# Summarizing Diverging String Sequences, with Applications to Chain-Letter Petitions

Patty Commins[1,2], David Liben-Nowell[1],
Tina Liu[1,3], and Kiran Tomlinson[1,4]

[1]Department of Computer Science, Carleton College
[2]Department of Mathematics, University of Minnesota
[3]Surescripts
[4]Department of Computer Science, Cornell University

CPM 2020

```
The US Congress has just authorized the President of
the US to go to war against Iraq. Please consider this
an urgent request.
UN Petition for Peace.
Stand for Peace.
Islam is not the Enemy.
War is NOT the Answer.
Today we are at a point of imbalance in the world and
are moving toward what may be the beginning of a
THIRD WORLD WAR.
If you are against this possibility, the UN is
gathering signatures in an effort to avoid a tragic
world event.
Please COPY (rather than Forward) this e-mail in a new
message, sign at the end of the list, and send it to
all the people whom you know. If you receive this list
with more than 500 names signed, please send a copy of
the message to:


usa@un.int <mailto:usa@un.int>
president@whitehouse.gov <mailto:president@whitehouse.gov>
```
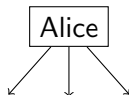
$\sim$ 3.5m emails

$\sim$ 170k signers

(Chierichetti, Kleinberg, & Liben-Nowell 2011)

Sent 20 February 2003, retrieved from G.W.B. Presidential Library

Please COPY (rather than Forward) this e-mail in a new
message, sign at the end of the list, and send it to
all the people whom you know.

Please COPY (rather than Forward) this e-mail in a new
message, sign at the end of the list, and send it to
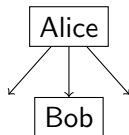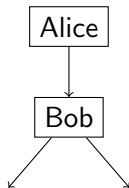all the people whom you know.

Alice

Please COPY (rather than Forward) this e-mail in a new
message, sign at the end of the list, and send it to
all the people whom you know.

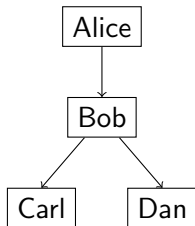Please COPY (rather than Forward) this e-mail in a new
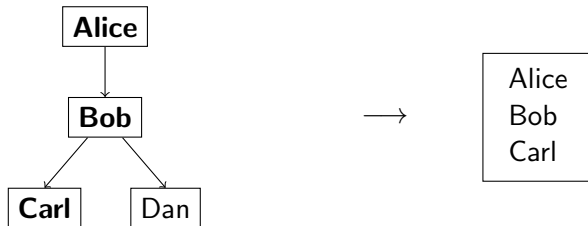message, sign at the end of the list, and send it to
all the people whom you know.

Please COPY (rather than Forward) this e-mail in a new
message, sign at the end of the list, and send it to
all the people whom you know.

Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

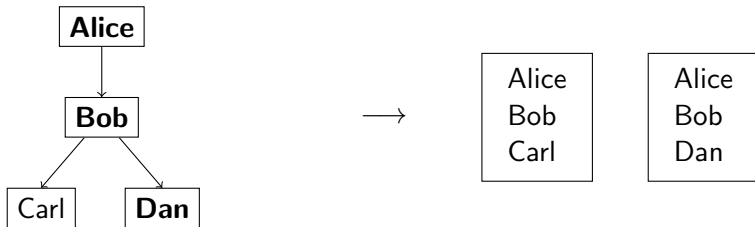Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

Please COPY (rather than Forward) this e-mail in a new message, sign at the end of the list, and send it to all the people whom you know.

| Alice |
| Bob |
| Carl |

| Alice |
| Bob |
| Dan |

**Central Question**

Can we reconstruct the propagation tree from signature lists?

