# Sequence alignment

## Multiple alignment scoring and dynamic programming

# Outline

- The multiple sequence alignment algorithm
- Scoring Multiple Sequence Alignments
- Dynamic programming

# What is multiple sequence alignment?

**Given**: three or more related biological sequences

**Do**: identify the subsets of positions across sequences that are truly related

*In other words*: find a simultaneous alignment of all input sequences such that the implied pairwise alignments identify the truly related positions between each pair of sequences

# An example multiple sequence alignment

```
                        10          20          30          40          50          60        7
Calb/1-357     - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Dhan/1-520     - M N Y T K L K S Y S A N A I S N I L P - I D R E T C G E L V D Y A L T L P T - - - D H E I E A H F L N L L G E S D E T S A F L T K F M S -
Kwal/1-512     - M G K E S A I S F G I K E I P H I I P - I D E D S A R Q L C E Q I L S D H G - Q E H D T I A Q K F L D I L G P E D A S L N F V L Q F N E -
Sklu/1-519     - M A K D E A I K Y A I N Q I P Q I L P - L E E K D V R E L V N Q V L T Q N G E H N S E G I A Q S F L D I L G H D D M S F E F V F M F N E -
Klac/1-498     - M T K E D A I E Y A I K E L P N I L P - L D T E Q I K D L C E Q T I K E G N - - N P E Q I A Q S F F D L L G Q D D S S V H F I F E F N E -
Scer/1-530     - M T R K Q A I D Y A I K Q V P Q I L P - L E E S D V K A L C E Q V L S T S S - D D P E Q I A S K F L E F L G H E D L S F E F V M K F N E -
Spar/1-530     - M T R Q Q A I D Y A T K K V P Q I L P - L E E S D V K A L C E Q V L S T T S - N N P E Q I A S K F L E F L G H E D L S F E F V M R F N E -
Smik/1-527     - M T R Q Q A I D Y A V K K V P Q I L P - L E E S D V K A L C E Q V L S T S S - S N P E Q I A S K F L E F L G H E D L S F E F V M M F N E -
Sbay/1-531     - M T R Q Q A I D Y A V K Q V P Q I L P - L E E S D V R A L C E Q V L S T S S - D N P E Q I A S K F L E F L G H E D L S F E F V M K F N E -
Scas/1-517     - M T K T Q A I Q Y A L T K V P E I L P - L E Q D D V K Q L C E N I I S - S S - H N P E A I A Q G F L D I L G H D D L S F E F V M K F N E -
Kpol/1-520     - M T R K D A I A Y A V K A I P E I L P - L E E Q D V K N L C D Q I L N T S N - N D F E L I A N E F L S M L G H E D L A F A F V V E F N R -
Cgla/1-532     - M T Q Q K A I D Y A I A T I P D I L P - L E A D E I R T L C D Q I I K S C N - G S P E Q I A E G F M G I L G Q E E L V F D F V I R F N E -
Ylip/1-455     M E K Y T V T E E Y A K D M V G R L L G G F D K E T V A Q L V D Q G M K K T D - - - P L E V H S Y F V E L L G E S E P V F R F V E E F N R -
Sjap/1-552     - M P K E S V E D W A I E K L K K L L A - L D N E T L T I L V H G L L D A P D - - - P E S T R E K F Y D W L G R S K A I E Q F V E E L L A I
```
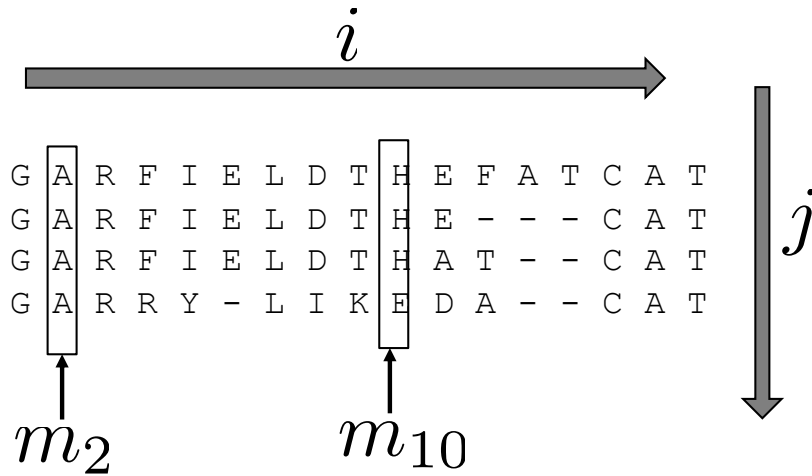
# Why multiple sequence alignment?

- Build phylogenetic trees (next module)
  - Determine evolutionary relationships between sequences
- A multiple sequence alignment can represent a family of proteins with similar function
  - Compare new sequence to a "family" of known proteins
  - For example the BLOCKS database used for BLOSUM contains several ungapped alignments for known protein families
- Discover common signatures or protein domains among a group of proteins
- Identify genetic variation among individuals of a population

# Some notation for multiple sequence alignments

- Let $m$ denote a Multiple Sequence Alignment
- $m_i$ is the $i^{th}$ column of the alignment $m$
- $m_i^j$ is the $i^{th}$ column and $j^{th}$ row
- $c_i^a$ count of residue $a$ in column $i$

# Example using notation

$i$

```
G A R F I E L D T H E F A T C A T
G A R F I E L D T H E - - - C A T
G A R F I E L D T H A T - - C A T
G A R R Y - L I K E D A - - C A T
```

$j$

$m_2$        $m_{10}$

$$m_3^1 = R$$

$$c_3^A = 0$$
$$c_2^A = 4$$
$$c_{10}^H = 3$$

# Scoring a Multiple Sequence Alignment (MSA)

- Key issue: how do we score a multiple sequence alignment?

