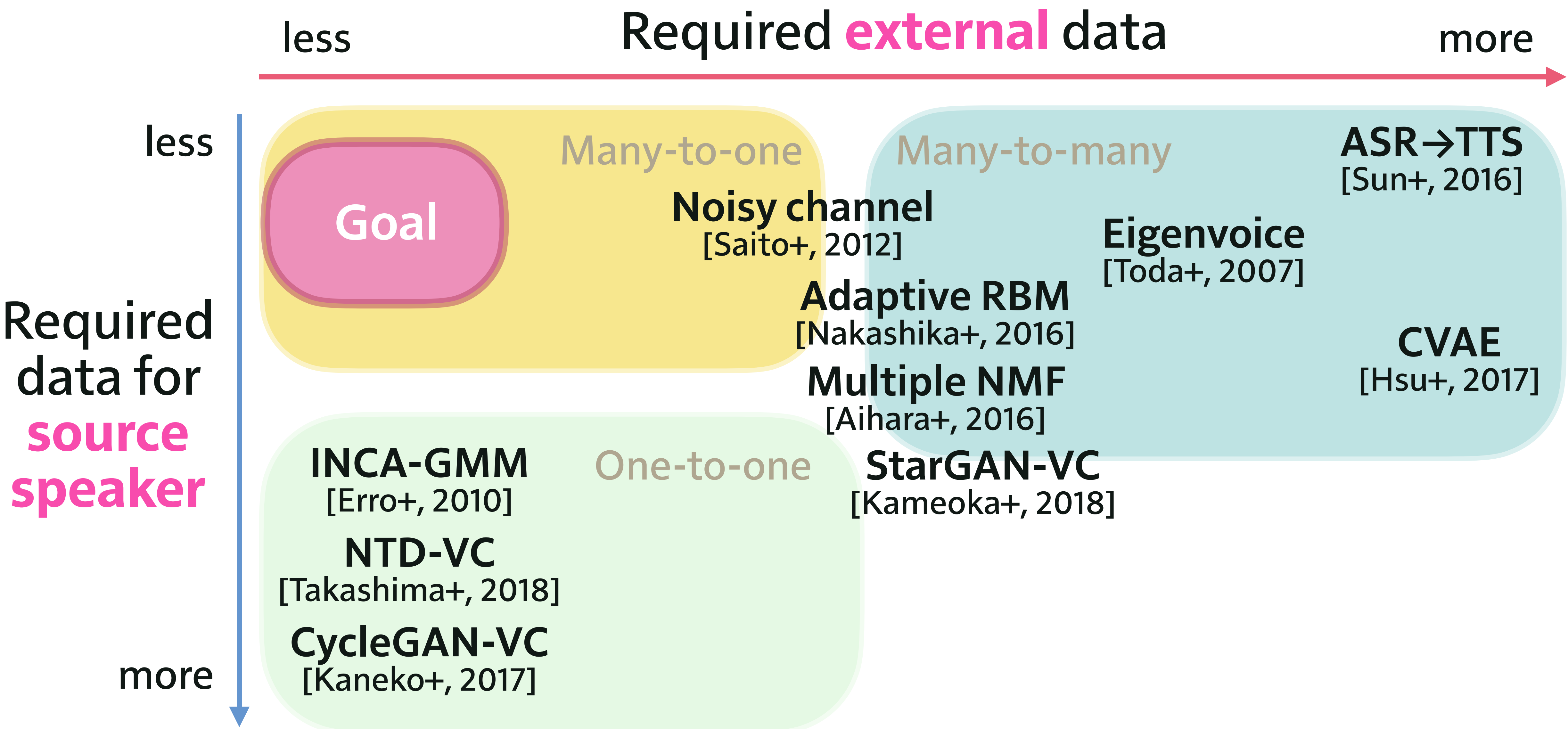# Nonparallel Training of Exemplar-based Voice Conversion System Using INCA-based Alignment Technique

