# Reflections|Projections 2020

Largest student run tech conference in the midwest, taking place September 21 through September 26
Virtual for the first time in 2020; featuring 20+ speakers and 5+ companies

Opportunity to hear about the bleeding edge on topics including
- Computer science for clinical methods (Tuesday @ 5)
- Ethical artificial intelligence (Wednesday @ 4, Friday @ 7)
- How to be successful in a startup (Monday @ 7, Wednesday @ 5)
- The essentials of product management (Thursday @ 5)
- Protecting privacy at a global scale (Friday @ 4)
- How to build accessible applications (Friday @ 6)
- And more! Check out our social media or website for more details

Network with top companies like Goldman Sachs, Grainger, IMC, Aechelon
Participate in a beginner-friendly, team-based AI competition with prizes
Swag pickup for attendees on campus
Opportunity to win free food!

Why wait? Register today at
**acmrp.org**

# Reflections | Projections



**Mark Ciaccio**

Senior Data Scientist @ AbbVie

**Tuesday, September 22 from 6-7PM**

reflectionsprojections.org

# CS 466
# Introduction to Bioinformatics
## Lecture 8

Mohammed El-Kebir

September 18, 2020

# Course Announcements

**Instructor:**

- Mohammed El-Kebir (melkebir)
- Office hours: Wednesdays, 3:15-4:15pm

**TAs:**

- Sarah Christensen (sac2) – Mondays, 3-4pm
- Wesley Wei Qian (weiqian3) – Fridays, 9-10am

Homework 2 released on 9/23 and due 10/1 by 11:59pm (CT)

Midterm will be 10/7, 7-10pm (CT)
Conflict midterm will be 10/8, 8-11am (CT).
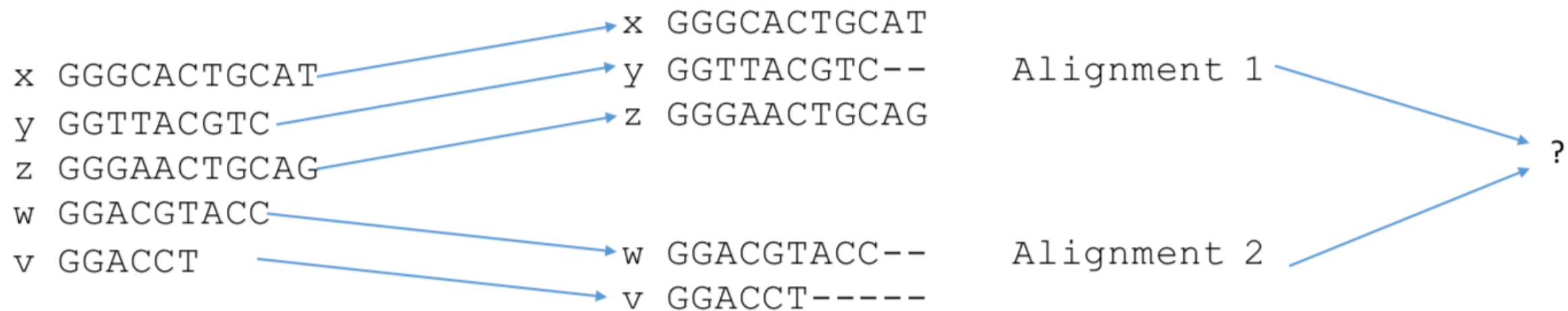You will need to sign up for conflict on Piazza.

# Outline

- Progressive alignment
  - Current methods
- Tree and star alignment

*[handwritten annotations in red: "heuristics → no guarantee on opt., or low bnd." and "approximation alg."]*

**Reading:**

- Material based on Chapter 14.6 in book "Algorithms on Strings, Trees and Sequences" by Dan Gusfield
- Lecture notes

# Heuristic: Iterative/Progressive Alignment

Iteratively add strings (or alignments) to existing alignment(s).

```
x GGGCACTGCAT              x GGGCACTGCAT
y GGTTACGTC                y GGTTACGTC--      Alignment 1
z GGGAACTGCAG              z GGGAACTGCAG
w GGACGTACC
v GGACCT                                                      ?
                           w GGACGTACC--      Alignment 2
                           v GGACCT-----
```

Issues:

1. How to merge alignments?
2. What order to use in merging strings/alignments?

# Profile Representation of Multiple Alignment

```
-  A  G  G  C  T  A  T  C  A  C  C  T  G
T  A  G  -  C  T  A  C  C  A  -  -  -  G
C  A  G  -  C  T  A  C  C  A  -  -  -  G
C  A  G  -  C  T  A  T  C  A  C  -  G  G
C  A  G  -  C  T  A  T  C  G  C  -  G  G
```

```
A        1              1        .8
C   .6         1          .4  1    .6 .2
G        1  .2                .2       .4  1
T   .2              1     .6              .2
-   .2         .8                   .4 .8 .4
```