$$\xleftarrow{?}$$

| Alice |
|-------|
| Bob |
| Carl |

| Alice |
|-------|
| Bob |
| Dan |

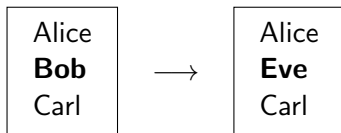# Challenge: Mutations

*People are bad at copy-paste.*

# Challenge: Mutations

*People are bad at copy-paste.*

1. Substitution

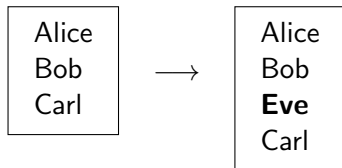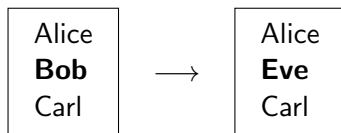| Alice | | Alice |
|-------|---------------|-------|
| **Bob** | $\longrightarrow$ | **Eve** |
| Carl | | Carl |

# Challenge: Mutations

*People are bad at copy-paste.*

1. Substitution

| Alice | | Alice |
|-------|-----|-------|
| **Bob** | $\longrightarrow$ | **Eve** |
| Carl | | Carl |

2. Insertion

| Alice | | Alice |
|-------|-----|-------|
| Bob | $\longrightarrow$ | Bob |
| Carl | | **Eve** |
| | | Carl |

# Challenge: Mutations

*People are bad at copy-paste.*

1. Substitution

   | Alice |              | Alice |
   | **Bob** | $\longrightarrow$ | **Eve** |
   | Carl |              | Carl |

2. Insertion

   | Alice |              | Alice |
   | Bob   | $\longrightarrow$ | Bob   |
   | Carl  |              | **Eve** |
   |       |              | Carl  |

3. Deletion

   | Alice |              | Alice |
   | **Bob** | $\longrightarrow$ | Carl  |
   | Carl  |              |       |

# Challenge: Mutations

*People are bad at copy-paste.*

1. Substitution

   | Alice |
   |-------|
   | **Bob** |
   | Carl |

   $\longrightarrow$

   | Alice |
   |-------|
   | **Eve** |
   | Carl |

2. Insertion

   | Alice |
   |-------|
   | Bob |
   | Carl |

   $\longrightarrow$

   | Alice |
   |-------|
   | Bob |
   | **Eve** |
   | Carl |

Character-level:
Carl $\to$ Car**o**l, Al**i**ce $\to$ Al**y**ce

3. Deletion

   | Alice |
   |-------|
   | **Bob** |
   | Carl |

   $\longrightarrow$

   | Alice |
   |-------|
   | Carl |

# Challenge: Mutations

*People are bad at copy-paste.*

1. Substitution

| Alice |
|-------|
| **Bob** |
| Carl |

$\longrightarrow$

| Alice |
|-------|
| **Eve** |
| Carl |

2. Insertion

| Alice |
|-------|
| Bob |
| Carl |

$\longrightarrow$

| Alice |
|-------|
| Bob |
| **Eve** |
| Carl |

Character-level:
Carl $\rightarrow$ Car**o**l, Al**i**ce $\rightarrow$ Al**y**ce

All present in the Iraq War petition
(Liben-Nowell & Kleinberg 2008)

3. Deletion

| Alice |
|-------|
| **Bob** |
| Carl |

$\longrightarrow$

| Alice |
|-------|
| Carl |

# Reconstruction with Mutations

| Alice | | Alice | | Alice |
|-------|--|-------|--|-------|
| Carl  | | Bob   | | Bob   |
| Eve   | | Carol | | Carl  |
|       | | Dan   | | Frank |

$\overset{?}{\longrightarrow}$

## Key chain letter features

# Reconstruction with Mutations

| Alice<br>Carl<br>Eve | Alice<br>Bob<br>Carol<br>Dan | Alice<br>Bob<br>Carl<br>Frank |
|---|---|---|

$\overset{?}{\longrightarrow}$

## Key chain letter features

1. One-ended growth

# Reconstruction with Mutations

| Alice | Alice | Alice |
|-------|-------|-------|
| Carl  | Bob   | Bob   |
| Eve   | Carol | Carl  |
|       | Dan   | Frank |

$\overset{?}{\longrightarrow}$

## Key chain letter features

1. One-ended growth
2. Divergence

# Reconstruction with Mutations

| Alice | Alice | Alice |
|-------|-------|-------|
| Carl  | Bob   | Bob   |
| Eve   | Carol | Carl  |
|       | Dan   | Frank |