**Hitoshi Suda**, Gaku Kotani, Daisuke Saito

The University of Tokyo

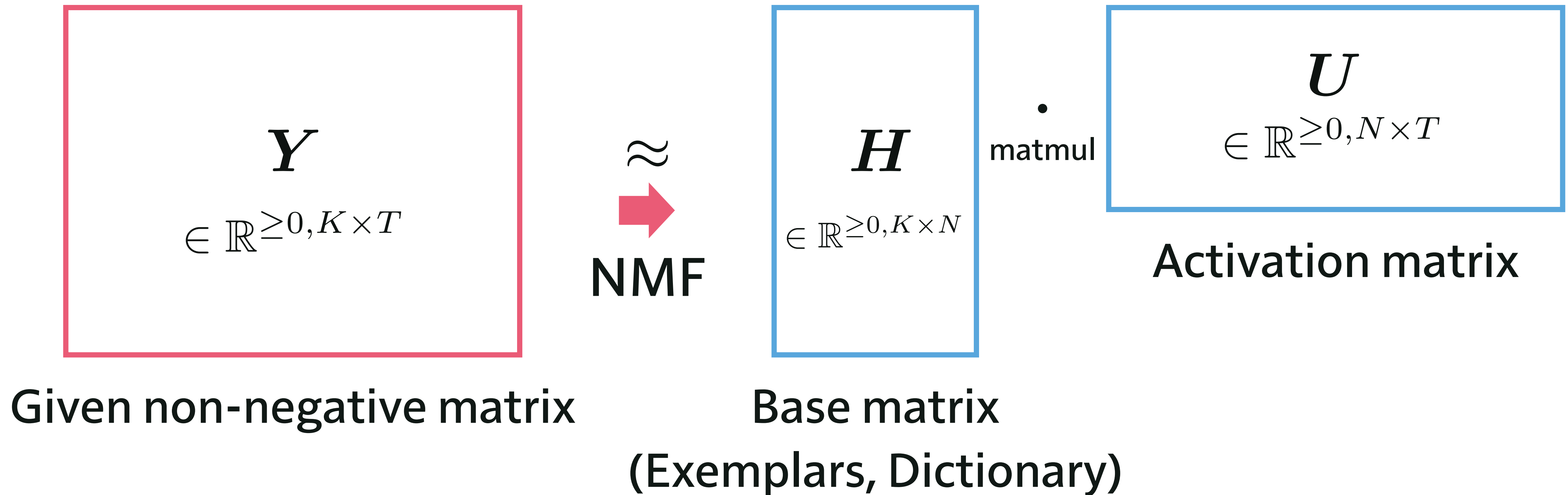INTERSPEECH 2020 Virtual Conference

# Outline

- **Baseline 1: NMF-based voice conversion**

- Baseline 2: INCA algorithm

- Proposed: Nonparallel training of NMF-based voice conversion

- Experiments

# Non-negative matrix factorization (NMF)

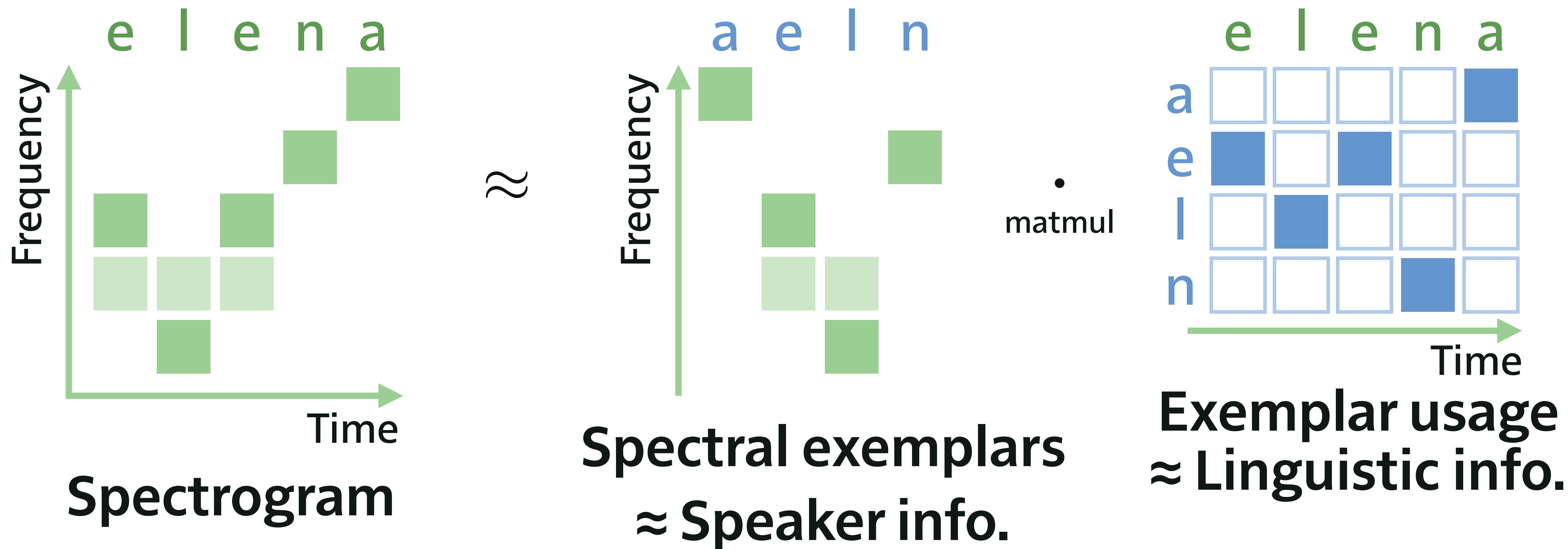[D. D. Lee+, 1999]

$$Y \in \mathbb{R}^{\geq 0, K \times T} \quad \approx \quad H \in \mathbb{R}^{\geq 0, K \times N} \quad \cdot \text{ matmul} \quad U \in \mathbb{R}^{\geq 0, N \times T}$$

NMF

Given non-negative matrix

Base matrix
(Exemplars, Dictionary)

