

Sequence alignment

Meaning and significance

Outline

- What is sequence alignment?
- What are the applications of alignment?
- What does it *mean* to *align* sequences?

What is sequence alignment?

Pattern matching

suffix trees, Burrows-Wheeler Transform,...

Database searching

BLAST

Optimization problem

Needleman-Wunsch, Smith-Waterman,...

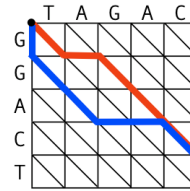
Statistical problem

Pair HMMs, TKF, Karlin-Altschul statistic...

...CATCG**ATGACTAT**CCG...
ATGACTGT

CATGCTTGCTGGCGTAAA

...
CATGCATGCTGCGTAC
CATGGTTGCTCACAAGTAC
CATGCTTGCTGGCGTAA
TACGTGCCTGACCTGCGTAC
CATGCCGAATGCTG
...



$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{\int_{\theta'} P(D|\theta')}$$

Applications of sequence alignment

- **Sequence assembly** – computing overlaps of reads with sequencing errors
- **Evolutionary functional analysis** – identifying evolutionarily-related sequences and conserved positions
- **Protein structure prediction** – use alignment of query protein to reference protein with known structure

DNA sequence edits

- Substitutions: **ACGA** \longrightarrow **AGGA**
- Insertions: **ACGA** \longrightarrow **ACCGGAGA**
- Deletions: **ACGGAGA** \longrightarrow **AGA**
- Transpositions: **ACGGAGA** \longrightarrow **AAGCGGA**
- Inversions: **ACGGAGA** \longrightarrow **ACTCCGA**

Alignment scales

- For proteins and short DNA sequences (gene scale) we will generally only consider
 - Substitutions: cause *mismatches* in alignments
 - Insertions/Deletions: cause *gaps* in alignments
- For long DNA sequences (genome scale) we will consider additional events
 - Transposition
 - Inversion
- In this course we will focus on the case of short sequences

What is a pairwise alignment?

- We will focus on *evolutionary* alignment
- matching of *homologous* positions in two sequences
- positions with no homologous pair are matched with a *space* ‘-’
- A group of consecutive spaces is a *gap*

