													
-7	0	-1	-2	0	0	2	5												
-7	-1	-2	-1	0	-1	1	3	5											
-7	-2	-2	0	-1	-3	0	1	2	6										
-4	-2	-3	-1	-3	-4	2	0	-1	3	7									
-4	-1	-2	-1	-3	-4	-1	-3	-3	1	1	6								
-7	-1	-1	-2	-2	-3	1	-1	-1	0	-2	2	5							
-6	-2	-1	-3	-2	-4	-3	-4	-4	-1	-4	-1	0	8						
-3	-2	0	-3	-1	-4	-2	-3	-3	-3	-4	-2	-2	1	6					
-7	-4	-3	-3	-3	-5	-4	-5	-4	-2	-3	-4	-4	3	1	5				
-2	-2	0	-2	0	-2	-3	-3	-3	-3	-3	-4	1	3	1	5				
-6	-3	-4	-5	-4	-5	-4	-7	-6	-6	-2	-4	-6	-1	0	0	-3	8		
-1	-3	-3	-6	-4	-6	-2	-5	-4	-5	-1	-6	-6	-4	-2	-3	-3	4	8	
-8	-2	-6	-7	-7	-8	-5	-8	-8	-6	-5	1	-5	-7	-7	-5	-8	-1	12	

See any trends or patterns?

# PAM = Point Accepted Mutation

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

- Margaret Dayhoff, 1978 – foundational work but still used, still OK
- Aligned closely-related proteins, counted amino acid substitutions at each position
- Family of matrices:
  - PAM1 = 1/100 AA changed
    - = short evolutionary span
    - = fairly intolerant of evolution
  - PAM250 = 250/100 AA have changed
    - = distant evolution span
    - = more tolerant of distant evolution

# BLOSUM-62 substitution matrix

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
9																			
-1	4																		
-1	1	5																	
-3	-1	-1	7																
0	1	0	-1	4															
-3	0	-2	-2	0	6														
-3	1	0	-2	-2	0	6													
-3	0	-1	-1	-2	-1	1	6												
-4	0	-1	-1	-1	-2	0	2	5											
-3	0	-1	-1	-1	-2	0	0	2	5										
-3	-1	-2	-2	-2	-2	1	-1	0	0	8									
-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5								
-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5						
-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4					
-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4				
-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	3	2	1	3	1	4			
-2	-2	-2	-4	-2	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			
-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	3	7		
-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	
																		11	

small polar

small nonpolar

polar or acidic

basic

large  
nonpolar

aromatic

# BLOSUM = BLOcks SUbstitution

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments  
Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

- Henikoff & Henikoff 1992, 1996, very widely used
- Larger sample size than PAM
- Looked at distantly-related *blocks*
- Larger numbers: more distantly-related proteins:  
BLOSUM45 captures greater evolutionary distance than  
BLOSUM90

# Recommended substitution matrices for BLAST

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

<b>Query length</b>	<b>Substitution matrix</b>	<b>Gap costs</b>
< 35	PAM-30	(9,1)
35 – 50	PAM-70	(10,1)
50 – 85	BLOSUM-80	(10,1)
> 85	BLOSUM-62	(10,1)

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# Indels = Insertions or Deletions

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

>Rat  
THISSEQUENCE

>Cow  
THISISASEQUENCE

align  
→

Rat	THIS	---	SEQUENCE
Cow	THIS	IS A	SEQUENCE
consensus	****	*****	*****

What happened here biologically?  
Hint: Two obvious possibilities...

# High gap insertion penalty versus low penalty

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment  
Scoring  
Alignments  
Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

(A)

Bovine PI-3Kinase p110a      LNWENPDIMSELLFQNNEIIFKNGDDLRQMLTQIIRIMENIWQNGDLRMLPYGCLSIGDCVGLIEVRNRSHTIMQIQCKGGKLGAL  
cAMP-dependent protein kinase --WENPAQNTAHLDFERIKTLGTGSFGRVMVKHMETGNHYAMKILDKQKVVKLQIEHTLNEKRILQAYNPFVLKLEFSKDNLY