Activation matrix

- NMF acquires $H$ and $U$
  by minimizing divergence $\mathcal{D}(Y|HU)$

**Spectrogram**

**Spectral exemplars ≈ Speaker info.**

**Exemplar usage ≈ Linguistic info.**

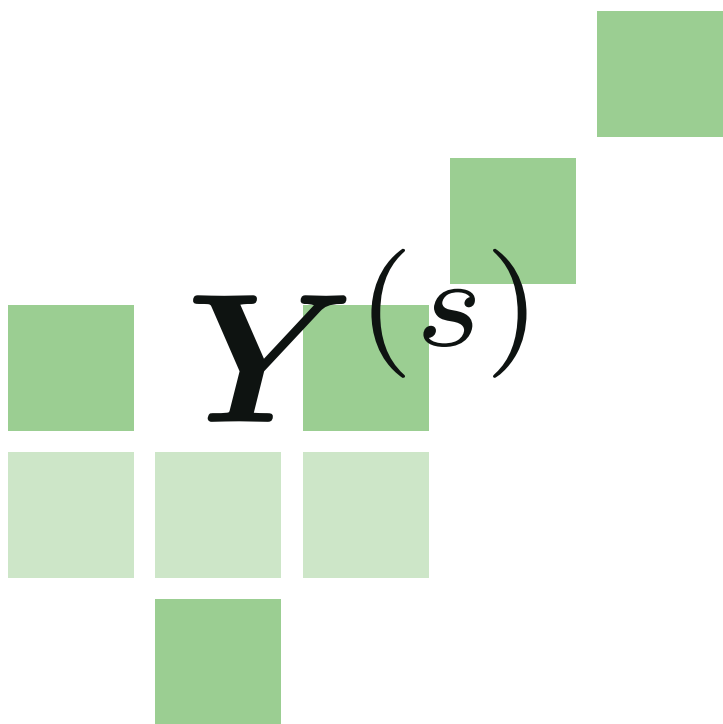- Supposing input matrix is spectrogram, **exemplars** contain **individuality** and **activity** contains **linguistic** information
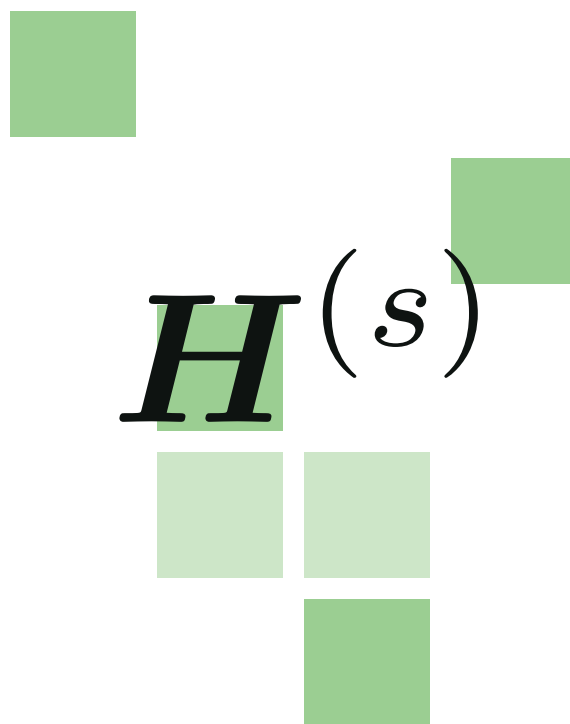
Training

[R. Takashima+, 2012]



**Source Speaker**

$$Y^{(s)} \approx H^{(s)} \cdot U$$

NMF

**Shared activity**

**Target Speaker**

$$Y^{(t)} \approx H^{(t)} \cdot U$$

NMF

**Aligned spectrogram**

**Parallel exemplars**

[R. Takashima+, 2012]

**Conversion**

**Source Speaker**

$$\boldsymbol{Y}^{(s)} \approx \boldsymbol{H}^{(s)} \cdot \boldsymbol{U}$$

NMF

**Input spectrogram**

**Trained exemplars**

$\boldsymbol{H}^{(t)} \cdot \boldsymbol{U}$

**Target Speaker**

**Estimated activity**

Copy

**Estimated spectrogram**

$= \quad \hat{\boldsymbol{Y}}^{(t)}$

matmul

# NMF-based Voice Conversion (NMF-VC)

[R. Takashima+, 2012]

- Uses **parallel exemplars as a conversion model**

  - NMF-VC is called "exemplar-based VC"