$\overset{?}{\longrightarrow}$

## Key chain letter features

1. One-ended growth
2. Divergence
3. Mutation with inheritance

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof
3. Efficient optimal solution for two lists

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof
3. Efficient optimal solution for two lists
4. Fixed-parameter tractable: poly-time algorithm for $O(1)$ lists

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof
3. Efficient optimal solution for two lists
4. Fixed-parameter tractable: poly-time algorithm for $O(1)$ lists
5. Fast heuristic for arbitrary number of lists

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof
3. Efficient optimal solution for two lists
4. Fixed-parameter tractable: poly-time algorithm for $O(1)$ lists
5. Fast heuristic for arbitrary number of lists
6. Experimental evaluation on synthetic data

# Summary of Contributions

1. Formal definition of chain letter reconstruction problem
2. NP-hardness proof*
3. Efficient optimal solution for two lists
4. Fixed-parameter tractable: poly-time algorithm for $O(1)$ lists*
5. Fast heuristic for arbitrary number of lists
6. Experimental evaluation on synthetic data

* see paper

- Chain letters
  - Iraq war petition tree structure (Liben-Nowell & Kleinberg 2008; Golub & Jackson 2010; Chierichetti, Liben-Nowell, & Kleinberg 2011)
  - Tree reconstruction from plea (Bennett, Li, & Ma 2003)

# Related Work

- Chain letters
  - Iraq war petition tree structure (Liben-Nowell & Kleinberg 2008; Golub & Jackson 2010; Chierichetti, Liben-Nowell, & Kleinberg 2011)
  - Tree reconstruction from plea (Bennett, Li, & Ma 2003)
- One-ended growth and divergence
  - Trie (De La Briandais 1959; Fredkin 1960)
  - Online conversations (Kumar, Mahdian, & McGlohon 2010)

## Related Work

- Chain letters
  - Iraq war petition tree structure (Liben-Nowell & Kleinberg 2008; Golub & Jackson 2010; Chierichetti, Liben-Nowell, & Kleinberg 2011)
  - Tree reconstruction from plea (Bennett, Li, & Ma 2003)
- One-ended growth and divergence
  - Trie (De La Briandais 1959; Fredkin 1960)
  - Online conversations (Kumar, Mahdian, & McGlohon 2010)
- Divergence and mutation
  - Molecular phylogenetics (Yang & Rannala 2012)
  - Stories; e.g., Little Red Riding Hood (Tehrani 2013)

# Outline

# Problem Definition, Informally

## DSSSP (Diverging String Sequence Summarization Problem)

Given diverging string sequences:

| Alice | | Alice | | Alice |
| Carl | | Bob | | Bob |
| Eve | | Carol | | Carl |
| | | Dan | | Frank |

# Problem Definition, Informally

## DSSSP (Diverging String Sequence Summarization Problem)

Given **diverging string sequences**:

| **Alice** | **Alice** | Alice |
|-----------|-----------|-------|
| **Carl**  | **Bob**   | Bob   |
| Eve       | **Carol** | Carl  |
|           | Dan       | Frank |

# Problem Definition, Informally

## DSSSP (Diverging String Sequence Summarization Problem)

Given diverging string sequences:

| Alice | | Alice | | Alice |
| Carl | | Bob | | Bob |
| Eve | | Carol | | Carl |
| | | Dan | | Frank |

Find best summary tree:

# Problem Definition, Informally

## DSSSP (Diverging String Sequence Summarization Problem)

Given diverging string sequences:

Find **best** summary tree:

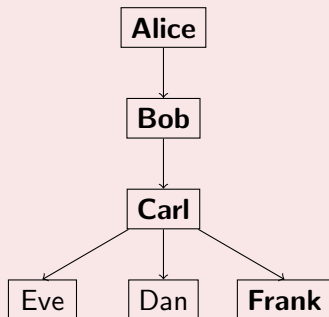| Alice | Alice | Alice |
| Carl | Bob | Bob |
| Eve | Carol | Carl |
| | Dan | Frank |

