

DNA1: (Last week)

Types of mutants
 Mutation, drift, selection
 Binomial for each
 Association studies χ^2 statistic
 Linked & causative alleles
 Alleles, Haplotypes, genotypes
 Computing the first genome,
 the second ...
 New technologies
 Random and systematic errors

1

DNA2: Aligning ancient diversity

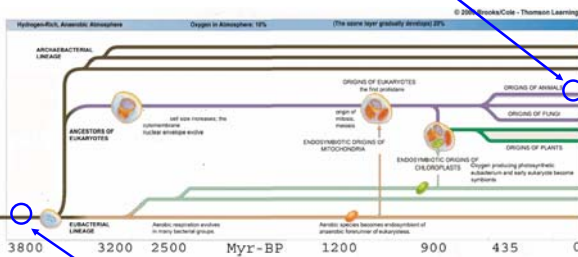
Comparing types of alignments & algorithms

Dynamic programming
 Multi-sequence alignment
 Space-time-accuracy tradeoffs
 Finding genes -- motif profiles
 Hidden Markov Model for CpG Islands

2

DNA2: diversity & continuity

DNA1: the last 5000 human generations



Integrate 1&2: ancient genetic code

http://www.colorado.edu/cpob/epob1030cornwall/fall_2001/Evolution_of_Life.gif

3

Applications of Dynamic Programming

- To sequence analysis
 - Shotgun sequence assembly
 - Multiple alignments
 - Dispersed & tandem repeats
 - Bird song alignments
 - Gene Expression time-warping
- 3D-structure alignment
- Through HMMs
 - RNA gene search & structure prediction
 - Distant protein homologies
 - Speech recognition

4

Alignments & Scores

Global (e.g. haplotype)

ACCACACA
 : : X X : : X :
 ACACCATA
 Score = $5(+1) + 3(-1) = 2$

Local (motif)

ACCACACA
 : : : :
 ACACCATA
 Score = $4(+1) = 4$

Suffix (shotgun assembly)

ACCACACA
 : : : :
 ACACCATA
 Score = $3(+1) = 3$

5

Increasingly complex (accurate) searches

Exact (StringSearch)	CGCG
Regular expression (PrositeSearch)	$CGN\{0-9\}CG = CGAACG$
Substitution matrix (BlastN)	$CGCG \approx CACG$
Profile matrix (PSI-blast)	$CGc(g/a) \approx CACG$
Gaps (Gap-Blast)	$CGCG \approx CGAACG$
Dynamic Programming (NW, SM)	$CGCG \approx CAGACG$

Hidden Markov Models (HMMER)



WU

6

"Hardness" of (multi-) sequence alignment

Align 2 sequences of length N allowing gaps.

```
ACCAC-ACA      ACCACACA
::x::x::x:    :xxxxxx:
AC-ACCATA ,   A-----CACCATA , etc.
```

2N gap positions, gap lengths of 0 to N each:
A naïve algorithm might scale by $O(N^{2N})$.
For $N = 3 \times 10^9$ this is rather large.

Now, what about $k > 2$ sequences?
or rearrangements other than gaps?

7

Testing search & classification algorithms

Separate Training set and Testing sets
Need databases of non-redundant sets.
Need evaluation criteria (programs)
Sensitivity and Specificity (false negatives & positives)

sensitivity ($\text{true_predicted}/\text{true}$)
specificity ($\text{true_predicted}/\text{all_predicted}$)

Where do training sets come from?
More expensive experiments: crystallography, genetics, biochemistry

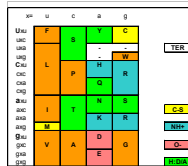
8

Comparisons of homology scores

Pearson WR Protein Sci 1995 Jun;4(6):1145-60
Comparison of methods for searching protein
sequence databases. Methods Enzymol 1996;266:227-58
Effective protein sequence comparison.

Algorithm: FASTA, Blastp, Blitz
Substitution matrix: PAM120, PAM250, BLOSUM50, BLOSUM62
Database: PIR, SWISS-PROT, GenPept

Switch to
protein
searches
when possible



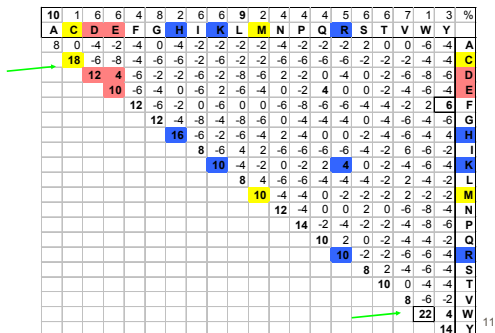
9

A Multiple Alignment of Immunoglobulins

```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCVSGTSFDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG--
VSLTCLVKGFYPSD--IAVEWESNG--
```

10

Scoring matrix based on large set of distantly related blocks: Blosum62



11

Scoring Functions and Alignments

Scoring function:

$\omega(\text{match}) = +1;$
 $\omega(\text{mismatch}) = -1;$
 $\omega(\text{indel}) = -2;$
 $\omega(\text{other}) = 0.$

} substitution matrix

Alignment score: sum of columns.

