

# CMSC423: Bioinformatic Algorithms, Databases and Tools

## Lecture 12

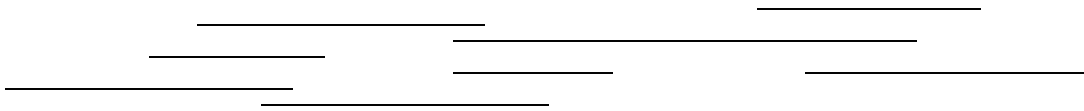
chaining algorithms  
multiple alignment

# Jobs

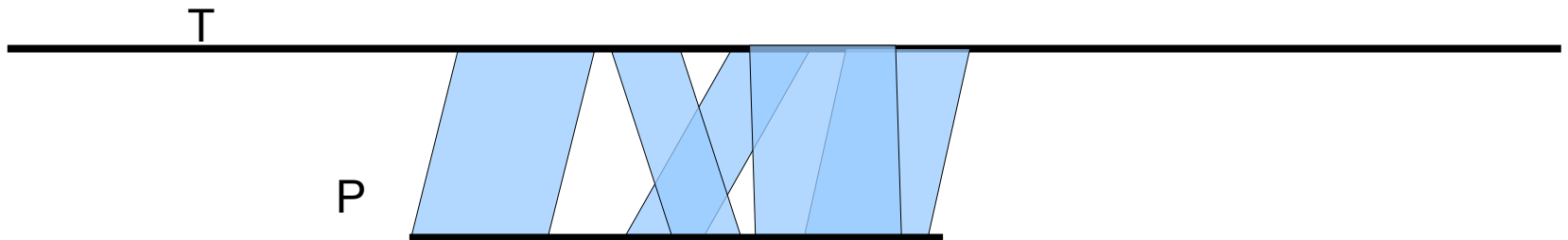
- Applied Predictive Technologies – looking for the best students – focus on databases (forwarded by Daniel Hackner) -not bioinformatics

# Chaining in 1-D

- Input: multiple overlapping intervals on a line
- Output: highest weight set of non-overlapping intervals
- Weight could be length of interval, or Smith-Waterman score, etc.

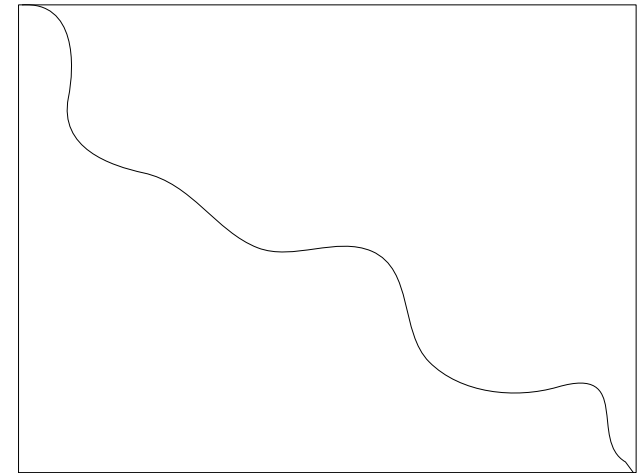
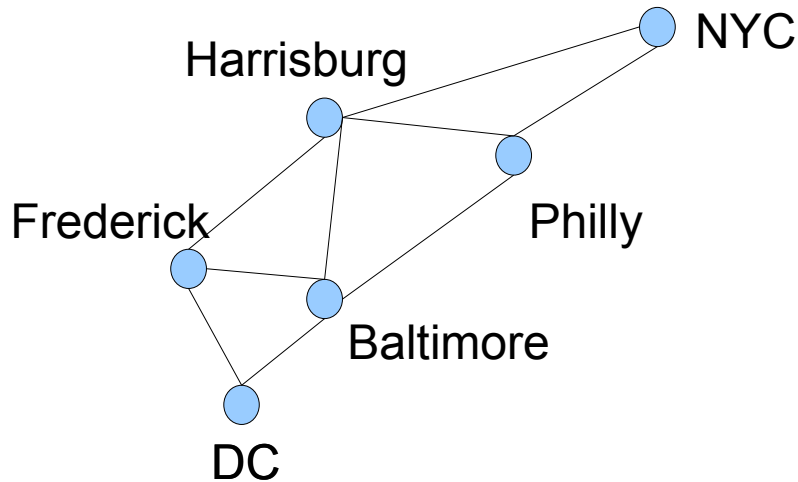


