

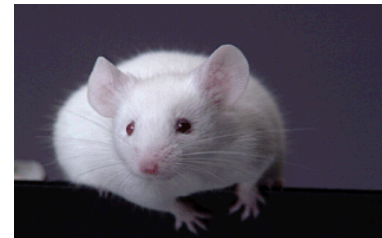
# Lecture outline

- Sequence alignment
  - Why do we need to align sequences?
  - Evolutionary relationships
- Gaps and scoring matrices
- Dynamic programming
  - Global alignment (Needleman & Wunsch)
  - Local alignment (Smith & Waterman)
- Database searches
  - BLAST
  - FASTA

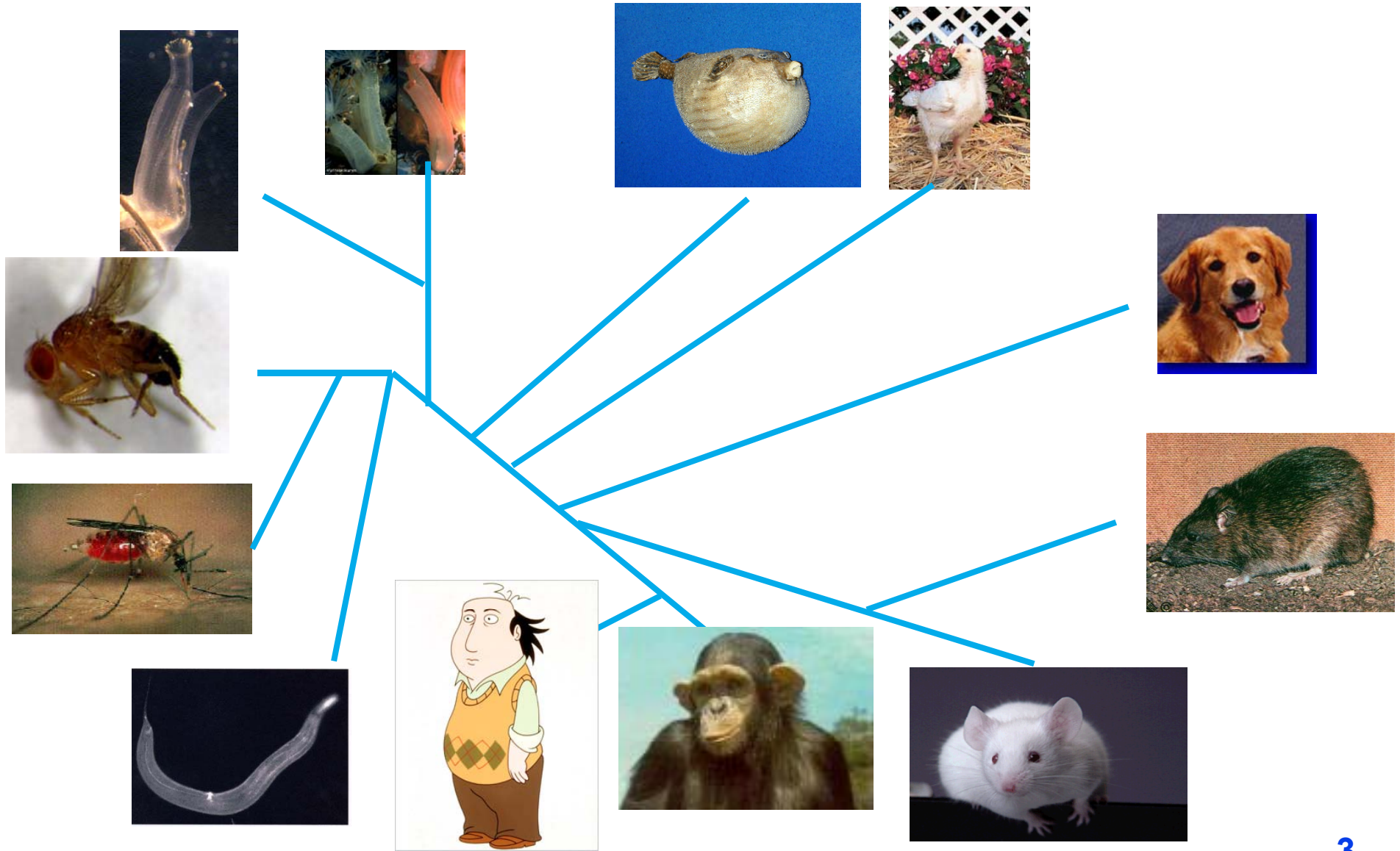
# Complete DNA Sequences



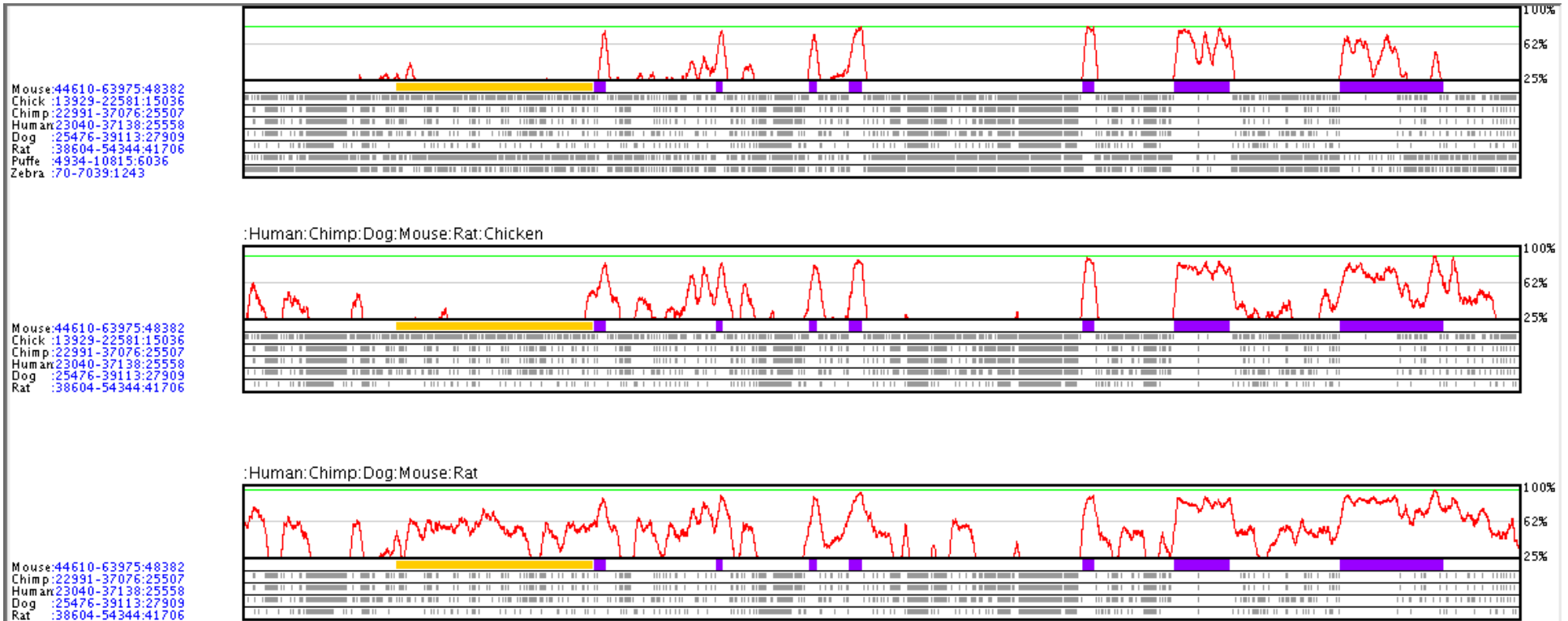
Whole genome  
sequencing projects  
for more than 2000  
species



# Evolution

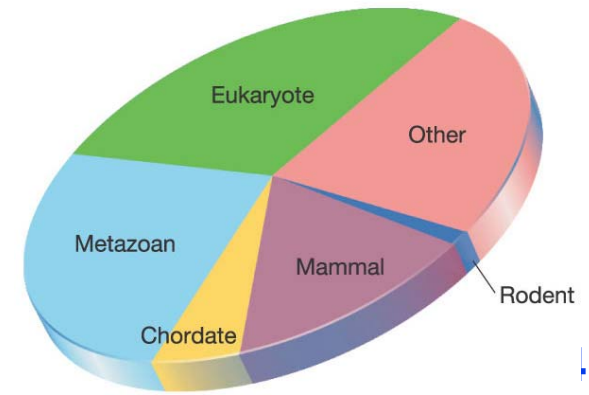


# Sequence conservation implies function



## Alignment is the key to

- Finding important regions
- Determining function
- Uncovering the evolutionary forces



# Sequence alignment

- Comparing DNA/protein sequences for
  - Similarity
  - Homology
- Prediction of function
- Construction of phylogeny
- Shotgun assembly
  - End-space-free alignment / overlap alignment
- Finding motifs

# Sequence Alignment

Procedure of comparing two (pairwise) or more (multiple) sequences by searching for a series of individual characters that are in the same order in the sequences

GCTAGTCAGATCTGACGCTA

| | | | | | | | | |

TGGTCACATCTGCCGC

# Sequence Alignment

Procedure of comparing two (pairwise) or more (multiple) sequences by searching for a series of individual characters that are in the same order in the sequences

```
VLS PADKTNVKA AWGKVG AHAGYEG
| | |   |   |   | | |   |   |
VLSEG DWQLVLHVWAKVEADVAGEG
```

