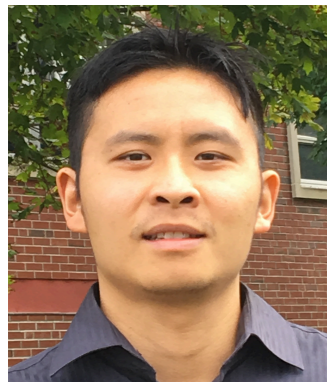# Optimal Transport-based Alignment of Learned Character Representations for String Similarity
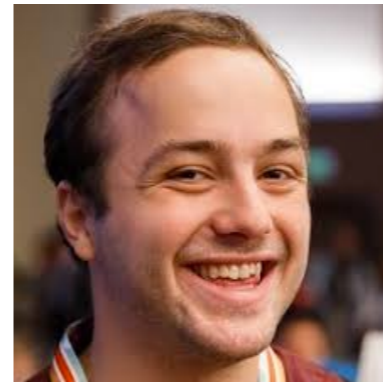
Derek Tam

Nicholas Monath

Ari Kobren
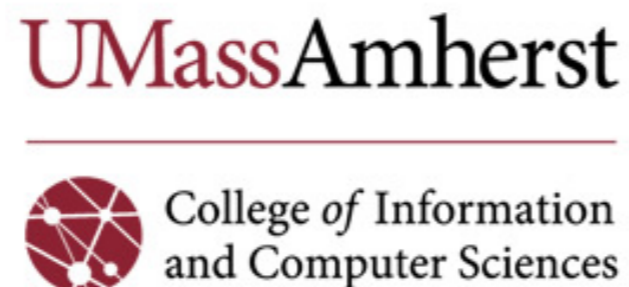
Aaron Traylor

Rajarshi Das

Andrew McCallum

# Record Linkage

US Patent Assignee Records

| Company Name | Location | Patent Title |
|---|---|---|
| Ethicon Surgery, Inc. | Somerville, NJ, US | Surgical Stapler Safety and Sequencing Mechanisms |
| Ethicon Endo Surgery | Somerville, NJ, US | Pneumatically Actuated Surgical Stapler Head |

# Coreference and Entity Linking

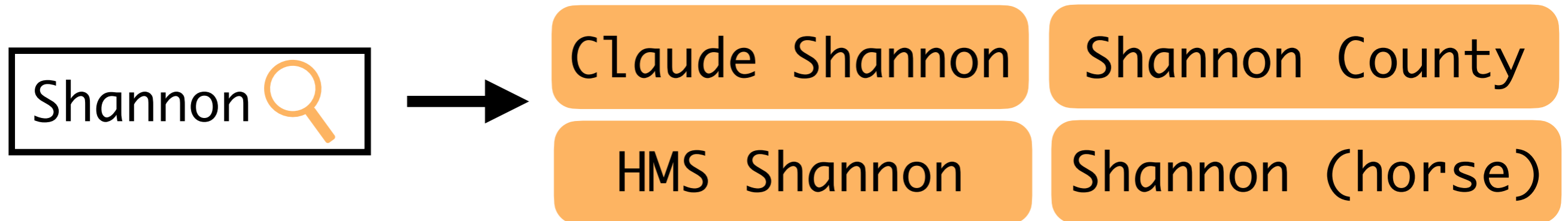Excited for these Grammys! Just a weird opening with **Tay Sway.**

**T-Swift** opens the #Grammys

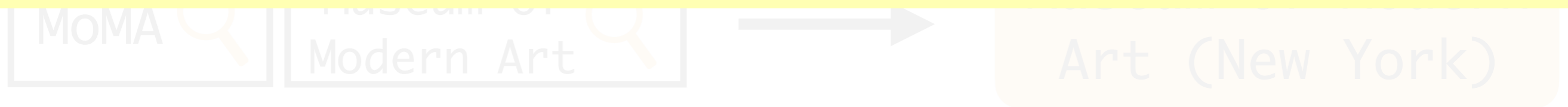Always get goosebumps before the #Grammys!!! **Taylor Swift** is on!

# Search

MoMA 🔍  Museum of Modern Art 🔍 → Museum of Modern Art (New York)

# Disambiguation

Shannon 🔍 →

Claude Shannon    Shannon County

HMS Shannon    Shannon (horse)

# Record Linkage

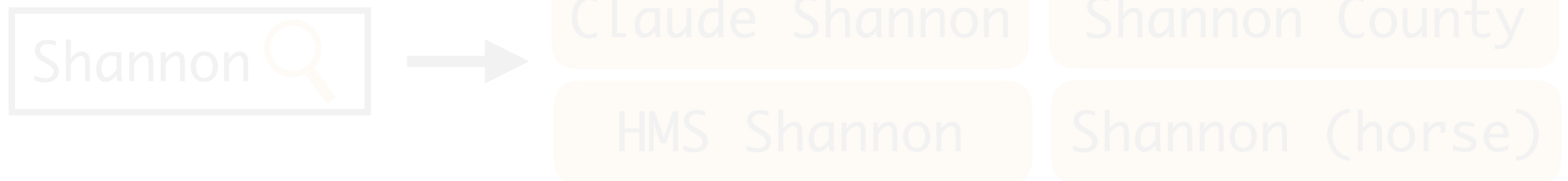| | Company Name | Location | Patent Title |
|---|---|---|---|
| US Patent Assignee Records | Ethicon Surgery, Inc. | Somerville, NJ, US | Surgical Stapler Safety and Sequencing Mechanisms |
| | Ethicon Endo Surgery | Somerville, NJ, US | Pneumatically Actuated Surgical Stapler Head |

# Coreference and Entity Linking

Excited for these
Grammys! Just a weird

T-Swift opens the

Always get goosebumps
before the #Grammys!

**Similarity of mention strings informs whether or not they refer to the same entity.**

MoMA 🔍   Museum of Modern Art 🔍   →   Art (New York)

# Disambiguation

Shannon 🔍   →   Claude Shannon   Shannon County

HMS Shannon   Shannon (horse)

# String Similarity for Entity Aliases

**Which strings can refer to the same entity?**
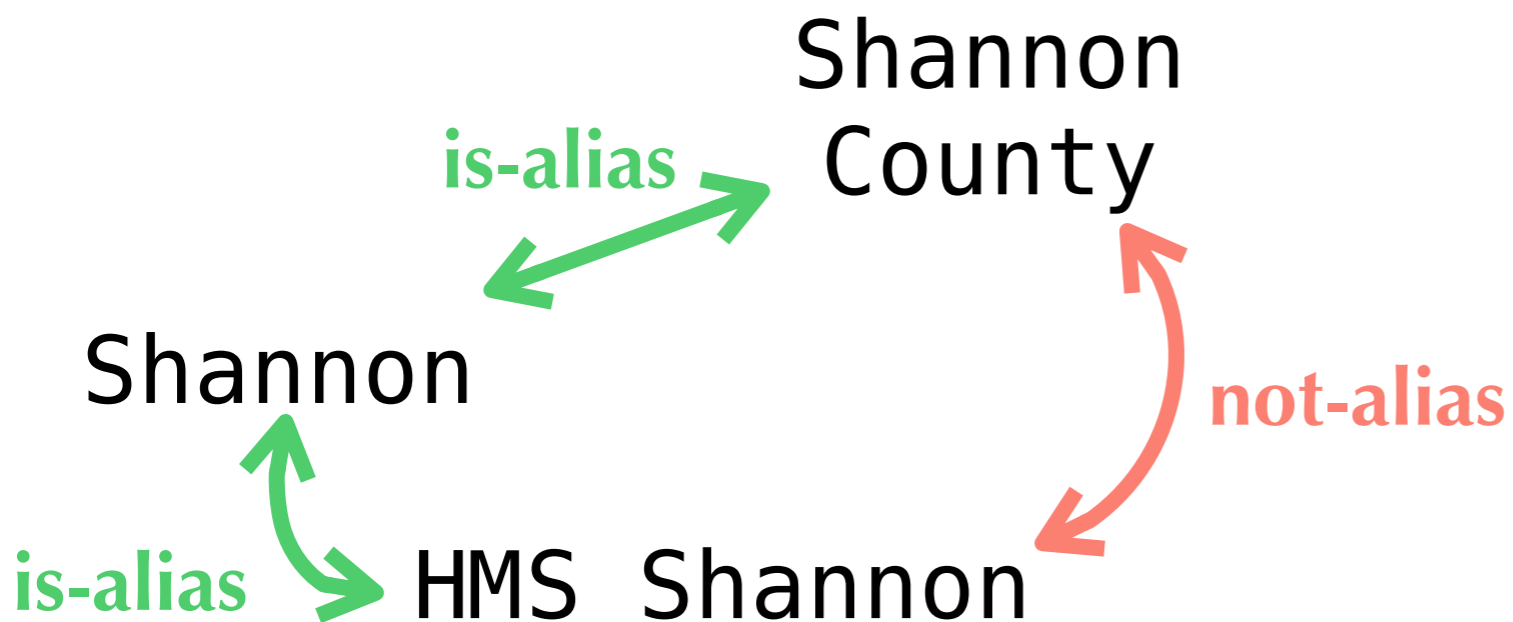
**Design similarity function *f***