Alice
Carl
Eve

Alice
Bob
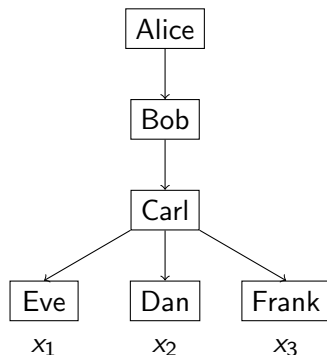Carol
Dan

Alice
Bob
Carl
Frank

### Accurate representation

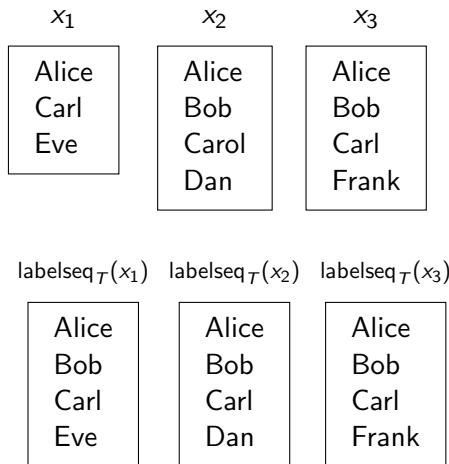# Competing Objectives

# Competing Objectives



11 / 26

# Competing Objectives

# Competing Objectives

# Competing Objectives



11 / 26

# Competing Objectives

# Measuring Representation Accuracy

$x_1$

| Alice |
| Carl |
| Eve |

$x_2$

| Alice |
| Bob |
| Carol |
| Dan |

$x_3$

| Alice |
| Bob |
| Carl |
| Frank |

labelseq$_T(x_1)$

| Alice |
| Bob |
| Carl |
| Eve |

labelseq$_T(x_2)$

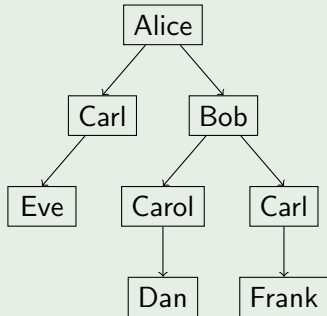| Alice |
| Bob |
| Carl |
| Dan |

labelseq$_T(x_3)$

| Alice |
| Bob |
| Carl |
| Frank |

### AED$(x, y)$

Allowed operations:

1. Insert string into $x$
2. Substitute string

Costs using Levenshtein ED

# Minimizing Redundancy

# Minimizing Redundancy
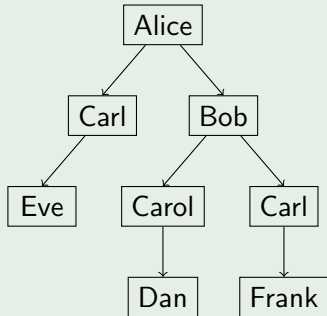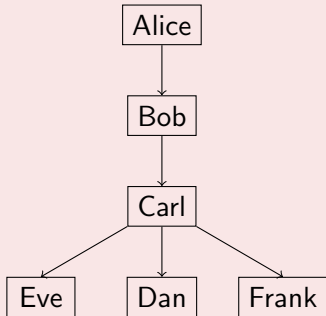


**Accurate representation**

Alice → Carl, Bob
Carl → Eve
Bob → Carol, Carl
Carol → Dan
Carl → Frank

8 nodes

**Minimal redundancy**

Alice → Bob → Carl → Eve, Dan, Frank

6 nodes

# Minimizing Redundancy



8 nodes            6 nodes

$\Rightarrow$ Cost $\lambda$ per node

# Problem Definition, Formally

## DSSSP

Given diverging string sequences $x_1, \ldots, x_m$ and node cost $\lambda$, find tree $T$ that minimizes

$$\mathrm{err}_\lambda(T) = \underbrace{\sum_{i=1}^{m} \mathrm{AED}(x_i, \mathrm{labelseq}_T(x_i))}_{\text{loss}} + \underbrace{\lambda \cdot |T|}_{\text{regularization}}$$