# Sequence Alignment

AGGCTATCACCTGACCTCCAGGCCGATGCCC  
TAGCTATCACGACCGCGGGTCGATTTGCCCGAC

**-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---**  
**TAG-CTATCAC--GACCGC--GGTCGATTTGCCCGAC**

## Definition

Given two strings  $x = x_1x_2\dots x_M$ ,  $y = y_1y_2\dots y_N$ ,

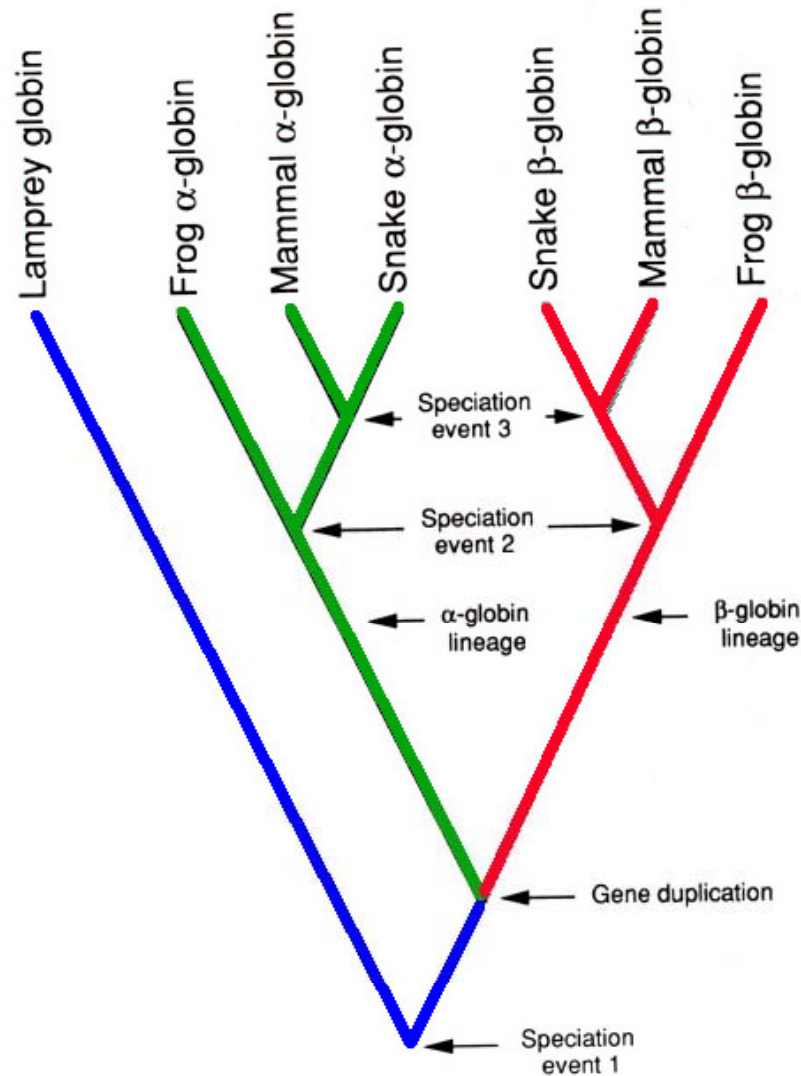
an alignment is an assignment of gaps to positions  $0, \dots, M$  in  $x$ , and  $0, \dots, N$  in  $y$ , so as to line up each letter in one sequence with either a letter, or a gap in the other sequence



# Homology

- Orthologs
  - Divergence follows speciation
  - Similarity can be used to construct phylogeny between species
- Paralogs
  - Divergence follows duplication
- Xenologs
- [Article on terminology](#)
- [ISMB tutorial on protein sequence comparison](#)

# Orthologs and paralogs



# Understanding evolutionary relationships

molecular

molecular

Nothing in biology makes sense except in the light of evolution



Dobzhansky, 1973

# Sources of variation

- Nucleotide substitution
  - Replication error
  - Chemical reaction
- Insertions or deletions (indels)
  - Unequal crossing over
  - Replication slippage
- Duplication
  - a single gene (complete gene duplication)
  - part of a gene (internal or partial gene duplication)
    - Domain duplication
    - Exon shuffling
  - part of a chromosome (partial polysomy)
  - an entire chromosome (aneuploidy or polysomy)
  - the whole genome (polyploidy)

# Differing rates of DNA evolution

- Functional/selective constraints (particular features of coding regions, particular features in 5' untranslated regions)
- Variation among different gene regions with different functions (different parts of a protein may evolve at different rates).
- Within proteins, variations are observed between
  - surface and interior amino acids in proteins (order of magnitude difference in rates in haemoglobins)
  - charged and non-charged amino acids
  - protein domains with different functions
  - regions which are strongly constrained to preserve particular functions and regions which are not
  - different types of proteins -- those with constrained interaction surfaces and those without

# Common assumptions

- All nucleotide sites change independently
- The substitution rate is constant over time and in different lineages
- The base composition is at equilibrium
- The conditional probabilities of nucleotide substitutions are the same for all sites, and do not change over time
- Most of these are not true in many cases...

# Lecture outline

- Sequence alignment
  - Why do we need to align sequences?
  - Evolutionary relationships
- Gaps and scoring matrices
- Dynamic programming
  - Global alignment (Needleman & Wunsch)
  - Local alignment (Smith & Waterman)
- Database searches
  - BLAST
  - FASTA

# A simple alignment

- Let us try to align two short nucleotide sequences:
  - AATCTATA and AAGATA
- Without considering any gaps (insertions/deletions) there are 3 possible ways to align these sequences

**AATCTATA**

**AAGATA**

**AATCTATA**

**AAGATA**

**AATCTATA**

**AAGATA**

- Which one is better?