✓ **Fast on GPU**, because NMF uses only basic matrix operations

✗ Requires parallel corpora

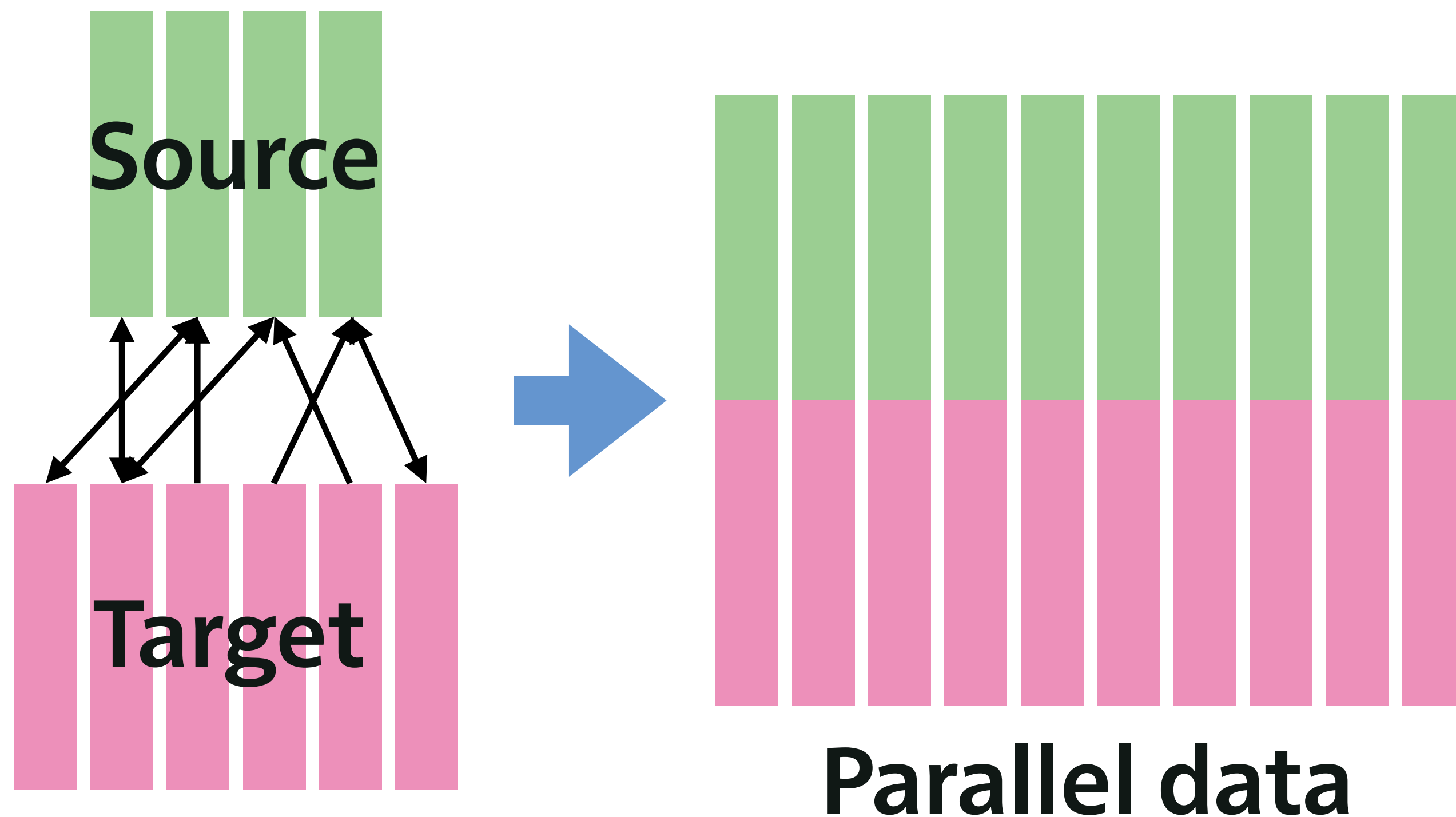✗ Degrades because of activation mismatch

# Outline

- Baseline 1: NMF-based voice conversion

- **Baseline 2: INCA algorithm**

- Proposed: Nonparallel training of NMF-based voice conversion
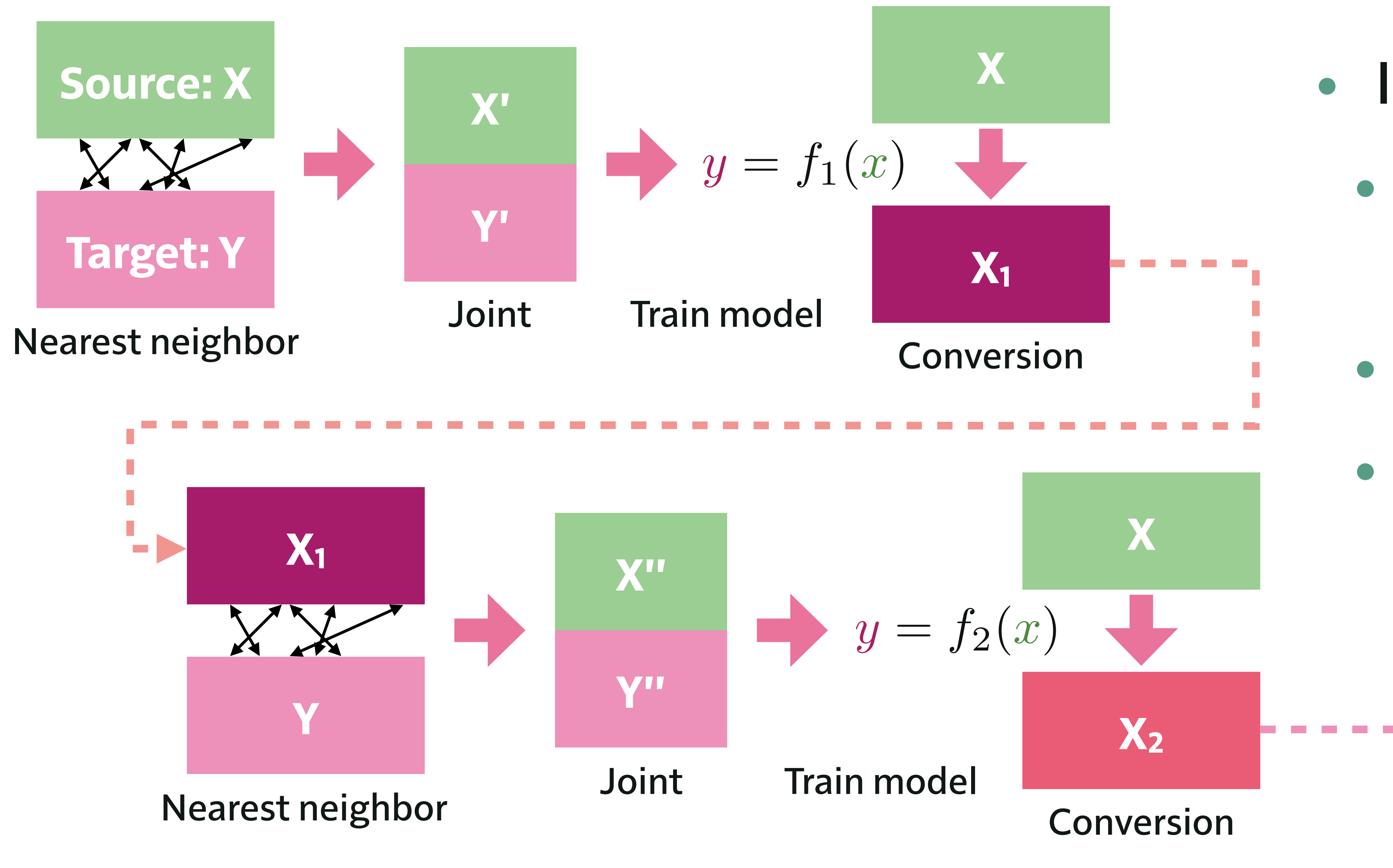
- Experiments

