# An Unsupervised Probability Model for Speech-to-Translation Alignment of Low-Resource Languages

Antonios Anastasopoulos[1], David Chiang[1], Long Duong[2]

[1] University of Notre Dame, USA
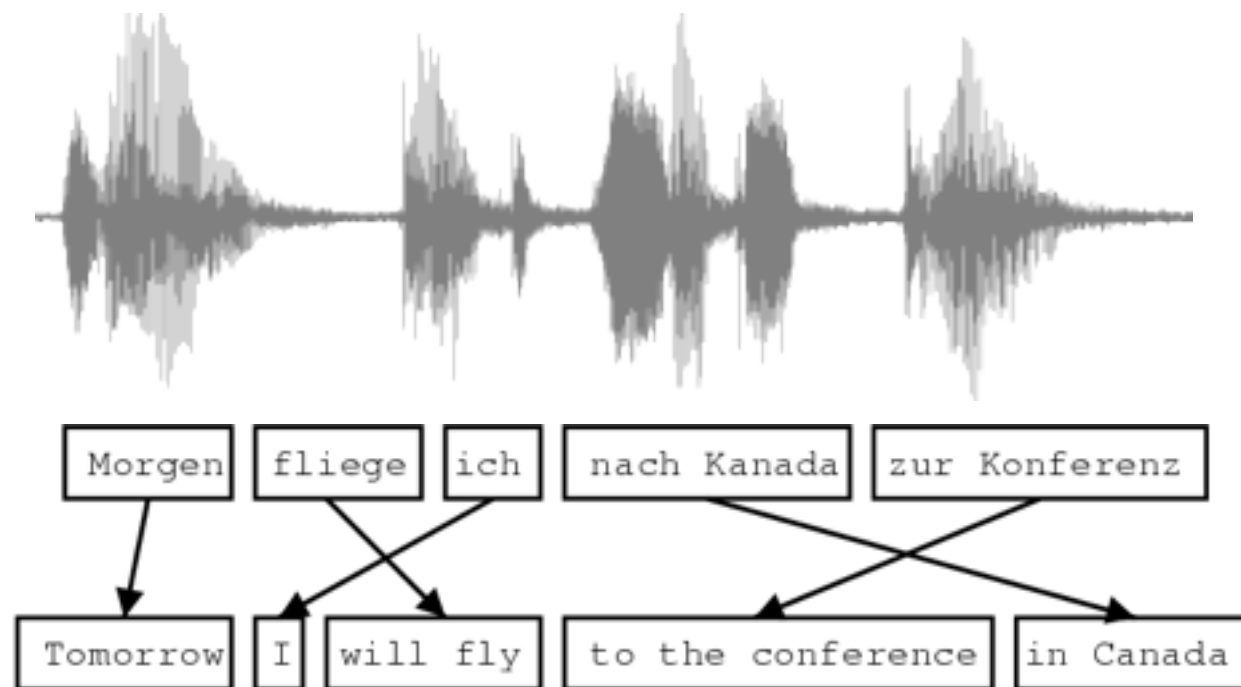[2] University of Melbourne, Australia

# Motivation

**Why Speech-based MT?**

90% of languages do not have a writing system

# Motivation
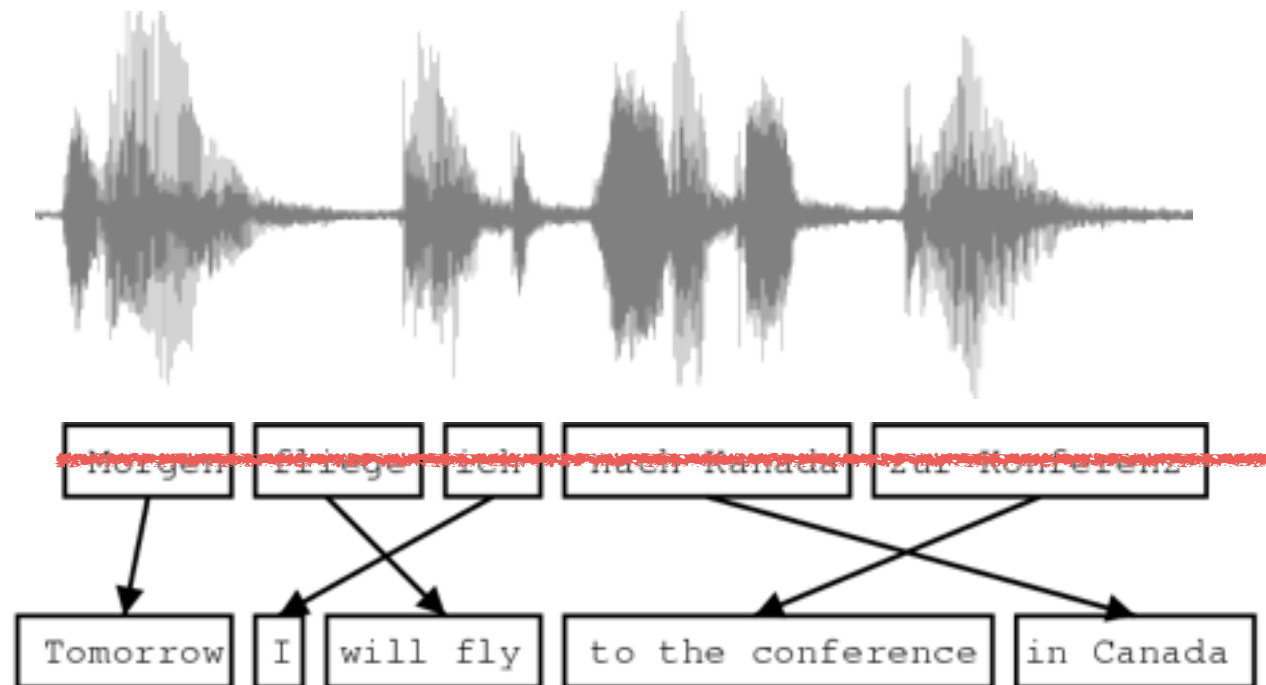
**Why Speech-based MT?**

90% of languages do not have a writing system

# Motivation

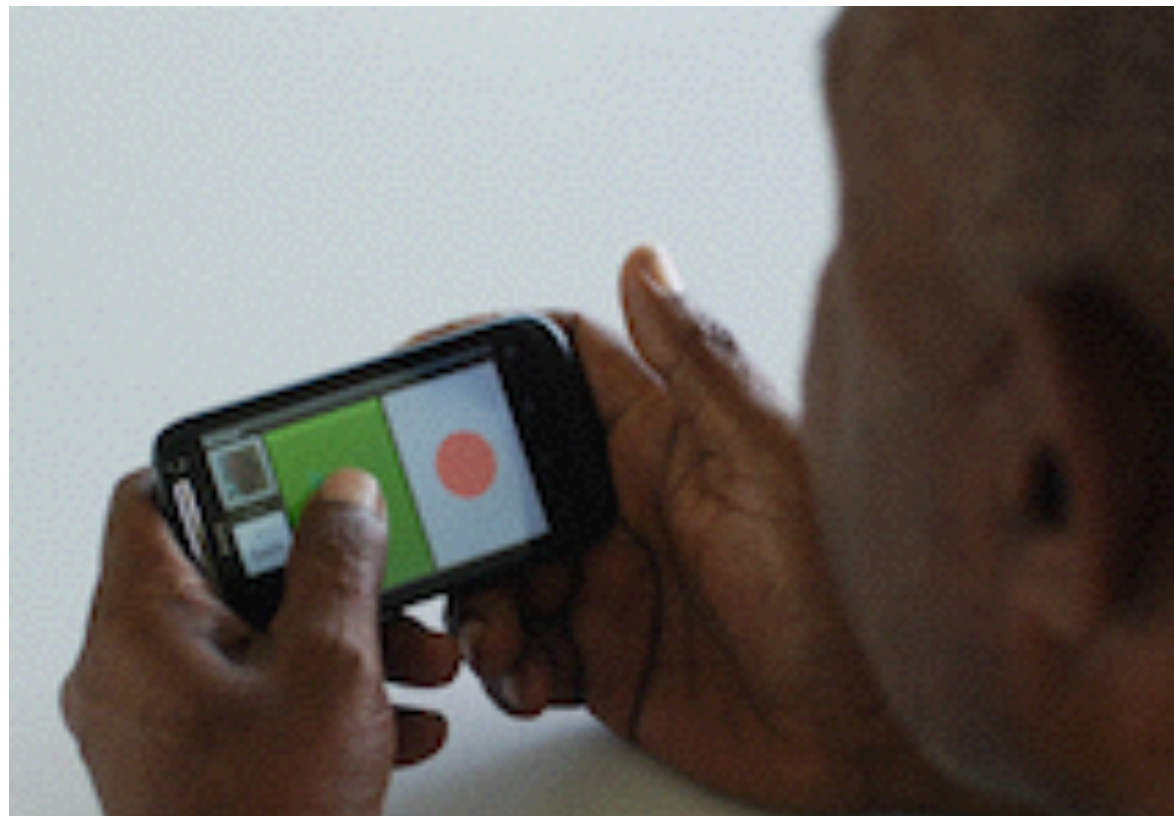**Why Speech-based MT?**

90% of languages do not have a writing system

# Motivation

**Endangered languages documentation**

Use speech with translations



Using the Aikuma (Bird 2010)
app to collect parallel speech
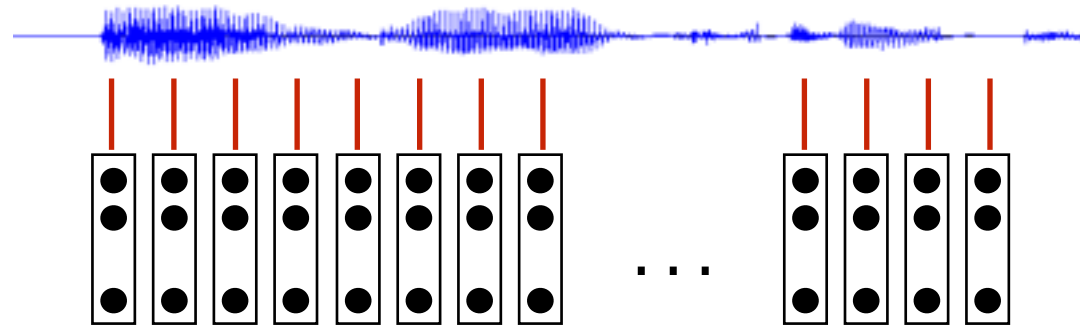
# Motivation

**Low-resource languages**

Utilize translations rather than transcriptions
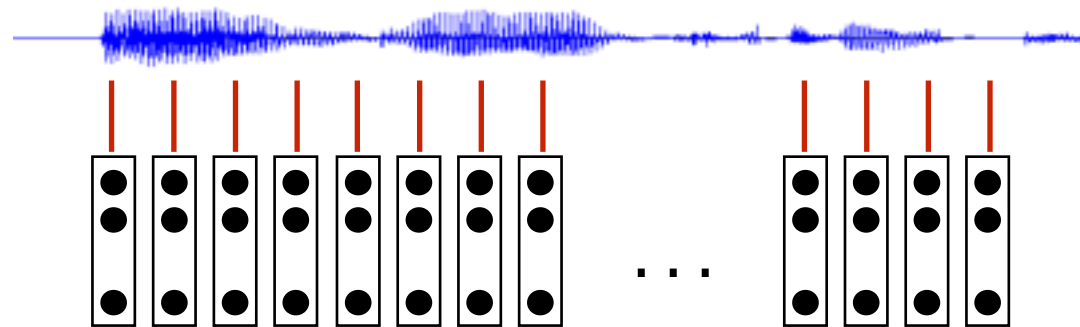
# Task Description

# Task Description

**Source side:** Frames of the speech signal

# Task Description

**Source side:** Frames of the speech signal
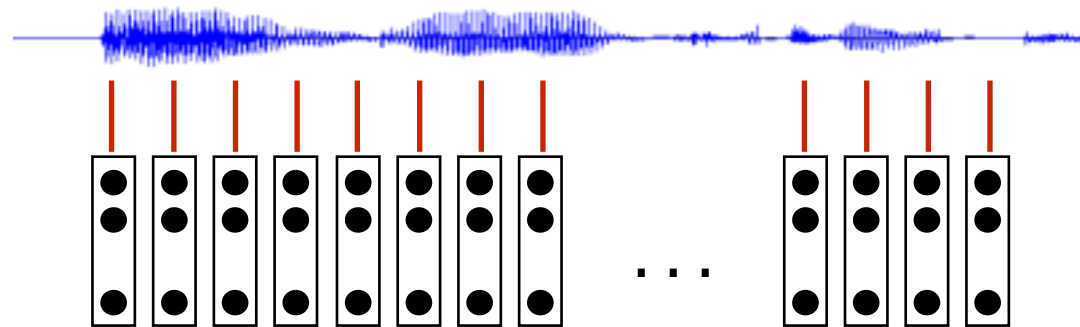
**Target side:** Translation text



a
little
bit
of
knowledge

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

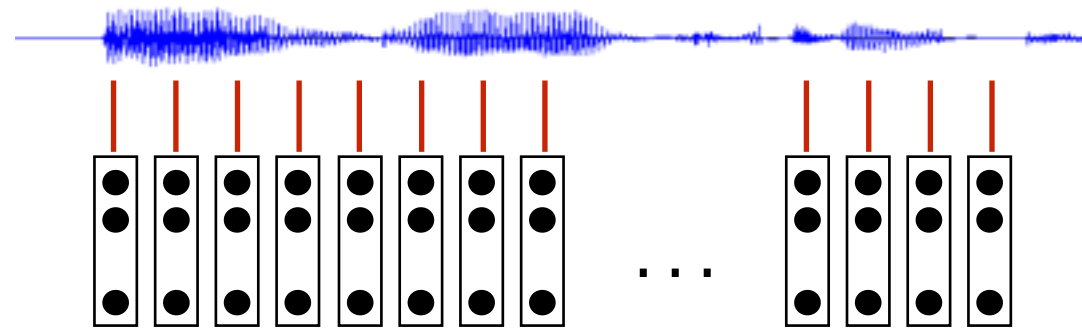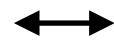**Task**: find best alignment between source and target side