Optimal alignment: maximum score.

12

Calculating Alignment Scores

(1) ATGA A-TGA
 :xx: : : :
 ACTA ACT-A

$$(1) \text{Score} = \omega\left(\begin{matrix} A \\ A \end{matrix}\right) + \omega\left(\begin{matrix} T \\ C \end{matrix}\right) + \omega\left(\begin{matrix} G \\ T \end{matrix}\right) + \omega\left(\begin{matrix} A \\ A \end{matrix}\right) = 1 - 1 - 1 + 1 = 0.$$

$$(2) \text{Score} = \omega\left(\begin{matrix} A \\ A \end{matrix}\right) + \omega\left(\begin{matrix} - \\ C \end{matrix}\right) + \omega\left(\begin{matrix} T \\ T \end{matrix}\right) + \omega\left(\begin{matrix} G \\ - \end{matrix}\right) + \omega\left(\begin{matrix} A \\ A \end{matrix}\right) = 1 - 2 + 1 - 2 + 1 = -1.$$

if $\omega(\text{indel}) = -1$, Score = $1 - 1 + 1 - 1 + 1 = +1$.

DNA2: Aligning ancient diversity

Comparing types of alignments & algorithms

Dynamic programming

Multi-sequence alignment

Space-time-accuracy tradeoffs

Finding genes -- motif profiles

Hidden Markov Model for CpG Islands

What is dynamic programming?

A dynamic programming algorithm solves every subsubproblem just once and then saves its answer in a table, avoiding the work of recomputing the answer every time the subsubproblem is encountered.

-- Cormen *et al.* "Introduction to Algorithms", The MIT Press.

Recursion of Optimal Global Alignments

$s\left(\begin{matrix} u \\ v \end{matrix}\right)$: optimal global alignment score of u and v .

$$s\left(\begin{matrix} ATGA \\ ACTA \end{matrix}\right) = \max \begin{cases} s\left(\begin{matrix} ATGA \\ ACT \end{matrix}\right) + \omega\left(\begin{matrix} - \\ A \end{matrix}\right); \\ s\left(\begin{matrix} ATGA \\ ACTA \end{matrix}\right) + \omega\left(\begin{matrix} A \\ - \end{matrix}\right); \\ s\left(\begin{matrix} ATG \\ ACTA \end{matrix}\right) + \omega\left(\begin{matrix} A \\ - \end{matrix}\right). \end{cases}$$

Recursion of Optimal Local Alignments

$s\left(\begin{matrix} u \\ v \end{matrix}\right)$: optimal local alignment score of u and v .

$$s\left(\begin{matrix} u_1u_2\dots u_i \\ v_1v_2\dots v_j \end{matrix}\right) = \max \begin{cases} s\left(\begin{matrix} u_1u_2\dots u_i \\ v_1v_2\dots v_{j-1} \end{matrix}\right) + \omega\left(\begin{matrix} - \\ v_j \end{matrix}\right); \\ s\left(\begin{matrix} u_1u_2\dots u_{i-1} \\ v_1v_2\dots v_j \end{matrix}\right) + \omega\left(\begin{matrix} u_i \\ v_j \end{matrix}\right); \\ s\left(\begin{matrix} u_1u_2\dots u_{i-1} \\ v_1v_2\dots v_j \end{matrix}\right) + \omega\left(\begin{matrix} u_i \\ - \end{matrix}\right); \\ 0. \end{cases}$$

Computing Row-by-Row

		A	C	T	A
	0	min	min	min	min
A	min	1	-1	-3	-5
T	min				
G	min				
A	min				

		A	C	T	A
	0	min	min	min	min
A	min	1	-1	-3	-5
T	min	-1	0	0	-2
G	min				
A	min				

		A	C	T	A
	0	min	min	min	min
A	min	1	-1	-3	-5
T	min	-1	0	0	-2
G	min	-3	-2	-1	-1
A	min				