# Outline

# Two sequences: dynamic programming

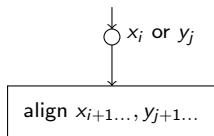Two string sequences $x, y$; align $x_{i...}$ and $y_{j...}$

$$\text{EDG}(i, j) = \min \left\{ \right.$$

# Two sequences: dynamic programming

Two string sequences $x, y$; align $x_{i\ldots}$ and $y_{j\ldots}$

$$\text{EDG}(i,j) = \min \begin{cases} \text{EDG}(i+1, j+1) + \lambda + \text{ED}(x_i, y_j) & \text{(substitution)} \\ \\ \\ \\ \end{cases}$$

substitution

$x_i$ or $y_j$
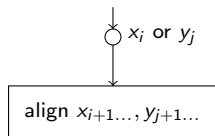
align $x_{i+1\ldots}, y_{j+1\ldots}$
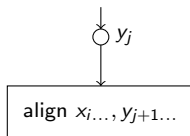
# Two sequences: dynamic programming

Two string sequences $x, y$; align $x_{i...}$ and $y_{j...}$

$$\text{EDG}(i,j) = \min \begin{cases} \text{EDG}(i+1, j+1) + \lambda + \text{ED}(x_i, y_j) & \text{(substitution)} \\ \text{EDG}(i, j+1) \quad\quad + \lambda + \text{ED}(\varepsilon, y_j) & \text{(insertion)} \\ \\ \end{cases}$$

substitution

$\downarrow$ $x_i$ or $y_j$

align $x_{i+1...}, y_{j+1...}$

insertion

$\downarrow$ $y_j$

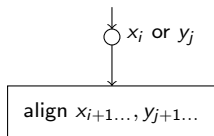align $x_{i...}, y_{j+1...}$
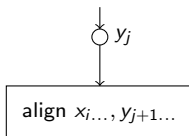
# Two sequences: dynamic programming

Two string sequences $x, y$; align $x_{i...}$ and $y_{j...}$

$$\text{EDG}(i,j) = \min \begin{cases} \text{EDG}(i+1, j+1) + \lambda + \text{ED}(x_i, y_j) & \text{(substitution)} \\ \text{EDG}(i, j+1) + \lambda + \text{ED}(\varepsilon, y_j) & \text{(insertion)} \\ \text{EDG}(i+1, j) + \lambda + \text{ED}(x_i, \varepsilon) & \text{(deletion)} \end{cases}$$
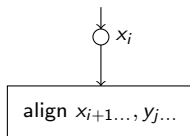
substitution

$x_i$ or $y_j$

align $x_{i+1...}, y_{j+1...}$

insertion

$y_j$

align $x_{i...}, y_{j+1...}$
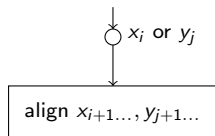
deletion

$x_i$

align $x_{i+1...}, y_{j...}$
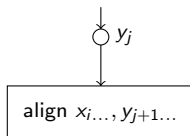
# Two sequences: dynamic programming

Two string sequences $x, y$; align $x_{i\ldots}$ and $y_{j\ldots}$

$$\text{EDG}(i,j) = \min \begin{cases} \text{EDG}(i+1, j+1) & + \lambda + \text{ED}(x_i, y_j) & \text{(substitution)} \\ \text{EDG}(i, j+1) & + \lambda + \text{ED}(\varepsilon, y_j) & \text{(insertion)} \\ \text{EDG}(i+1, j) & + \lambda + \text{ED}(x_i, \varepsilon) & \text{(deletion)} \\ \lambda(|x| - i + 1) & + \lambda(|y| - j + 1) & \text{(give up)} \end{cases}$$
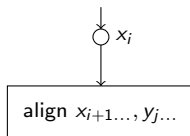
# Two sequences: dynamic programming