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDKNKGEIYDAAIDLFTRSCAGYCVATFILGIGDRHNSNIMVKDGGQLFHIDFGHFLDHKKKK  
cAMP-dependent protein kinase MVMEYVPGEMFSHLRRIGFSEHARFYAQIVLHSDLDKLPKENLQQQGYIQVTDFGFAKRVKGRTWXCGTPEYLAP

Bovine PI-3Kinase p110a      LIVSKGQAQECTKTREFERFQEMCYKAYLAIRQHANLFINLFSMMLGSGMPELQSFDIAYIRKTLABDK  
cAMP-dependent protein kinase EIILSKGYNKAVDWALGVLIYEMAGYPPFFADQPIQEYKIVSGKVRFSHFSDLKDLRNLQVDLTKRFGNLKGVNDIKHG

Bovine PI-3Kinase p110a      WTTKMDWIFHTIKQHALN-----  
cAMP-dependent protein kinase ATTDWIAIYQRKVEAPFIPKFGPGDTSNFDYEEEIRVXIEKCGKESEFEF

(B)

Bovine PI-3Kinase p110a      LNWENPDIMSELLFQNNEIIFKNGDDLRQMLTQIIRIMENIWQNGDLRMLPYGCLSIGDCVGLIEVRNRSHTIMQIQCKGGKLGAL  
cAMP-dependent protein kinase ?WENPAQNTAHLDFERIKTLGTGSFGRVMVKH--ETGNHYAMKILDKQKVVKLQIEHTLNEKRILQAYNPFVLKLEFSKDNLY

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDKNKGEIYDAAIDLFTRSCAGYCVATFILGIGDRHNSNIMVKDGGQLFHIDFGHFLDHKKKK  
cAMP-dependent protein kinase -SNLYMVMEYVPGGEMFSHLRRIGFSEHARFYAQIVLHSDLDIYRDLKPENLIDQQQGYIQVTDFGFAKRVKGRTWXCGT

Bovine PI-3Kinase p110a      QDFL--IVSKGQAQECTKTREFERFQEMCYKAYLAIRQHANLFINLFSMMLGSGMPELQSFDIAYIRKTLABD  
cAMP-dependent protein kinase PEYLAPEILSKGYNKAVDWALGVLIYEMAGYPPFFADQPIYEKIVSGKVRF--SHFSDLKDLRNLQVDLTKR--FGNLKN

Bovine PI-3Kinase p110a      QMNDAHHGGTTKMDWI-----  
cAMP-dependent protein kinase GVNDIKNHKMFATTDWIAIYQRKVEAPFIPKFGPGDTSNFDYEEEIRVXIEKCGKESEFEF

# We can include gaps in the similarity score

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments  
Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

Example protein alignment scoring system:

- ① Lookup BLOSUM-62 score at each position.
- ② Add these scores.
- ③ Subtract 10 for each indel position.

A better protein alignment scoring system:

- ① Lookup BLOSUM-62 score at each position.
- ② Add these scores
- ③ Subtract 10 for each new string of gaps
- ④ Subtract 2 for each position in a gap.

Hard question:

Why does the second scoring system better reflect evolution?

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# Multi-domain proteins and alignment

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

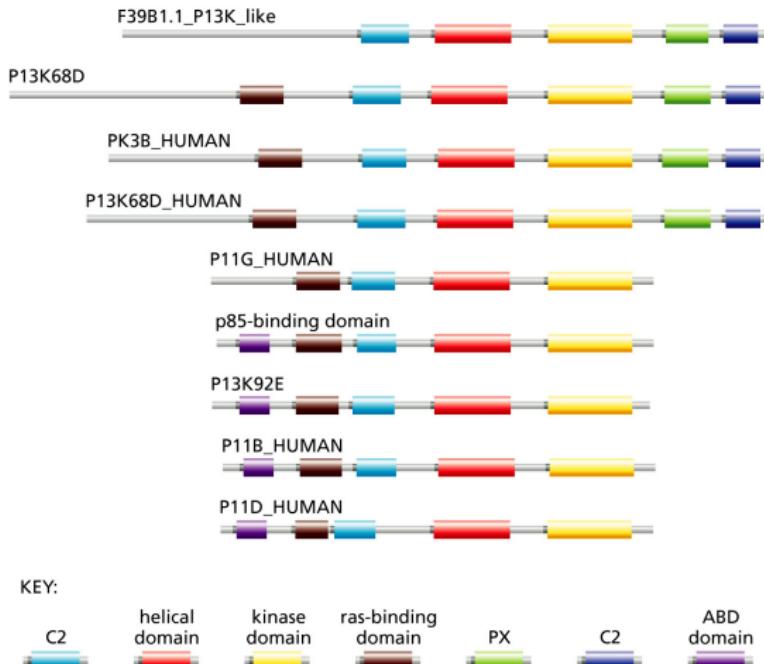
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class