		A	C	T	A
	0	min	min	min	min
A	min	1	-1	-3	-5
T	min	-1	0	0	-2
G	min	-3	-2	-1	-1
A	min	-5	-4	-3	0

min = -10⁹⁹

Traceback Optimal Global Alignment

		A	C	T	A
	0	min	min	min	min
A	min	1	-1	-3	-5
T	min	-1	0	-1	-2
G	min	-3	-2	-1	-1
A	min	-5	-4	-3	0

$\begin{pmatrix} A & G & T & A \\ \vdots & \times & \times & \vdots \\ A & T & C & A \end{pmatrix}$
 19

Local and Global Alignments

		A	C	C	A	C	A	C	A
	0	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	1	0	1
C	0	0	2	1	0	2	0	2	0
A	0	1	0	1	2	0	3	1	3
C	0	0	2	1	0	3	1	4	2
C	0	0	0	3	2	1	2	2	3
A	0	1	0	1	4	2	2	1	3
T	0	0	0	0	2	3	1	1	1
A	0	1	0	0	1	1	4	2	2

	0	m	m	m	m	m	m	m	m
A	m	1	-1	-3	-5	-7	-9	-11	-13
C	m	-1	2	0	-2	-4	-6	-8	-10
A	m	-3	0	1	1	-1	-3	-5	-7
C	m	-5	-2	1	0	2	0	-2	-4
C	m	-7	-4	-1	0	1	1	1	-1
A	m	-9	-6	-3	0	-1	2	0	2
T	m	-11	-8	-5	-2	-1	0	1	0
A	m	-13	-10	-7	-4	-1	0	-1	2

Time and Space Complexity of Computing Alignments

For two sequences $u=u_1u_2\dots u_n$ and $v=v_1v_2\dots v_m$, finding the optimal alignment takes $O(mn)$ time and $O(mn)$ space.

An $O(1)$ -time operation: one comparison, three multiplication steps, computing an entry in the alignment table...

An $O(1)$ -space memory: one byte, a data structure of two floating points, an entry in the alignment table...

Time and Space Problems

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

Time & Space Improvement for w-band Global Alignments

- Two sequences differ by at most w bps ($w \ll n$).
- w-band algorithm: $O(wn)$ time and space.
- Example: $w=3$.

		A	C	C	A	C	A	C	A
	0	m	m	m					
A	m	1	-1	-3	-5				
C	m	-1	2	0	-2	-4			
A	m	-3	0	1	1	-1	-3		
C	m	-5	-2	1	0	2	0	-2	
C	m	-7	-4	-1	0	1	1	1	-1
A	m	-9	-6	-3	0	-1	2	0	2
T	m	-11	-8	-5	-2	-1	0	1	0
A	m	-13	-10	-7	-4	-1	0	-1	2

Summary

- Dynamic programming
- Statistical interpretation of alignments
- Computing optimal global alignment
- Computing optimal local alignment
- Time and space complexity
- Improvement of time and space
- Scoring functions

DNA2: Aligning ancient diversity

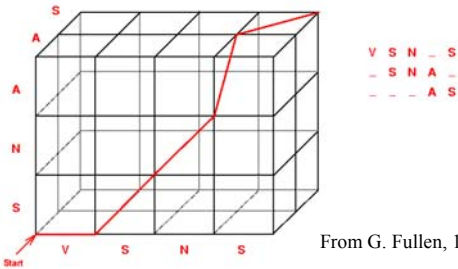
Comparing types of alignments & algorithms
 Dynamic programming
Multi-sequence alignment
 Space-time-accuracy tradeoffs
 Finding genes -- motif profiles
 Hidden Markov Model for CpG Islands

A Multiple Alignment of Immunoglobulins

```

VTISCTGSSSNIGAG-NHVKNWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCVSGTSEDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPSD--IAVEWESNG--
    
```

A multiple alignment \Leftrightarrow Dynamic programming on a hyperlattice



From G. Fullen, 1996.

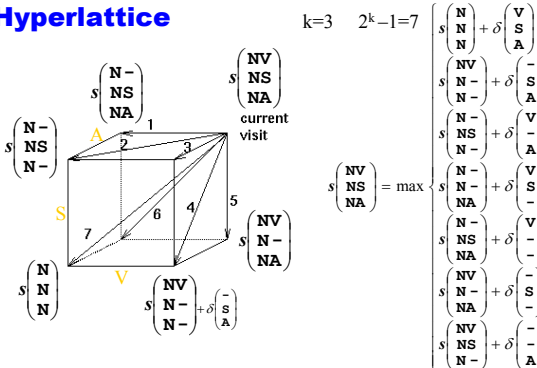
Multiple Alignment vs Pairwise Alignment