$f(\texttt{string1}, \texttt{string2})$    <span style="color:green">**high similarity**</span>    if ***can refer*** to the same entity

$f(\texttt{string1}, \texttt{string2})$    <span style="color:salmon">**low similarity**</span>    if ***cannot refer*** to the same entity

Shannon County

**is-alias**

Shannon

**not-alias**

**is-alias**    HMS Shannon

**Designed to inform coreference decisions**

# Classic Approaches

Similarity determined by number and type of edits

Music in Chile
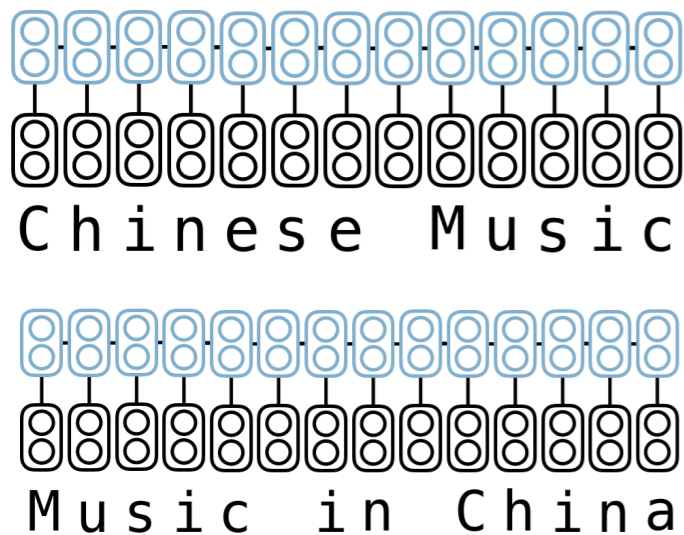
**# Edits = 2**

Music in China

**# Edits = 12**

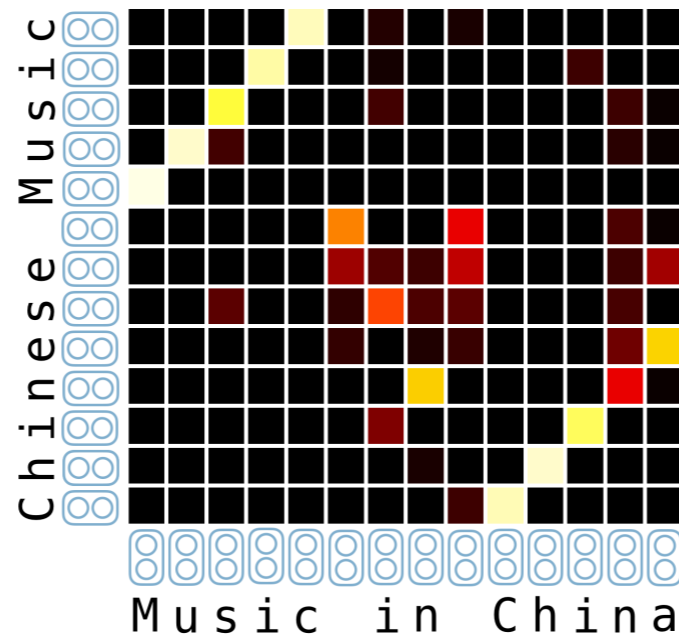Chinese Music

***Character edits alone insufficient!***

# STANCE

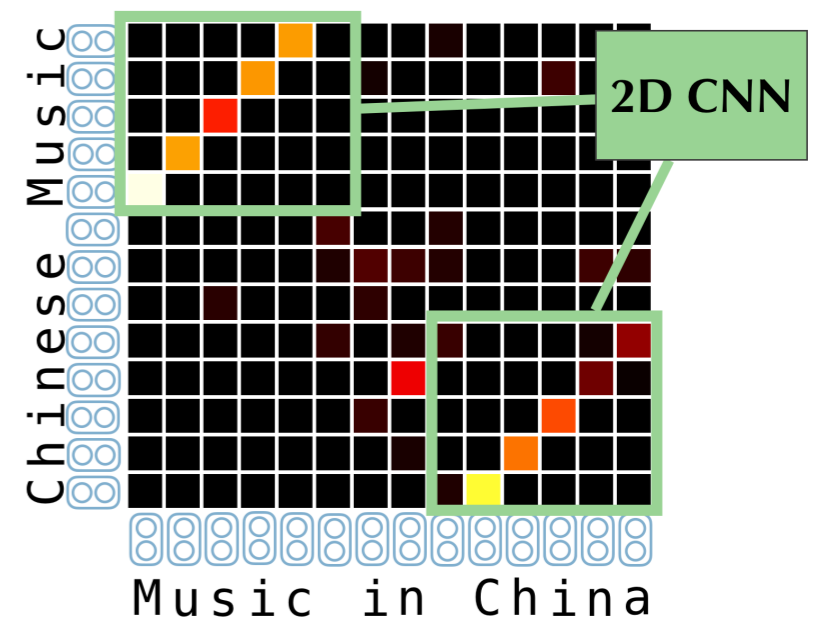## Similarity of Transport Aligned Neural Character Encodings



**Character Embeddings**

Chinese Music

Music in China

**Optimal Transport based Alignment**

Chinese Music

Music in China

**CNN Scoring Function**

Chinese Music

2D CNN

Music in China

# STANCE

## Similarity of Transport Aligned Neural Character Encodings



**Character Embeddings**

Chinese Music

Music in China

Optimal Transport based Alignment

Chinese Music

Music in China

CNN Scoring Function

Chinese Music

2D CNN

Music in China

# Character Representations

Encode with RNN, Measure Pairwise Similarities

# Character Representations

Encode with RNN, Measure Pairwise Similarities



High Similarity

Low Similarity

**Repeated characters may suffer from spurious high similarities**
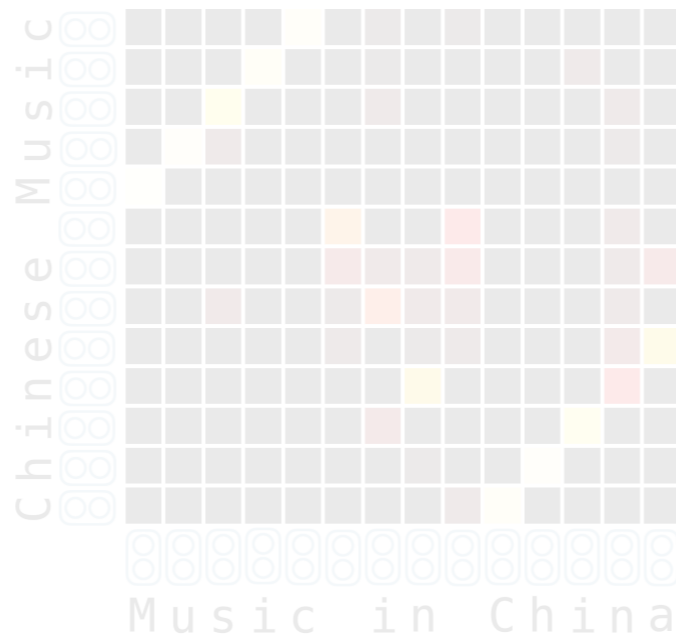
$h_i = RNN(x_i, h_{i-1})$