a
little
bit
of
knowledge

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

**Task:** find best alignment between source and target side



a ↔

little
bit
of
knowledge

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

**Task**: find best alignment between source and target side
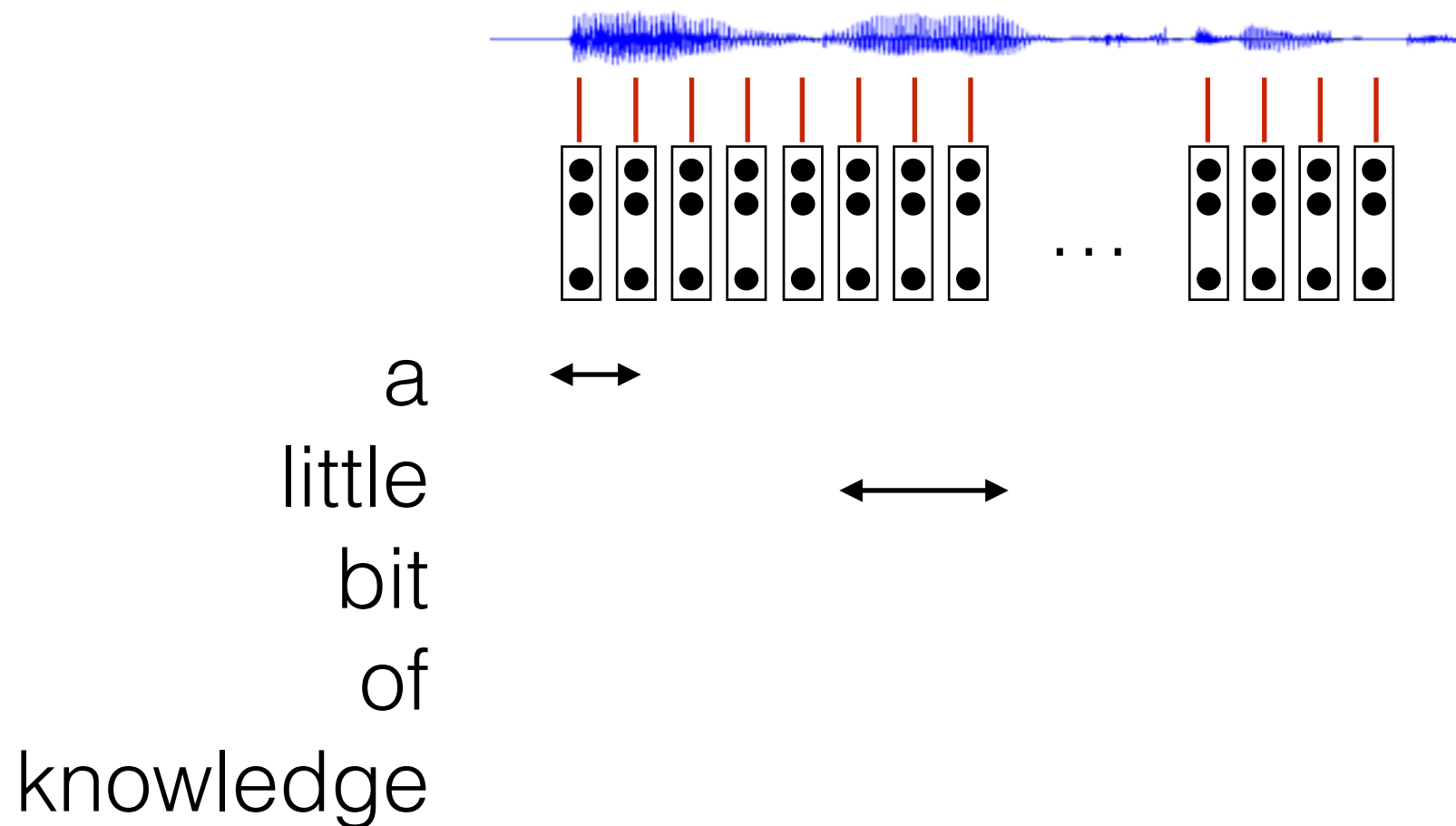


a

little

bit

of

knowledge

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text
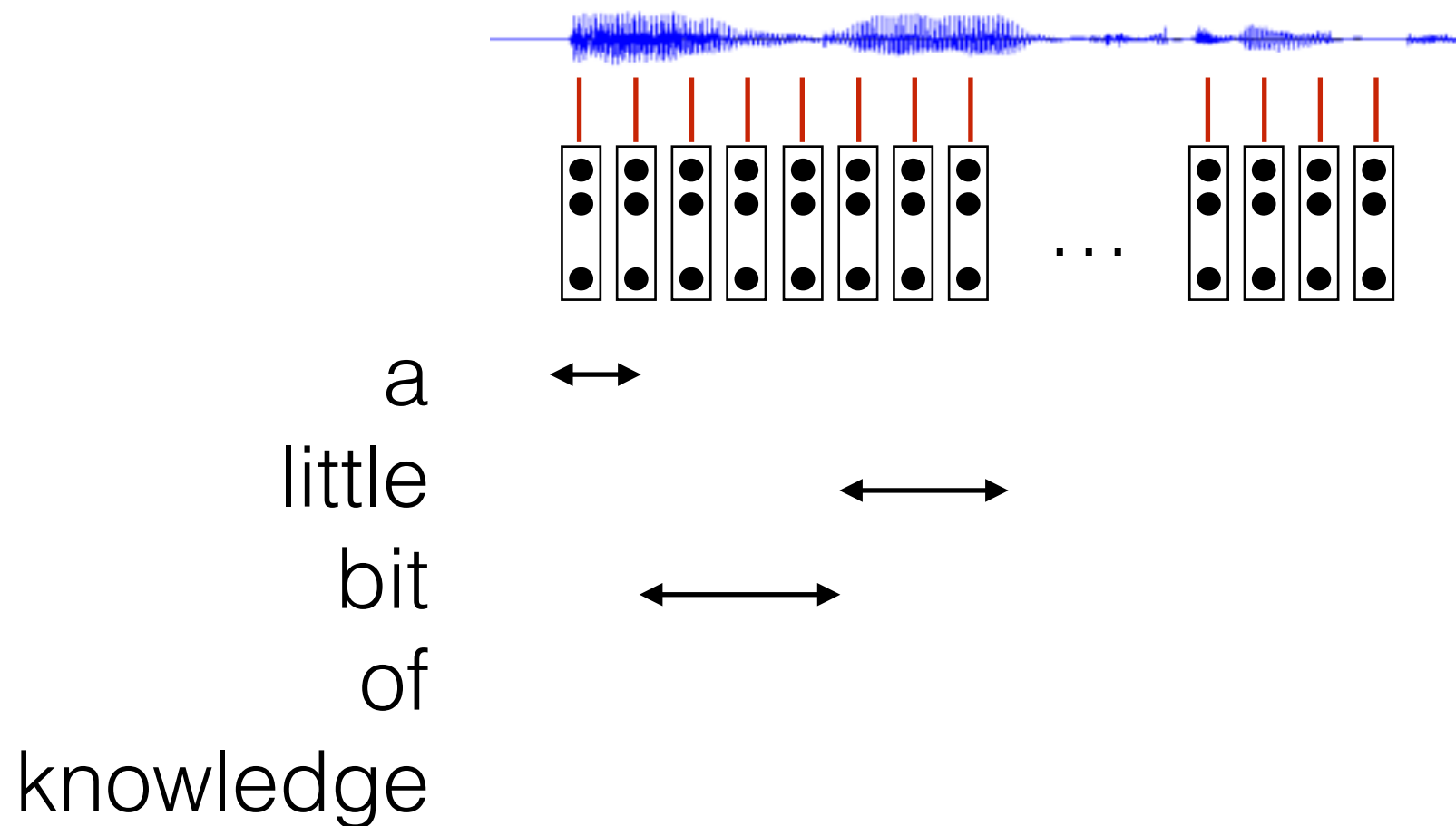
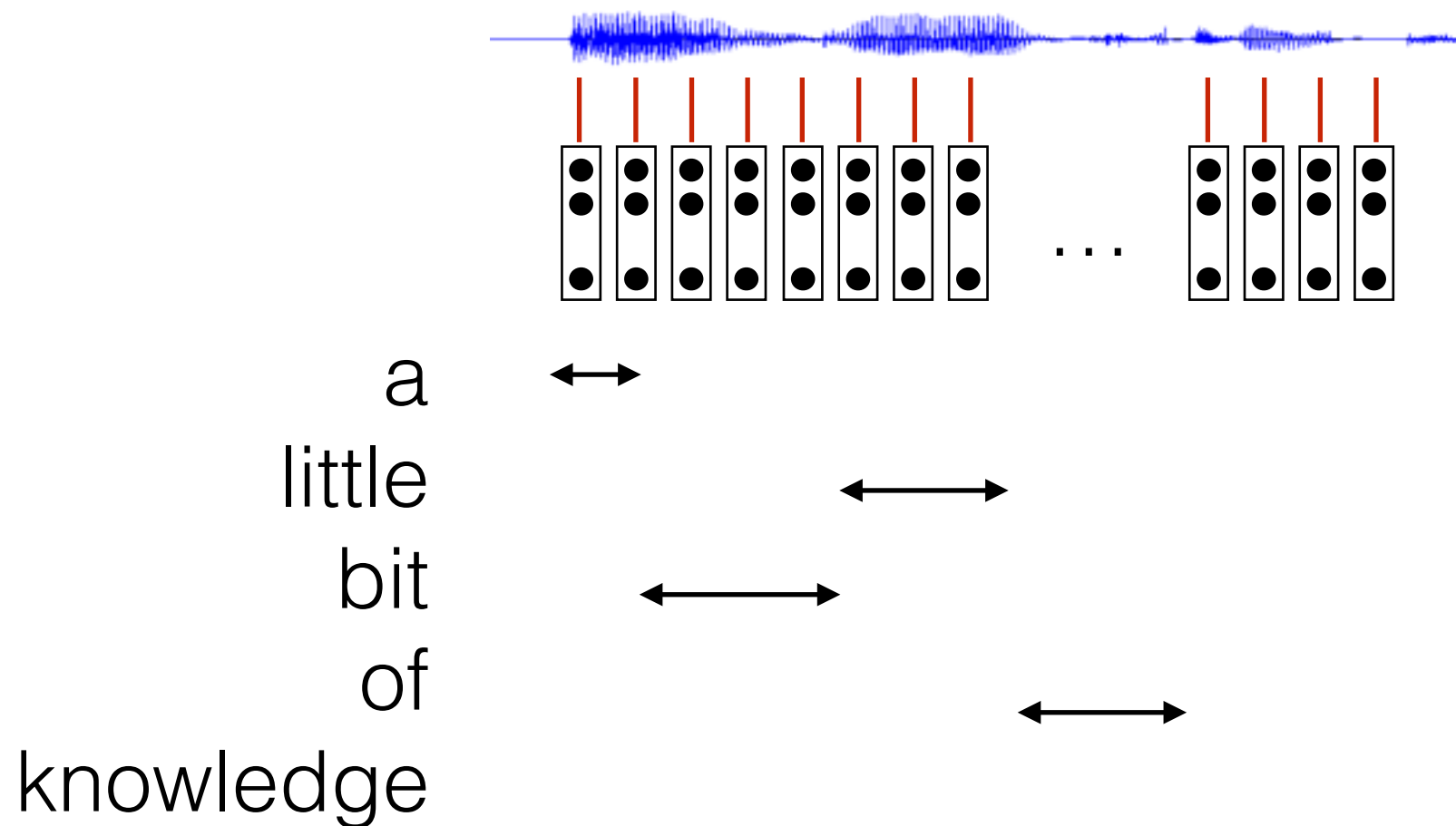**Task**: find best alignment between source and target side

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

**Task**: find best alignment between source and target side



a ↔

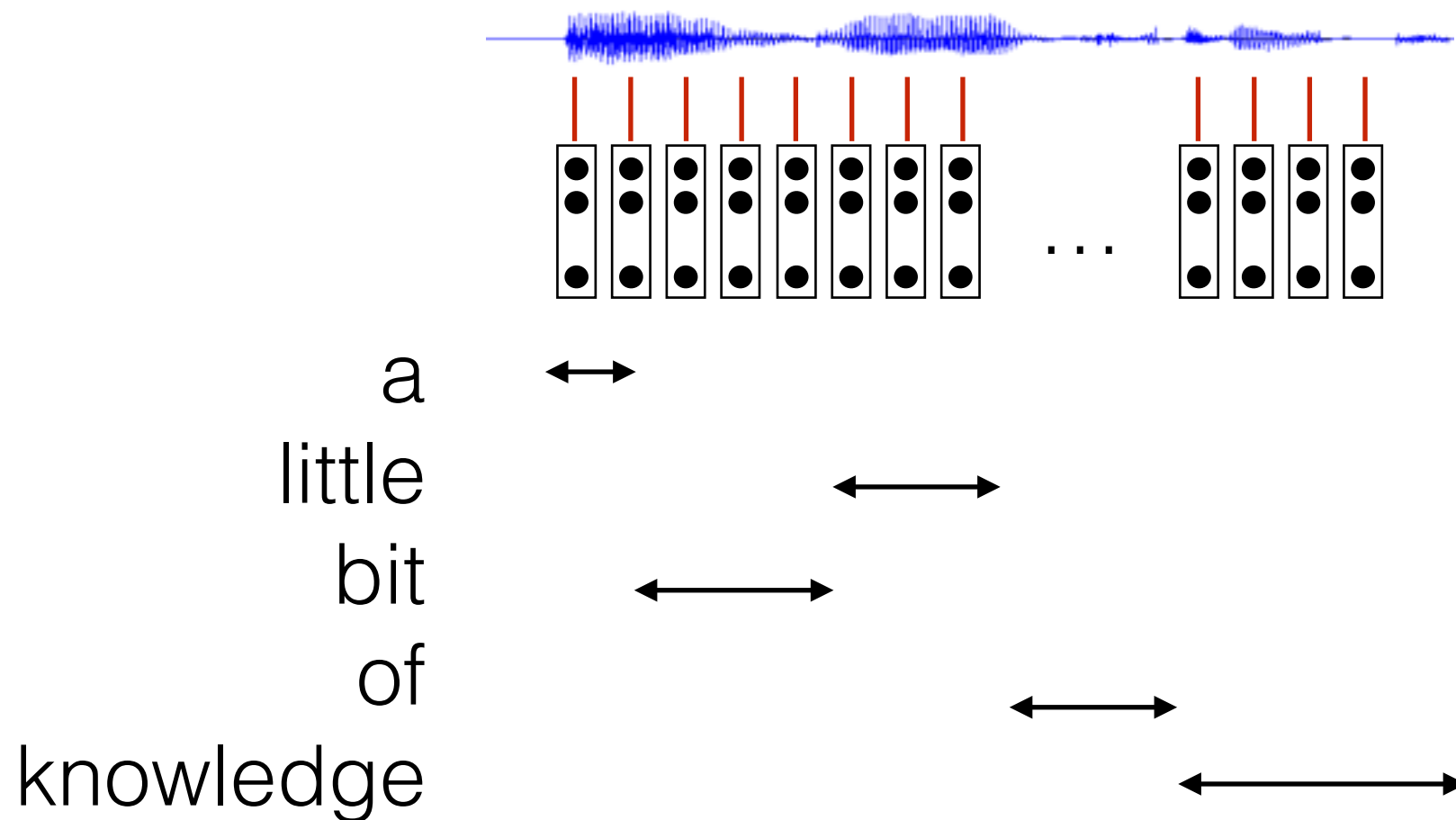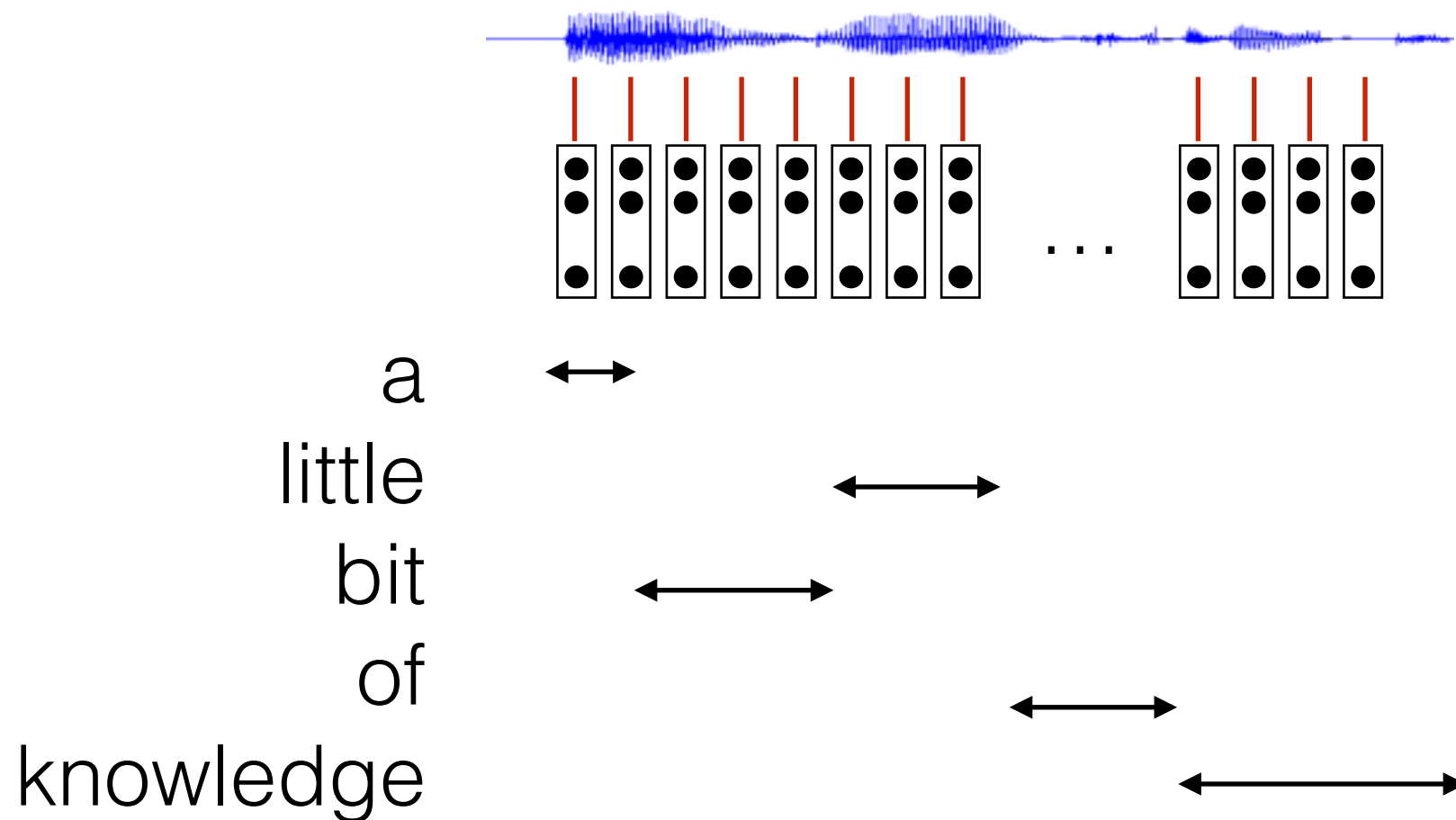little ⟷