Two string sequences $x, y$; align $x_{i\ldots}$ and $y_{j\ldots}$

$$\text{EDG}(i,j) = \min \begin{cases} \text{EDG}(i+1, j+1) & + \lambda + \text{ED}(x_i, y_j) & \text{(substitution)} \\ \text{EDG}(i, j+1) & + \lambda + \text{ED}(\varepsilon, y_j) & \text{(insertion)} \\ \text{EDG}(i+1, j) & + \lambda + \text{ED}(x_i, \varepsilon) & \text{(deletion)} \\ \lambda(|x| - i + 1) & + \lambda(|y| - j + 1) & \text{(give up)} \end{cases}$$
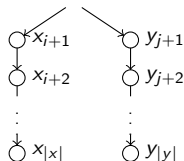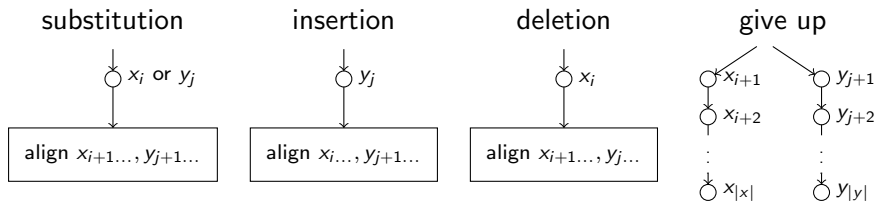


substitution — $x_i$ or $y_j$ — align $x_{i+1\ldots}, y_{j+1\ldots}$

insertion — $y_j$ — align $x_{i\ldots}, y_{j+1\ldots}$

deletion — $x_i$ — align $x_{i+1\ldots}, y_{j\ldots}$

give up — $x_{i+1}$ $x_{i+2}$ $\vdots$ $x_{|x|}$ — $y_{j+1}$ $y_{j+2}$ $\vdots$ $y_{|y|}$

## Theorem

*This produces an optimal two-sequence DSSSP solution.*

# Algorithm for more sequences

## Theorem

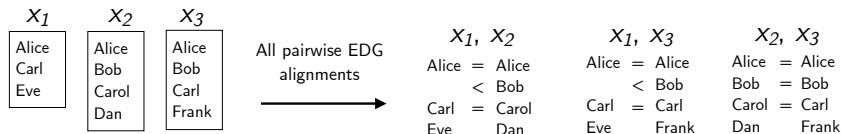*DSSSP is NP-hard with an unbounded number of sequences.*

# Algorithm for more sequences

## Theorem

*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

Repeatedly merge pair of sequences that diverges last

## Theorem

*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

Repeatedly merge pair of sequences that diverges last
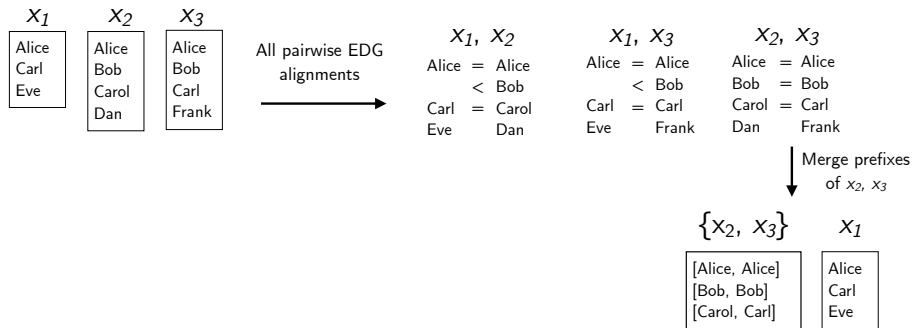
| $x_1$ | $x_2$ | $x_3$ |
|---|---|---|
| Alice | Alice | Alice |
| Carl | Bob | Bob |
| Eve | Carol | Carl |
| | Dan | Frank |

# Algorithm for more sequences

## Theorem

*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

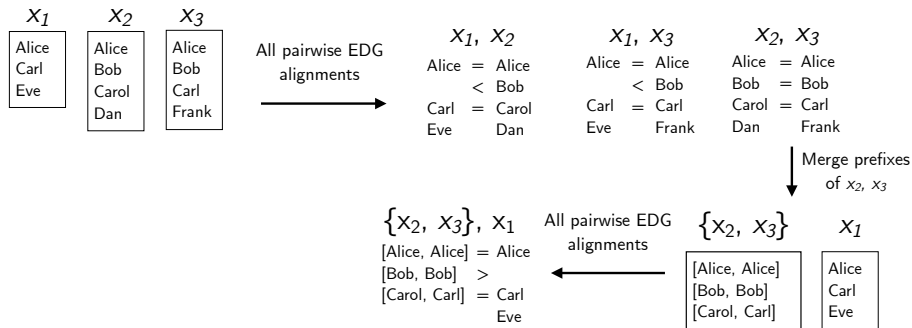Repeatedly merge pair of sequences that diverges last