H

X

Music in China

Chinese Music
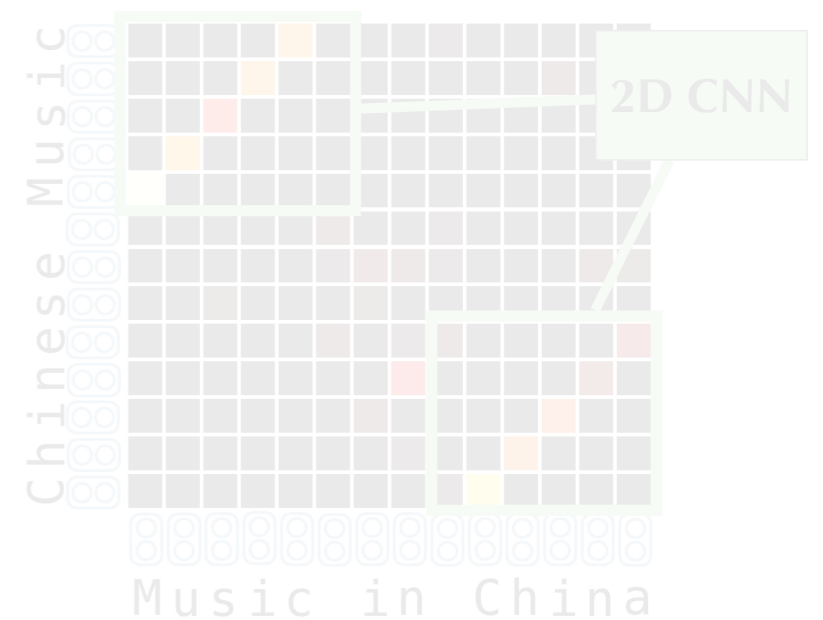
# STANCE

## Similarity of Transport Aligned Neural Character Encodings
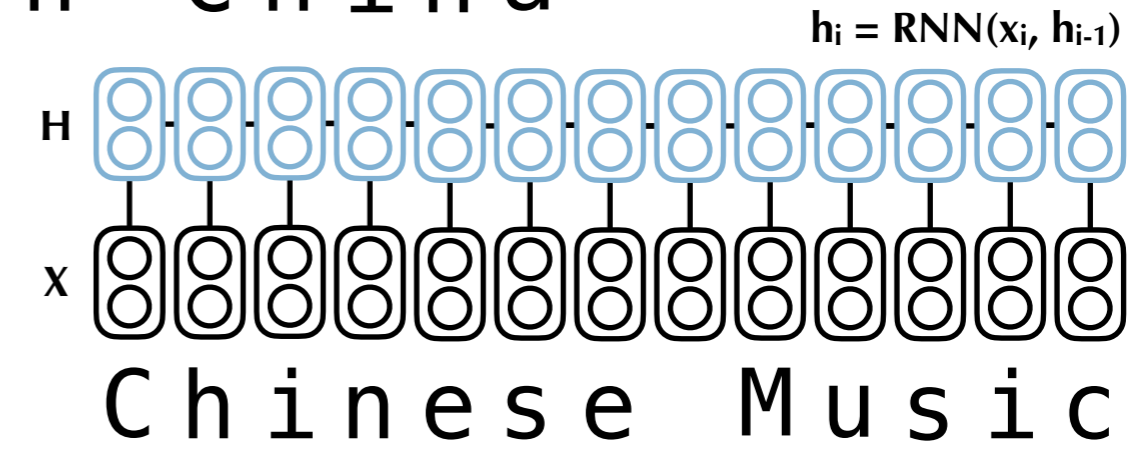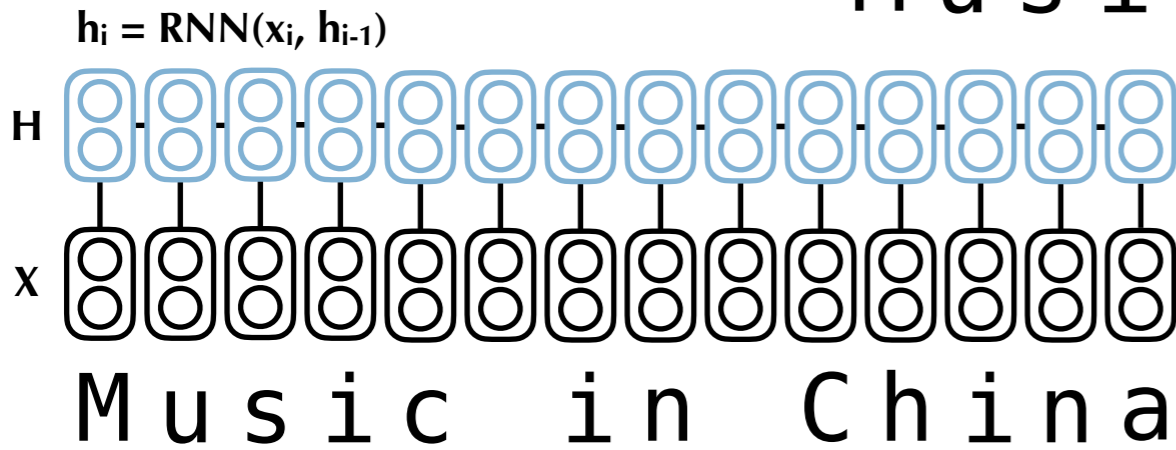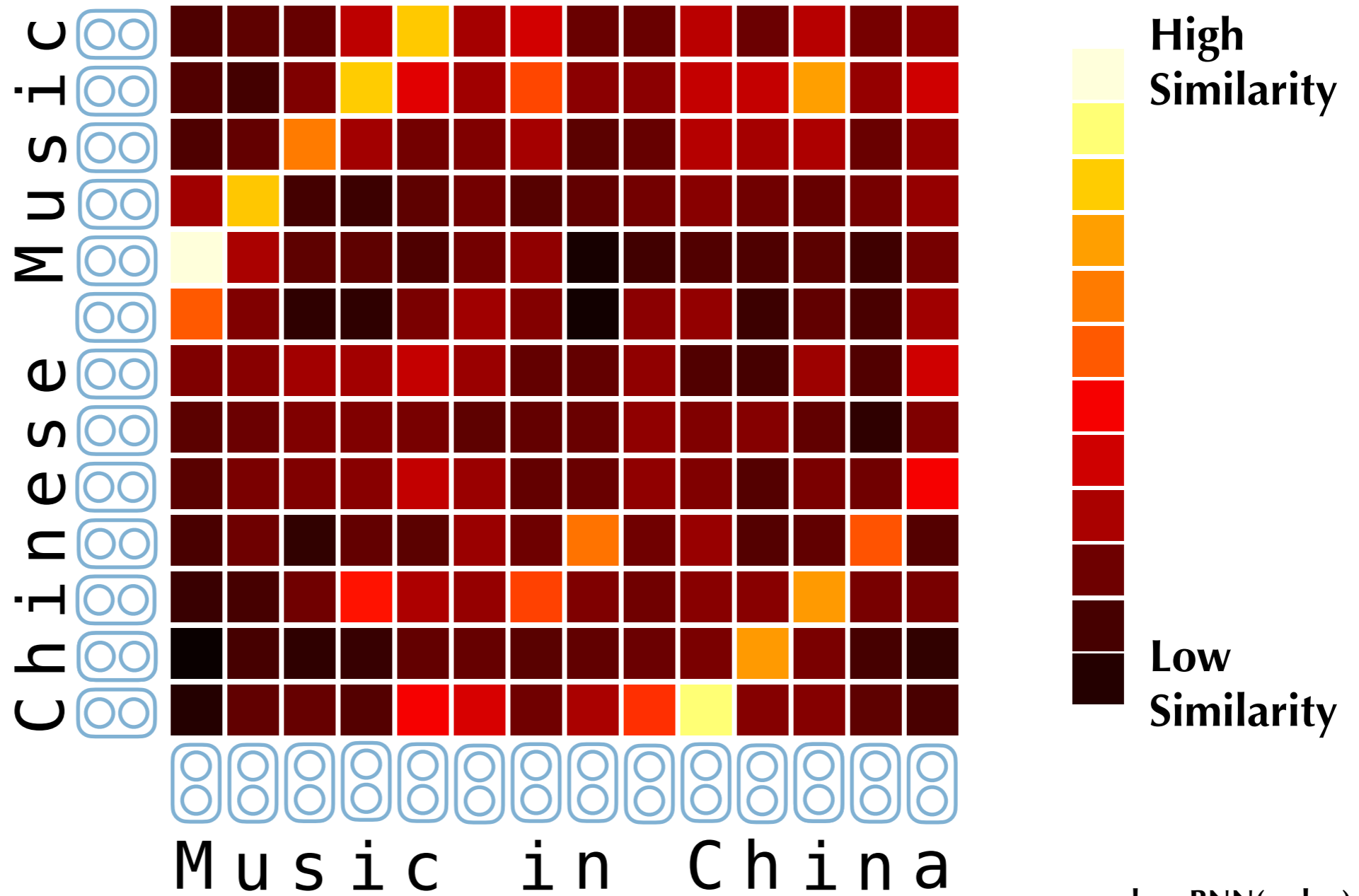