bit ⟷

of ⟷

knowledge

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

**Task**: find best alignment between source and target side

# Task Description

**Source side:** Frames of the speech signal

**Target side:** Translation text

**Task**: find best alignment between source and target side



**Our method outperforms both baselines**
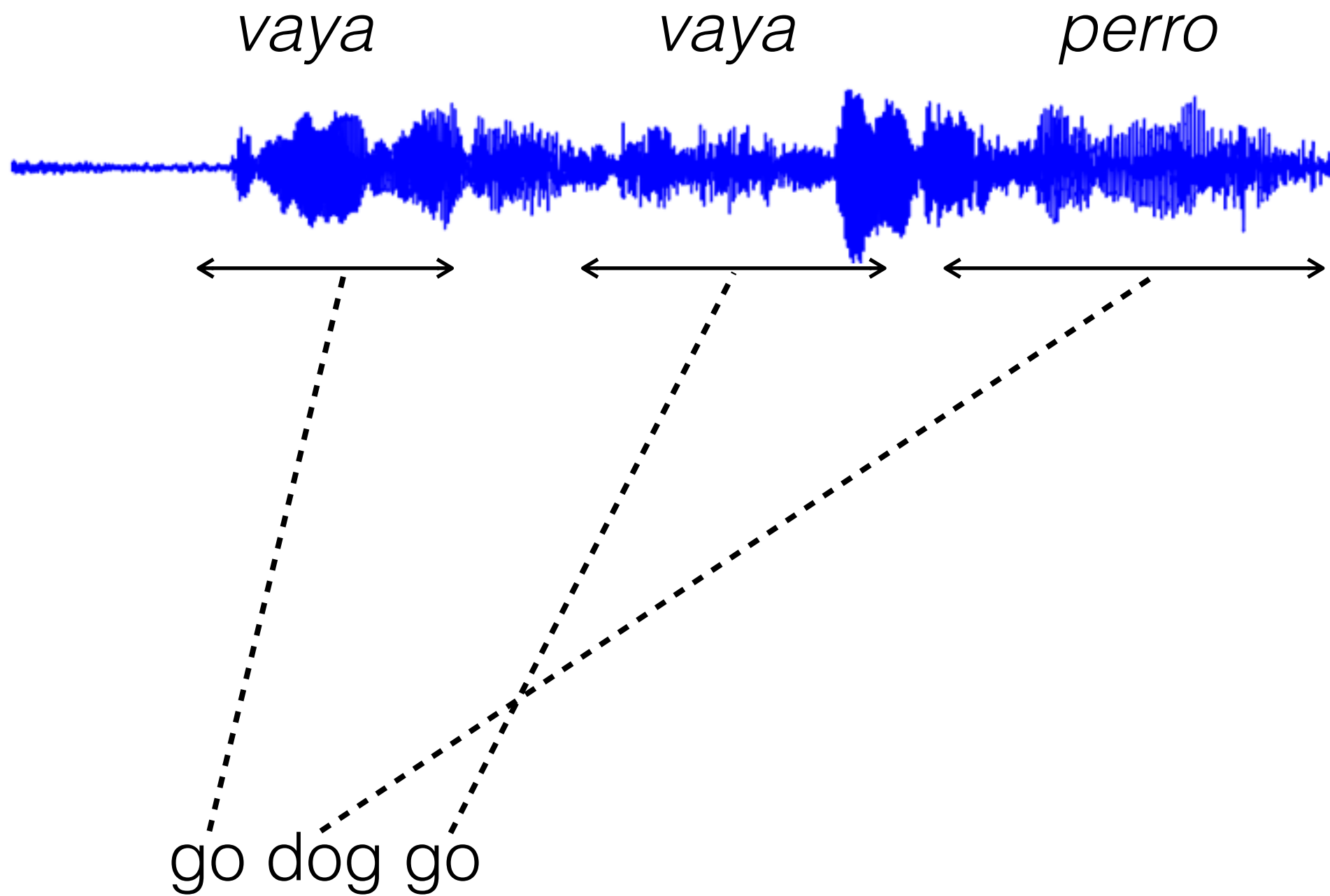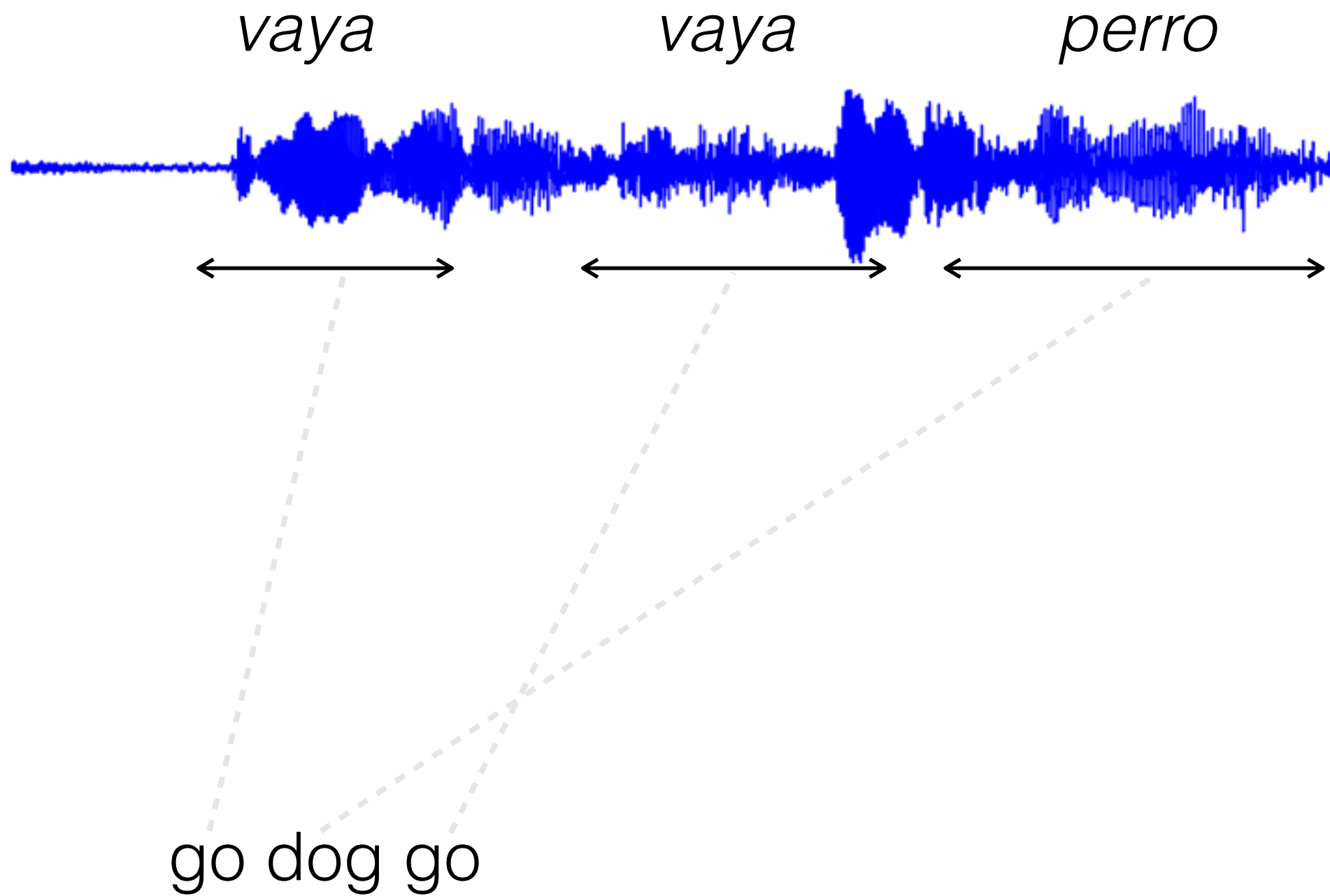
# The big picture