AT	
AT	
AT	
A-	
-T	
AT	A-
AT	-T

Optimal Multiple Alignment

Non-Optimal Pairwise Alignment

Computing a Node on Hyperlattice



Challenges of Optimal Multiple Alignments

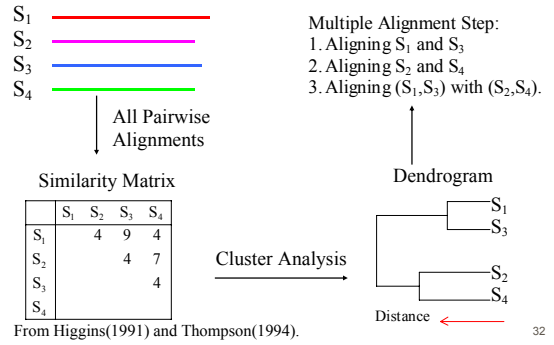
- Space complexity (hyperlattice size): $O(n^k)$ for k sequences each n long.
- Computing a hyperlattice node: $O(2^k)$.
- Time complexity: $O(2^k n^k)$.
- Find the optimal solution is exponential in k (non-polynomial, NP-hard).

Methods and Heuristics for Optimal Multiple Alignments

- Optimal: dynamic programming
Pruning the hyperlattice (MSA)
- Heuristics:
 - tree alignments(ClustalW)
 - star alignments
 - sampling (Gibbs) (discussed in RNA2)
 - local profiling with iteration (PSI-Blast, ...)

31

ClustalW: Progressive Multiple Alignment



32

Star Alignments

$s_1 = \text{ATTGCCATT}$
 $s_2 = \text{ATGGCCATT}$
 $s_3 = \text{ATCCAATTTT}$
 $s_4 = \text{ATCTTCTT}$
 $s_5 = \text{ACTGACC}$

Similarity Matrix

	s_2	s_3	s_4	s_5
s_1	7	-2	0	-3
s_2		-2	0	-4
s_3			0	-7
s_4				-3

Find the Central Sequence s_1

Multiple Alignment

ATTGCCATT--
 ATGGCCATT--
 ATC-CAATTTT
 ATCTTC-TT--
 ACTGACC----
 AT*GCCATTTT

Combine into Multiple Alignment

Pairwise Alignment

ATTGCCATT
 ATGGCCATT
 ATTGCCATT--
 ATC-CAATTTT
 ATTGCCATT
 ATCTTC-TT
 ATTGCCATT
 ACTGACC

33

DNA2: Aligning ancient diversity

- Comparing types of alignments & algorithms
- Dynamic programming
- Multi-sequence alignment
- Space-time-accuracy tradeoffs
- Finding genes -- motif profiles
- Hidden Markov Model for CpG Islands

34

Accurately finding genes & their edges

What is distinctive ?

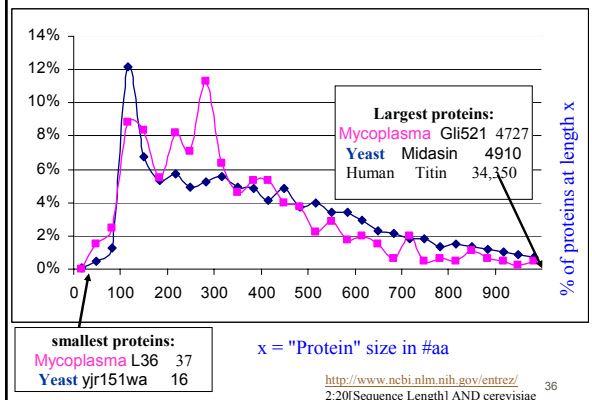
- Promoters & CGs islands
- Preferred codons
- RNA splice signals
- Frame across splices
- Inter-species conservation
- cDNA for splice edges

Failure to find edges?

- Variety & combinations
- Tiny proteins (& RNAs)
- Alternatives & weak motifs
- Alternatives
- Gene too close or distant
- Rare transcript

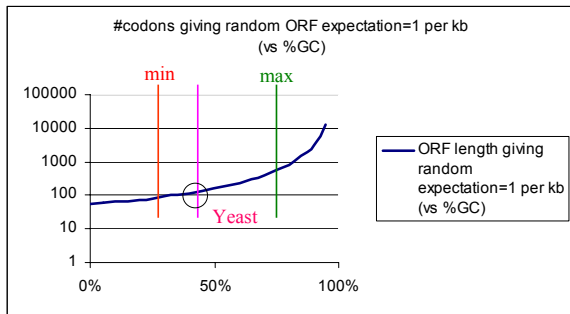
35

Annotated "Protein" Sizes



36

Predicting small proteins (ORFs)