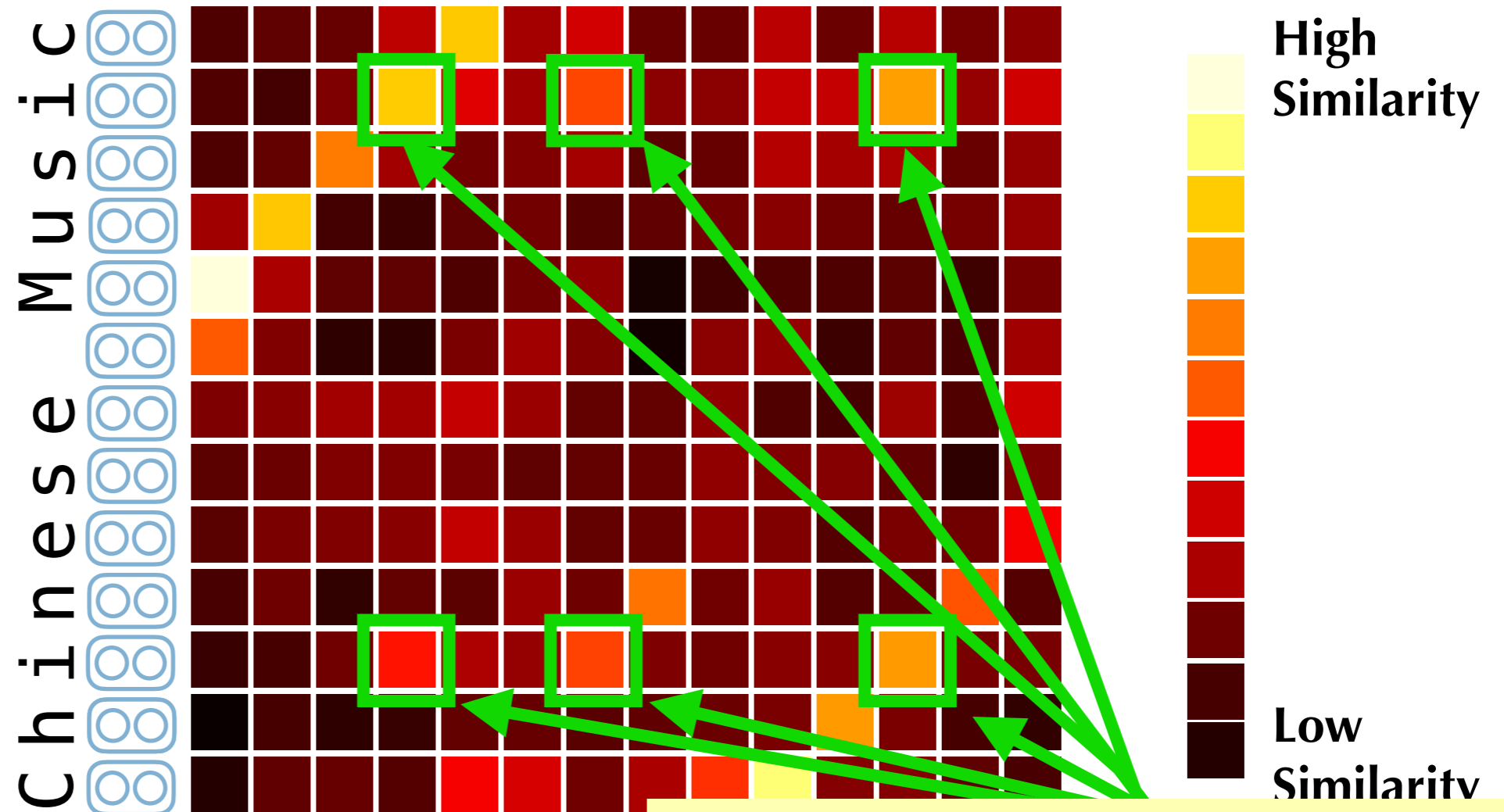Character Embeddings

Optimal Transport based Alignment

CNN Scoring Function

# Optimal Transport-based Alignment

Similarity Matrix



High Similarity

Low Similarity

**Repeated characters may suffer from spurious high similarities**

Chinese Music

Music in China

# Optimal Transport-based Alignment

**Each character aligned to closest character(s) in other string**



Similarity Matrix

High Similarity

Low Similarity

**Keep "good" alignments**

**Remove "bad" alignments**

Chinese Music

Music in China

# Alignment as Optimal Transport

Chinese Music

Music in China

**The amount of transported mass indicates degree of alignment.**

# Alignment as Optimal Transport



**The amount of transported mass indicates degree of alignment.**

# Alignment as Optimal Transport



**The amount of transported mass indicates degree of alignment.**

# Alignment as Optimal Transport

The amount of transported mass indicates degree of alignment.



Similarity Matrix

High Similarity

Low Similarity

Cost of transport inversely proportional to similarity.

# Alignment as Optimal Transport

The amount of transported mass indicates degree of alignment.