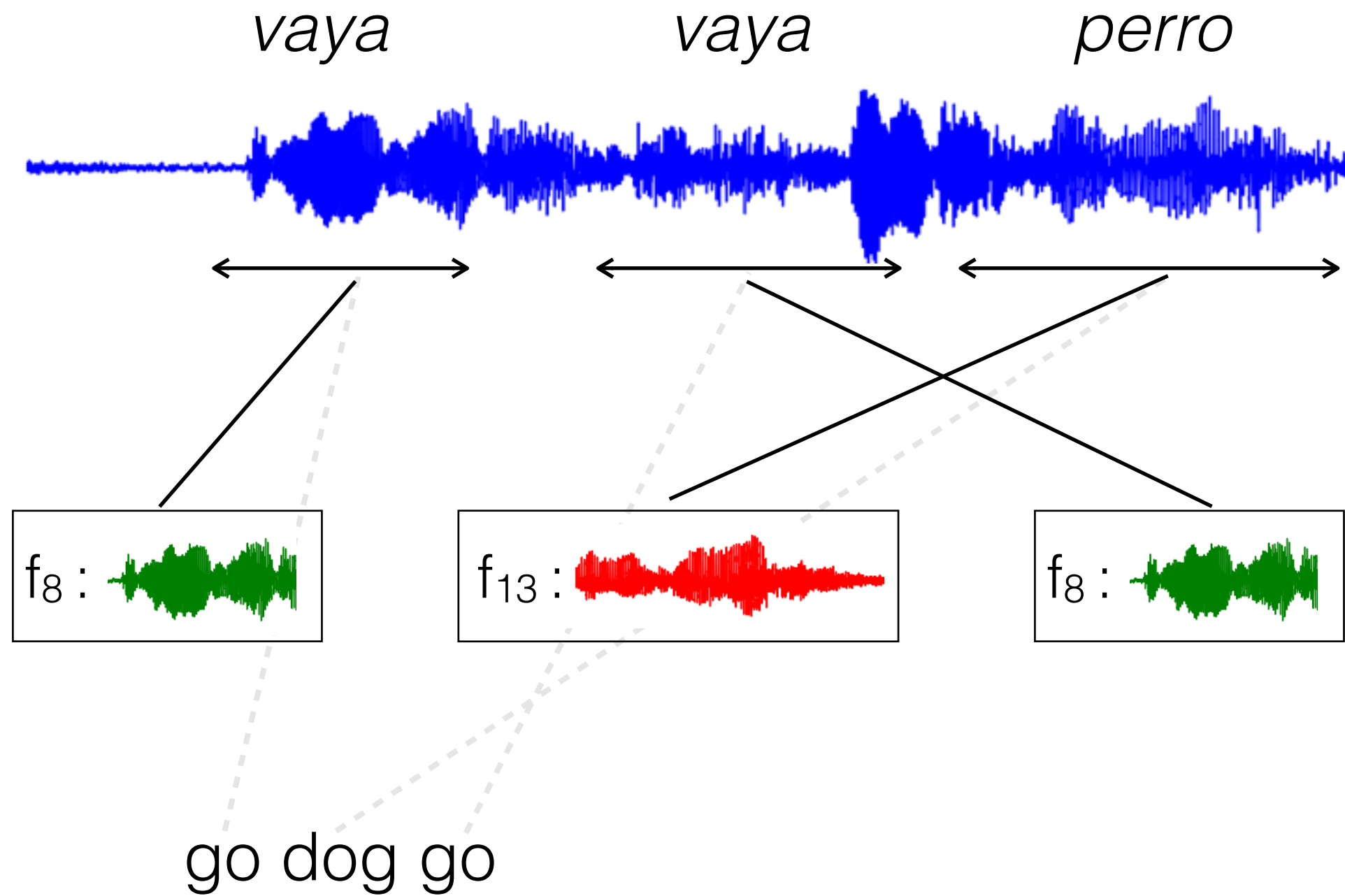# The big picture

*vaya*  *vaya*  *perro*



go dog go

# The big picture

# The big picture



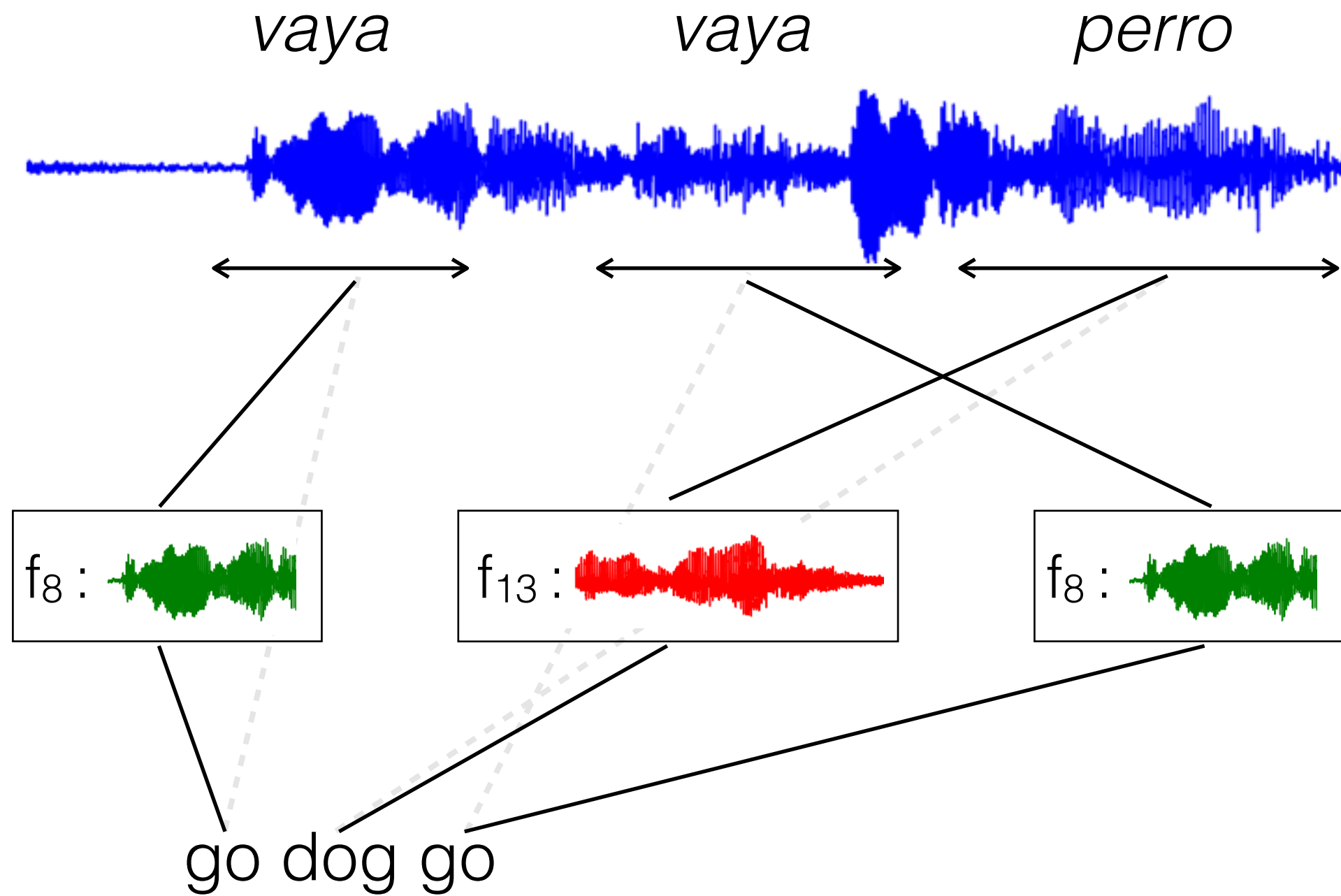*vaya*     *vaya*     *perro*

go dog go

# The big picture

# The big picture

# Distortion model

Vaya    vaya    pero

Controls the reordering of the target words
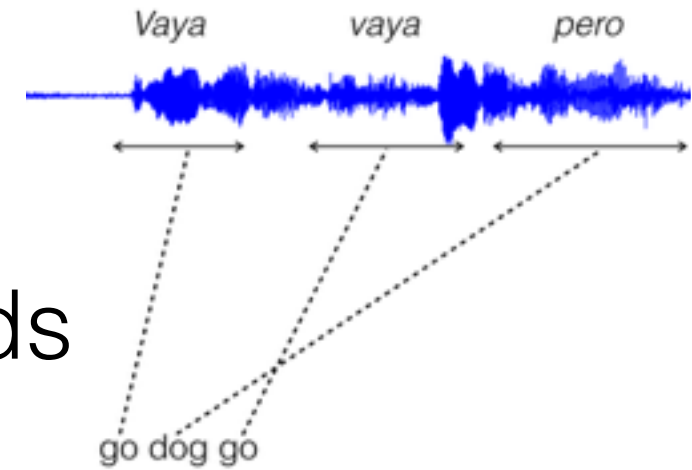  - based on fast-align [Dyer et al.]
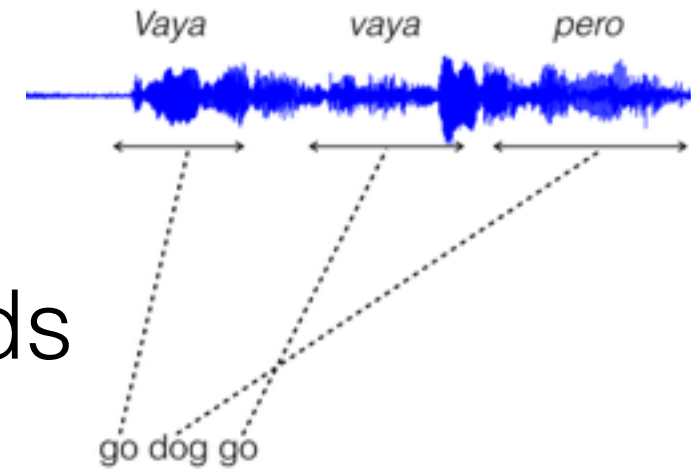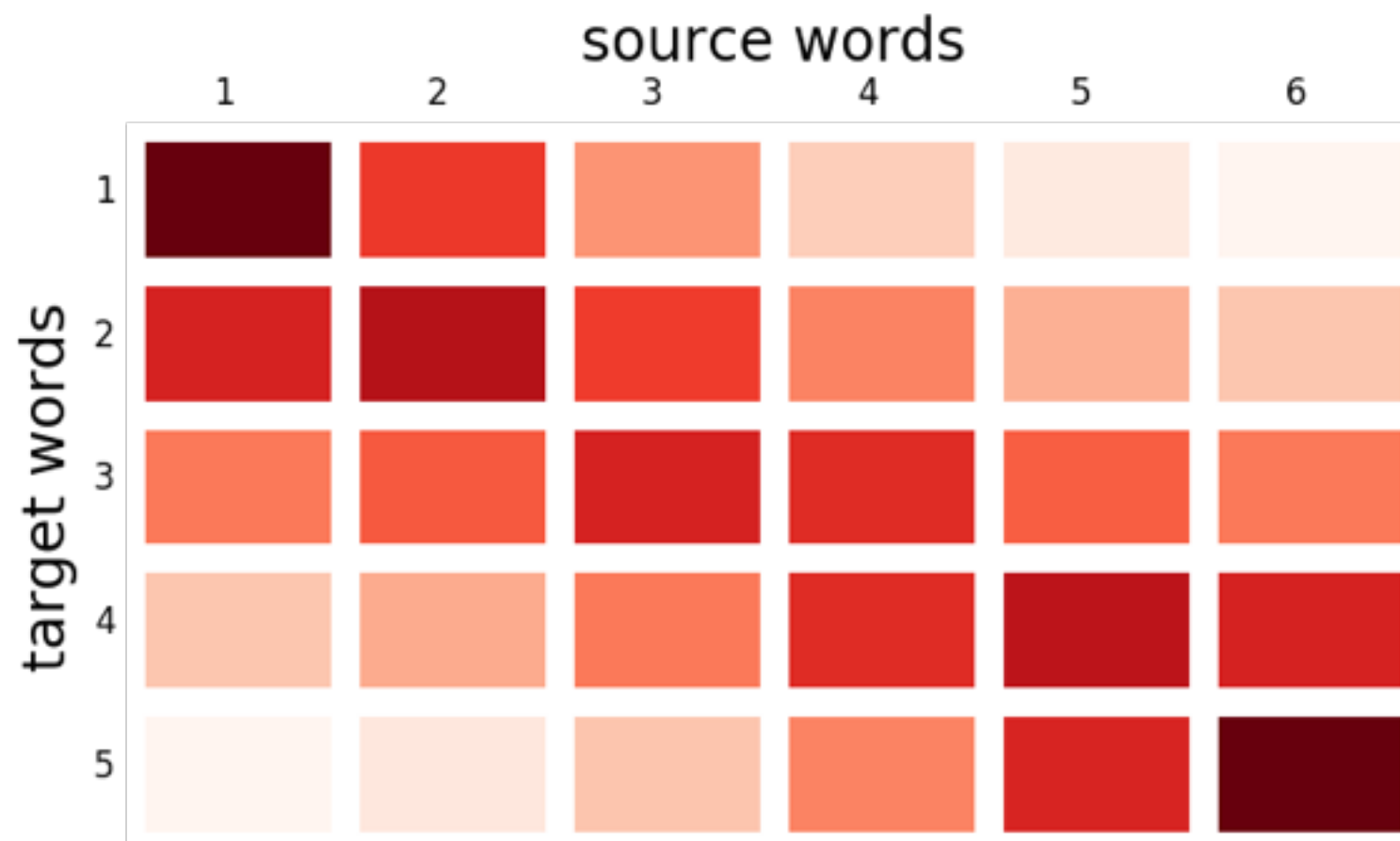
go dog go

# Distortion model

Controls the reordering of the target words
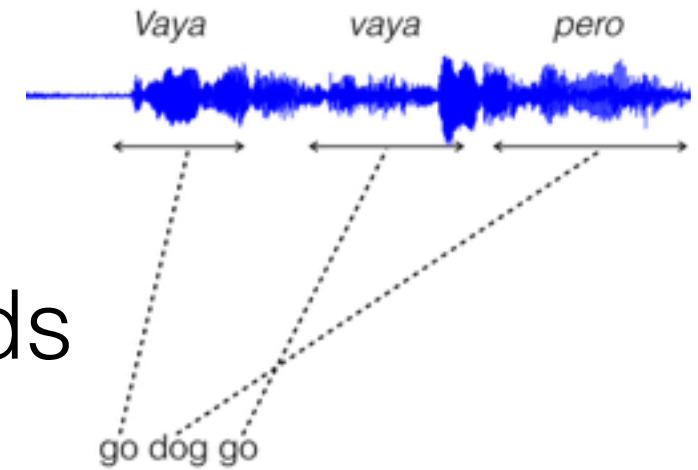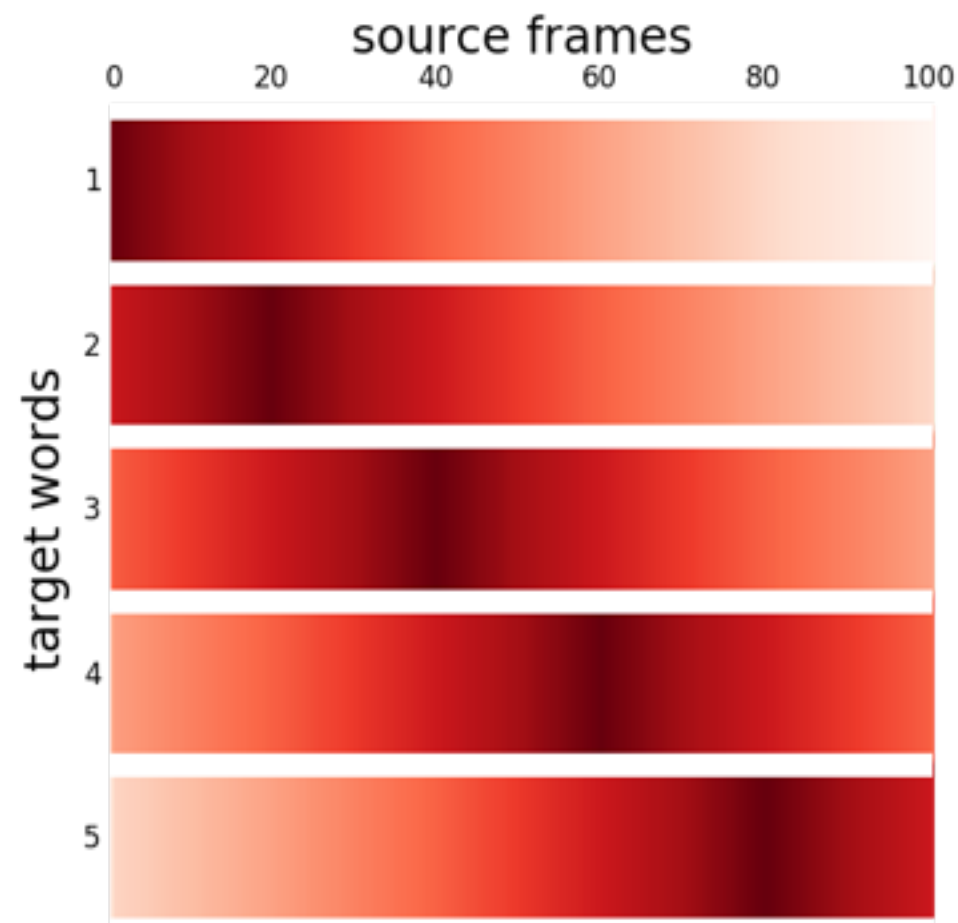 - based on fast-align [Dyer et al.]
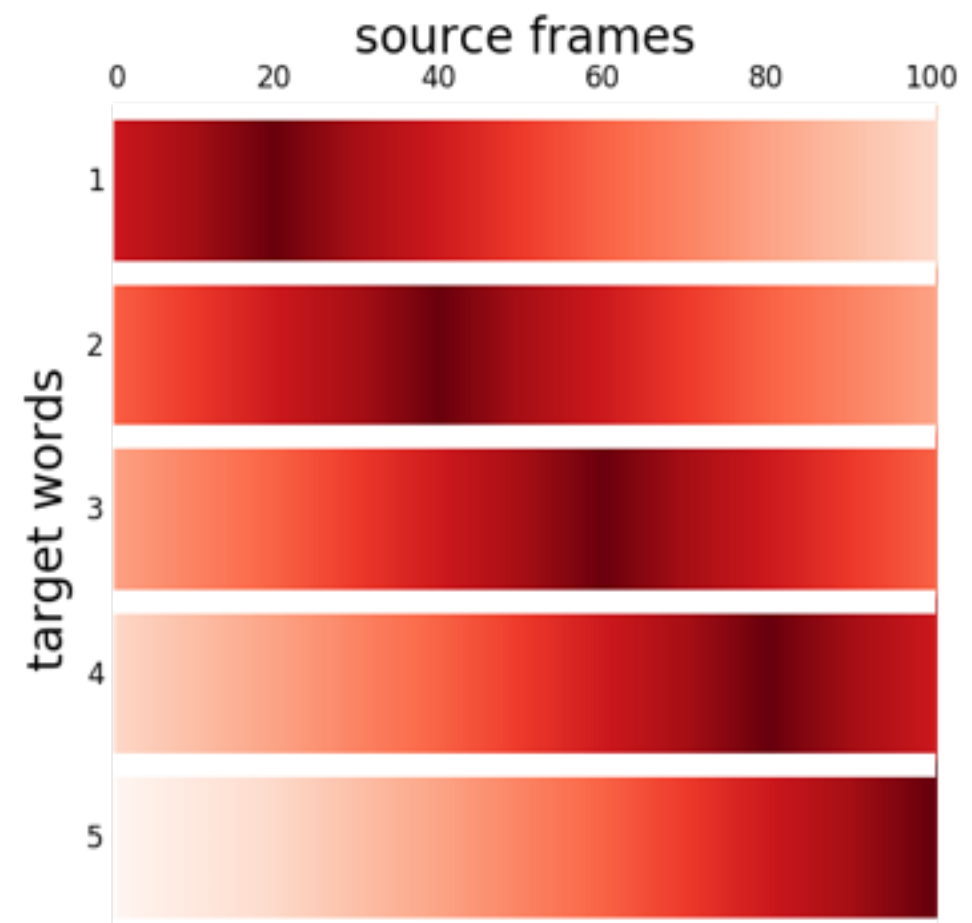
Original

# Distortion model

Controls the reordering of the target words
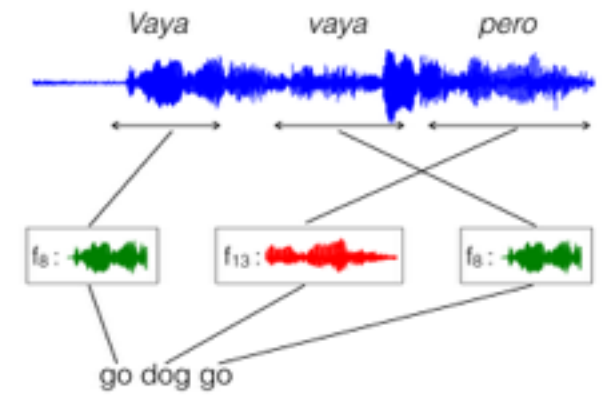- based on fast-align [Dyer et al.]