A **profile** $P = [p_{i,j}]$ is a $(|\Sigma| + 1) \times l$ matrix,
where $p_{i,j}$ is the frequency of $i$-th letter in $j$-th position of alignment

# Aligning String to Profile

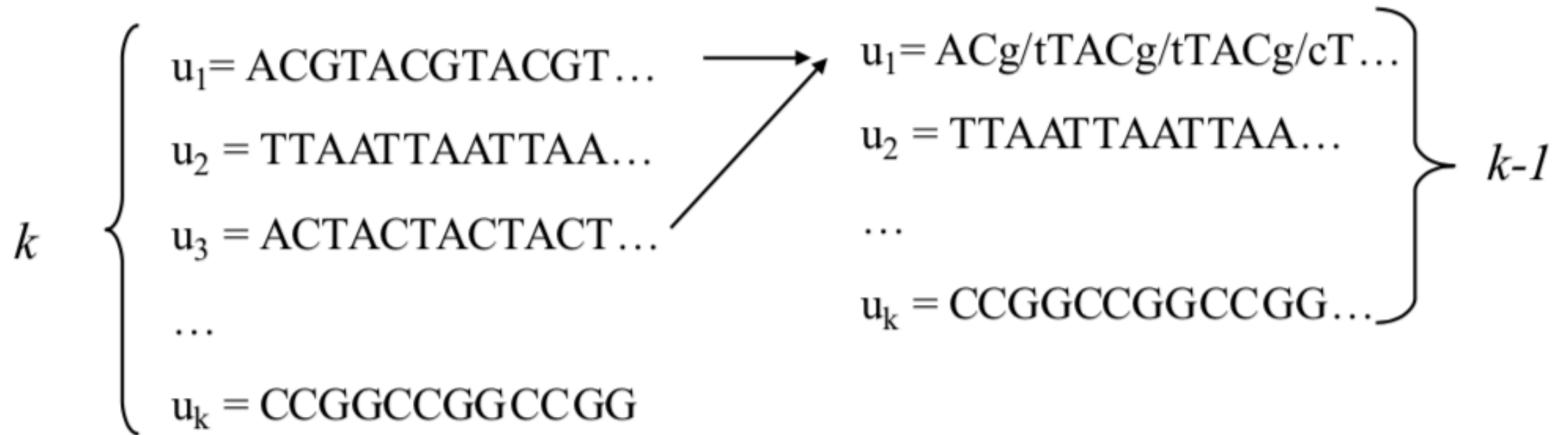$$\tau(x, j) = \sum_{y \in \Sigma \cup \{-\}} p_{y,j} \cdot \delta(x, y)$$

*(handwritten: profile, $v_i$, "profile = weighted average")*

$$s[i,j] = \max \begin{cases} 0, & \text{if } i = 0 \text{ and } j = 0, \\ s[i-1, j] + \delta(v_i, -), & \text{if } i > 0, \quad \textbf{Insert space in profile} \\ s[i, j-1] + \tau(-, j), & \text{if } j > 0, \quad \textbf{Insert space in string} \\ s[i-1, j-1] + \tau(v_i, j), & \text{if } i > 0 \text{ and } j > 0. \end{cases}$$

*(handwritten annotations: column of profile; column index of string)*

- $s[i,j]$ is optimal alignment of $v_1, \dots, v_i$ and first $j$ columns of $P$
- $\delta(x, y)$ is score for aligning characters $x$ and $y$
- $\tau(x, j)$ is score for aligning character $x$ and column $j$ of $P$

# Progressive Multiple Alignment: Greedy Algorithm

Choose most similar pair among *k* input strings, combine into a profile. This reduces the original problem to alignment of *k-1* sequences to a profile. Repeat.

$$k \begin{cases} u_1 = \text{ACGTACGTACGT}\ldots \\ u_2 = \text{TTAATTAATTAA}\ldots \\ u_3 = \text{ACTACTACTACT}\ldots \\ \ldots \\ u_k = \text{CCGGCCGGCCGG} \end{cases} \longrightarrow \begin{cases} u_1 = \text{ACg/tTACg/tTACg/cT}\ldots \\ u_2 = \text{TTAATTAATTAA}\ldots \\ \ldots \\ u_k = \text{CCGGCCGGCCGG}\ldots \end{cases} k\text{-}1$$

# Outline

- Progressive alignment
  - Current methods
- Tree and star alignment

**Reading:**

- Material based on Chapter 14.6 in book "Algorithms on Strings, Trees and Sequences" by Dan Gusfield