# What is a good alignment?

AGGCTAGTT, AGCGAAGTTT

AGGCTAGTT-                      6 matches, 3 mismatches, 1 gap  
AGCGAAGTTT

AGGCTA-GTT-                      7 matches, 1 mismatch, 3 gaps  
AG-CGAAGTTT

AGGC-TA-GTT-                      7 matches, 0 mismatches, 5 gaps  
AG-CG-AAGTTT

# Scoring the alignments

- We need to have a scoring mechanism to evaluate alignments
  - match score
  - mismatch score

- We can have the total score as:

$$\sum_{i=1}^n \text{match or mismatch score at position } i$$

- For the simple example, assume a match score of 1 and a mismatch score of 0:

**AATCTATA**

**AATCTATA**

**AATCTATA**

**AAGATA**

**AAGATA**

**AAGATA**

4

1

3

# Good alignments require gaps

- Maximal consecutive run of spaces in alignment
  - Matching mRNA (cDNA) to DNA
  - Shortening of DNA/protein sequences
  - Slippage during replication
  - Unequal crossing-over during meiosis
  - ...
- We need to have a scoring function that considers gaps also

# Simple alignment with gaps

- Considering gapped alignments vastly increases the number of possible alignments:

**AATCTATA**

**AATCTATA**

**AATCTATA**

more?

**AAG-AT-A**

**AA-G-ATA**

**AA--GATA**

1

3

3

- If gap penalty is -1 what will be the new scores?

# Scoring Function

- Sequence edits:

- Mutations
- Insertions
- Deletions

AGGCCTC

AGGACTC

AGGGCCTC

AGG . CTC

Alternative definition:

***minimal edit distance***

*“Given two strings x, y, find minimum # of edits (insertions, deletions, mutations) to transform one string to the other”*

## Scoring Function:

Match: +m

Mismatch: -s

Gap: -d

Score  $F = (\# \text{ matches}) \times m - (\# \text{ mismatches}) \times s - (\# \text{ gaps}) \times d$

# More complicated gap penalties

- Nature favors small number of long gaps compared to large number of short gaps.
- How do we adjust our scoring scheme to account for this fact above?

By having different gap opening and gap extension penalties.

- Choices of gap penalties
  - Linear
  - Affine
    - Gap open penalty
    - Gap extension penalty
  - Arbitrary

# Score matrix

- Instead of having a single match/mismatch score for every pair of nucleotides or amino acids, consider chemical, physical, evolutionary relationships:
  - E.g.
    - alanine vs. valine or alanine vs. lysine? Alanine and valine are both small and hydrophobic, but lysine is large and charged.
    - which substitutions occur more in nature?
- Assign scores to each pair of symbol
  - Higher score means more similarity

Table 1 - The log odds matrix for 250 PAMs (multiplied by 10)

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2	-2	0	0	-4	1	-1	-1	-1	-2	-1	0	1	0	-2	1	1	0	-6	-3
C		12	-5	-5	-4	-3	-3	-2	-5	-6	-5	-4	-3	-5	-4	0	-2	-2	-8	0
D			4	3	-6	1	1	-2	0	-4	-3	2	-1	2	-1	0	0	-2	-7	-4
E				4	-5	0	1	-2	0	-3	-2	1	-1	2	-1	0	0	-2	-7	-4
F					9	-5	-2	1	-5	2	0	-4	-5	-5	-4	-3	-3	-1	0	7
G						5	-2	-3	-2	-4	-3	0	-1	-1	-3	1	0	-1	-7	-5
H							6	-2	0	-2	-2	2	0	3	2	-1	-1	-2	-3	0
I								5	-2	2	2	-2	-2	-2	-2	-1	0	4	-5	-1
K									5	-3	0	1	-1	1	3	0	0	-2	-3	-4
L										6	4	-3	-3	-2	-3	-3	-2	2	-2	-1
M											6	-2	-2	-1	0	-2	-1	2	-4	-2
N												2	-1	1	0	1	0	-2	-4	-2
P													6	0	0	1	0	-1	-6	-5
Q														4	1	-1	-1	-2	-5	-4
R															6	0	-1	-2	2	-4
S																2	1	-1	-2	-3
T																	3	0	-5	-3
V																		4	-6	-2
W																			17	0
Y																				10



Table 2 - The log odds matrix for BLOSUM 62

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

## Major Differences between PAM and BLOSUM

PAM	BLOSUM
Built from global alignments	Built from local alignments
Built from small amount of Data	Built from vast amount of Data
Counting is based on minimum replacement or maximum parsimony	Counting based on groups of related sequences counted as one
Perform better for finding global alignments and remote homologs	Better for finding local alignments
Higher PAM series means more divergence	Lower BLOSUM series means more divergence