# INCA algorithm

[D. Erro+, 2010]

- **I**terative combination of a **N**earest neighbor search step and a **C**onversion step **A**lignment method

- Algorithm to acquire **alignment of nonparallel corpora**



Source

Target

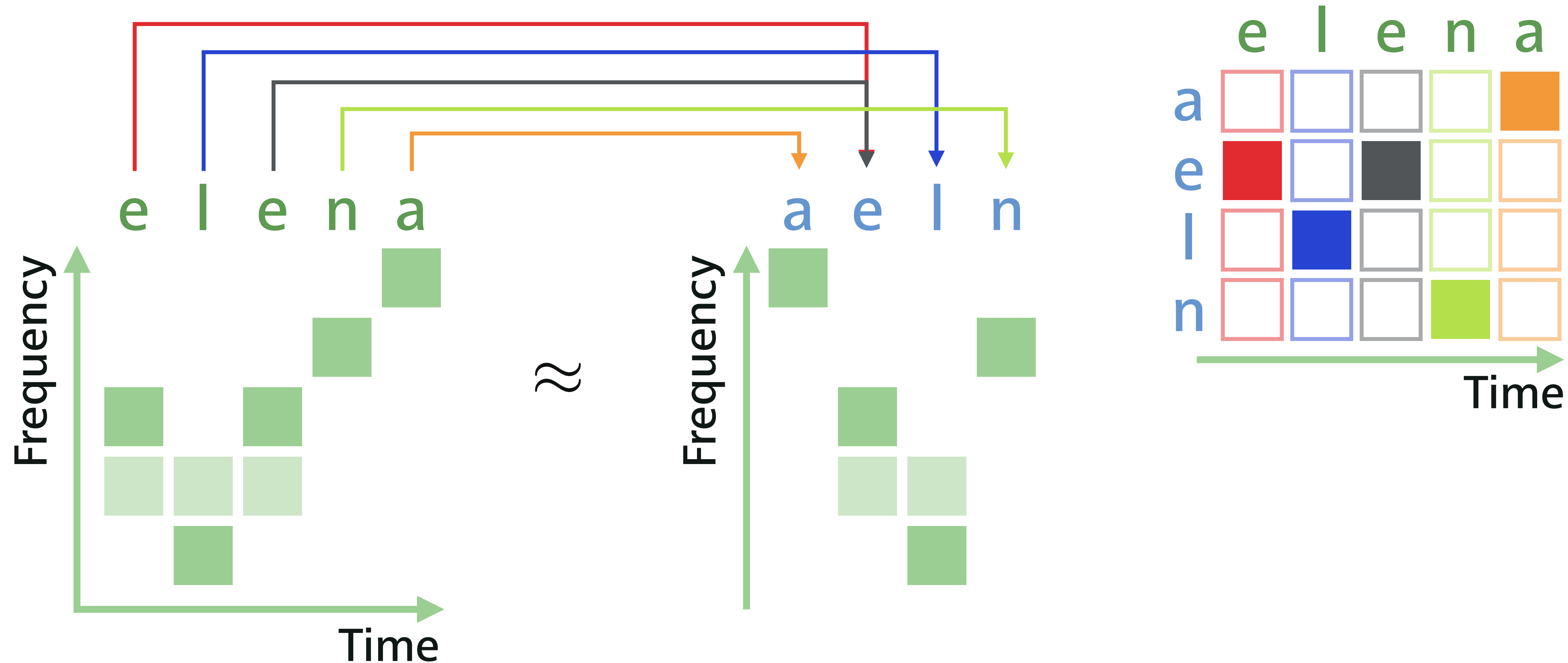Parallel data

[D. Erro+, 2010]