`C h i n e s e   M u s i c`

`M u s i c   i n   C h i n a`

Cost of transport inversely proportional to similarity.
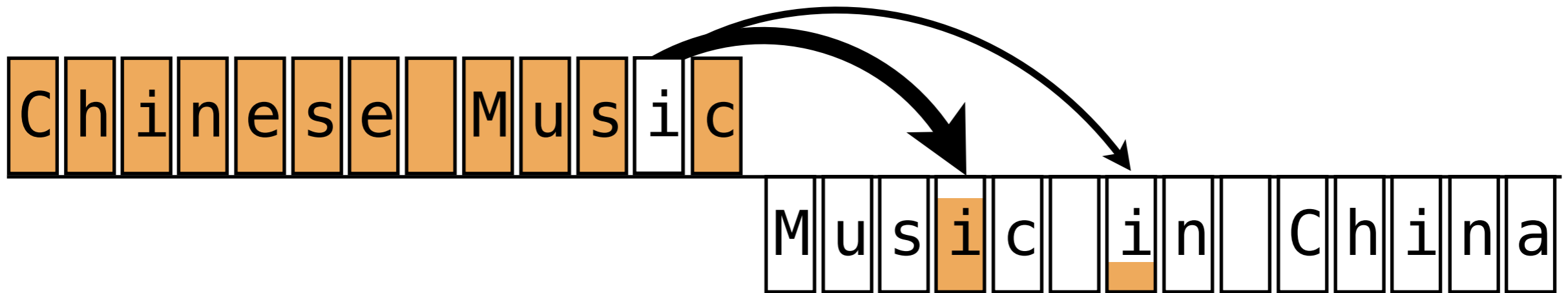
To transport:

$$\text{mass}(M) = 1/\text{StringLength}$$

To receive:

$$\text{mass}(M) = 1/\text{StringLength}$$

**Characters have fixed amount of mass to transport (or receive).**

**All characters must transport (or receive) entire mass.**

# Alignment as Optimal Transport

The amount of transported mass indicates degree of alignment.

C h i n e s e   M u s i c

M u s i c   i n   C h i n a

Cost of transport inversely proportional to similarity.

**AlignmentCost(S1,S2) =**

$$\sum_{i \in S1} \sum_{j \in S2} T_{i \to j} \, Cost(i, j)$$

Characters in S1
(e.g., "Chinese Music")

Characters in S2
(e.g., "Music in China")

How much of i is transported to j

Inversely proportional to similarity of i & j

# Alignment as Optimal Transport

The amount of transported mass indicates degree of alignment.



Chinese Music

Music in China
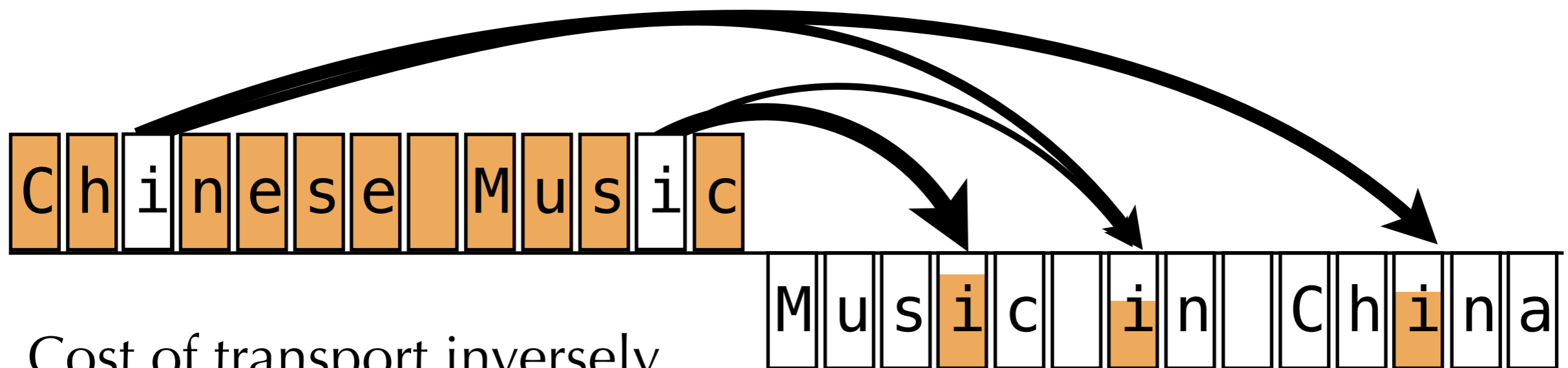
**Find minimum cost alignment between characters of the two strings**

$$\sum_{i \in S1} \sum_{j \in S2} T_{i \to j} \; Cost(i, j)$$

Characters in S1
(e.g., "Chinese Music")

Characters in S2
(e.g., "Music in China")

How much of i is
transported to j

Inversely proportional
to similarity of i & j

# Optimal Transport-based Alignment

Minimum cost soft alignment btw characters of the two strings



$$\sum \blacksquare = \text{mass}(\boxed{\texttt{i}})$$

Transport Matrix (**T**)

**How *aligned* the i in Chinese is to the i in China**

More Aligned

Less Aligned

**Sparsity in transport removes spurious matches**

$$\sum \blacksquare = \text{mass}(\boxed{\texttt{i}})$$

**Solved with Sinkhorn's Algorithm: Efficient and differentiable**

# Optimal Transport-based Alignment

Re-weight similarity by transport matrix

# STANCE

## Similarity of Transport Aligned Neural Character Encodings

**Character Embeddings**

**Optimal Transport based Alignment**

**CNN Scoring Function**



2D CNN

# CNN Scoring Function

Capture patterns of sequential alignment between characters.



**Similarity x Transport**

Chinese Music

Music in China

2D Convolution    Pooling    Output

**Entire Model Trained with Ranking-based Objective**

# Experimental Results

## Task 1: Alias Detection

## Task 2: Cross Document Coreference

## Qualitative Analysis & Ablation Study

# Alias Detection

**Aliases - Two strings that *can* refer to the same entity**

**Given a *query* string, *rank* candidate aliases.**

*Query*

Peace Agreement

*Candidates*

Peace Treaty

Peace Pact

Lease Agreement

Peacekeeping Troops

*Ranking*

Peace Treaty

Peace Pact

Peacekeeping Troops

Lease Agreement

# Datasets

**Built 5 datasets for alias detection from open KBs**

**Wikipedia**

Irish music **is-alias** Irish Folk

# Datasets

**Built 5 datasets for alias detection from open KBs**

**Wikipedia**  **Wikipedia-People**

Queen Elizabeth II **is-alias**
Queen Elizabeth the Second

# Datasets

**Built 5 datasets for alias detection from open KBs**

**Wikipedia**   **Wikipedia-People**   **Music Artist**

Red Hot Chili Peppers **is-alias** RHCP

# Datasets

**Built 5 datasets for alias detection from open KBs**

**Wikipedia** **Wikipedia-People** **Music Artist**

**Patent Assignee**

The Proctor & Gamble Company **is-alias** Proctor and Gamble

# Datasets

**Built 5 datasets for alias detection from open KBs**

**Wikipedia** **Wikipedia-People** **Music Artist**

**Patent Assignee** **Disease**

black water fever **is-alias** hemolytic malaria

# Alias Detection Experiments

Compare STANCE to 8 baseline methods including:

## Alignment Methods

- Levenshtein Similarity

- Learned Dynamic Time Warping - LDTW (Cuturi et al. 2017)

## Neural Methods

- Deep Conflation Model - DCM (Gan et al. 2017)

# Alias Detection Experiments

Compare STANCE to 8 baseline methods including:

## Alignment Methods

- Levenshtein Similarity

- Learned Dynamic Time Warping - LDTW (Cuturi et al. 2017)

## Neural Methods

- Deep Conflation Model - DCM (Gan et al. 2017)

# Alias Detection Experiments

Compare STANCE to 8 baseline methods including:

## Alignment Methods

- Levenshtein Similarity

- Learned Dynamic Time Warping - LDTW (Cuturi et al. 2017)

## Neural Methods

- Deep Conflation Model - DCM (Gan et al. 2017)