# Local versus global alignment

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## • Local alignment

- Best-known algorithm: BLAST family
- Looks for **any** regions of similarity
- Ignores regions of dissimilarity
- Useful for fishing expedition
- Great for finding protein domains

## • Global alignment

- Best-known algorithm: CLUSTALW
- Aligns **entire** sequence
- Can produce poor results in multi-domain proteins
- May need manual adjustment after alignment
- Very useful with protein families, e.g., hemoglobin

# Example of local versus global alignment

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

(A) local (BLAST)

PI3-kinase DRHNSNIMVKDDGQLFHI DFG  
cAMP PK DLKPENNLIDQQGYIQT DFG

(B) global (CLUSTALW)

PI3-kinase HQLGNLR--L**E**ECR**I**--MSSAKRPLWLWNWENP**D**IMSELLFQNNE**I**FKNGDDLRQ**D**M**L**  
cAMP PK GNAAAAKKG**X**E**Q**ESV**K**EFLAKAKED**F**LKKWENPAQNTAH**D**Q**F**ERIK**T**LTG**G**S**F**R**V****M**  
10 20 30 40 50

PI3-kinase L**Q**I**T**IR**I**ME--N**I**W**N**Q**N**Q**G****L**D**L**R**M**P**Y**G**C**L**S**I**G**D**C**V**G****L**E**V**V**R**N**S**H**T****I****M****D**-I**R**C**K****G****G****L****K****G****A**  
cAMP PK ---V**K**H**M**E**T**G**N****H****Y****A****M****K****I****L**D**K****Q****K****V****V**K-----L**K****Q****I****E****T****H****L****N****E****K****R****I****L****Q****A****V****N****F****P****L****V****K****L**  
60 70 80 90 100 110  
60 70 80 90 100

PI3-kinase Q**F**N**S**H**T****I****L**H**Q**W**L**K**D**K**N**K**G****E****I****Y****D****A****---**I**D****L**F**T**R**S****C****A****G****C**V**A****T****F****I****L****G****I****G****D****R****H****N****S****N****I****M****V****K****D****-****D**  
cAMP PK S**F**K**D**N**S****N****L****Y****M****V****E****Y****P****G****E****M****F****S****H****L****R****I****G****F****E****P****A****R****Y****A****Q****I****V****L****T****E****Y****L****H****S****L****D****I****Y****R****D**  
110 120 130 140 150 160

PI3-kinase G**Q**L**F**H**I****D****G****H****F**L**D****H****K****K****K****F****G****Y****K****R****E****R****V****P**-----F**V****L**T**Q****D****F****L**---I**V****I****S****K****G****A****Q****E****C****T****K****T****R****E****E**  
cAMP PK P**E**N**L****I****D****Q****G****Y****I****---**Q**V****T**D**F****G****FA****K**-R**V**K**G****R****T****W****X****L****G****T****P****E****Y****L****A****P****E****I****I****L****S****K****G****N****K****A****D****W****W****A****G**  
170 180 190 200 210 220

PI3-kinase R**F**-Q**E****M****---**Y**K****A****L****A****I****R****Q****H****A****N****L****F****I****N****L****S****M****L****G****S****G****P****E****L****Q****S****F****D****I****A****Y****I****R****K****T****A****L****D****K****T****E****Q****E****A**  
cAMP PK V**L****I****Y****E****M****A****G****Y****P****F****F****A****-**D**O****P****I****Q****I****Y****E****K****I****V****S****G****K****V****--**F**P****S****H****F****S****S****D****L****K****D****L****R****N****L****Q****V****D****L****T****K****R****--**  
230 240 250 260 270