Modification

Span Start

Span End

# Clustering model



Assuming a *"prototype"* for each cluster

# Clustering model

Assuming a *"prototype"* for each cluster

# Clustering model

Assuming a *"prototype"* for each cluster

# Clustering model

Assuming a *"prototype"* for each cluster



$f_8$ ⟶

# Clustering model

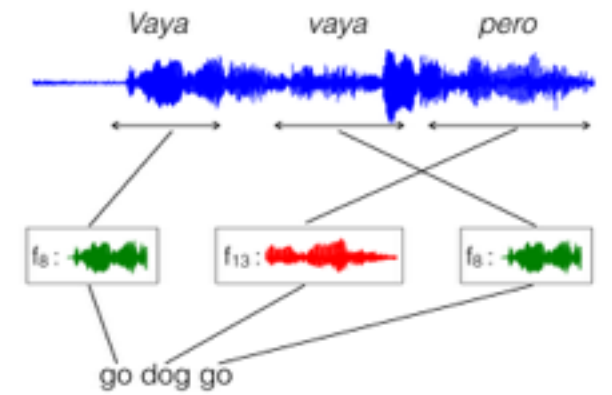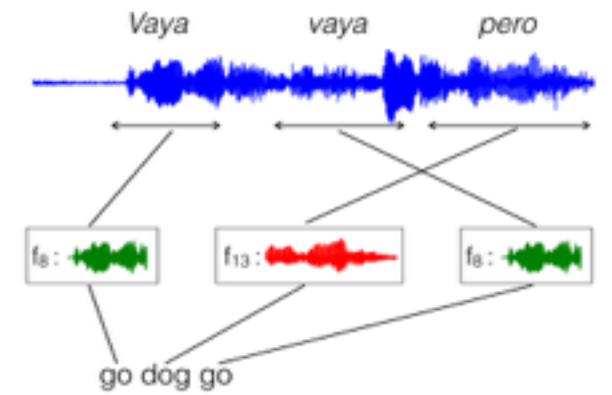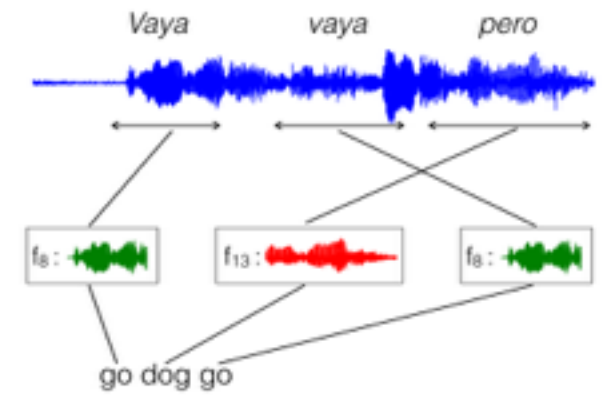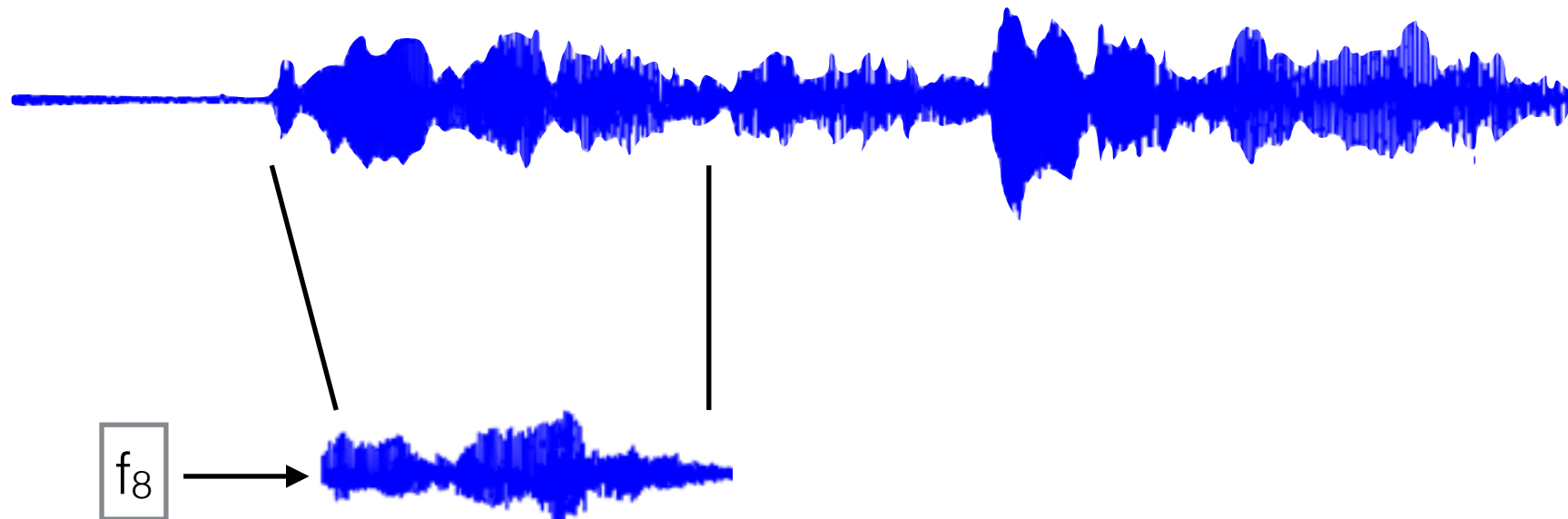Assuming a *"prototype"* for each cluster
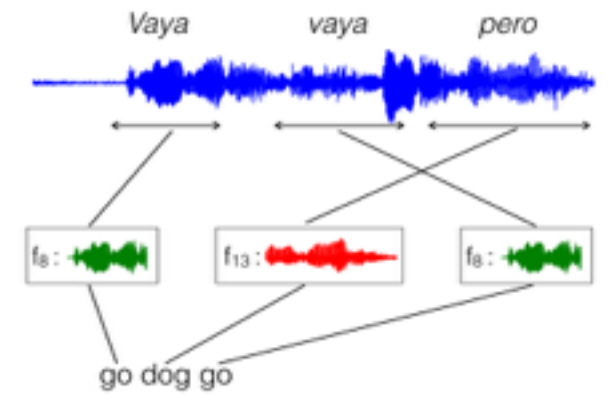


$f_8$

# Clustering model

Assuming a *"prototype"* for each cluster

# Clustering model

Assuming a *"prototype"* for each cluster

# Clustering model

Assuming a *"prototype"* for each cluster

$f_8$

# Clustering model

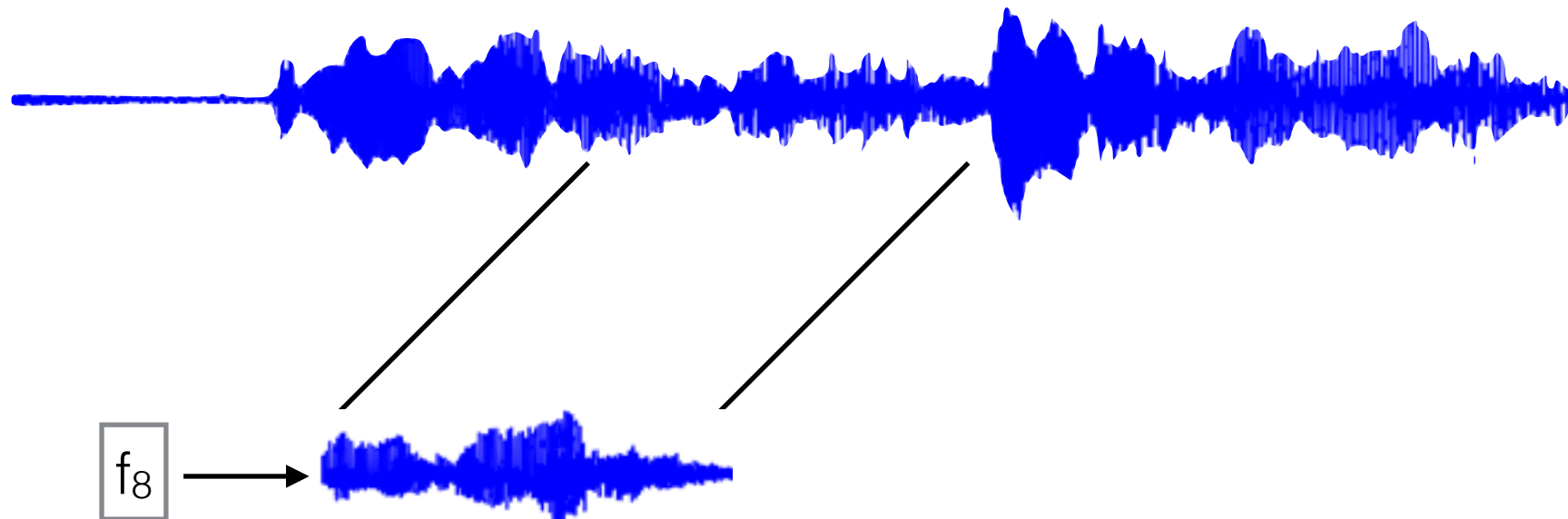Assuming a *"prototype"* for each cluster



$f_8$

# Training

Expectation-Maximization

# Training

Expectation-Maximization



*Vaya*     *vaya*     *pero*

go dog go

# Training

Expectation-Maximization

Initialize spans and clusters

Vaya            vaya            pero

go dog go

# Training

Expectation-Maximization

Initialize spans and clusters

# Training

Expectation-Maximization

Initialize spans and clusters

# Training

Expectation-Maximization

Initialize spans and clusters

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
  - Re-estimate prototypes

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
  - Re-estimate prototypes

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
    - Re-estimate prototypes

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
    - Re-estimate prototypes
- E step:

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
  - Re-estimate prototypes
- E step:
  - Assign cluster and align

Vaya        vaya        pero

$f_8$

go dog go

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
    - Re-estimate prototypes
- E step:
    - Assign cluster and align

# Training

Expectation-Maximization

Initialize spans and clusters
- M step:
    - Re-estimate prototypes
- E step:
    - Assign cluster and align

# Training



Expectation-Maximization

Initialize spans and clusters
- M step:
    - Re-estimate prototypes
- E step:
    - Assign cluster and align
    - We restrict the search space:
        - voice activity detection
        - phone boundary detection [Khanaga et al.]

# Experiments

| Language Pair | Dataset | Number of utterances |
|---|---|---|
| Griko - Italian | [Lekakou et al] | 330 |
| Spanish - English | CALLHOME (sample) | 2k |
| | CALLHOME (all) | 17k |
| | Fisher | 143k |

# Baselines

# Baselines

- Naive:
    - frames/word ~ #characters
    - along the diagonal

# Baselines

- Naive:
  - frames/word ~ #characters
  - along the diagonal

*Austin is great*

# Baselines

*Austin is great*

- Naive:
  - frames/word ~ #characters
  - along the diagonal

- Neural [Duong et al]:
  - DNN optimised for direct translation of speech
  - convert attention mechanism weights to alignments

# Results

Alignment F-score

# Results

## Alignment F-score



naive     neural     ours

F-score axis: 0, 10, 20, 30, 40, 50, 60

Griko-Italian    Callhome (2k)    Callhome (17k)    Fisher (143k)

# Results

## Alignment F-score

# Results

## Alignment F-score



| | naive | neural | ours |
|---|---|---|---|
| Griko-Italian | 46.7 | 27.0 | |
| Callhome (2k) | 35.8 | 26.4 | |
| Callhome (17k) | 35.7 | 29.1 | |
| Fisher (143k) | 27.8 | 26.2 | |