- Lecture notes
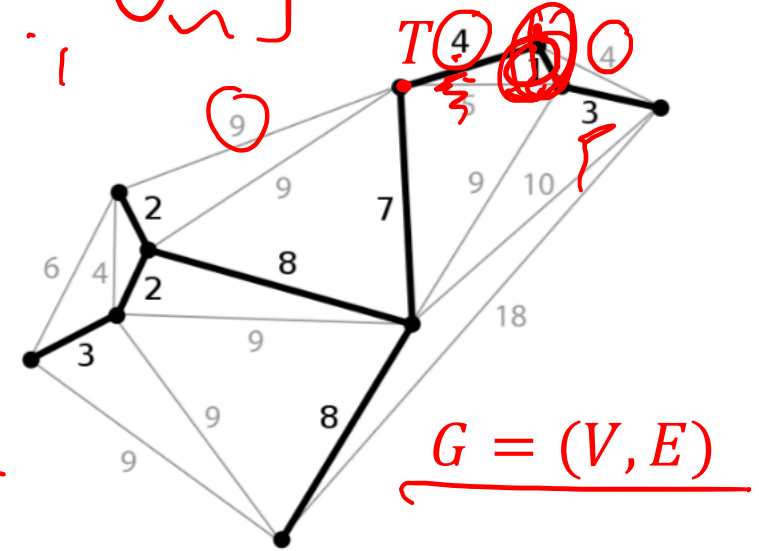
# Progressive Alignment – Feng and Doolittle (1987)

1. Compute pairwise sequence alignments of $n$ sequences

2. Generate complete graph $G = (V, E)$ with edge weights $w$ :
$E \to \mathbb{R}$

3. Compute a (rooted) minimum spanning tree $T$ of $G$

4. Perform sequence-sequence, sequence-alignment and alignment-alignment alignment to construct MSA according to guide tree $T$ (from most similar to least similar)



$$V = \{ T_1, \ldots, T_n \}$$

guide tree

$G = (V, E)$

Minimum spanning tree is a tree $T$ spanning all vertices of $G$ with minimum total weight

'Once a gap, always a gap'

# Progressive Alignment – ClustalW (1994)

*weighted*

- Widely used alignment method by Thompson, Higgins and Gibson (1994)
- W stands for weighted:
  - Input sequences are weighted to compensate for biased representation
  - Different substitution matrices depending on expected similarity in guide tree (BLOSUM80 for closely related sequences, and BLOSUM50 for distant sequences )
  - Position-specific gap-open and gap-extend penalties depending on context (hydrophobic vs. hydrophilic)

*affine gap penalties*

**Three steps:**

1. Construct pairwise alignments
2. Build guide tree $T$ using neighbor joining*
3. Progressive profile alignment guided by $T$

# ClustalW – Step 2: Guide Tree

## Create Guide Tree using the similarity matrix

("cluster" distances.  Details to come…)

|       | $V_1$ | $V_2$ | $V_3$ | $V_4$ |
|-------|-------|-------|-------|-------|
| $V_1$ | –     |       |       |       |
| $V_2$ | .17   | –     |       |       |
| $V_3$ | .87   | .28   | –     |       |
| $V_4$ | .59   | .33   | .62   | –     |

*(handwritten annotations: "distance not similarity", "neighbor joining", "additive tree")*

*(tree diagram showing $V_1$, $V_3$, $V_4$, $V_2$)*

ClustalW uses the neighbor-joining method

Guide tree roughly reflects evolutionary relationships

Calculate:

$V_{1,3}$ $\quad$ = alignment $(v_1, v_3)$

$V_{1,3,4}$ $\quad$ = alignment$((v_{1,3}), v_4)$

$V_{1,2,3,4}$ $\quad$ = alignment$((v_{1,3,4}), v_2)$

# ClustalW – Step 3: Progressive Alignment

- Start by aligning the two most similar sequences

- Following the guide tree, add in the next sequences, aligning to the existing alignment

- Insert gaps as necessary

```
FOS_RAT       PEEMSVTS-LDLTGGLPEATTPESEEAFTLPLLNDPEPK-PSLEPVKNISNMELKAEPFD
FOS_MOUSE     PEEMSVAS-LDLTGGLPEASTPESEEAFTLPLLNDPEPK-PSLEPVKSISNVELKAEPFD
FOS_CHICK     SEELAAATALDLG----APSPAAAEEAFALPLMTEAPPAVPPKEPSG--SGLELKAEPFD
FOSB_MOUSE    PGPGPLAEVRDLPG-----STSAKEDGFGWLLPPPPPPP---------------LPFQ
FOSB_HUMAN    PGPGPLAEVRDLPG-----SAPAKEDGFSWLLPPPPPPP---------------LPFQ
              .      . :    **  .       :..   *:.*    *    . *         **:
```

Dots and stars show how well-conserved a column is.