PI3-kinase L**E****Y****F****M****K****R****M****N****D****A****H****G****G****W****T****K****M****D****W****I****---**-----F**H****T****I****K****Q****H****A****L****N****---**  
cAMP PK F**G****N****L****K****N****G****V****N****D****I****K****N****H****M****F****A****T****D****W****I****A****I****Y****Q****R****K****V****E****A****P****I****P****K****F****K****G****P****G****D****T****S****N****F****D****D****Y****E****E****E****E****I****R****V****X****I**  
290 300 310 320 330 340 350

Critical kinase residues:

- DLKPEN
- DFG

Local alignment  
found them

Global alignment  
didn't find them

When would we want  
local alignment?

When would we want  
global alignment?

# Pairwise alignment versus multiple alignment

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

Pairwise alignment: Aligns two sequences (BLAST)

AQP1nuc. SEQ	GCTGTCA-TGTATATCATCGCC
AQP2nuc. SEQ	GCTGCCT-TCTATGTTGGCTGCC

Multiple alignment: Aligns  $\geq$  two sequences (CLUSTALW)

AQP1nuc. SEQ	GCTGTCA-TGTATATCATCGCC
AQP2nuc. SEQ	GCTGCCT-TCTATGTTGGCTGCC
AQP3nuc. SEQ	GCTGCCCATCTACACACTGGCA
AQP4nuc. SEQ	TCTGTCT-TCTACATCACTGCG
AQP5nuc. SEQ	GCTGTCT-TCTACGTGGCAGGCC

# Multiple sequence alignment can discover conserved and critical residues

More of the big picture

Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps  
Types of Alignment

Summary

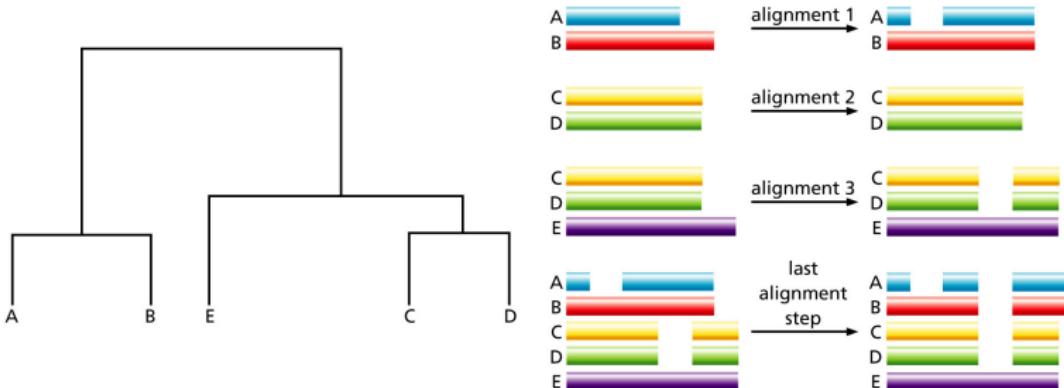
Reading for next class

(A)	p110 $\alpha$ cAMP-kinase	TFILGIG <b>D</b> RHNS <b>N</b> IMVKDDG-QLFHI <b>D</b> FGHFLDHKKKKFGYKRERVPFVLT--QDFLIVI 142 QIVLT <b>F</b> EYLHSLDLIY <b>R</b> DLK <b>P</b> EN <b>L</b> IDQQGYIQVT <b>D</b> FGFAKRVKGRTWLCGTPEYLAPE 179
(B)	p110 $\beta$ p110 $\delta$ p110 $\alpha$ p110 $\gamma$ p110_dicti cAMP-kinase	SYVLGIG----- <b>D</b> RHSD <b>N</b> INV <b>K</b> KT <b>G</b> QLFHI <b>D</b> FGHILGNFKSKFGIKRERVPFILT 136 TYVLGIG----- <b>D</b> RHSD <b>N</b> IMIRE <b>S</b> GQLFHI <b>D</b> FGHFLGNFKTKFGINRERVPFILT 136 TFILGIG----- <b>D</b> RHNS <b>N</b> IMVKDD <b>G</b> QLFHI <b>D</b> FGHFLDHKKKKFGYKRERVPFVLT 135 TFVLGIG----- <b>D</b> RHND <b>N</b> IMITET <b>G</b> NLFHI <b>D</b> FGHILGNYKSFLGINKERVPFVLT 135 TYVLGIG----- <b>D</b> RHND <b>N</b> LMV <b>T</b> KGG <b>R</b> LFHI <b>D</b> FGHFLGNYKKFGFKRERAPFVFT 135 QIVLT <b>F</b> EYLHSLDLIY <b>R</b> DLK <b>P</b> EN <b>L</b> IDQQGYIQVT <b>D</b> FGFAKRVKGRTWLCG--TPEYLA 177