CA--GATTCGAAT

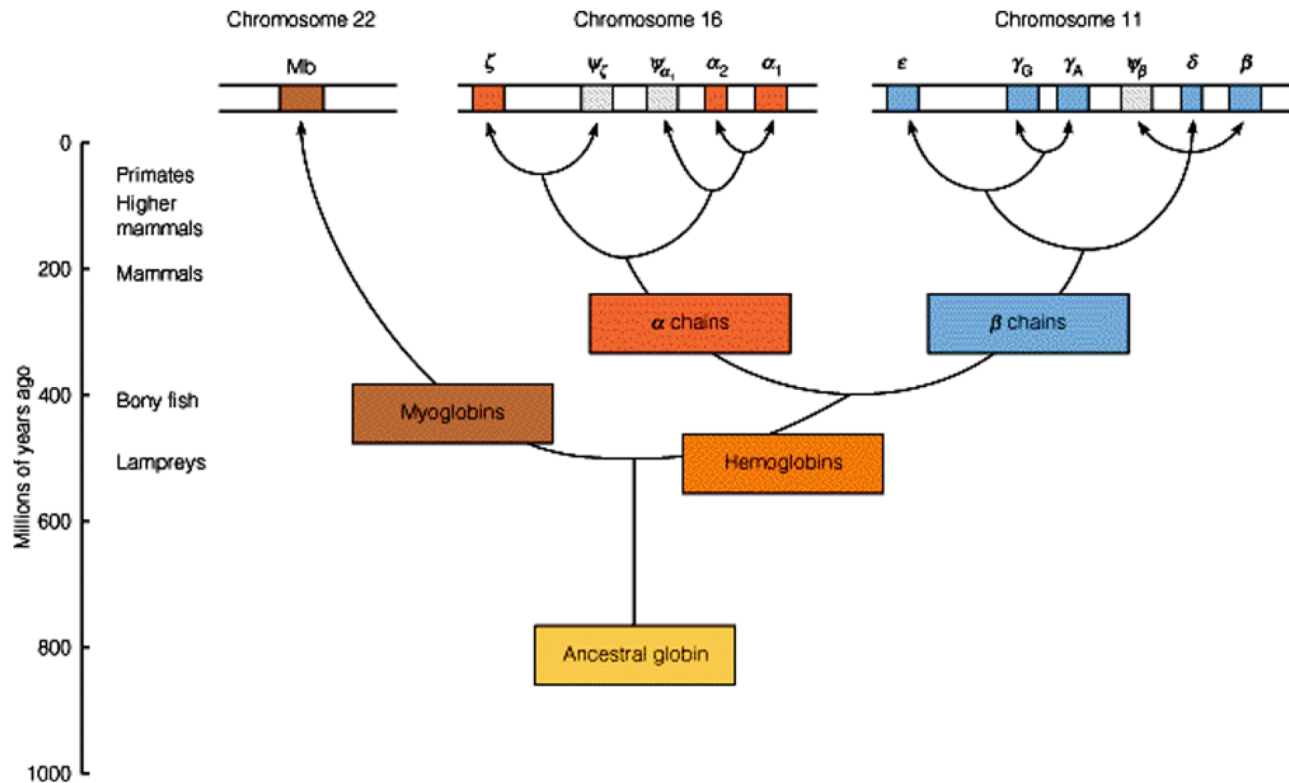
CGCCGATT---AT


gap

The Role of Homology

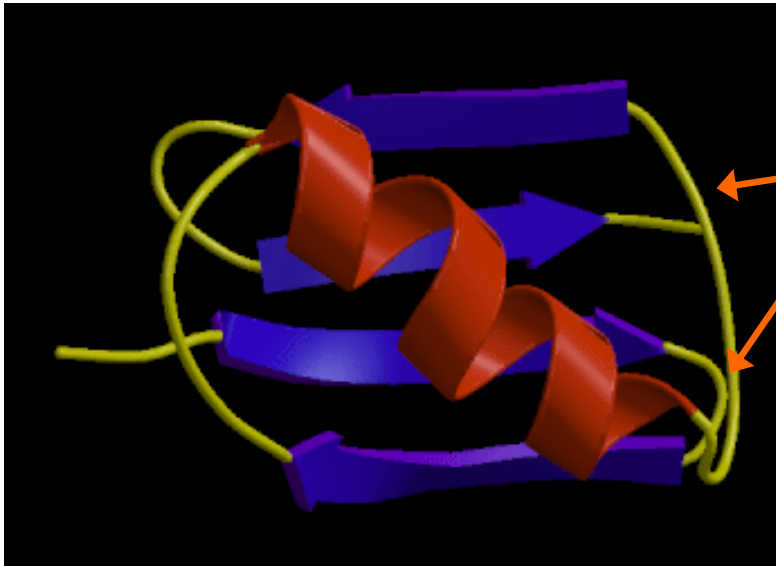
- *character*: some feature of an organism (could be molecular, structural, behavioral, etc.)
- *homology*: the relationship of two characters that have descended from a common ancestor
- homologous characters tend to be similar due to their common ancestry and evolutionary pressures
- thus we often infer homology from similarity
- thus we can sometimes infer structure/function from sequence similarity

Homology Example: Evolution of the Globins



Insertions/Deletions and Protein Structure

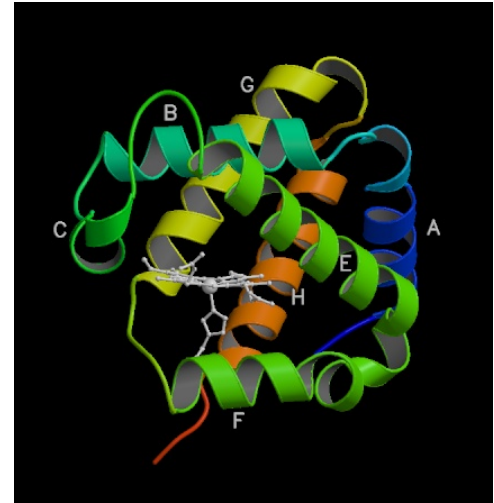
- Why is it that two “similar” sequences may have large insertions/deletions?
 - some insertions and deletions may not significantly affect the structure of a protein



loop structures: insertions/deletions here not so significant

Example Alignment: Globins

- figure at right shows prototypical structure of globins
- figure below shows part of alignment for 8 globins (-'s indicate gaps)



	A0	A4	A8	A12	B1	B6	B14	C2	CD1	CD4	
	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	
Hb_a	-----VL	SPADK	TNVKA	AWGKV	G-----	HAGEY	GAEAL	ERMFL	SFP	TTKTYF	PHF
Hb_b	-----VHL	TPEEK	SAVTAL	WGKV-----	NVDE	VGGEAL	GRLLV	VYPWT	QRRF	FESF	
Mb_SW	-----VL	SEGEW	QLVLH	VWAKVEA-----	DVAGH	GQDIL	IRLFK	SHPET	LEKF	DRF	
LegHb	-----GAL	TESQA	ALVKSS	WEEFN-----	NIPKH	THRFF	FILVLE	IAPAA	KDLF	SFL	
BacHb	-----LDQ	QTINI	IKATVP	VLKEHG-----	V-TIT	TTYF	KNLFA	KHPEV	RPLF	---	
SeaHb	GGTLAI	QAQGD	LTAAK	KIVRKT	WHQLMR---	NKTSF	VTDVF	IRIFAY	DPSAQ	NKFP	QM
AscHb	-----	ANKTR	ELCMK	SLEHAK	VDTSN	EARQD	GIDLY	KHMFEN	YPP	PLRKY	FKS-
Eryt.	-----L	SADQI	ISTVQ	ASF	DKVKG-----	DPVG	IILYA	VFKAD	PSI	MAKF	TQF

Summary

- Meaning of the term "sequence alignment"
- The task of evolutionary pairwise alignment
 - Homology
- Representations of alignments
- Sequence edits