# MUSCLE (Edgar, 2004)

Multiple Sequence Comparison by Log-Expectation
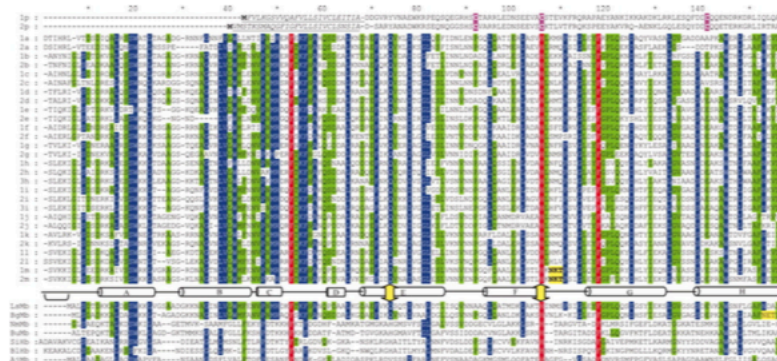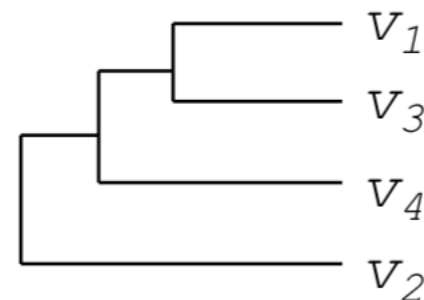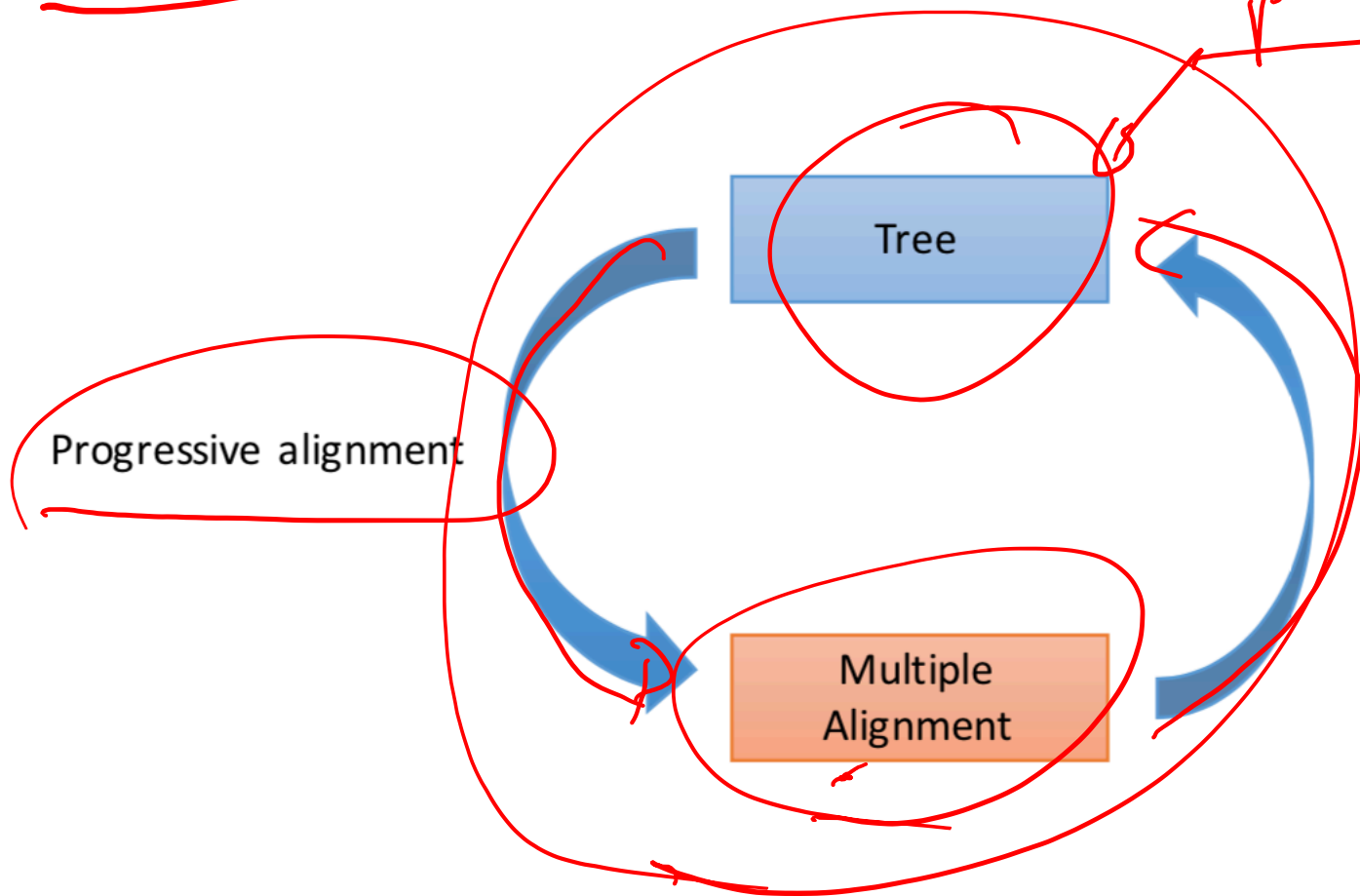
## Three phases:

1. Draft progressive alignment: fast heuristic

2. Improved progressive: use tree derived in phase 1

3. Refinement of MSA
   - Remove sequence from MSA and realign to profile of remaining sequences
   - Repeat until convergence

# Progressive MSA

phylogenetic tree



Tree

Multiple Alignment

Progressive alignment

$v_1$

$v_3$

$v_4$

$v_2$

Circularity!
Ideally, want to derive alignment and tree simultaneously → Hard

# Outline

- Progressive alignment  *≥ messy, heuristics*
  - Current methods

- Tree and star alignment  *← ("clean")*

**Reading:**

- Material based on Chapter 14.6 in book "Algorithms on Strings, Trees and Sequences" by Dan Gusfield

- Lecture notes

# Multiple Sequence Alignment Problem w/ SP-Score

**MSA-SP problem:** Given strings $\mathbf{v}_1, \ldots, \mathbf{v}_k$ and scoring function $\delta: (\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \to \mathbb{R}$, find multiple sequence alignment $\mathcal{M}^*$ with **<span style="color:red">maximum</span>** value of SP-score$(\mathcal{M}^*) = \sum_{i=1}^{k} \sum_{j=i+1}^{k} S(\mathbf{v}_i, \mathbf{v}_j)$ where $S(\mathbf{v}_i, \mathbf{v}_j)$ is the score of the induced pairwise alignment of $(\mathbf{v}_i, \mathbf{v}_j)$ in $\mathcal{M}^*$ using $\delta$

**Weighted SP-Edit Distance problem:** Given strings $\mathbf{v}_1, \ldots, \mathbf{v}_k$ and cost function $\delta: (\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \to \mathbb{R}$, find multiple sequence alignment $\mathcal{M}^*$ with **<span style="color:red">minimum</span>** value of SP-score$(\mathcal{M}^*) = \sum_{i=1}^{k} \sum_{j=i+1}^{k} S(\mathbf{v}_i, \mathbf{v}_j)$ where $S(\mathbf{v}_i, \mathbf{v}_j)$ is the cost of the induced pairwise alignment of $(\mathbf{v}_i, \mathbf{v}_j)$ in $\mathcal{M}^*$ using $\delta$

# Tree Alignment

$D(\bar{v}_1, \bar{v}_3) = 1$

$D(\bar{v}_i, \bar{v}_j)$

$=$ opt/min cost

of aligning

$\bar{v}_i$ and $\bar{v}_j$

$\left( \left( \left( (\bar{v}_1, \bar{v}_3) \, \bar{v}_2 \right), \bar{v}_4 \right), \bar{v}_5 \right) \, T$

Input: $\bar{v}_1, \dots, \bar{v}_k$

Let $T$ be a tree with $k$ vertices labeled uniquely by $\{\bar{v}_1, \dots, \bar{v}_k\}$.

$\begin{cases} A \, X \, X \, Z \\ A \, X \, \sim \, Z \end{cases}$

$d(\bar{v}_1, \bar{v}_3) = 1$

A multiple alignment $A$ of $\bar{v}_1, \dots, \bar{v}_n$ is consistent with $T$ if induced pairwise alignment of $\bar{v}_i$ and $\bar{v}_j$ has cost $D(\bar{v}_i, \bar{v}_j)$ for all edges $(v_i, v_j)$ of $T$

AXZ  $\bar{v}_1$

AXZ  $\bar{v}_2$

$\bar{v}_3$ AXXZ

a)

$\bar{v}_i =$

$\bar{v}_4$ AYZ

$\bar{v}_5$ AYXYZ

b)

| 3 | A | X | X | _ | Z |
|---|---|---|---|---|---|
| 1 | A | X | _ | _ | Z |
| 2 | A | _ | X | _ | Z |
| 4 | A | Y | _ | _ | Z |
| 5 | A | Y | X | X | Z |

$\Big\} A$

**Figure 14.6:** a. A tree with its nodes labeled by a (multi)set of strings, b. A multiple alignment of those strings that is consistent with the tree. The $T$ pairwise scoring scheme scores a zero for each match and a one for each mismatch or space opposite a character. The reader can verify that each of the four induced alignments specified by an edge of the tree has a score equal to its respective optimal distance. However, the induced alignment of two strings which do not label adjacent nodes may have a score greater than their optimal pairwise distance.

**Thm 1** (Gusfield). Let $T$ be a tree whose nodes are labeled uniquely by $\bar{v}_1, \ldots, \bar{v}_k$. We can compute an alignment $A(T)$ of $\bar{v}_1, \ldots, \bar{v}_k$ that is consistent with $T$ in $O(k^2 n^2)$ time.