What's going on here biologically?

# CLUSTALW operates in several steps

- More of the big picture
  - Sequence Alignment
    - Principles of Sequence Alignment
    - Scoring Alignments
    - Substitution Matrices
    - Inserting Gaps
    - Types of Alignment
  - Summary
  - Reading for next class
- ① Compute alignment scores between all sequence pairs
  - ② Take highest scoring pair, create a consensus from the pair
  - ③ Take next highest scoring pair, create a consensus sequence
  - ④ Repeat until all sequences are merged into a single consensus.



# Alignment programs produce different results

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## (A) structural/functional alignment from BAliBase

```
1csy SHEKMPWFHGKISREEESEQIVLIGSKTNGKFLIRARD--NNGSYALCLLHEGKVLHYRIDKDGTGKLSIPEGK-KFDTLWQLVEHYSYKA-----DGLLRVL-TVPCQK
1gri EMKHPWFFGKIPRAKAAEML-SQRHHDGAFLIRESES-APGDFSLSVKFGNDVQHFKVLRDGAGKYFL-WVV-FKNSLNELVDYHRSTS-S-VSRNQQIFLIRDIE&VPQQ-
1aya ---MRRWFHNPITGVEAEENLLL-TRGV-DGSFLARPSKS-NPGDFTLSVRNGAVTHIKIQNT--TDGYDLYGGEKFA-TLAELVQYYMEHHGQLKEKNGDVIEL-KYPLN-
2pna -LQDAEWYWGDISREEVENEKLDT--ADGTFLVRDASTKMHGDYTLTLRKGGNNKLKIFHRDGKY-GFSDPL-TFNSVELINHYRNES-LAQYNPKLDVKL-LYPVS-
1bfi HHDEKTNWVGSSNRNKAENLLRGK--RGDTFLVRESS--KQGCYACSVVVDGEVKHCVINKTATG-YGFAEPYNLYSSLKELVLHYQHTS-LVQHNDSLNVTL-AYPVYA
```

## (B) DIALIGN multiple sequence alignment

```
1csy SHEKMPWFHGKISREEESEQIVLIGSKTNGKFLIRAR-DN--NGSYALCLLHEGKVLHYRIDKDGTGKLSIPEGK-KFTDLWQLVEHYSYKA-----DGLLRVL-TVPCQK
1gri EMKHPWFFGKIPRAKAAEML-SQRHHDGAFLIRESES-A-PGDFSLSVKFGNDVQHFKVLRDGAGKYFL-WVV-K-FNSLNELVDYHRSTS-S-VSRNQQIFLIRDIE&VPQQ-
1aya M---R-RWFHNPITGVEAEENLLL-TRGV-DGSFLARPSKS-N-PGDFSLSVRFNGAVTHIKIANTGDYDLYG-GEK-FATLAELVQYYMEHHGQLKEKNGDVIEL-KYPLN-
2pna LQDAEWYWGDISREEVENEKL-RDTA-DGTFLVRDA-STKMHGDYTLTLRKGGNNKLKIFHRDGKYGFSD-PLT-FNSVVELINHYRN-E-SLAGYNPKLDVKL-LYPVS-
1bfi HHDEKTNWVGSSNRNKAENLL-RGKR-DGTFLVRRES-SK--QGCYACSVVVDGEVKHCVINKTATG-YGFAEPYNLYSSLKELVLHYQHTS-SLVQHNDSLNVTL-AYPVYA
```