- Iteration of
- Nearest neighbor
- Train model
- Conversion

# INCA algorithm

✓ Applicable to any parallel VC frameworks

✓ Easy to implement


✗ Requires as much data as parallel VC systems

- VC system does not require much data for source speaker

  - **VC system is a generator of target speaker**

# Outline

- Baseline 1: NMF-based voice conversion

- Baseline 2: INCA algorithm

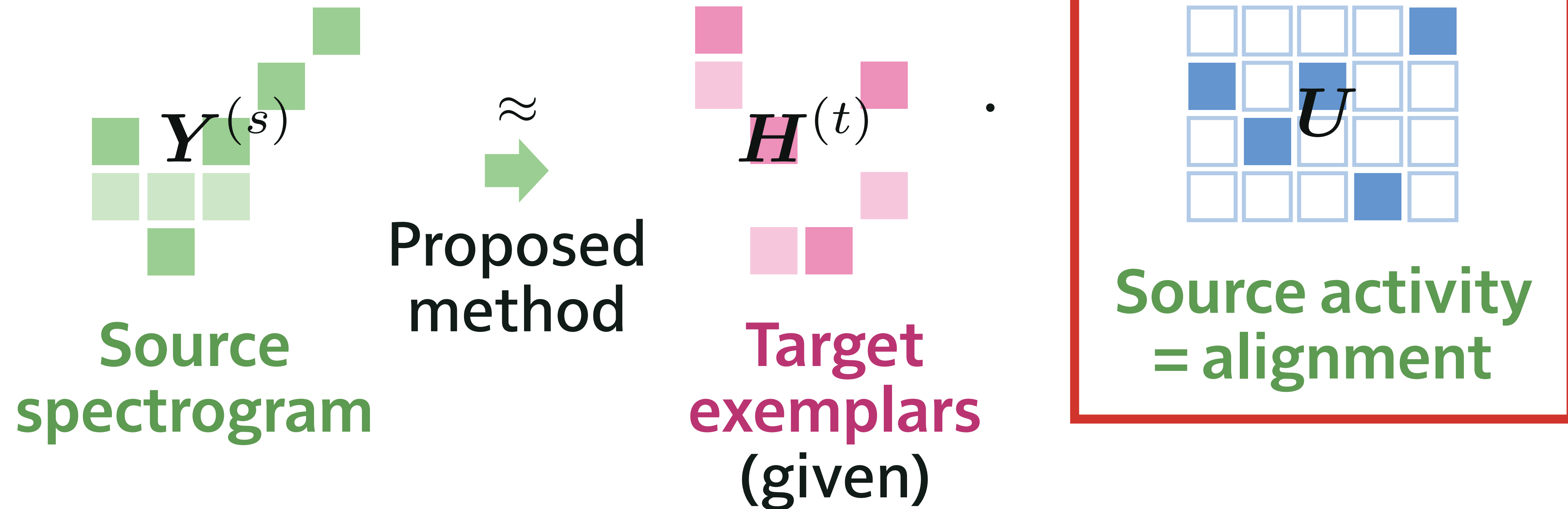- **Proposed: Nonparallel training of NMF-based voice conversion**

- Experiments

- NMF is equivalent to
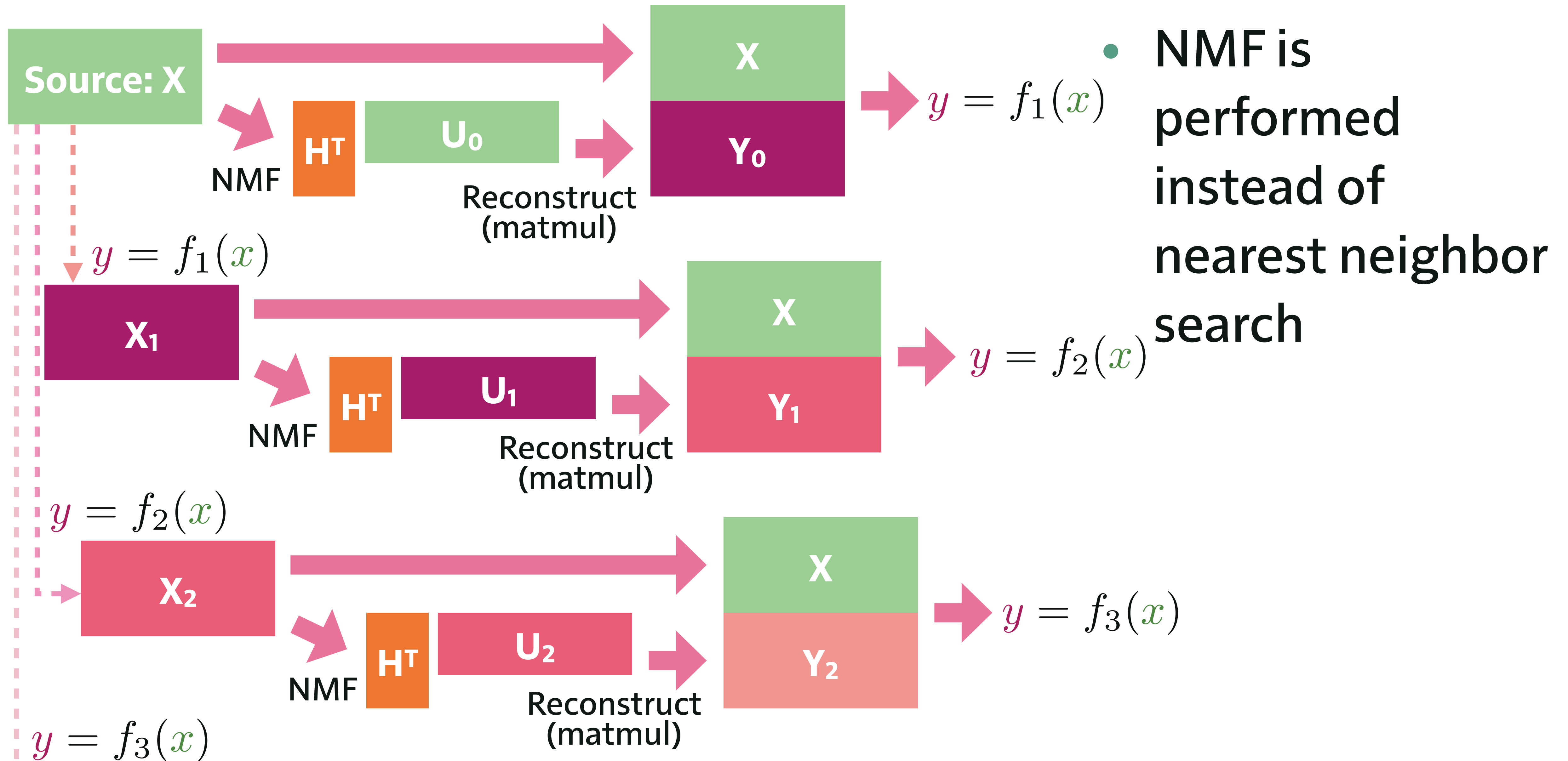  **alignment from spectrogram to exemplars**