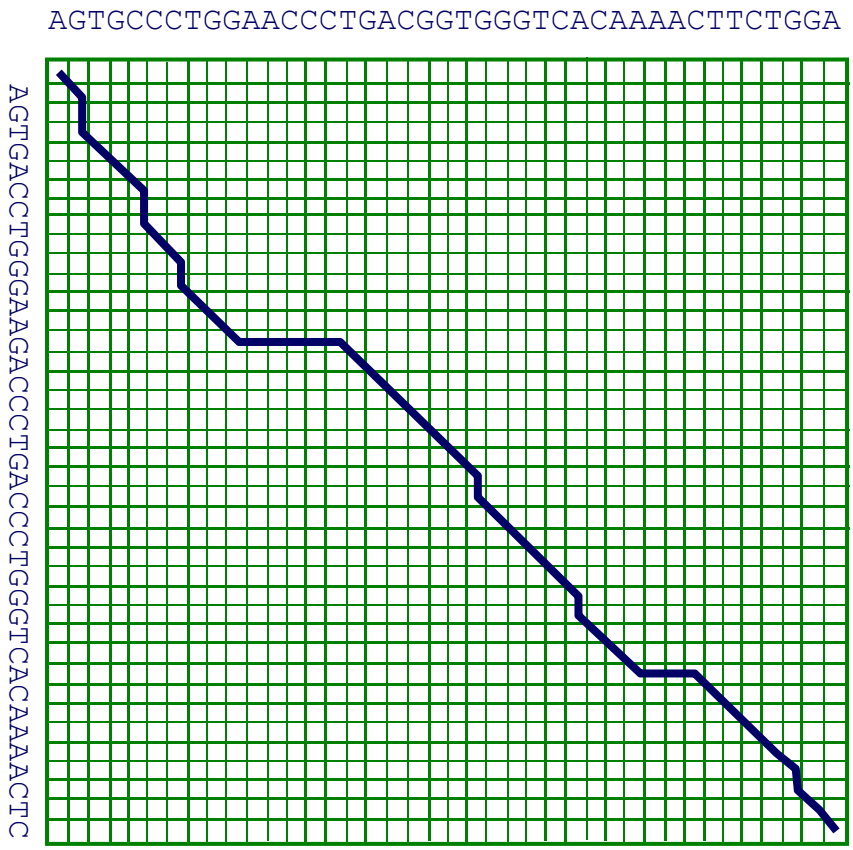
# Typical score matrix

- DNA
  - Match = +1
  - Mismatch = -3
  - Gap penalty = -5
  - Gap extension penalty = -2
- Protein sequences
  - Blossum62 matrix
  - Gap open penalty = -11
  - Gap extension = -1

# Lecture outline

- Sequence alignment
  - Why do we need to align sequences?
  - Evolutionary relationships
- Gaps and scoring matrices
- **Dynamic programming**
  - Global alignment (Needleman & Wunsch)
  - Local alignment (Smith & Waterman)
- Database searches
  - BLAST
  - FASTA

# How do we compute the best alignment?



Too many possible alignments:

$$\gg 2^N$$

(exercise)

# Alignment is additive

Observation:

The score of aligning

$x_1 \dots x_M$

$y_1 \dots y_N$

is additive

Say that  
aligns to

$x_1 \dots x_i$

$x_{i+1} \dots x_M$

$y_1 \dots y_j$

$y_{j+1} \dots y_N$

The two scores add up:

$$F(x[1:M], y[1:N]) = F(x[1:i], y[1:j]) + F(x[i+1:M], y[j+1:N])$$

# Types of alignment

- Global (Needleman & Wunsch)
  - Strings of similar size
    - Genes with a similar structure
    - Larger regions with a preserved order (syntenic regions)
- Local (Smith & Waterman)
  - Finding similar regions among
    - Dissimilar regions
    - Sequences of different lengths

# Dynamic programming

- Instead of evaluating every possible alignment, we can create a table of partial scores by breaking the alignment problem into subproblems.
- Consider two sequences CACGA and CGA
  - we have three possibilities for the first position of the alignment

First position	Score	Remaining seqs.
C	+1	ACGA
C		GA
-	-1	CACGA
C		GA
C	-1	ACGA
-		CGA



# Example

score(H,P) = -2, gap penalty=-8 (linear)

	-	H	E	A	G	A	W	G	H	E	E
-	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
P	-8	-2									
A	-16										
W	-24										
H	-32										
E	-40										
A	-48										
E	-56										

# The DP recurrence relation

- $s(a,b)$  = score of aligning  $a$  and  $b$
- $F(i,j)$  = optimal similarity of  $A(1:i)$  and  $B(1:j)$
- Recurrence relation
  - $F(i,0) = \sum s(A(k),-), 0 \leq k \leq i$
  - $F(0,j) = \sum s(-,B(k)), 0 \leq k \leq j$
  - $F(i,j) = \max [F(i,j-1) + s(-,B(j)),$   
 $F(i-1,j) + s(A(i),-),$   
 $F(i-1,j-1) + s(A(i),B(j))]$
  - Assume linear gap penalty

# Example contd.

$\text{score}(E,P) = 0$ ,  $\text{score}(E,A) = -1$ ,  $\text{score}(H,A) = -2$

	-	H	E	A	G	A	W	G	H	E	E
-	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
P	-8	-2	-8								
A	-16	-10	-3								
W	-24										
H	-32										
E	-40										
A	-48										
E	-56										

Optimal alignment:

HEAGAWGHE - E  
 - P - - AW - HEAE

		H	E	A	G	A	W	G	H	E	E
	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
P	-8	-2	-8	-16	-24	-33	-42	-49	-57	-65	-73
A	-16	-10	-3	-4	-12	-19	-28	-36	-44	-52	-60
W	-24	-18	-11	-6	-7	-15	-4	-12	-21	-29	-37
H	-32	-14	-18	-13	-8	-9	-12	-6	-2	-11	-19
E	-40	-22	-8	-16	-16	-9	-12	-14	-6	4	-5
A	-48	-30	-16	-3	-11	-11	-12	-12	-14	-4	2
E	-56	-38	-24	-11	-6	-12	-14	-15	-12	-8	2

The value in the final cell is the best score for the alignment

# More examples

- Sequence alignment applet at:

<http://www.iro.umontreal.ca/~casagran/baba.html>

# Semi-global alignment

- In Needleman&Wunsch DP algorithm the gap penalty is assessed regardless of whether gaps are located internally or at the terminal ends.
- Terminal gaps may not be biologically significant

**AATCTATA**

**--TCT---**

- Treat terminal gaps differently than internal gaps → semi-global alignment
- What modifications should be made to the original DP?