# Alias Detection Experiments

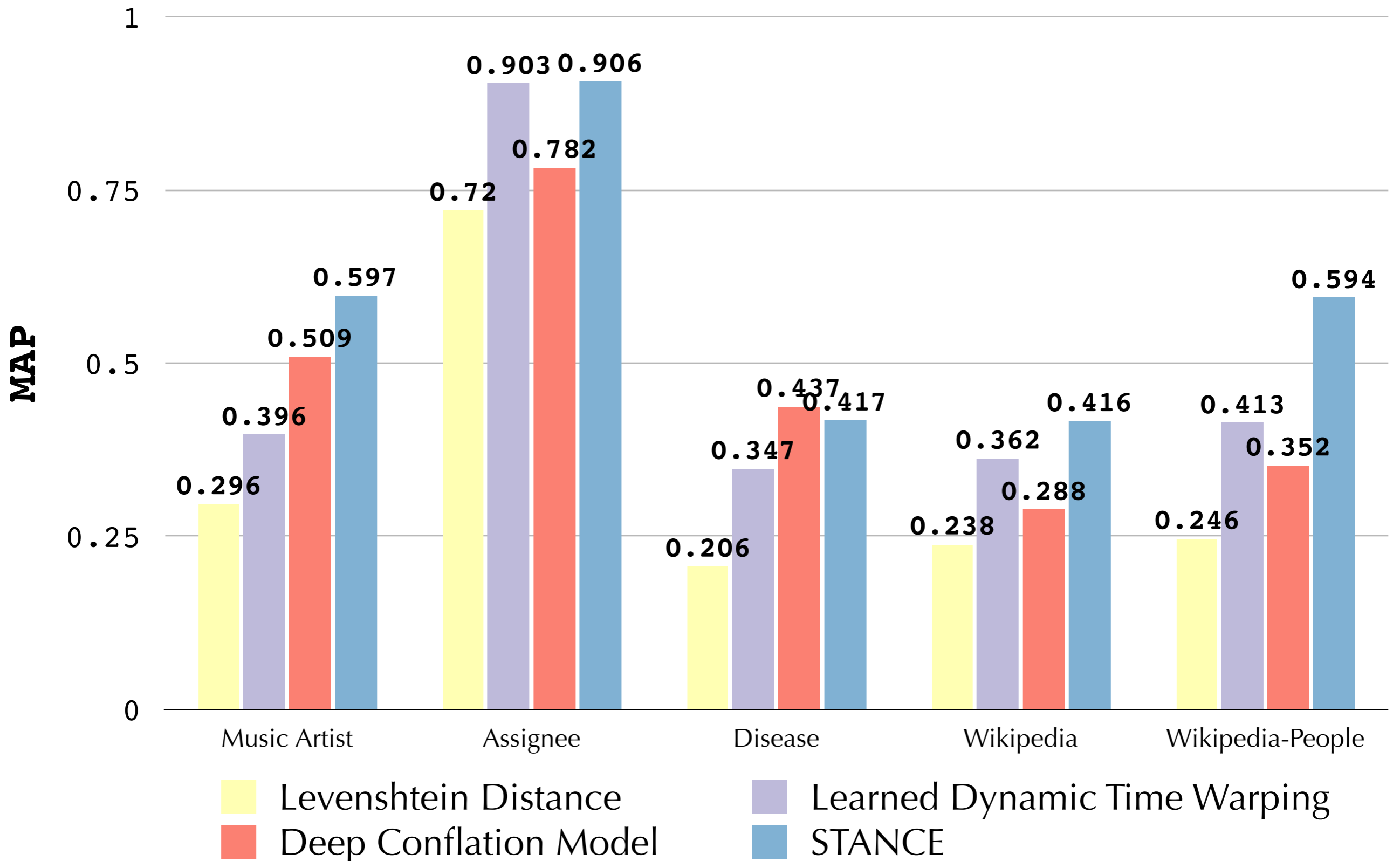Compare STANCE to 8 baseline methods including:

## Alignment Methods

- Levenshtein Similarity

- Learned Dynamic Time Warping - LDTW (Cuturi et al. 2017)

## Neural Methods

- Deep Conflation Model - DCM (Gan et al. 2017)

# Alias Detection - Mean Average Precision (MAP)

**MAP**

| | Music Artist | Assignee | Disease | Wikipedia | Wikipedia-People |
|---|---|---|---|---|---|
| Levenshtein Distance | 0.296 | 0.72 | 0.206 | 0.238 | 0.246 |
| Learned Dynamic Time Warping | 0.396 | 0.903 | 0.347 | 0.362 | 0.413 |
| Deep Conflation Model | 0.509 | 0.782 | 0.437 | 0.288 | 0.352 |
| STANCE | 0.597 | 0.906 | 0.417 | 0.416 | 0.594 |

Legend: Levenshtein Distance, Learned Dynamic Time Warping, Deep Conflation Model, STANCE

# Experimental Results

Task 1: Alias Detection

## Task 2: Cross Document Coreference

Qualitative Analysis & Ablation Study

# Cross-Document Coreference

## Twitter at the Grammy's Dataset (Dredze et al, 2016)

### 4577 Mentions, 273 Entities

Excited for these Grammys! Just a weird opening with **Tay Sway**.

**T-Swift** opens the #Grammys

Always get goosebumps before the #Grammys!!! **Taylor Swift** is on!

**Taylor,** what happened, this is madness. #grammys

# Cross-Document Coreference

## Twitter at the Grammy's Dataset (Dredze et al, 2016)
## 4577 Mentions, 273 Entities

Excited for these Grammys! Just a weird opening with **Tay Sway**.

**T-Swift** opens the #Grammys

Always get goosebumps before the #Grammys!!! **Taylor Swift** is on!

**Taylor,** what happened, this is madness. #grammys

**LL Cool J** has swag for days. No better person to host the #Grammys!

**El-El Cool John.** #Grammy

**LL Cool James** just mispronounced @edsheeran's name AGAIN at the #Grammys!

# Cross-Document Coreference

**Our approach**

Average-Linkage
**Hierarchical
Agglomerative
Clustering**



Use **pre-trained STANCE** model on
Wikipedia-People as **pairwise similarity function**.

Tune **threshold** to cut tree for
**predicting entities** on dev set

# Cross-Document Coreference Performance



**B³ F1**

- 77.2 — Green et al (2012) (Spelling Only)
- 72.3 — Andrews et al (2014) (Spelling Only)
- 79.7 — Green et al (2012) (w Context)
- 72.1 — Andrews et al (2014) (with Context)
- 72.3 — Andrews et al (2014) (with Context & Time)
- 82.5 — STANCE

Baseline Results from
Dredze et al (2016)

# Cross-Document Coreference

## Twitter at the Grammy's Dataset (Dredze et al, 2016)

Excited for these Grammys! Just a weird opening with **Tay Sway.**

**T-Swift** opens the #Grammys

Always get goosebumps before the #Grammys!!! **Taylor Swift** is on!

**Taylor,** what happened, this is madness. #grammys

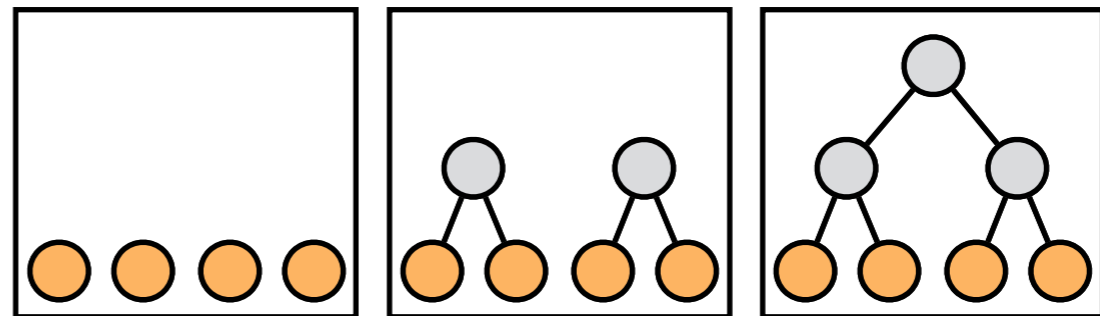**LL Cool J** has swag for days. No better person to host the #Grammys!