- Usually, we assume that *columns* of an alignment are independent

$$Score(m) = G(m) + \sum_i S(m_i)$$

gap function

score of $i^{th}$ column

- For now, we will simplify the score by assuming a linear gap penalty

$$Score(m) = \sum_i S(m_i)$$

8

# Gap penalty (G)

- We will use a simple linear gap penalty function
  - Penalty for a space: s
- Let S(a,b) denote the cost of substituting a by b.
- Linear gap penalty can be incorporated into the substitution matrix
  - S(a,-) = S(-, a) = s
  - S(-,-)=0

# Two common ways of scoring a multiple alignment

- Entropy based scores

- Sum of pairs

# Entropy of a distribution

- A measure of uncertainty of an outcome
- For a discrete distribution $P(X)$, where $X$ takes k values $x_1, .. x_k$ it is defined as

$$H(X) = -\sum_{i=1}^{k} P(x_i)\log P(x_i)$$

- Entropy is greatest when we are most uncertain, that is, for a uniform distribution
- Entropy is least when we are most certain, e.g. deterministic event

# Entropy in extreme cases

# Score of a column: Entropy based

- Score of the $i^{th}$ column of alignment $m$ is

$$S(m_i) = -\sum_a c_i^a \log(p_{ia})$$

$p_{ia}$ : Probability of character $a$ in column $i$
$c_i^a$ : Number of occurrences of $a$ in column $i$

- This has an entropy-based interpretation
  - Let $X_i$ be a random variable representing a character in column i
  - Consider each entry of column i to be observations of $X_i$ across multiple independent experiments
  - We estimate $P(X_i = a)$ by $p_{ia} = \dfrac{c_i^a}{n}$
  - Column score is proportional to the entropy of $X_i$

13

# Scoring an alignment: Entropy based score

- High entropy: More uniform distribution/more variability of characters
- Low entropy: Less uniform distribution/less variability of characters

$$S(m_i) = -\sum_a c_i^a \log(p_{ia})$$

# Scoring of a column: Sum of Pairs

- Compute the sum of the pairwise scores

$$S(m_i) = \sum_{k<l} s(m_i^k, m_i^l)$$

Iterate over all pairs of rows in the column

$s(m_i^k, m_i^l)$ Substitution score from a substitution/match matrix such as BLOSUM or PAM

# Dynamic Programming (DP) for global multiple sequence alignment

- Assume columns are independent
  - Score of alignment is sum of column scores
- Generalization of methods for pairwise alignment
  - consider k-dimensional matrix for k sequences (instead of 2-dimensional matrix)
  - each matrix element represents alignment score for k prefixes (instead of 2 prefixes)

# Notation for DP

- Assume we have $k$ sequences $x^1, \cdots, x^k$
- $i_1$ denotes the length of the prefix for sequence 1
- $i_2$ denotes the length of the prefix for sequence 2
- …
- $i_k$ denotes the length of the prefix for sequence $k$
- $x^k_{i_k}$ denotes the character at $i_k$ position of sequence $x^k$
- $F$: $k$-dimensional matrix where

$$F(i_1, i_2, \cdots, i_k)$$

denotes the score of the best alignment of the $i_1$, $i_2$.. $i_k$ prefixes of the sequences

# Recall the DP for the pairwise alignment

$$F(i_1, i_2) = \max \begin{cases} F(i_1 - 1, i_2 - 1) + S(x^1_{i_1}, x^2_{i_2}) \\ F(i_1, i_2 - 1) + S(-, x^2_{i_2}) \\ F(i_1 - 1, i_2) + S(x^1_{i_1}, -) \end{cases}$$

# DP for Multiple sequence alignment

$$F(i_1, \cdots, i_k) = \max \begin{cases} F(i_1 - 1, \cdots, i_k - 1) + S(x^1_{i_1}, \cdots, x^k_{i_k}) \\ F(i_1, i_2 - 1, \cdots, i_k - 1) + S(-, x^2_{i_2}, \cdots, x^k_{i_k}) \\ F(i_1 - 1, i_2, \cdots, i_k - 1) + S(x^1_{i_1}, -, \cdots, x^k_{i_k}) \\ \vdots \\ F(i_1, i_2 - 1, \cdots, i_k) + S(-, x^2_{i_2}, \cdots, -) \\ \vdots \end{cases}$$

max score of alignment for the $k$ prefixes

How many items do we need to maximize over?    $2^k - 1$

# DP algorithm is too expensive

- For $k$ sequences each of length $n$
  - Space complexity: $O(n^k)$
  - Time complexity: $O(n^k 2^k)$

# Summary

- Multiple alignment task and applications
- Two scoring functions
  - Entropy-based
  - Sum of pairs
- Inefficient dynamic programming extension