## (C) ClustalW multiple sequence alignment

```
1csy SHEKMPWFHGKISREEESEQIVLIGSKTNGKFLIRARDN--NGSYALCLLHEGKVLHYRIDKDGTGKLSIPEGKKF-TDLWQLVEHYSYK-----ADGLLRVL-TVPCQK
1gri EMKHPWFFGKIPRAKAAEML-SQRHHDGAFLIRESES-A-PGDFSLSVKFGNDVQHFKVLRDGAGKYFL-WVVK-FNSLNELVDYHRSTS-S-VSRNQQIFLIRDIE&VPQQ-
1aya ---MRRWFHNPITGVEAEENLLL-TRGV-DGSFLARPSKS-N-PGDFFTLSVRNGAVTHIKIQNT-GDYYDLYGGEKFA-TLAELVQYYMEHHGQLKEKNGDVIEL-KYPLN-
2pna -LQDAEWYWGDISREEVN-EKLRDTADGTFLVRDASTKMHGDYTLTLRKGGNNKLKIFHRDGKY-GFSDPL-TFNSVELINHYRNES-LAQYNPKLDVKL-LYPVS-
1bfi HHDEKTNWVGSSNRNKAENL--NLLRGKRDGTFLVRESSK--QGCYACSVVVDGEVKHCVINKTATG-YGFAEPYNLYSSLKELVLHYQHTS-LVQHNDSLNVTLAYPVA
```

## (D) divide-and-conquer multiple sequence alignment

```
1csy SHEKMPWFHGKISREEESEQIVLIGSKTNGKFLIR-A-RDNN-GSYALCLLHEGKVLHYRIDKDGTGKLSIPEGK-KFTDLWQLVEH-Y-S-----KADGLLRVL-L-TVPCQK
1gri EMKHPWFFGKIPRAKAAEMLS-SQRHHDGAFLIRE-SESAPGDFSLSVKFGNDVQHFKVLRDGAGKYFL-WVVK-FNSLNELVDYH-RSTS-VSRNQQIFLIRDIE&VPQQ-
1aya ---MRRWFHNPITGVEAEENLLL-TRGV-DGSFLARPSKS-KSNPGDFTLSVRNGAVTHIKIQNTGDYDLYGGEK-FATLAELVQYYMEHHGQLKEKNGDVIEL-KYPLN-
2pna -LQDAEWYWGDISREEVENEKL-RDTADGTFLVRDASTKMHGDYTLTLRKGGNNKLKIFHRDGKY-GFSDPL-FNSVVELINHYRNES-LAQNPKLDVKL-LYPVS-
1bfi HHDEKTNWVGSSNRNKAENL--RGKRDGTFLVRE-SSKQ-GCYACSVVVDGEVKHCVINKTATG-YGFAEPYNLYSSLKELVLHYQHTS-SLVQHNDSLNVTLAYPVA
```

If the inputs are the same, how can this happen?

What's going on here?

# Why can't programs find the best alignment?

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

Two kinds of computational problems:

① Tractable:

- Runs in a “reasonable” amount of time
- Can find best solution

② Intractable:

- Requires huge amount of computer time to find best solution → **would take longer than age of universe.**
- Famous examples:
  - Traveling salesman problem
  - Multiple sequence alignment

Q: So what do we do?

A: Try for an *approximate answer* in a reasonable amount of time

→ **heuristic solution** → most of bioinformatics.

What are the practical implications?

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment

Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps

Types of  
Alignment

Summary

Reading for  
next class

## ① More of the big picture

## ② Sequence Alignment

Principles of Sequence Alignment

Scoring Alignments

Substitution Matrices

Inserting Gaps

Types of Alignment

## ③ Summary

## ④ Reading for next class

# Reading for next class

More of the  
big picture

Sequence  
Alignment

Principles of  
Sequence  
Alignment  
Scoring  
Alignments

Substitution  
Matrices

Inserting Gaps  
Types of  
Alignment

Summary

Reading for  
next class

## Chapter 3

Section “Pairwise Alignment and Limits of Detection”  
to  
End of Chapter

Pages	Notes
94–112	Read

## Chapter 4

Section “Introduction”  
to  
Section “Stand-alone BLAST”

Pages	Notes
121–135	Read