# Local sequence alignment

- Suppose, we have a long DNA sequence (e.g., 4000 bp) and we want to compare it with the complete yeast genome (12.5M bp).
- What if only a portion of our query, say 200 bp length, has strong similarity to a gene in yeast.
  - Can we find this 200 bp portion using (semi) global alignment?

Probably not. Because, we are trying to align the complete 4000 bp sequence, thus a random alignment may get a better score than the one that aligns 200 bp portion to the similar gene in yeast.

# The local alignment problem

Given two strings

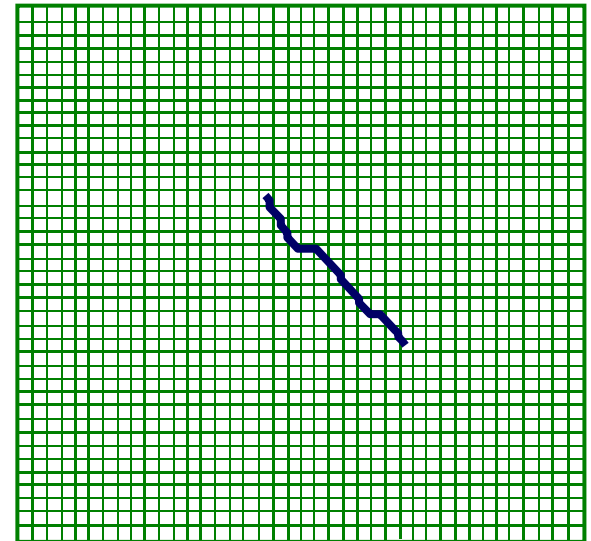
$$x = x_1 \dots x_M,$$

$$y = y_1 \dots y_N$$

Find substrings  $x'$ ,  $y'$  whose similarity  
(optimal global alignment value)  
is maximum

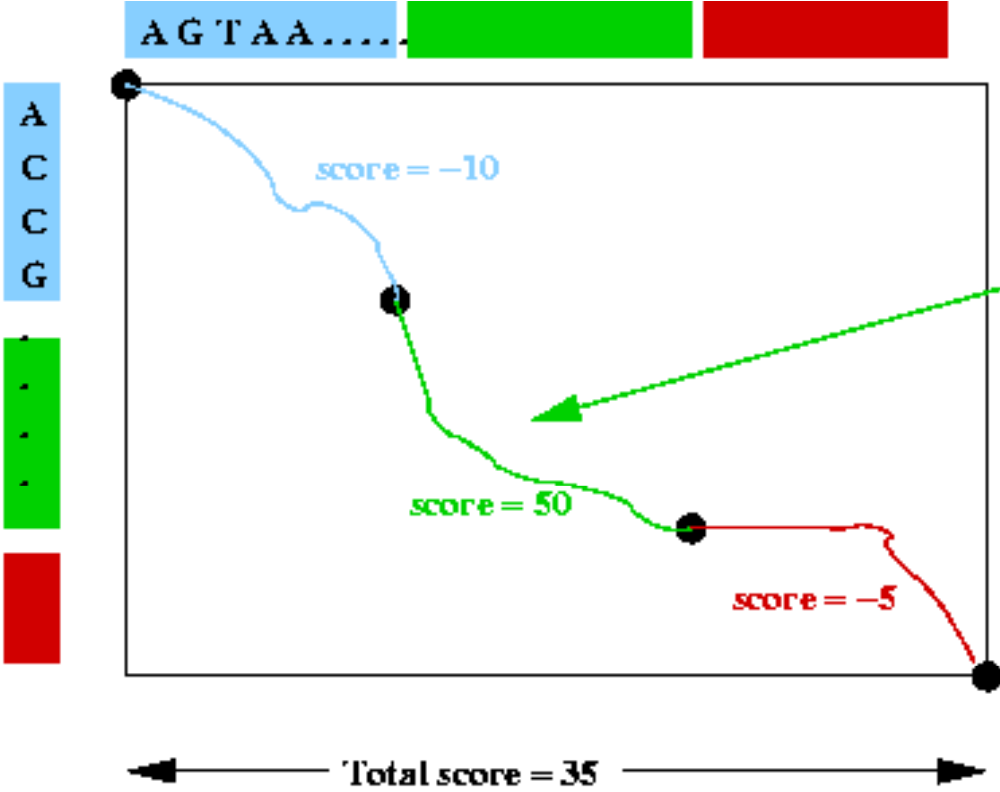
$x = \text{aaaacc} \boxed{\text{cccggg}} \text{gtta}$

$y = \text{tt} \boxed{\text{cccggg}} \text{aaccaacc}$





# Local alignment



local alignment  
may have higher  
score than overall  
global alignment

# Local sequence alignment (Smith-Waterman)

- $F(i,j)$  = optimal local similarity among suffixes of  $A(1:i)$  and  $B(1:j)$
- Recurrence relation
  - $F(i,0) = 0$
  - $F(0,j) = 0$
  - $F(i,j) = \max [0,$ 
    - $F(i,j-1) + s(-,B(j)),$
    - $F(i-1,j) + s(A(i),-),$
    - $F(i-1,j-1) + s(A(i),B(j))]$
  - Assume linear gap model