- Rationale? The pattern can have multiple inconsistent exact matches to the text – we want to pick a longest consistent set



# Path “planning” and dynamic programming

- One intuitive way to think about dynamic programming
  - similar to finding shortest path between two points
  - at each “point” ask – what are all possible ways to get here?
  - pick the best (shortest, fastest, etc.)

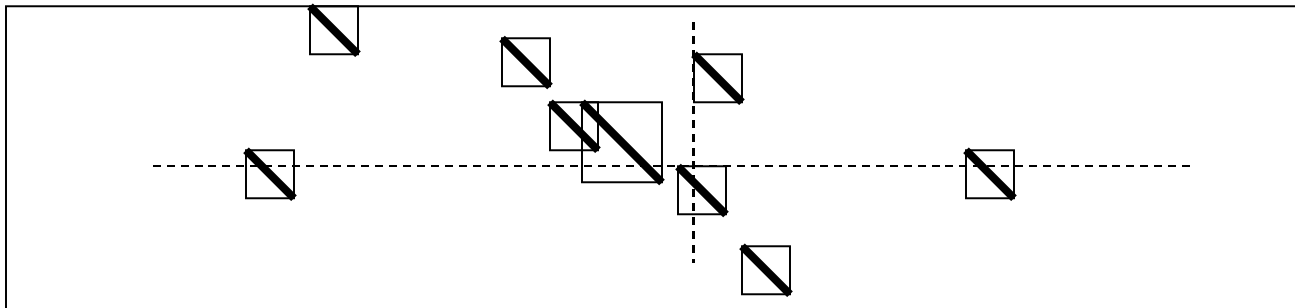


# Chaining in 1D

- Sort the endpoints (starts, ends) of the intervals
- For every interval  $j$ , store  $V[j]$  – best score of a chain ending in  $j$
- $MAX$  – store highest  $V[j]$  seen sofar
- Process endpoints in increasing order of  $x$  coordinate
- If we encounter left end (start) of interval  $j$ 
  - $V[j] = \text{weight}(j) + MAX$
- If we encounter right end (end) of interval  $j$ 
  - $MAX = \max\{V[j], MAX\}$
- Running time?

# Chaining in 2-D

- Easy to do in  $O(n^2)$  ( $n$  - # of intervals)
- View alignments as "boxes"
- All boxes in a chain must follow each other in a "diagonal" order, i.e. the range of the x coordinates and y coordinates of any two boxes in a chain cannot overlap
- Similar to 1-D approach except at each step we must check if current box can extend any of the previously built chains
- $V[j] = \max_{\text{all previous boxes } k} \{V[k] + \text{weight}(j)\}$
- More complex algorithm leads to  $O(n \log n)$  running time



# Multiple sequence alignment

# Multiple sequence alignment

- Simultaneously identify relationship between multiple sequences

```

HBB_HUMAN      FFESFGDLSTPDAVMGNPKVKAHGKKVL-----GAFSDGLAHL DNLKGTFF
HBB_HORSE     FFDSFGDLSNPGAVMG NPKVKAHGKKVL-----HSFGEGVHHL DNLKGTFF
HBA_HUMAN     YFPHF-DLS-----HGSAQVKGHGK KVA-----DALTNAVAHVDDMPNAL
HBA_HORSE     YFPHF-DLS-----HGSAQVKAHGK KVG-----DALTLAVGHLDDLPGAL
MYG_PHYCA     KFDREFKHLKTEAEMKASEDLKKHGVTVL-----TALGAILKKKGHHEAEL
GLB5_PETMA    FFPKFKGLTTADQLKKSADVRWHAERII-----NAVNDAVASMD DTEKMS
LGB2_LUPLU    LFSFLKGTSEVP--QNNPELQAHAGKVF KLVYEAAIQLQVTGVVVVTDATL
  