$\leftarrow$ max length is $n$

**Proof.** WLOG, assume that vertices $\bar{v}_1, \ldots, \bar{v}_k$ are arranged s.t. $\bar{v}_1, \ldots, \bar{v}_i$ indices a connected subtree of $T$. Prove inductively that adding each string $\bar{v}_i$ while maintaining consistency can be done in $O(in^2)$ time.

$$\sum_{i=1}^{k} O(in^2) = O(n^2) \sum_{i=1}^{k} O(i) = O(n^2 \cdot O(k^2))$$
$$= O(k^2 n^2).$$

**Base case:** $i=2$. $\bar{v}_2$ and $\bar{v}_2$, compute pairwise alignment in $\boxed{O(n^2)}$ time. By definition this alignment will be consistent with

$$T[\{\bar{v}_1, \bar{v}_2\}]$$



$$\bar{v}_1$$
$$\bar{v}_2$$
$$O(i n^2)$$

resulting alignment is consistent with $T[\{\bar{v}_1, \bar{v}_2\}]$ in time $O(n^2)$.

**Step:** $i > 2$. **IH:** We are given tree $T'$ and $A(T')$ that is consistent with strings $\bar{v}_1, \ldots, \bar{v}_{i-1}$.

Consider $\boxed{\bar{v}_i}$. By definition, there is a vertex among $\{\bar{v}_1, \ldots, \bar{v}_{i-1}\}$ adjacent to $\bar{v}_i$. Let's call this vertex $\bar{v}_j$. Extract gapped seq. $\tilde{v}_j$ corresponding to $\bar{v}_j$. We align $\tilde{v}_j$ and $\bar{v}_i$ with the added rule that $\boxed{\delta(-,-) = 0}$ i.e. Two opposing gaps have a cost of $0$.

Let $\tilde{v}_j'$ and $\tilde{v}_i$ be the two resulting gapped sequences. If this new alignment did not insert new gaps into $\tilde{v}_j$ then we add $\tilde{v}_i$ to $A(T')$ ( $\tilde{v}_j' = \tilde{v}_j$ ). The result is a multiple alignment with one more string, induced cost of $\tilde{v}_j' = \tilde{v}_j$ and $\tilde{v}_c$ $\boxed{\text{equals } D(\bar{v}_j, \bar{v}_i)}$

However, if opt. alignment of $\tilde{v}_j$ and $\bar{v}_c$ inserted a gap in $\tilde{v}_j$ ( i.e. $\tilde{v}_j' \neq \tilde{v}_j$ ), say b/w characters $\ell$ and $\ell+1$ .

23

We insert a gap b/w columns $\ell$ and $\ell+1$ in each sequence of the previous mult. alignment $A(T')$. Observe that the induced costs of edges of $T'$ remains unchanged.
$$\bar{v}_1, \ldots, \bar{v}_{i-1}.$$
On the other hand cost of induced pairwise alignment of $\bar{v}_j$ and $\bar{v}_i$ equals $D(\bar{v}_j, \bar{v}_i)$.

In both cases, the new alignment is consistent with the tree $T'$ extended with vertex $\bar{v}_i$.

For the running time, observe that $A(T')$ is composed of $(i-1)$ supported seqs. each of length at most $(i-1)n$. Worst case $\tilde{v}_j$ has length $(i-1)n = \underline{O(in)}$ while $\bar{v}_i$ has length $\leq n$ $O(n)$. We can compute alignment b/w $\tilde{v}_j$ and $\bar{v}_i$ in $\underline{O(in^2)}$ time

Input $\bar{v}_1, \ldots, \bar{v}_k \in \Sigma^*$ ; $\boxed{\delta} : (\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \mapsto \mathbb{R}$

Goal : Find multiple alignment $A$ with SP score
$SP(A)$ that comes with approximation guarantee

$A \qquad A^*$

$\downarrow$

$SP(A) \geqslant SP(A^*)$

$\dfrac{SP(A) = 10}{SP(A^*) = 5} \leq 2$

$$\dfrac{SP(A)}{SP(A^*)} \leq \dfrac{2}{?}$$

**Def 2.** A cost function $\delta$ satisfies the triangle inequality if

$$\delta(x,z) \leq \delta(x,y) + \delta(y,z)$$

$$1 \qquad \leq \qquad 1 \qquad + \quad 1$$

$$\forall x,y,z \in \Sigma \cup \{-\}$$

classical edit dist.

|   | $x$ | $y$ | $z$ |
|---|-----|-----|-----|
| $x$ | ⓪ | 1 | 1 |
| $y$ | 1 | ⓪ | 1 |
| $z$ | 1 | 1 | ⓪ |

27

$$\Rightarrow x = y = z$$

$$\delta(x,z) \leq \delta(x,y) + \delta(y,z)$$

$$\delta(x,x) \leq \delta(x,x) + \delta(x,x)$$

$$0 \leq \boxed{0} + \boxed{0}$$

$D(\bar{v}_i, \bar{v}_j)$ opt cost of aligning $\bar{v}_i$ and $\bar{v}_j$
(in isolation).