# Algorithm for more sequences

## Theorem

*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

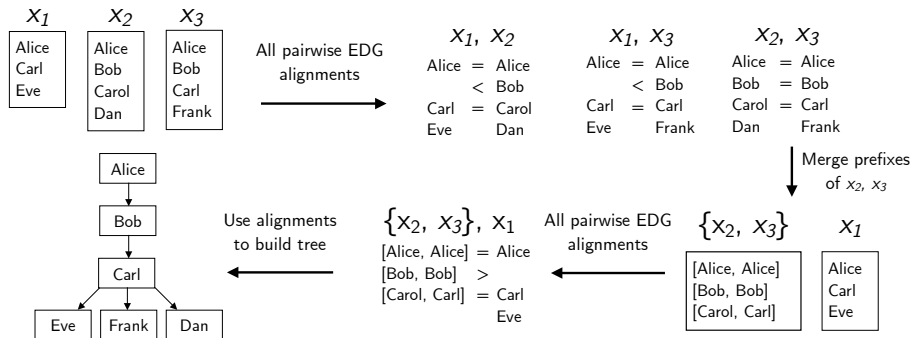Repeatedly merge pair of sequences that diverges last

# Algorithm for more sequences

## Theorem
*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

Repeatedly merge pair of sequences that diverges last

# Algorithm for more sequences

## Theorem

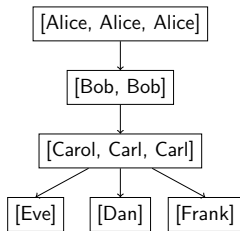*DSSSP is NP-hard with an unbounded number of sequences.*

## Idea: progressive alignment (Feng & Doolittle 1987)

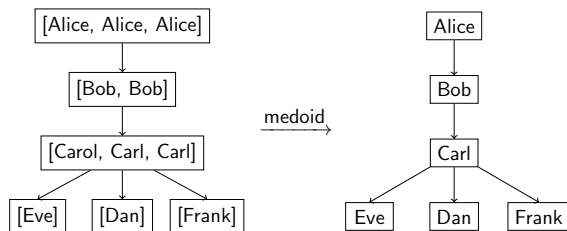Repeatedly merge pair of sequences that diverges last

1. Labeling the final tree

# Algorithm details
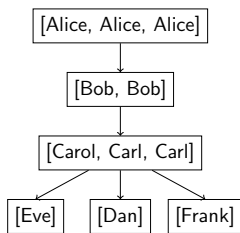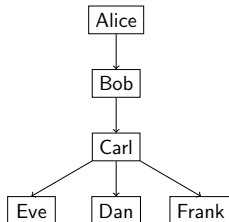
1. Labeling the final tree

# Algorithm details

1. Labeling the final tree



[Alice, Alice, Alice] → [Bob, Bob] → [Carol, Carl, Carl] → [Eve], [Dan], [Frank]
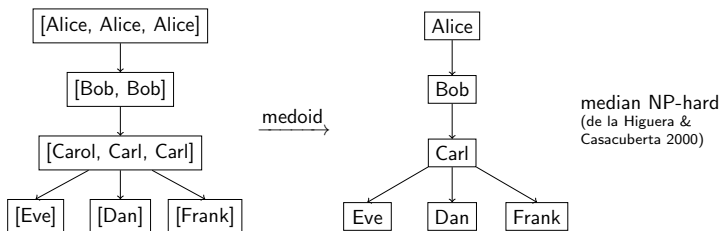
medoid