# Results



Alignment F-score

# Results

Alignment Precision

# Results

## Alignment Precision



naive ■    neural ■    ours ■

Precision
60
50
40
30
20
10
0

Griko-Italian    Callhome (2k)    Callhome (17k)    Fisher (143k)
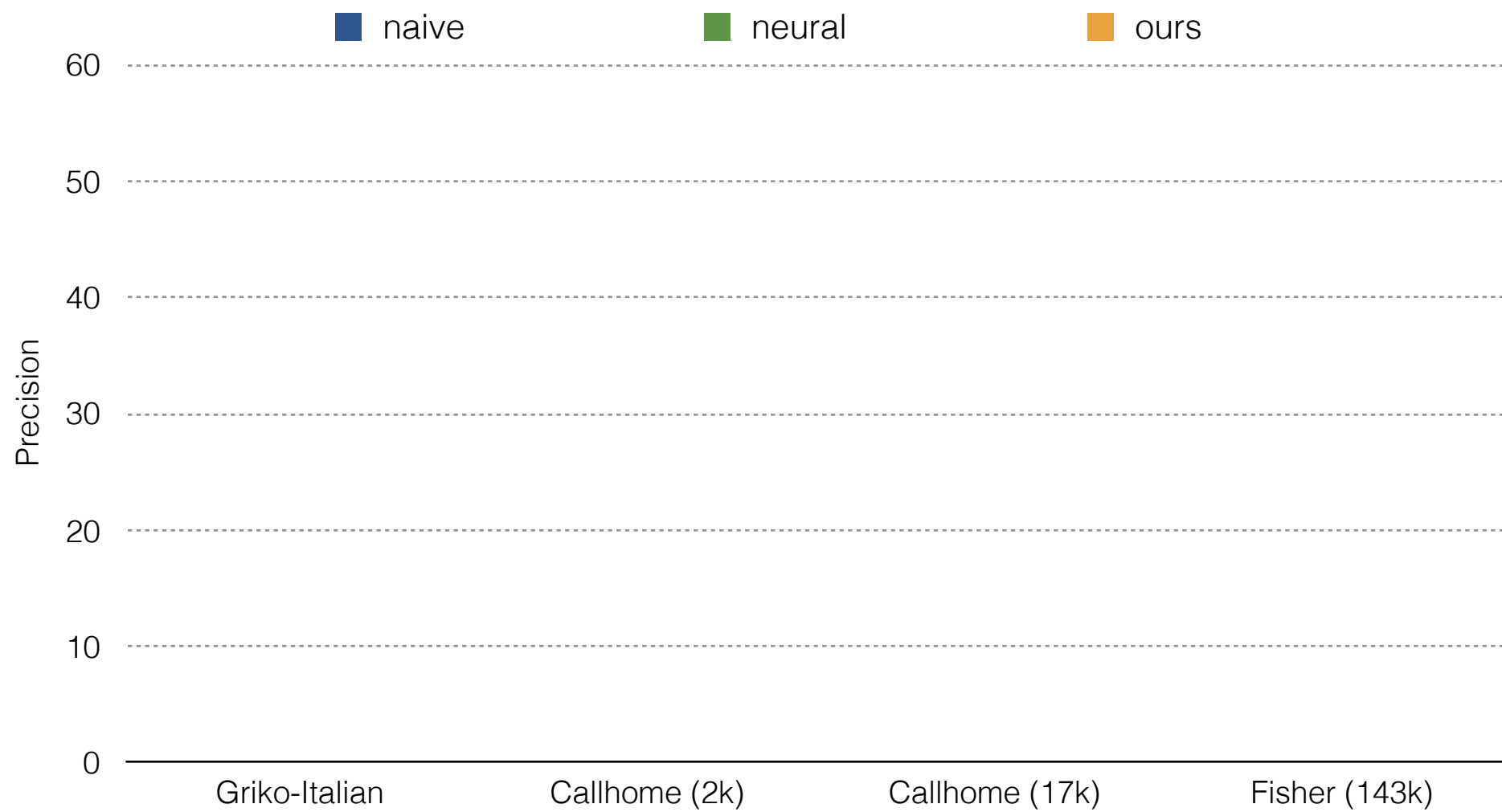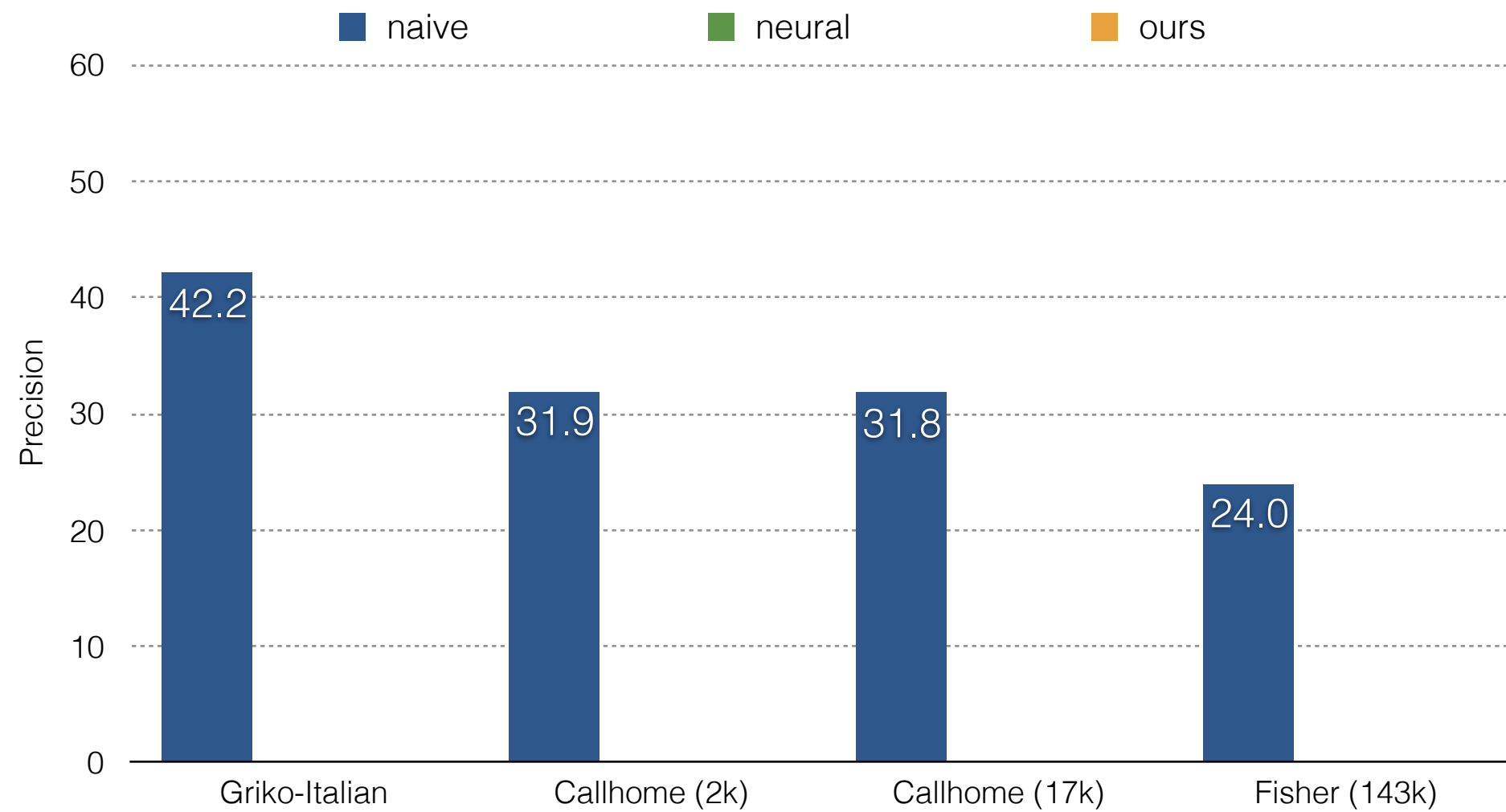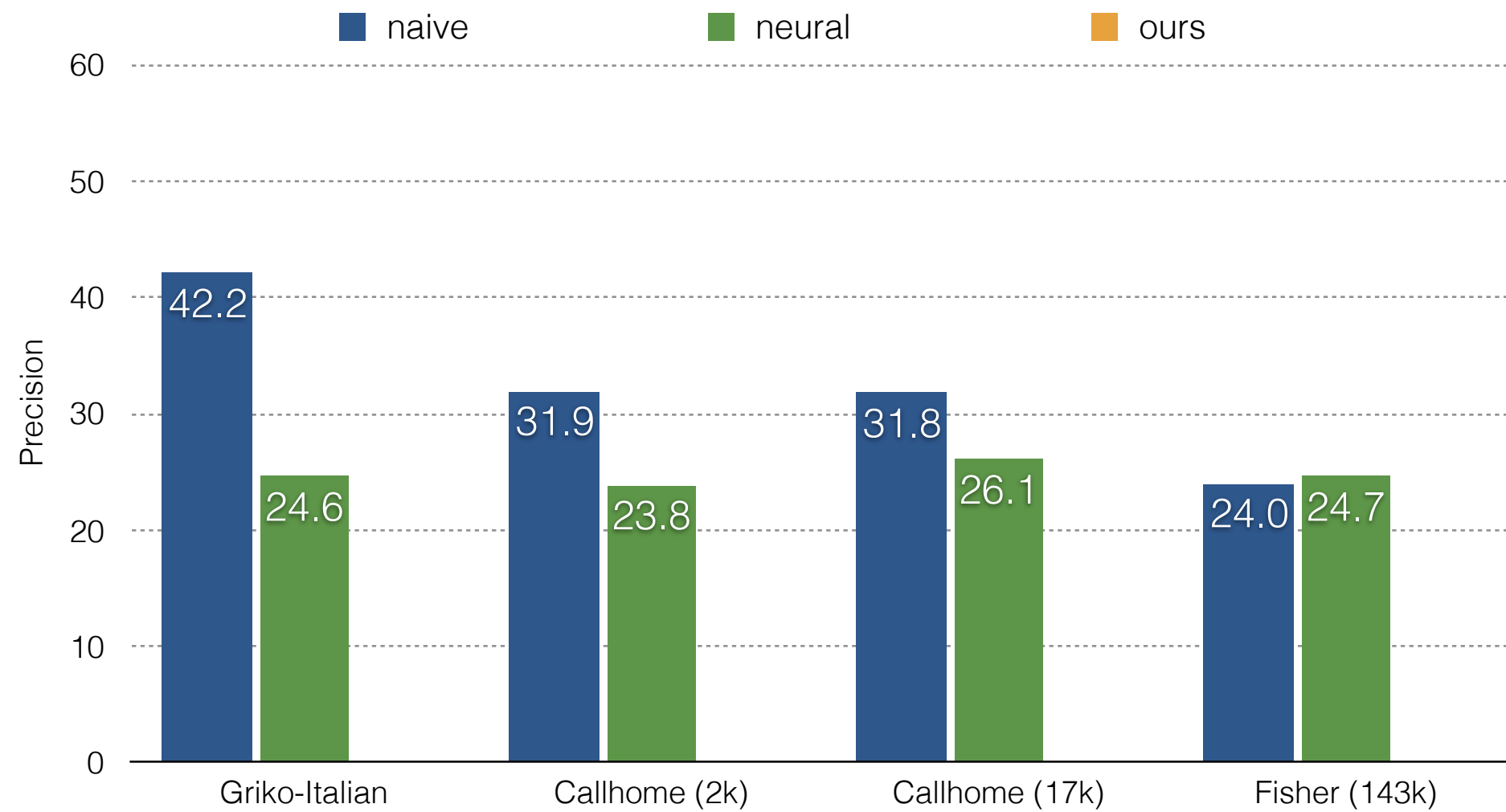
# Results



Alignment Precision

# Results



Alignment Precision

# Results

## Alignment Precision



Legend: ■ naive  ■ neural  ■ ours

| | naive | neural | ours |
|---|---|---|---|
| Griko-Italian | 42.2 | 24.6 | 56.6 |
| Callhome (2k) | 31.9 | 23.8 | 38.8 |
| Callhome (17k) | 31.8 | 26.1 | 38.4 |
| Fisher (143k) | 24.0 | 24.7 | 33.3 |

# Results



Word-level F-score

# Example



| | | | | |
|---|---|---|---|---|
| **Griko:** | ìcha | na | aforàso | to tsomì |
| **Gold:** | dovevo | | comprare | il pane |
| **Ours:** | dovevo | comprare | il | pane — 82.3 |
| **Naive:** | dovevo | comprare | il | pane — 72.4 |
| **Attention:** | dovevo | comprare il | | pane — 38.3 |

F-score

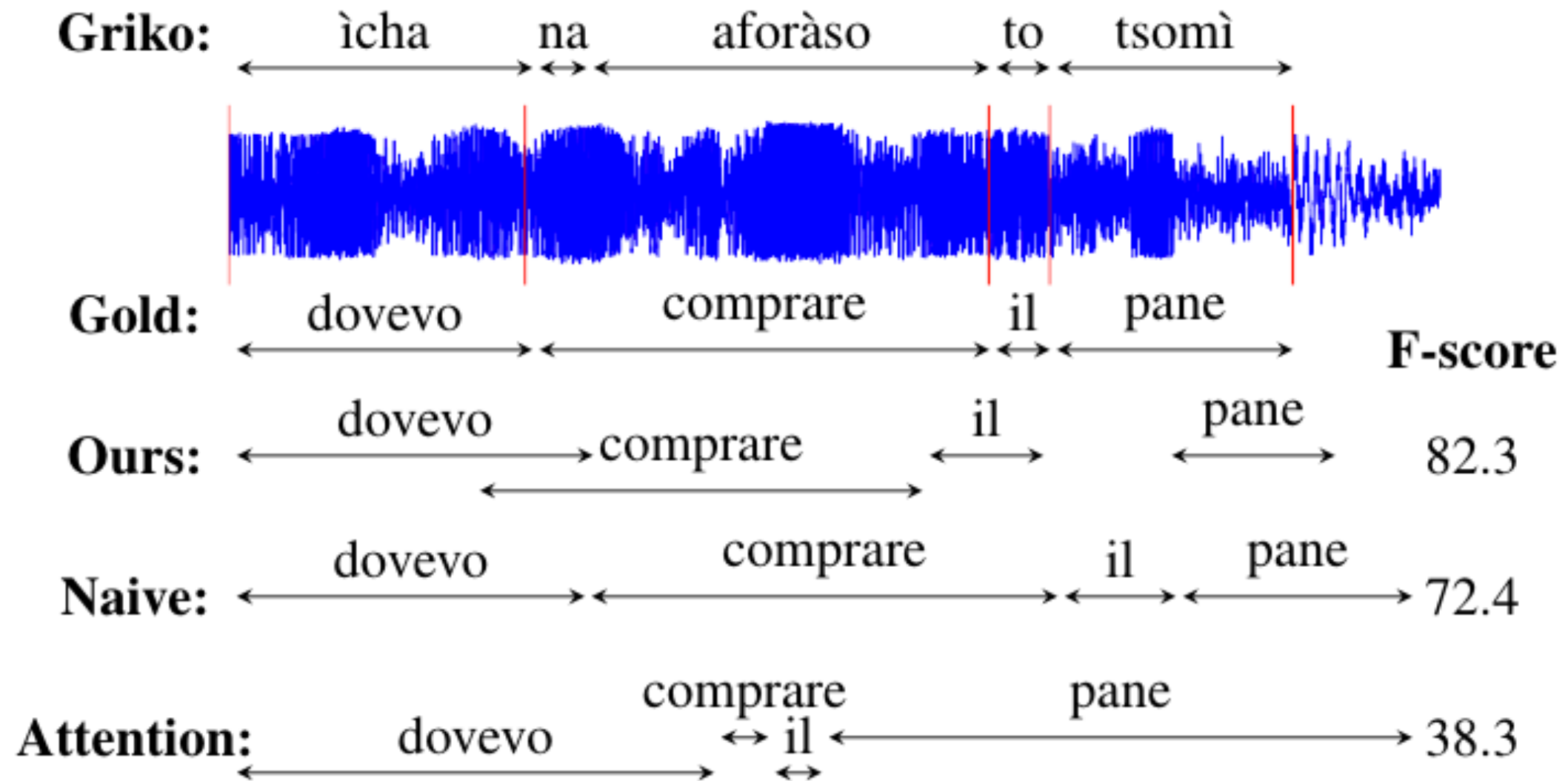# Conclusion

Alignment model

Extension of IBM-2 with fast-align for speech-to-translation

k-means clustering with DTW and DBA

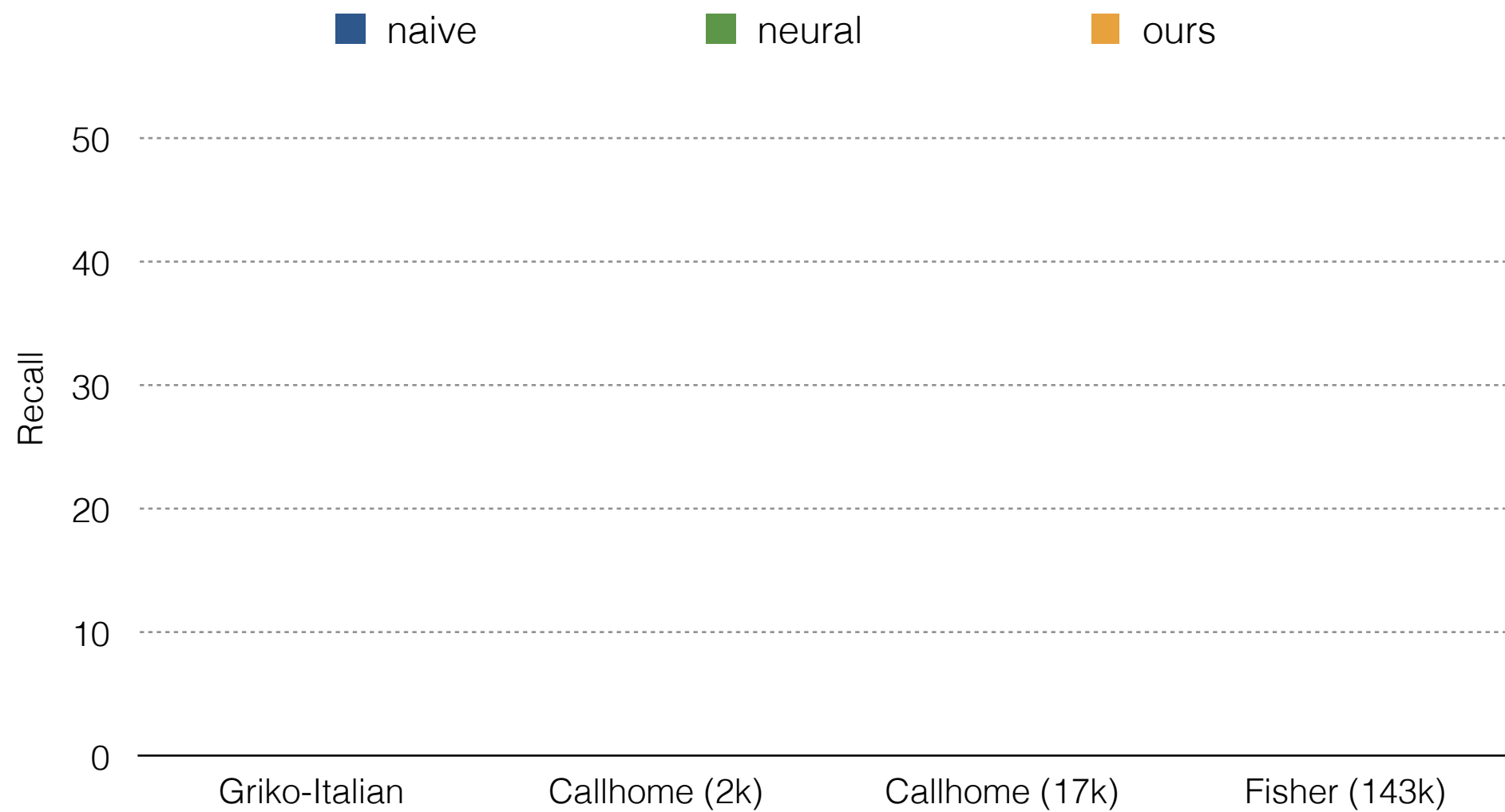Improvements in F-score and particularly Precision