Given a mult. alignment $A$

$d(\bar{v}_i(A), \bar{v}_j(A))$ cost of induced $\nearrow$ by $A$

pairwise alignment of $\bar{v}_i$ and $\bar{v}_j$

$$D(\bar{v}_i, \bar{v}_j) \leq d(\bar{v}_i(A), \bar{v}_j(A))$$

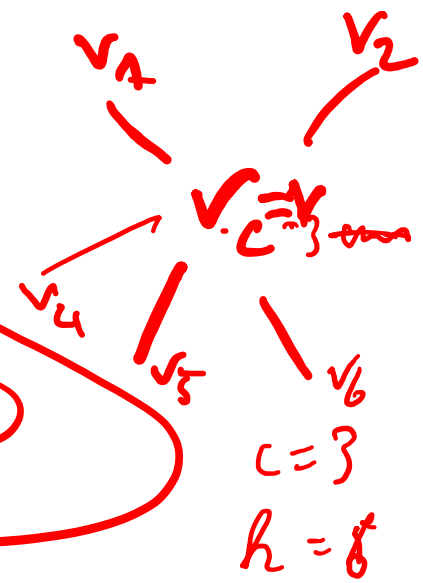$$\forall \text{ mult. } A, \bar{v}_i, \bar{v}_j.$$

**Def. 3** Given strings $\bar{v}_1, \ldots, \bar{v}_k$, the <u>center</u> <u>string</u> $\bar{v}_c$ is the input string ($c \in [k]$) that minimizes

$$\sum_{i=1}^{k} D(\bar{v}_c, \bar{v}_i)$$

$$D(\bar{v}_c, \bar{v}_c) = 0$$

$c = 3$
$k = 8$

The <u>center star</u> is a <u>star tree</u> of $k$ nodes with center node labeled by $\bar{v}_c$ and each of remaining $k-1$ nodes labeled by $\{\bar{v}_1, \ldots, \bar{v}_k\} \setminus \{\bar{v}_c\}$.

Construct alignment $A_c$ of $\bar{v}_1, \ldots, \bar{v}_k$ consistent with center star $T_c$.

① Identify $\bar{v}_c$
    — compute all pairwise alignments $O(k^2 n^2)$ time
    — identify $\bar{v}_c$    $O(k^2)$ time
    — gives you center star $T_c$

② Construct $A_c$ consistent with $T_c$
     $O(k^2 n^2)$ time.
       ↳ with SP cost $d(A_c)$

Thm.

$$\frac{d(A_c)}{d(A^*)} \leq \frac{2(k-1)}{k} \leq \boxed{2}$$

holds for
any cost function
satisfying
triangle ineq.

$k = 3$      $\frac{2 \cdot (3-1)}{3} = \frac{4}{3}$
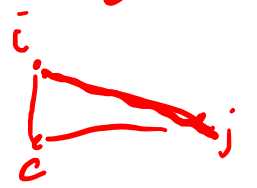
$k = 4$      $\frac{2(4-1)}{4} = \frac{6}{4}$

$\downarrow$

$k = \infty$          $\leq 2$

**Lemma 1.** Let $\delta$ be a cost function that satisfies the triangle ineq. and let $\bar{v}_c$ be the center string of a center star alignment $\underline{A_c}$ of input $\bar{v}_1, \ldots, \bar{v}_k$. Then, for any input string $\bar{v}_i$ and $\bar{v}_j$ it holds, that

$$d\left(\bar{v}_i\left(A_c\right), \bar{v}_j\left(A_c\right)\right) \leq d\left(\bar{v}_i\left(A_c\right), \bar{v}_c\left(A_c\right)\right) + d\left(\bar{v}_c\left(A_c\right), \bar{v}_j\left(A_c\right)\right)$$

$$\bigcirc\!\!\!= \quad D\left(\bar{v}_i, \bar{v}_c\right) + D\left(\bar{v}_c, \bar{v}_j\right)$$

Consider any column of $Ac$.

$\bar{v}_i(A_c)$

$\bar{v}_i(A_c)$

$v_j(A_c)$

$\begin{pmatrix} x \\ y \\ z \end{pmatrix}$ $\Sigma \cup \{-\}$

$$\boxed{\delta(x,z) \leq \delta(x,y) + \delta(y,z)}$$

**Thm 2.** $\dfrac{d(A_c)}{d(A^*)} \leq \dfrac{2(k-1)}{k}$

$d(A_c)$ SP score of $A_c$

$\sum\limits_{i=1}^{k} \sum\limits_{j=i+1}^{k} d(\bar{v}_i(A_c), \bar{v}_j(A_c))$

**Proof.** $f(A_c) = \sum\limits_{\substack{(i,j) \in [k]^2 \\ i \neq j}} d\left(\bar{v}_i(A_c), \bar{v}_j(A_c)\right)$

$f(A^*) = \sum\limits_{\substack{(i,j) \in [k]^2 \\ i \neq j}} d\left(\bar{v}_i(A^*), \tilde{v}_j(A^*)\right)$