- Acquires activation of NMF-VC using INCA technique

$$Y^{(s)} \approx H^{(t)} \cdot U$$

**Source spectrogram**

Proposed method

**Target exemplars** (given)

**Source activity = alignment**

- Activity is irrelevant to speaker of exemplars

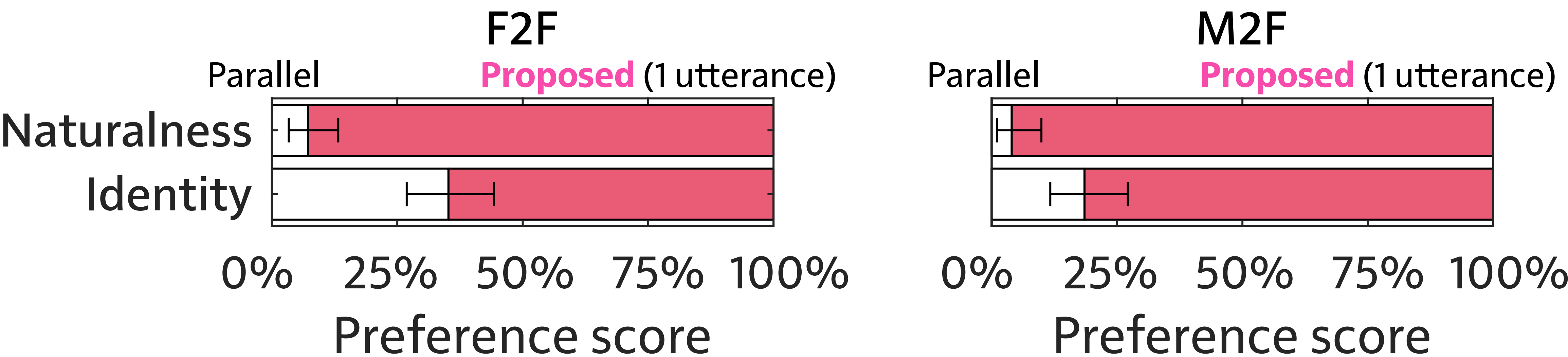- NMF is performed instead of nearest neighbor search

# Outline

- Baseline 1: NMF-based voice conversion

- Baseline 2: INCA algorithm

- Proposed: Nonparallel training of NMF-based voice conversion

- **Experiments**

# Experimental setups

- Methods for comparison

  - CycleGAN-VC [T. Kaneko+, 2019], Parallel NMF-VC [R. Takashima+, 2012]

- Dataset: Japanese versatile speech corpus [S. Takamichi+, 2019]

- Speaker pair: JVS066→JVS010 (F2F), JVS054→JVS010 (M2F)

- Number of sentences: 60 for target speaker, 1 or 10 for source speakers, 20 for test

  - Same 60 sentences are used for source speakers in parallel NMF-VC

- Analysis / Synthesis: WORLD [M. Morise+, 2016]

- Target of decomposition: Amplitude spectrograms of WORLD spectral envelopes

- Dictionary size (number of exemplars): 200

- Conditions of subjective experiments (≥ 25 subjects)