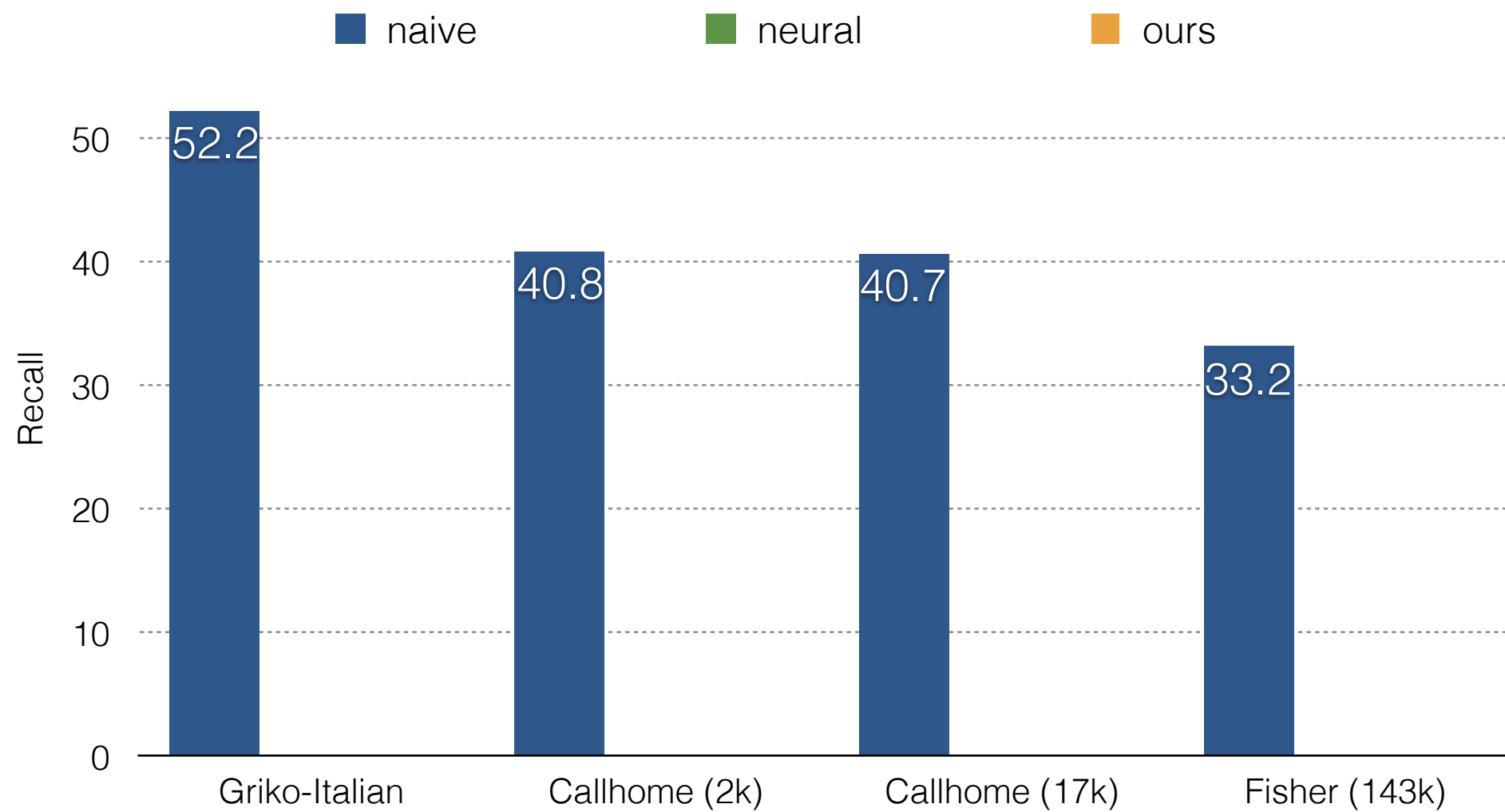**https://bitbucket.org/ndnlp/speech2translation**

# Results

Alignment Recall

# Results

## Alignment Recall

■ naive    ■ neural    ■ ours

Recall

50
40
30
20
10
0

Griko-Italian    Callhome (2k)    Callhome (17k)    Fisher (143k)

# Results

## Alignment Recall



Legend: ■ naive  ■ neural  ■ ours

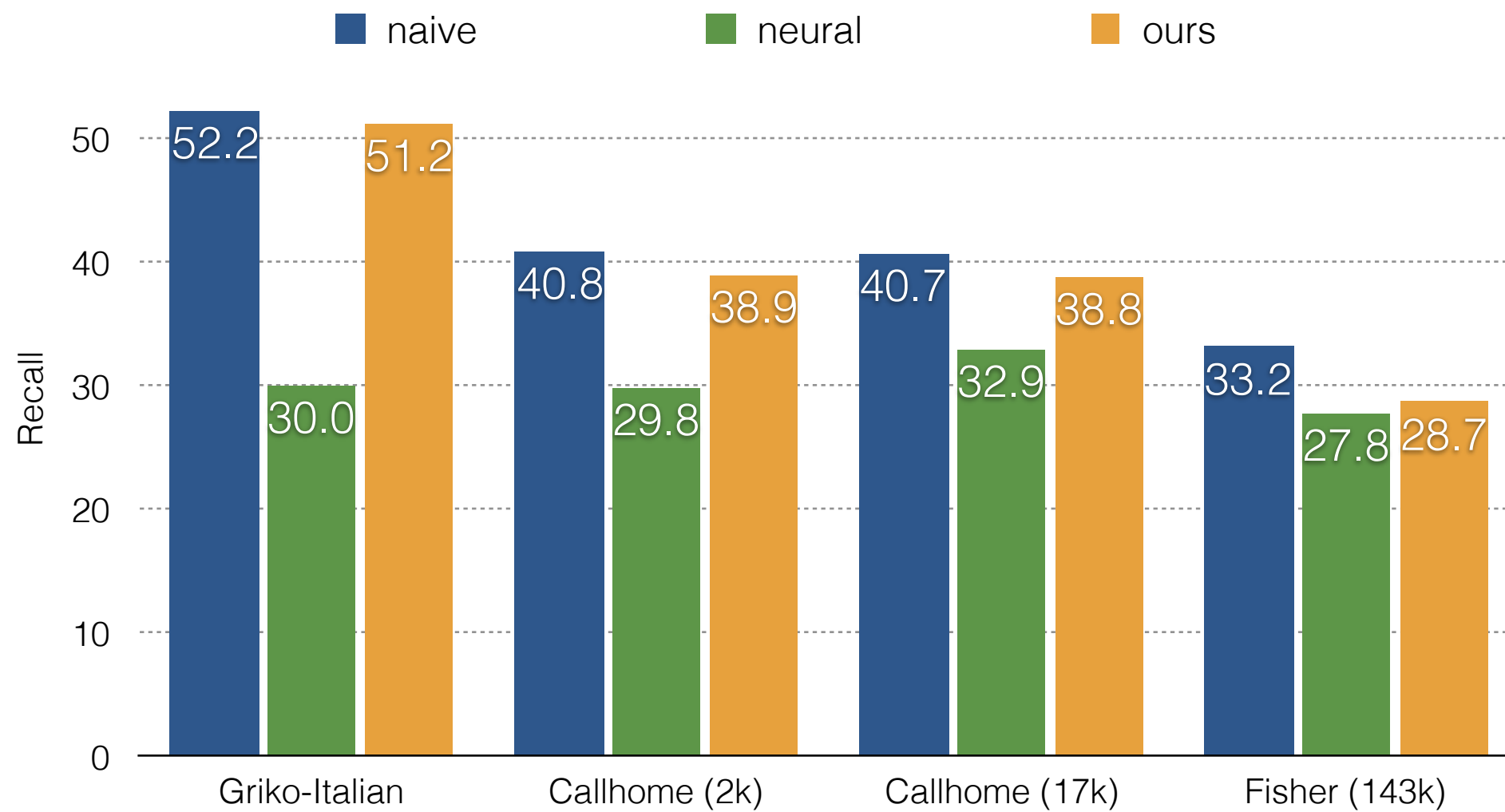| | Griko-Italian | Callhome (2k) | Callhome (17k) | Fisher (143k) |
|---|---|---|---|---|
| naive | 52.2 | 40.8 | 40.7 | 33.2 |

# Results

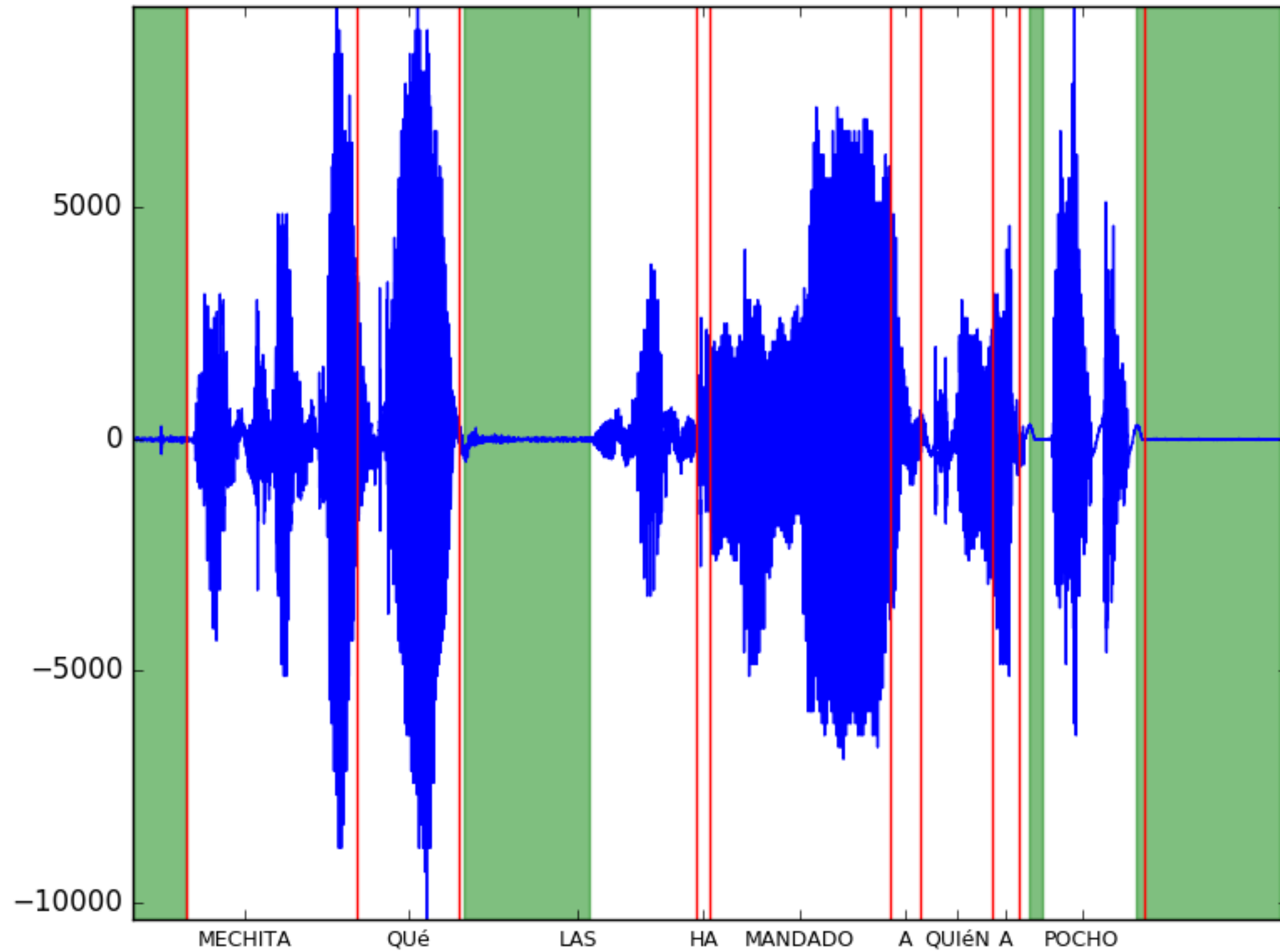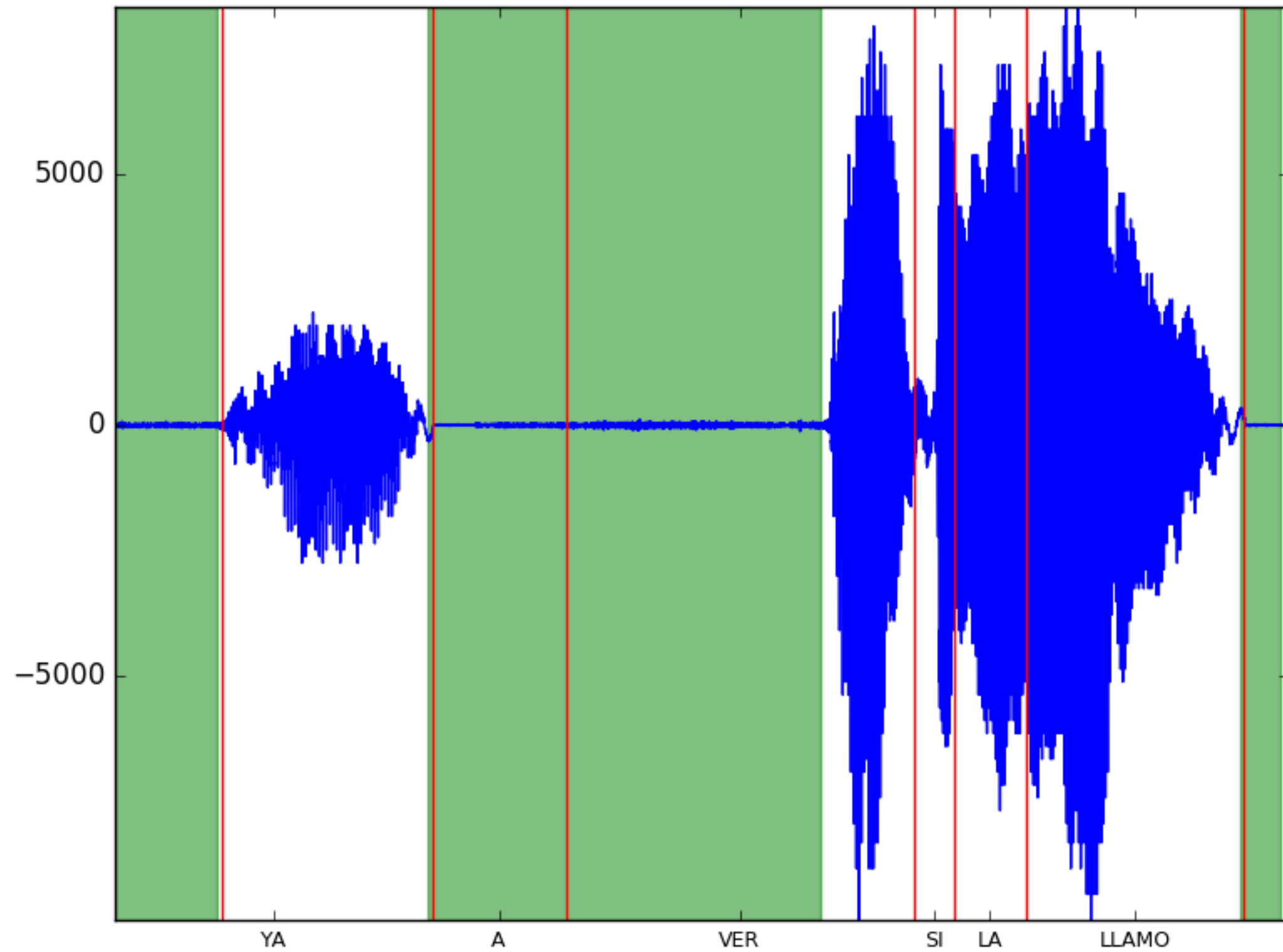## Alignment Recall