Clearly, $f(A_c) = 2 \cdot d(A_c)$          $f(A^*) = 2 \cdot d(A^*)$

$$F(A_c) = \sum_{\substack{(i,j) \in [h]^2 \\ i \neq j}} d\left(\bar{v}_i(A_c), \bar{v}_j(A_c)\right)$$

$$\leq \{triangle \ ineq.\} \ \{lemma\}$$

$$\sum_{\substack{(i,j) \in [h]^2 \\ i \neq j}} \left[ d\left(\bar{v}_i(A_c), \bar{v}_c(A_c)\right) + d\left(\bar{v}_c(A_c), \bar{v}_j(A_c)\right) \right]$$

$$= \sum_{\substack{i,j \in [h]^2 \\ i \neq j}} \left[ D\left(\bar{v}_i, \bar{v}_c\right) + D\left(\bar{v}_c, \bar{v}_j\right) \right]$$

Observe that for any fixed $j$ each of the terms $D(\bar{v}_c, \bar{v}_j)$ and $D(\bar{v}_j, \bar{v}_c)$ show up $h-1$ times. Moreover, $D(\bar{v}_c, \bar{v}_j) = D(\bar{v}_j, \bar{v}_c)$

$$F(A_c) \leq \sum_{\substack{i,j \in [h]^2 \\ i \neq j}} \left[ D(\bar{v}_i, \bar{v}_c) + D(\bar{v}_c, \bar{v}_j) \right]$$

$$= 2(h-1) \sum_{j=1}^{h} D(\bar{v}_c, \bar{v}_j).$$

$$f(A^*) = \sum_{\substack{(i,j) \in [h]^2 \\ i \neq j}} d\left(\bar{v}_i(A^*), \bar{v}_j(A^*)\right)$$

$$\geq \sum_{\substack{(i,j) \in [h]^2 \\ i \neq j}} D\left(\bar{v}_i, \bar{v}_j\right)$$

$$= \sum_{i=1}^{h} \sum_{j=1}^{h} D\left(\bar{v}_i, \bar{v}_j\right)$$

$D(\bar{v}_i, \bar{v}_i) = 0$

cost of a star with center $\bar{v}_i$

Recall that $\overline{v_c}$ corresponds to the star with smallest cost.

$$f^*(A^*) \geq \sum_{i=1}^{h} \boxed{\sum_{j=1}^{h} D(\overline{v_i}, \overline{v_j})}$$

star centered at $\overline{v_i}$

$$\geq h \sum_{j=1}^{h} \boxed{D(\overline{v_c}, \overline{v_j})},$$

$$\boxed{\sum_{j=1}^{h} D(\overline{v_i}, \overline{v_j})} \geq \sum_{j=1}^{h} D(\overline{v_c}, \overline{v_j})$$

$\hookrightarrow$ center star

$$F(A_c) \leq 2(k-1) \sum_{j=1}^{k} D(\bar{v}_c, \bar{v}_j)$$

$$F(A^*) \geq k \sum_{j=1}^{k} D(\bar{v}_c, \bar{v}_j)$$

$$\frac{d(A_c)}{d(A^*)} = \frac{F(A_c)}{F(A^*)} \leq \frac{2(k-1) \sum_{j=1}^{k} D(\bar{v}_c, \bar{v}_j)}{k \sum_{j=1}^{k} D(\bar{v}_c, \bar{v}_j)}$$

$$\leq \lim_{k \to \infty} \frac{2(k-1)}{k} = 2$$

# Summary

1. Optimal pairwise alignment by dynamic programming in $O(n^2)$ time

2. Optimal multiple alignment with SP-score by dynamic programming in $O(k^2 2^k n^k)$ time

3. Multiple alignment with SP-score is NP-hard (Jiang and Wang, 1994)

4. Carrillo-Lipman enables us to decide whether alignment passes through a vertex $(i_1, i_2, i_3)$ for $k = 3$ sequences (generalizes to $k > 3$)

5. Progressive alignment methods are widely used, but come with no theoretical bounds on alignment quality

6. Star alignment gives 2-approximation algorithm

# History

- 1975 Sankoff
  Formulated MSA problem and gave dynamic programming solution
- 1988 Carrillo-Lipman
  Branch and Bound approach for MSA
- 1990 Feng-Doolittle
  Progressive alignment
- 1993 Gusfield
  Star alignment: 2-approximation algorithm
- 1994 Jiang and Wang
  MSA with SP-score is NP-hard
- 1994 Thompson-Higgins-Gibson: ClustalW
  Most popular multiple alignment program
- 2000 Notredam-Higgins-Heringa: T-coffee
  Use library of pairwise alignments
- 2004 Edgar: MUSCLE
  Refinement