Small coding regions

Mutations in domain II of 23 S rRNA facilitate translation of a 23 S rRNA-encoded **pentapeptide** conferring erythromycin resistance. Dam et al. 1996 J Mol Biol 259:1-6

Trp (W) leader peptide, 14 codons:
MKAIFVLK**GW**RTS **STOP**

Phe (F) leader peptide, 15 codons:
MK**H**IPFFFAFFFT**F**P **STOP**

His (H) leader peptide, 16 codons:
MTRVQ**F**KHHHHHHHPD **STOP**

<http://www.andrew.cmu.edu/~berget/Education/attenuation/atten.html>

Other examples in proteomics lectures

38

Motif Matrices

```
a a t g
c a t g
g a t g
t g t g
```

```
a 1 3 0 0
c 1 0 0 0
g 1 1 0 4
t 1 0 4 0
```

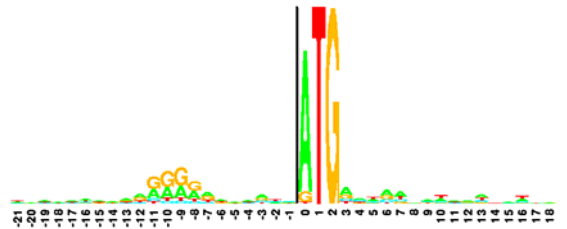
Align and calculate frequencies.
Note: Higher order correlations lost.

39

Protein starts

[GeneMark](#)

1055 E. coli Ribosome binding sites listed in the Miller book



40

Motif Matrices

```
a a t g 1+3+4+4 = 12
c a t g 1+3+4+4 = 12
g a t g 1+3+4+4 = 12
t g t g 1+1+4+4 = 10
```

```
a 1 3 0 0
c 1 0 0 0
g 1 1 0 4
t 1 0 4 0
```

Align and calculate frequencies.
Note: Higher order correlations lost.

Score test sets:
a c c c 1+0+0+0 = 1

41

DNA2: Aligning ancient diversity

Comparing types of alignments & algorithms

Dynamic programming

Multi-sequence alignment

Space-time accuracy tradeoffs

Finding genes-- motif profiles

Hidden Markov Model for CpG Islands

42

Why probabilistic models in sequence analysis?

- **Recognition** - Is this sequence a protein start?
- **Discrimination** - Is this protein more like a hemoglobin or a myoglobin?
- **Database search** - What are all of sequences in SwissProt that look like a serine protease?

43

A Basic idea

Assign a number to every possible sequence such that

$$\sum_s P(s|M) = 1$$

$P(s|M)$ is a probability of sequence s given a model M .

44

Sequence recognition

Recognition question - What is the probability that the sequence s is from the start site model M ?

$$P(M|s) = P(M) * P(s|M) / P(s)$$

(Bayes' theorem)

$P(M)$ and $P(s)$ are prior probabilities and $P(M|s)$ is posterior probability.

45

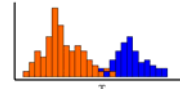
Database search

■ N = null model (random bases or AAs)

■ Report all sequences with

$$\log P(s|M) - \log P(s|N) > \log P(N) - \log P(M)$$

■ Example, say $\alpha\beta$ hydrolase fold is rare in the database, about 10 in 10,000,000. The threshold is 20 bits. If considering 0.05 as a significant level, then the threshold is $20 + 4.4 = 24.4$ bits.



46

Plausible sources of mono, di, tri, & tetra- nucleotide biases

C rare due to lack of uracil glycosylase (cytidine deamination)
 TT rare due to lack of UV repair enzymes.
 CG rare due to 5methylCG to TG transitions (cytidine deamination)
 AGG rare due to low abundance of the corresponding Arg-tRNA.
 CTAG rare in bacteria due to error-prone "repair" of CTAGG to C*CAGG.
 AAAAA excess due to polyA pseudogenes and/or polymerase slippage.

AmAcid	Codon	Number	/1000	Fraction
Arg	AGG	3363.00	1.93	0.03
Arg	AGA	5345.00	3.07	0.06
Arg	CGG	10558.00	6.06	0.11
Arg	CGA	6853.00	3.94	0.07
Arg	CGT	34601.00	19.87	0.36
Arg	CGC	36362.00	20.88	0.37

<ftp://sanger.otago.ac.nz/pub/Tranterm/Data/codons/bct/Esccol.cod>

47

CpG Island + in a ocean of - First order Hidden Markov Model

MM=16, HMM= 64 transition probabilities (adjacent bp)