# Example

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Linear gap model

Gap = -1

Match = 4

Mismatch = -2

	-	E	Q	L	L	K	A	L	E	F	K	L
-												
K												
V												
L												
E												
F												
G												
Y												

# Example

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Linear gap model

Gap = -1

Match = 4

Mismatch = -2

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0											
V	0											
L	0											
E	0											
F	0											
G	0											
Y	0											

# Example

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Linear gap model

Gap = -1

Match = 4

Mismatch = -2

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
E	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Y	0	1	0	0	0	0	0	2	7	12	11	10

# Example

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Linear gap model

Gap = -1

Match = 4

Mismatch = -2

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
E	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Y	0	1	0	0	0	0	0	2	7	12	11	10

# Example

## Alignment

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Q: K A - L E F  
P: K - V L E F

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
E	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Y	0	1	0	0	0	0	0	2	7	12	11	10

# Example

## Alignment

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Q: K - A L E F  
P: K V - L E F

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
E	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Y	0	1	0	0	0	0	0	2	7	12	11	10



# Example

## Alignment

Q: E Q L L K A L E F K L  
P: K V L E F G Y

Q: K A L E F  
P: K V L E F

	-	E	Q	L	L	K	A	L	E	F	K	L
-	0	0	0	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	4	3	2	1	0	4	3
V	0	0	0	0	0	3	2	1	0	0	3	2
L	0	0	0	4	4	3	2	6	5	4	3	7
E	0	4	3	3	3	2	1	5	10	9	8	7
F	0	3	2	2	2	1	0	4	9	14	13	12
G	0	2	1	1	1	0	0	3	8	13	12	11
Y	0	1	0	0	0	0	0	2	7	12	11	10

# Another Example

Find the local alignment between:

Q: G C T G G A A G G C A T

P: G C A G A G C A C G

Linear gap model

Gap = -4

Match = +5

Mismatch = -4

Q

P

	--	G	C	T	G	G	A	A	G	G	C	A	T
--	0	0	0	0	0	0	0	0	0	0	0	0	0
G	0												
C	0												
A	0												
G	0												
A	0												
G	0												
C	0												
A	0												
C	0												
G	0												

# Another Example

Q's subsequence: G A A G - G C A

P's subsequence: G C A G A G C A

Q

		--	G	C	T	G	G	A	A	G	G	C	A	T	
P	--		0	0	0	0	0	0	0	0	0	0	0	0	0
	G	0	5	1	0	5	5	1	0	5	5	1	0	0	
	C	0	1	10	6	2	1	1	0	1	1	10	6	2	
	A	0	0	6	6	2	0	6	6	2	0	6	15	11	
	G	0	5	2	2	11	7	3	2	11	7	3	11	11	
	A	0	1	1	0	7	7	11	8	7	7	3	8	7	
	G	0	5	1	0	5	11	7	7	13	12	8	4	4	
	C	0	0	10	6	2	7	7	3	9	8	17	13	9	
	A	0	0	6	6	2	3	11	12	8	5	13	22	18	
	C	0	0	5	2	2	0	7	8	8	4	18	18	18	
G	0	5	1	1	7	7	5	4	13	13	14	14	14		

# Local vs. Global alignment

PIR Entry				Similarity Score		
				Global		Local
				End Penalty	No End Penalty	
HBHU	vs	HBHU	Hemoglobin beta-chain—human	725	725	725
		HAHU	Hemoglobin alpha-chain—human	314	320	322
		MYHU	Myoglobin—Human	121	164	166
		GPYL	Leghemoglobin—Yellow lupin	8	28	43
		LZCH	Lysozyme precursor—Chicken	-107	16	32
		NRBO	Pancreatic ribonuclease—Bovine	-124	16	31
		CCHU	Cytochrome c—Human	-160	10	26
MCHU	vs	MCHU	Calmodulin—Human	671	671	671
		TPHUCS	Troponin C, skeletal muscle	395	430	438
		PVPK2	Parvalbumin beta—Pike	-57	103	115
		CIHUH	Calpain heavy chain—Human	-2085	89	100
		AQJFNV	Aequorin precursor—Jelly fish	-65	48	76
		KLSWM	Calcium binding protein—Scallop	-89	45	52
QRHULD	vs	EGMSMG	Epidermal growth factor precursor	-591	475	655

# Complexity

- $O(mn)$  time
- $O(mn)$  space
  - $O(\max(m,n))$  if only distance value is needed
- More complicated “divide-and-conquer” algorithm that doubles time complexity and uses  $O(\min(m,n))$  space [Hirschberg, JACM 1977]

# Time and space bottlenecks

- Comparing two one-megabase genomes.
- Space:
  - An entry: 4 bytes;
  - Table:  $4 * 10^6 * 10^6 = 4 \text{ T bytes}$  memory.
- Time:
  - 1000 MHz CPU: 1M entries/second;
  - $10^{12}$  entries: 1M seconds = 10 days.

# Lecture outline

- Sequence alignment
  - Why do we need to align sequences?
  - Evolutionary relationships
- Gaps and scoring matrices
- Dynamic programming
  - Global alignment (Needleman & Wunsch)
  - Local alignment (Smith & Waterman)
- Database searches
  - BLAST
  - FASTA

# BLAST

- Basic Local Alignment Search Tool
  - Altschul *et al.* 1990,1994,1997
- Heuristic method for local alignment
- Designed specifically for database searches
- Idea: good alignments contain short lengths of exact matches



# Steps of BLAST

1. Filter low complexity regions (optional)
2. *Query words* of length 3 (for proteins) or 11 (for DNA) are created from query sequence using a sliding window

```
MEFPGLGSLGTSEPLPQFVDPALVSS
MEF
  EFP
    FPG
      PGL
        GLG
```

# Steps of BLAST

3. Scan each database sequence for an exact match to *query words*\*. Each match is a *seed* for an ungapped alignment.

\*Blast actually uses a list of *high scoring words* created from words *similar* to query words.