**El-El Cool John.** #Grammy

**LL Cool James** just mispronounced @edsheeran's name AGAIN at the #Grammys!

**Name variation more informative than context**

# Experimental Results

**Task 1: Alias Detection**

**Task 2: Cross Document Coreference**

**Qualitative Analysis & Ablation Study**

# Qualitative Analysis

**Query:** Boom Microphones

**Nearest Neighbors:**

**STANCE**

Boom mike

Boom mics

**LDTW**

Open Microphone

Shotgun Microphone

**DCM**

Open Microphone

Condensor Microphone

# Qualitative Analysis
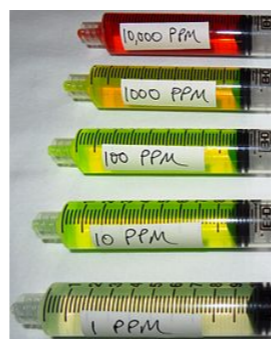
Query:
RPM

Nearest Neighbors:

**STANCE**

RPM Weekly
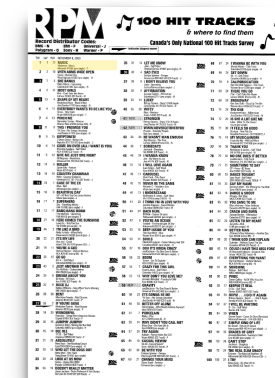


Randle Patrick McMurphy



**LDTW**

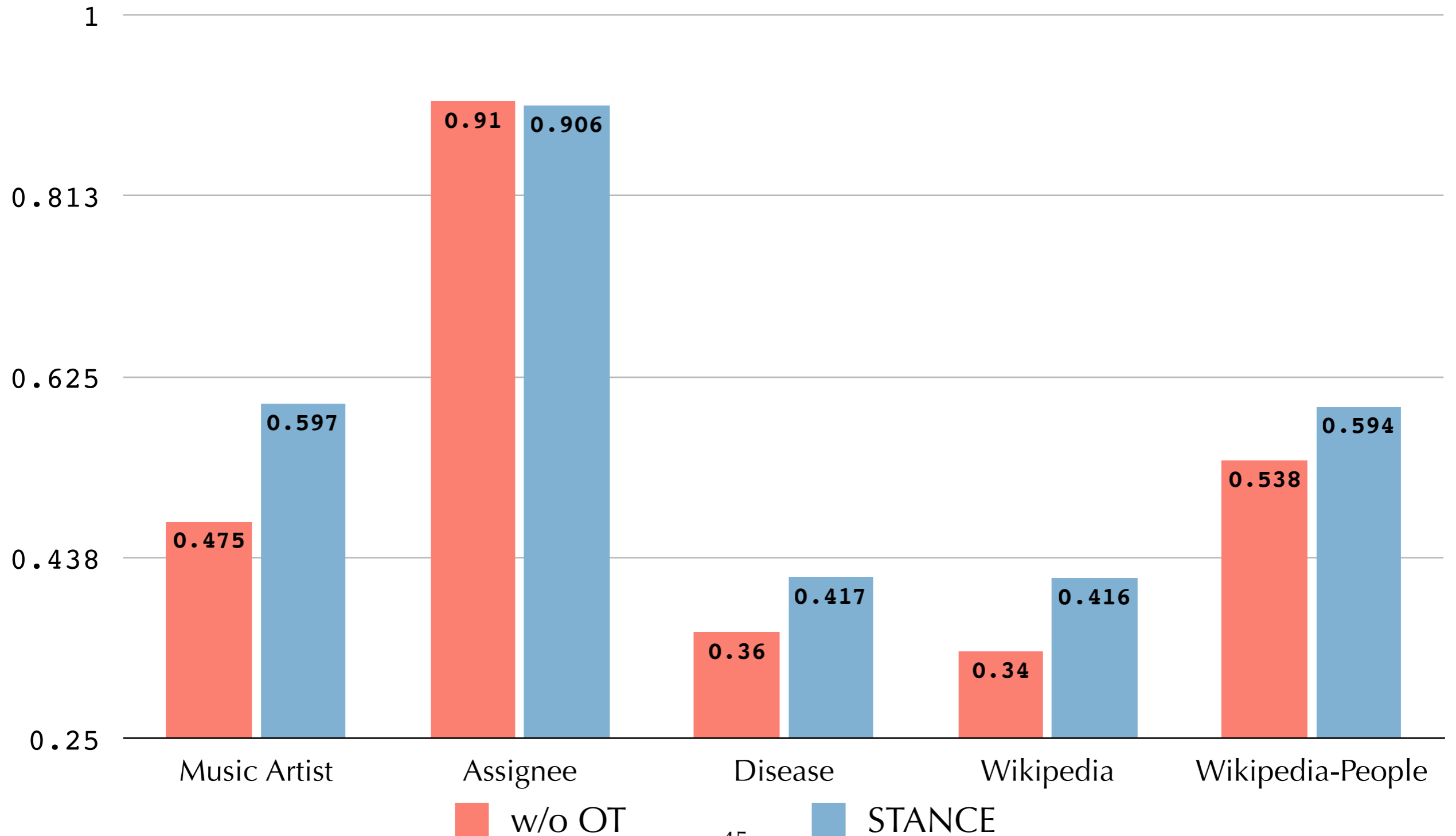PPM



RPM Alternative 30



**DCM**

RP1



PRM

# Impact of Optimal Transport in STANCE

## OT component improves results on 4 of 5 datasets.

# Benefit of OT - Noise Reduction

**Query:** Saath Saath Banayenger Ek Aashi

**Non-Alias Candidate**: Teen Bahuraaniyaan

**Significant number of repeated characters and character bigrams**

## Similarity Matrix - w/o OT

# Benefit of OT - Noise Reduction

**Query:** Saath Saath Banayenger Ek Aashi

**Non-Alias Candidate**: Teen Bahuraaniyaan

**Significant number of repeated characters and character bigrams**
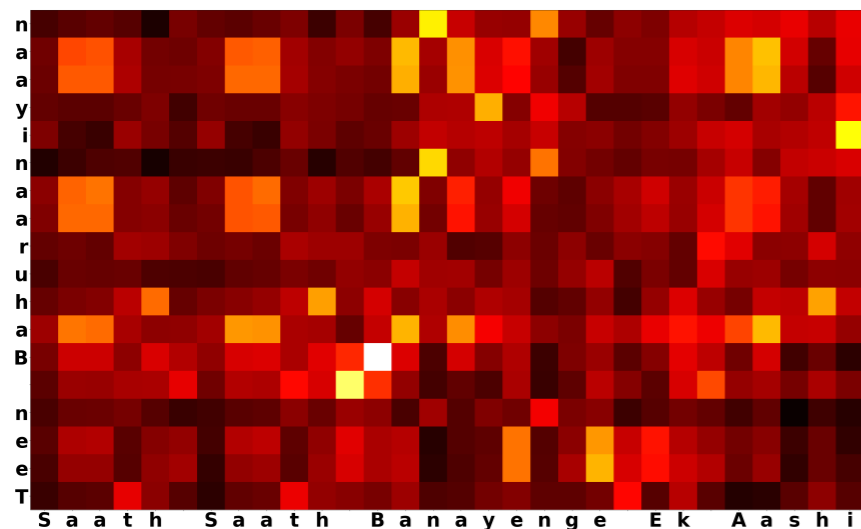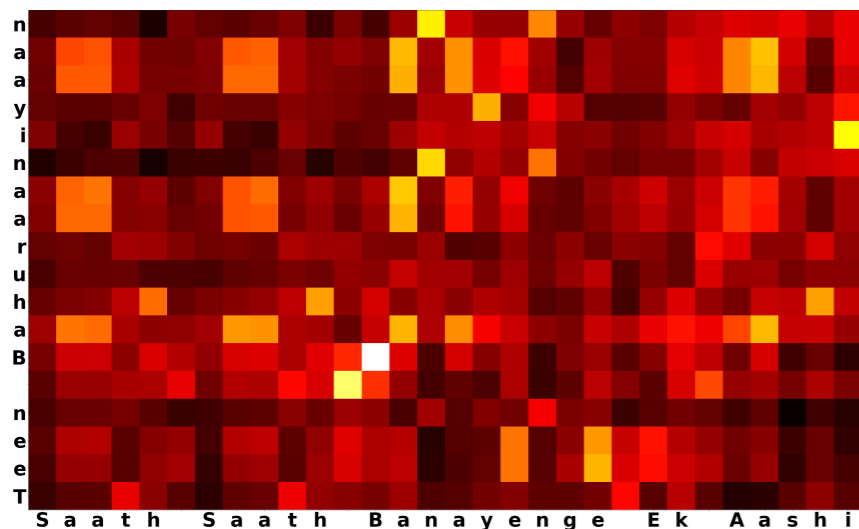
**Similarity Matrix - w/o OT**