```

\* : . . . . . \* . : . . :

- Note: multiple alignment implies (not necessarily optimal) pairwise alignment between the individual sequences

```

HBA_HUMAN     YFPHF-DLS-----HGSAQVKGHGK KVA-----DALTNAVAHVDDMPNAL
HBA_HORSE     YFPHF-DLS-----HGSAQVKAHGK KVG-----DALTLAVGHLDDLPGAL
  
```



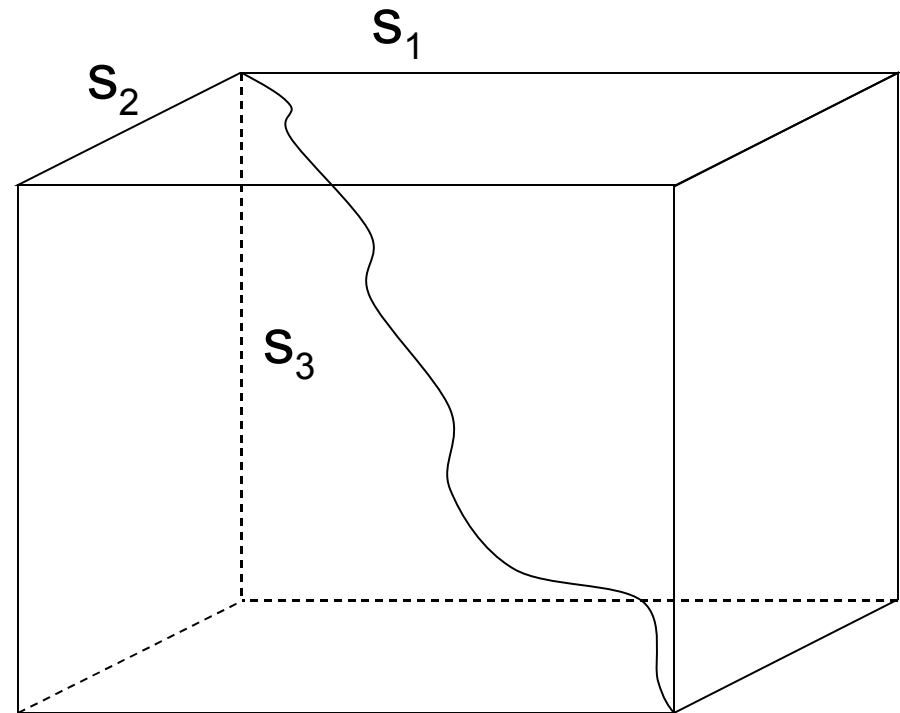
# Multiple alignment – formal definition

- $M$  – multiple sequence alignment for  $s_1, \dots, s_k$
- $D(s_i, s_j)$  – optimal score of alignment between  $s_i, s_j$
- $d(s_i, s_j)$  – score of alignment btwn  $s_i, s_j$  induced by  $M$
- score of  $M$   $d(M) = \sum_{\text{all pairs } s_i, s_j} d(s_i, s_j)$
- also called sum-of-pairs
  
- Optimal multiple alignment minimizes  $d(M)$
  
- Computing optimal  $d(M)$  is NP hard
- Note: in multiple alignment we think of "distance" rather than "similarity"

# But....here's a solution

- Dynamic programming solution. e.g. 3 sequences
- $\text{Score}(i, j, k)$  – optimal alignment between  $s1[1..i]$ ,  $s2[1..j]$ ,  $s3[1..k]$  – do DP as usual

- $s(i,j,k) = \max \{$   
     $s(i-1, j-1, k-1) +$   
     $\text{match}(s1[i], s2[j], s3[k]),$   
    ...