# Steps of BLAST

4. (*Original BLAST*) extend matching words to the left and right using ungapped alignments. Extension continues as long as score increases or stays same. This is a HSP (high scoring pair).

(*BLAST2*) Matches along the same diagonal (think dot plot) within a distance  $A$  of each other are joined and then the longer sequence extended as before. Need at least two contiguous hits for extension.

# Steps of BLAST

5. Using a cutoff score  $S$ , keep only the extended matches that have a score at least  $S$ .
6. Determine statistical significance of each remaining match.
7.
  - (*Original BLAST*) Only ungapped alignments; sometimes combined together
  - (*BLAST2*) Extend the HSPs using gapped alignment

# Summarizing BLAST

- One of the few algorithms to make it as a verb
  - Blast(v): to run a BLAST search against a sequence database
- Extension is the most time-consuming step
- BLAST2 reported to be 3 times faster than the original version at same quality

# Example BLAST run

- [BLAST website](#)

# Steps of FASTA

1. Find k-tups in the two sequences (k=1-2 for proteins, 4-6 for DNA sequences)
2. Select top 10 scoring “local diagonals” with matches and mismatches but no gaps.
  - a. For proteins, each k-tup found is scored using the PAM250 matrix
  - b. For DNA, use the number of k-tups found
  - c. Penalize intervening regions of mismatches

# Finding k-tups

```
position 1 2 3 4 5 6 7 8 9 10 11
protein 1 n c s p t a . . . . .
protein 2 . . . . . a c s p r k
```

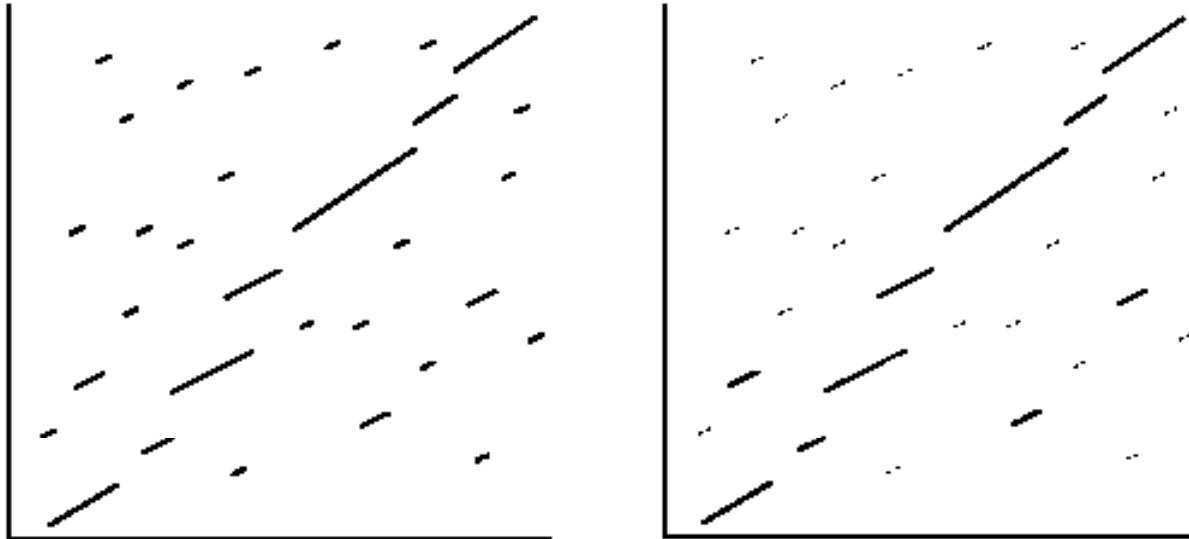
amino acid	position in		offset
	protein A	protein B	pos A - posB
a	6	6	0
c	2	7	-5
k	-	11	
n	1	-	
p	4	9	-5
r	-	10	
s	3	8	-5
t	5	-	

Note the common offset for the 3 amino acids c,s and p  
A possible alignment is thus quickly found -

```
protein 1 n c s p t a
           | | |
protein 2 a c s p r k
```



# FASTA



Finding 10 best diagonal runs

# FASTA

3. Rescan top 10 diagonals (representing alignments), score with PAM250 (proteins) or DNA scoring matrix. Trim off the ends of the regions to achieve highest scores.
4. Join regions that are consistent with gapped alignments. (maximal weighted paths in a graph).

# FASTA

5. After finding the best initial region (step 3), FASTA performs a DP global alignment in a gap centered on the best initial region.

# Summarizing FASTA

- Statistics based on histograms on values of intermediate and final scores.
- Begins with exact matches unlike BLAST
- Less of a statistical basis for comparison
- Quality and complexity similar to BLAST

# History of sequence searching

- 1970: NW
- 1980: SW
- 1985: FASTA
- 1989: BLAST
- 1997: BLAST2