# Benefit of OT - Noise Reduction

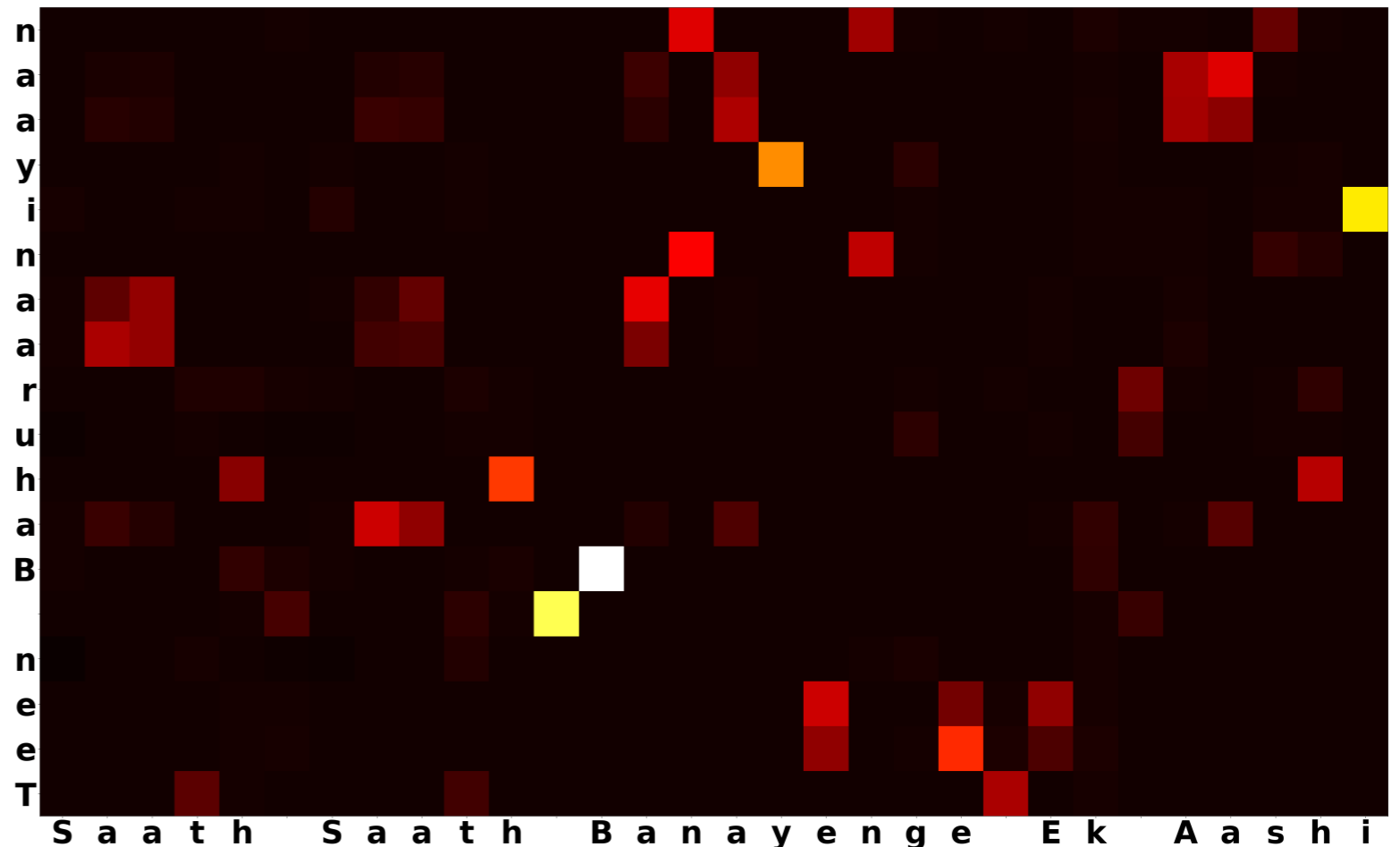**Query:** Saath Saath Banayenger Ek Aashi

**Non-Alias Candidate:** Teen Bahuraaniyaan

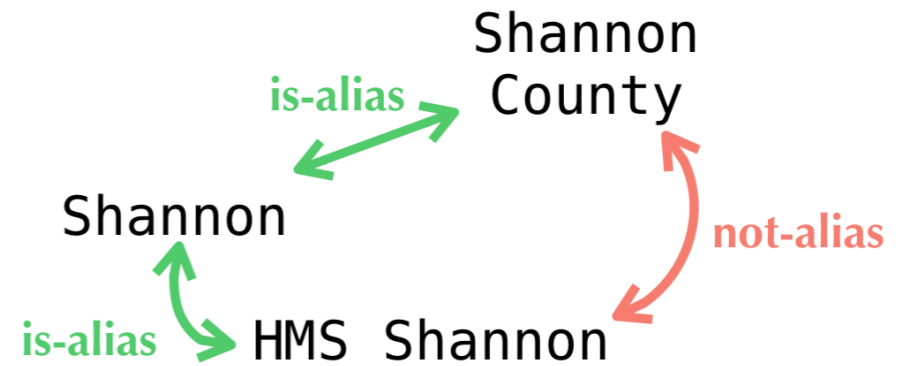## Significant number of repeated characters and character bigrams

### Similarity Matrix - w/o OT

### Similarity Matrix - STANCE

# Summary

**Learned String Similarity**



Shannon County

is-alias

Shannon

not-alias

is-alias

HMS Shannon

## STANCE

**S**imilarity of **T**ransport **A**ligned **N**eural **C**haracter **E**ncodings



**Character Embeddings**

Chinese Music

Music in China

**Optimal Transport based Alignment**

Chinese Music

Music in China

**CNN Scoring Function**

Chinese Music

2D CNN

Music in China

**New Datasets and Results**



MAP

1

0.75

0.5

0.25

0

Music Artist    Assignee    Disease    Wikipedia    Wikipedia-People

# Thanks! Questions?

## Code: https://github.com/iesl/stance

### Similarity Matrix

**High Similarity**

**Low Similarity**

Chinese Music

Music in China

Chinese Music

Music in China

**Green et al (2012) (Spelling Only)**
**Andrews et al (2014) (Spelling Only)**
**Green et al (2012) (w Context)**
**Andrews et al (2014) (with Context)**
**Andrews et al (2014) (with Context & Time)**
**STANCE**

77.2  72.3  79.7  72.1  72.3  82.5

*Query*

**Peace Agreement**

*Candidates*

**Peace Treaty**    **Peace Pact**

Lease Agreement    Peacekeeping Troops

22

*Ranking*

**Peace Treaty**

**Peace Pact**

**Peacekeeping Troops**

**Lease Agreement**

Excited for these Grammys! Just a weird opening with **Tay Sway**.

**T-Swift** opens the #Grammys

Always get goosebumps before the #Grammys!!! **Taylor Swift** is on!

**Taylor,** what happened, this is madness. #grammys

**LL Cool J** has swag for days. No better person to host the #Grammys!

**El-El Cool John.** #Grammy

**LL Cool James** just mispronounced @edsheeran's name AGAIN at the #Grammys!

50