# But... it's expensive

- 3 sequences – need to fill in the cube  $O(n^3)$
- k sequences – k-dimensional cube  $O(n^k)$  time/space
- There are tricks that can help – similar to AI techniques for reducing the search space
- Basic idea – if we can estimate optimal score, we can prune the search space.
- Note – these are just heuristics – not guaranteed to work faster

# Alternative – approximation algorithm

- Can we efficiently compute a multiple alignment with a score that's not too bad?
- The Star method:
  - build all  $k^2$  pairwise alignments ( $O(k^2n^2)$ )
  - pick sequence  $sc$  that is closest to all other sequences:  
 $\sum_{s_i} D(sc, s_i)$  is minimal over all choices of  $sc$
  - iteratively align each sequence to  $sc$
- Theorem: sum-of-pairs score of star alignment is at most twice as big as optimal multiple alignment score

# Iterative alignment

SC YFPHF~~DLS~~HGSAQVKAHGKKVGDAL~~T~~LAVGHLDDLP~~GAL~~

- Take sequences  $s_i$  in order:

- align  $s_1$  with  $sc$  - results in gaps being inserted in both sequences

SC YFPHF~~DLS~~HGSAQVKAHGKKVGDAL~~T~~LAVGHLDDLP~~GAL~~

S1 YFPHF~~DLS~~HG-AQVKG--KKVADAL~~T~~NAVAHVDDMP~~NAL~~

- align  $s_2$  with  $sc$  - if gaps must be inserted – insert in previously aligned sequences

SC YFPHF-DLS-----HGSAQVKAHGKKVG-----DAL~~T~~LAVAH~~L~~DDLP~~GAL~~

S1 YFPHF-DLS-----HG-AQVKG-GKKVA-----DAL~~T~~NAVAHVDDMP~~NAL~~

S2 FFPKFKGLTTADQLKKSADV~~R~~WHAERII-----NAV~~N~~DAVASMDDTEK~~M~~S

- and so on (note: if gaps coincide with previously introduced gaps no need to change previously aligned sequences)

SC YFPHF-DLS-----HGSAQVKAHGKKVG-----DAL~~T~~LAVAH~~L~~DDLP~~GAL~~

S1 YFPHF-DLS-----HG-AQVKG-GKKVA-----DAL~~T~~NAVAHVDDMP~~NAL~~

S2 FFPKFKGLTTADQLKKSADV~~R~~WHAERII-----NAV~~N~~DAVASMDDTEK~~M~~S

S3 LFSFLKGTSEVP--QNNPELQAHAGKVF~~K~~LVYEAAIQ~~L~~QVTG~~V~~VVTDAT~~L~~

# Theorem proof

- Theorem: star alignment is 2-optimal
- Assumption: distances obey triangle inequality

$$\text{OPT} = \sum_{s_i, s_j} d^*(s_i, s_j) \geq \sum_{s_i, s_j} D(s_i, s_j) \geq k \sum_{s_i} D(s_i, s_c)$$

$$\begin{aligned} \text{STAR} &= \sum_{s_i, s_j} d(s_i, s_j) \leq \sum_{s_i, s_j} (D(s_i, s_c) + D(s_j, s_c)) \quad \# \text{ triangle ineq.} \\ &= \sum_{s_j, s_j} D(s_j, s_c) + \sum_{s_j, s_j} D(s_i, s_c) \\ &= 2k \sum_{s_i} D(s_i, s_c) \end{aligned}$$

$$\Rightarrow \text{STAR/OPT} \leq 2 \quad \text{Q.E.D}$$

note:  $\sum_{s_i} D(s_i, s_c)$  – is score optimized by choice of  $s_c$

$d^*(s_i, s_j)$  – score of alignment btwn  $s_i, s_j$  within optimal alignment

$d(s_i, s_j)$  – score of alignment btwn  $s_i, s_j$  within star alignment

$D(s_i, s_j)$  – score of optimal alignment btwn  $s_i, s_j$