Alice → Bob → Carl → Eve, Dan, Frank

median NP-hard
(de la Higuera & Casacuberta 2000)
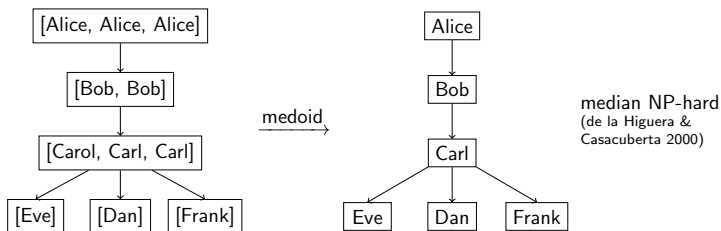
# Algorithm details

1. Labeling the final tree



2. Generalizing EDG to sequences of *lists* of strings

# Algorithm details

1. Labeling the final tree



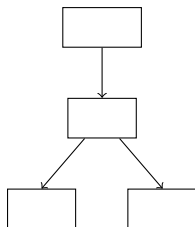2. Generalizing EDG to sequences of *lists* of strings

   Substitution cost for lists $A, B$:

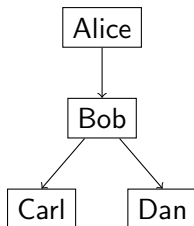   $$\mathcal{C}(A, B) := (\text{AED error if we merge } A, B) - (\text{AED error if we don't})$$
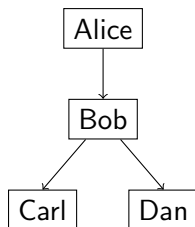
# Outline

1. Run branching process (Watson & Galton 1875)

# Generating synthetic data

1. Run branching process (Watson & Galton 1875)
2. Label with random strings

1. Run branching process (Watson & Galton 1875)
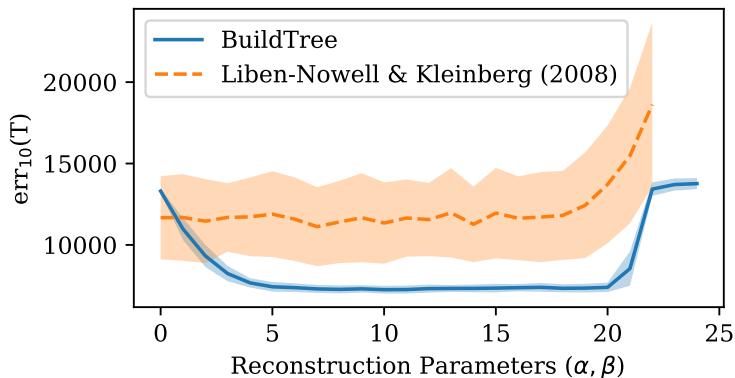2. Label with random strings
3. Simulate noisy propagation down the tree

15 sequences, 500 trials

100 sequences, 8 trials

# Approximate comparison with true tree

15 sequences, 500 trials

# Outline

# Takeaways and open questions

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem

2. NP-hard in general, but dynamic programming solution for two sequences and poly-time algorithm for $O(1)$ sequences

# Takeaways and open questions

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem

2. NP-hard in general, but dynamic programming solution for two sequences and poly-time algorithm for $O(1)$ sequences

3. Efficient heuristic for more sequences

# Takeaways and open questions

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem
2. NP-hard in general, but dynamic programming solution for two sequences and poly-time algorithm for $O(1)$ sequences
3. Efficient heuristic for more sequences

## Open questions

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem
2. NP-hard in general, but dynamic programming solution for two sequences and poly-time algorithm for $O(1)$ sequences
3. Efficient heuristic for more sequences

## Open questions

1. Approximation algorithm: bounding topological error seems hard

# Takeaways and open questions

## Takeaways

1. Chain letter petitions exhibit one-ended growth, divergence, and mutation: intriguing reconstruction problem
2. NP-hard in general, but dynamic programming solution for two sequences and poly-time algorithm for $O(1)$ sequences
3. Efficient heuristic for more sequences

## Open questions

1. Approximation algorithm: bounding topological error seems hard
2. Efficient algorithms for small $\lambda$

# Acknowledgment

Thanks to:

- Jon Kleinberg
- Anna Johnson
- Hailey Jones
- Dave Musicant

- Layla Oesper
- Anna Rafferty
- Ethan Somes