  - A/B preference tests for naturalness

  - ABX tests for speaker similarity

- Comparison with parallel NMF-VC

F2F

Parallel **Proposed** (1 utterance)
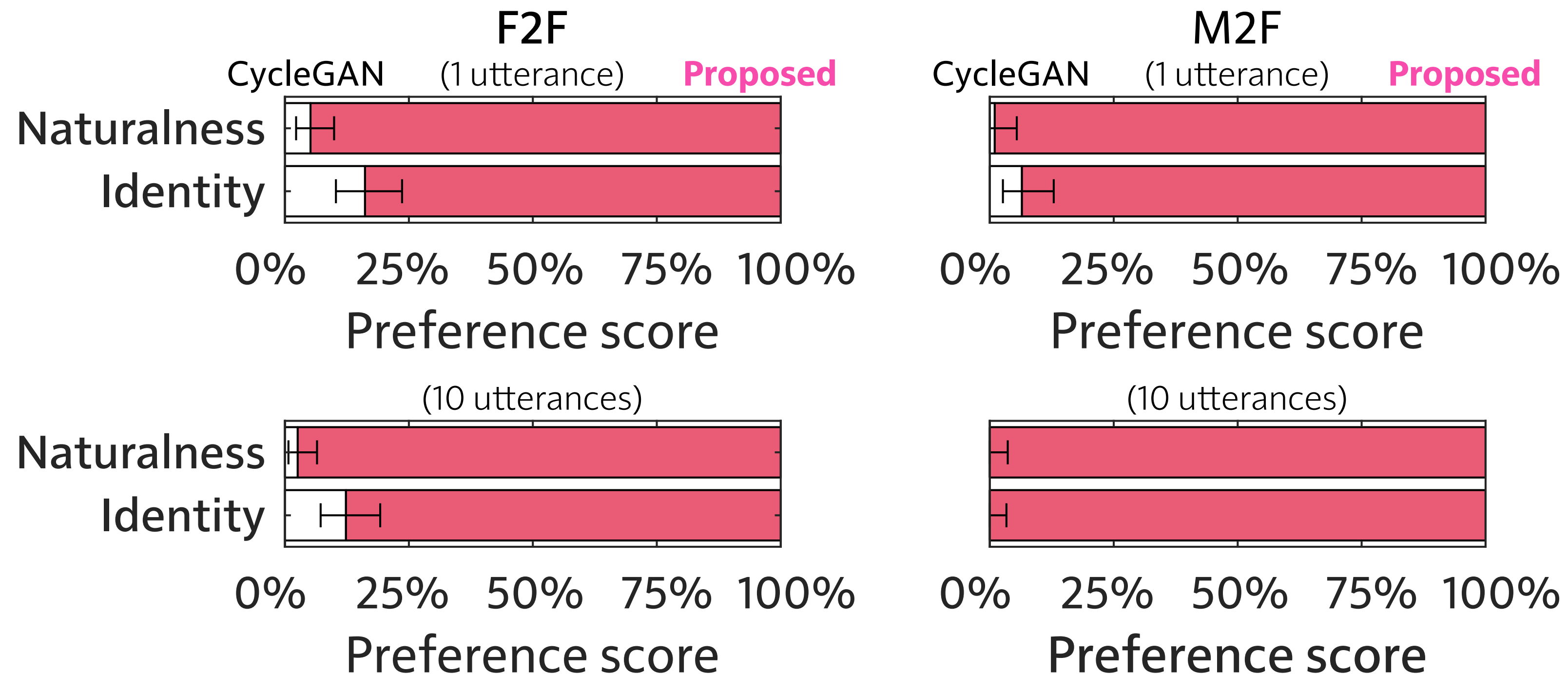
Naturalness
Identity

0%  25%  50%  75%  100%

Preference score

M2F

Parallel **Proposed** (1 utterance)

Naturalness
Identity

0%  25%  50%  75%  100%

Preference score

- Proposed framework outperformed conventional parallel NMF-VC framework
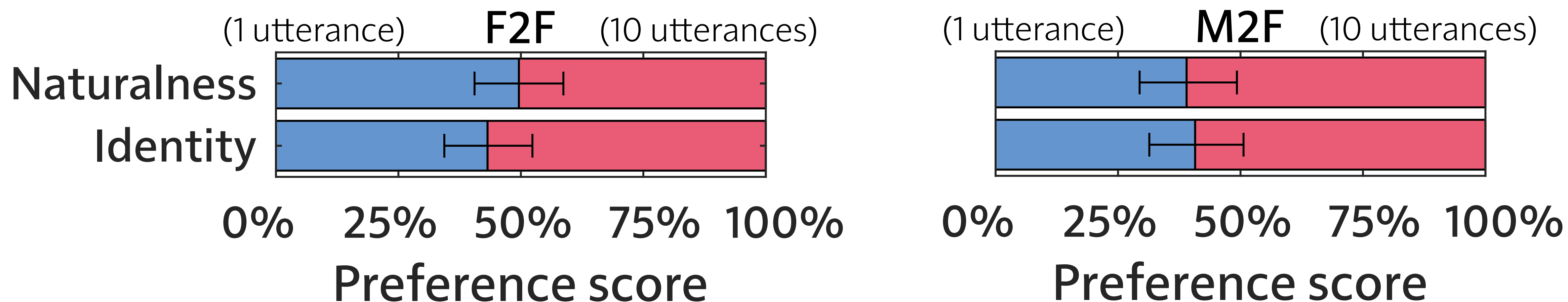
⊢⊣ : 95% confidential

- ## Comparison with CycleGAN-VC



- ## Proposed framework efficiently used a small source speakers' data

⊢——⊣ : 95% confidential

# Results of subjective experiments

- Effects of an amount of source speakers' data



- More **source** speakers' data provided more speaker similarity about **target** speaker

# Audio examples

| | Parallel NMF-VC | CycleGAN-VC | Proposed |
|---|---|---|---|
| Target (Female) | | 🔊 | |
| F2F Source | | 🔊 | |
| F2F Converted | 🔊 | 🔊 | 🔊 |
| M2F Source | | 🔊 | |
| M2F Converted | 🔊 | 🔊 | 🔊 |

- 1 sentence is used for training each source speaker

- Examples are available online (see last slide)

# Conclusions

- We introduce **nonparallel training method of NMF-VC** inspired by the INCA algorithm

- The method achieved the goal with the small amount of the training data for source speakers

## Future Works

- Higher quality conversion

- Inter-language conversion

https://www.gavo.t.u-tokyo.ac.jp/~hitoshi/nmfvc2020interspeech/