# Results

## Alignment Recall

# Example

# Example

# Proper model

Deficient:
$$p(\mathbf{e}, \mathbf{a}, \mathbf{b}, \mathbf{f} \mid \boldsymbol{\phi}) = p(l) \prod_{i=1}^{l} u(f_i) \times$$
$$s(a_i, b_i \mid f_i, \boldsymbol{\phi}) \times$$
$$\delta(a_i, b_i \mid i, l, |\boldsymbol{\phi}|) \times$$
$$t(e_i \mid f_i).$$

# Proper model

Proper:
$$p(\mathbf{e}, \mathbf{a}, \mathbf{b}, \mathbf{f} \mid \boldsymbol{\phi}) = p(l) \prod_{i=1}^{l} \delta(a_i, b_i \mid i, l, |\boldsymbol{\phi}|) \times$$
$$s(f_i \mid a_i, b_i, \boldsymbol{\phi}) \times$$
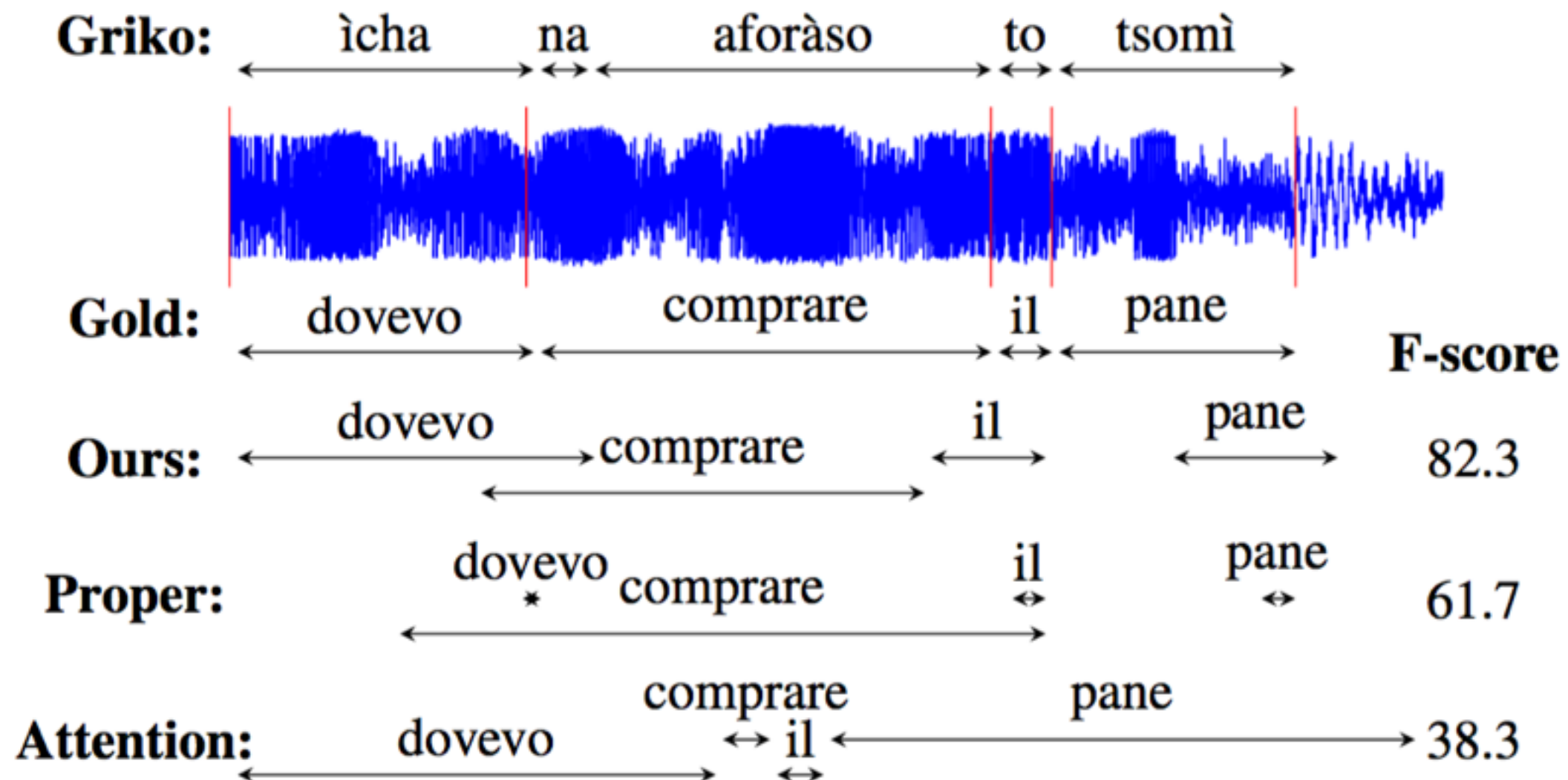$$t(e_i \mid f_i).$$

# Proper model

Proper:
$$p(\mathbf{e}, \mathbf{a}, \mathbf{b}, \mathbf{f} \mid \boldsymbol{\phi}) = p(l) \prod_{i=1}^{l} \delta(a_i, b_i \mid i, l, |\boldsymbol{\phi}|) \times$$
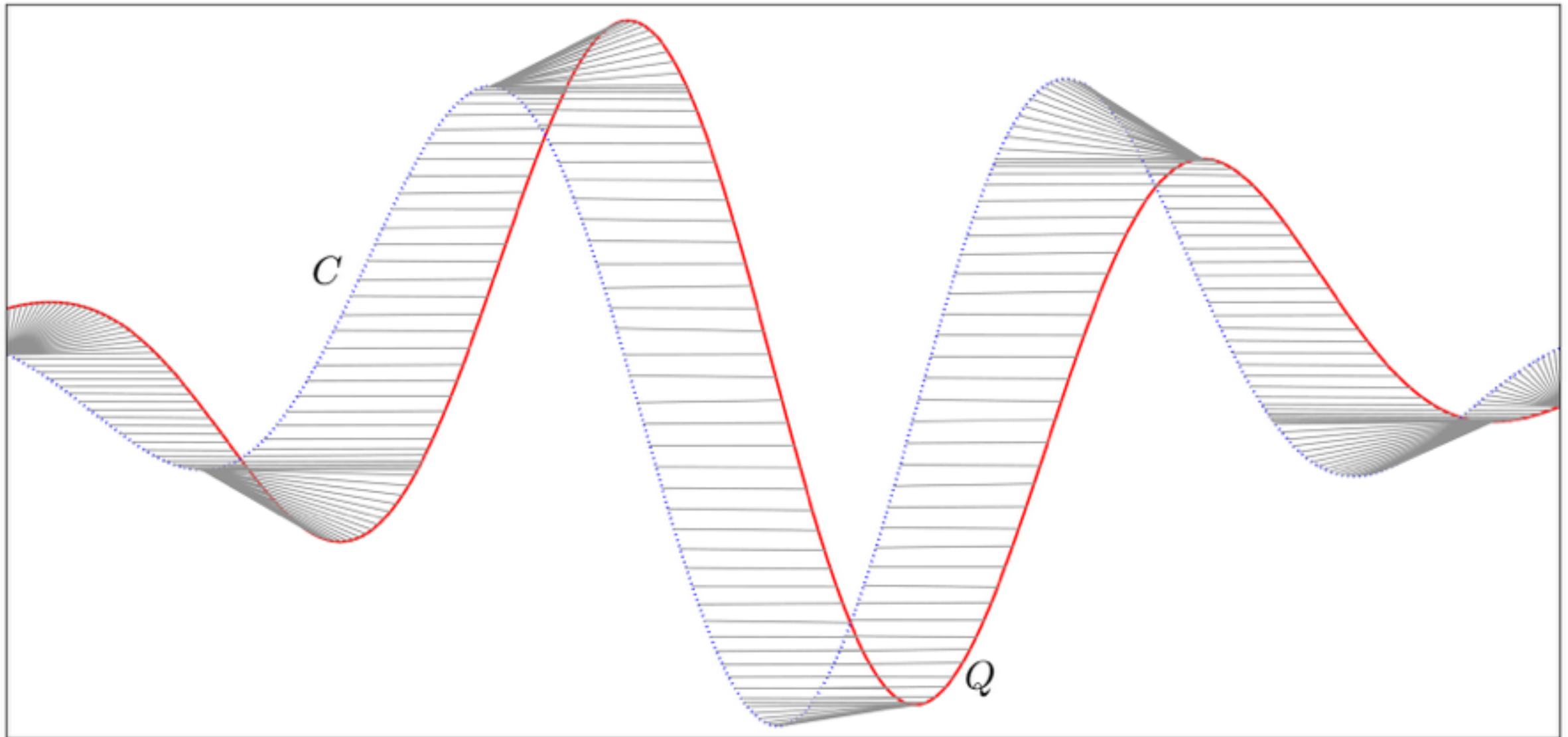$$s(f_i \mid a_i, b_i, \boldsymbol{\phi}) \times$$
$$t(e_i \mid f_i).$$

The proper model performs much worse.
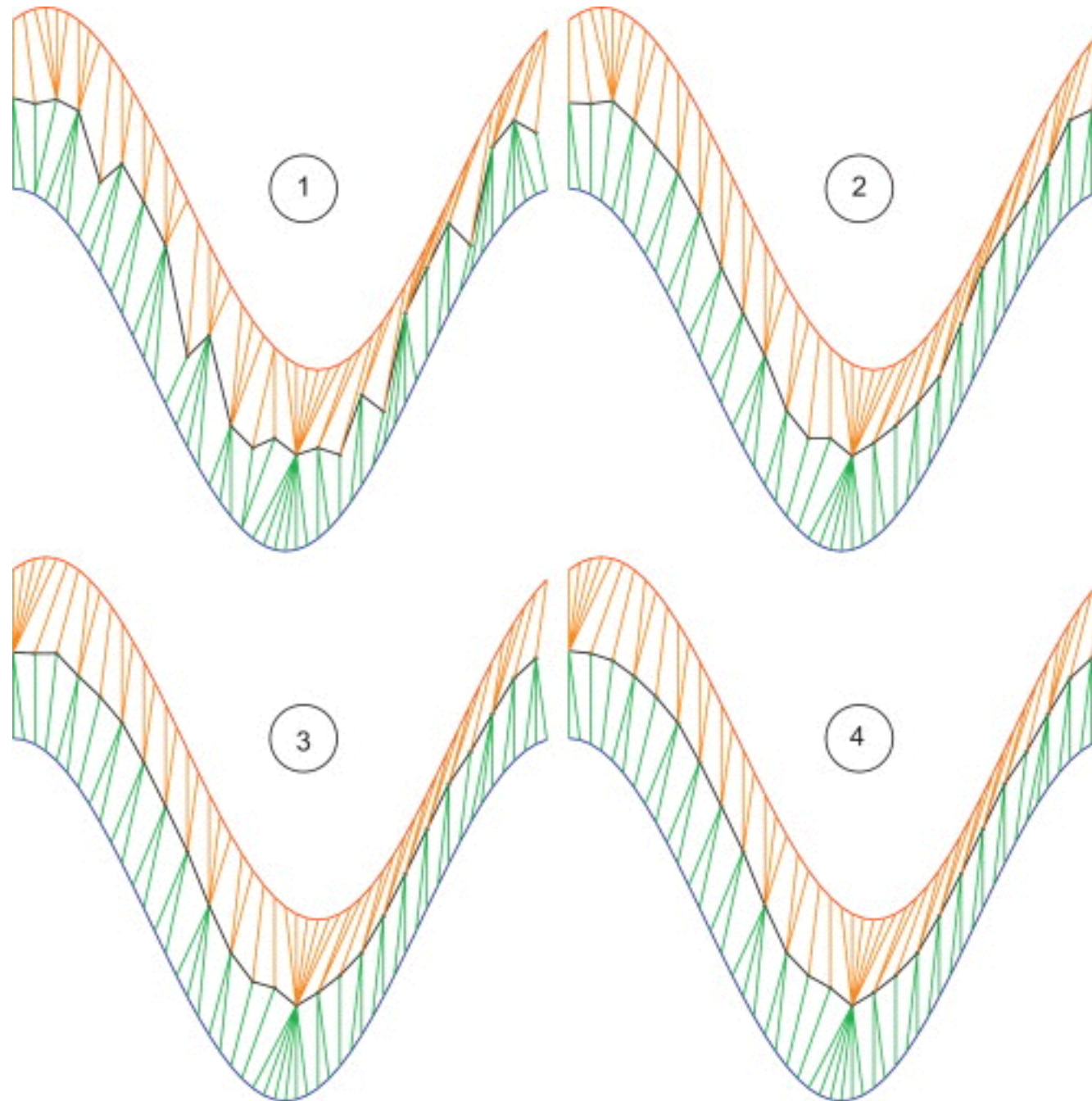-It favours too long or too short spans

# Example

# Background: DTW and DBA

Dynamic Time Warping (DTW)

# Background: DTW and DBA
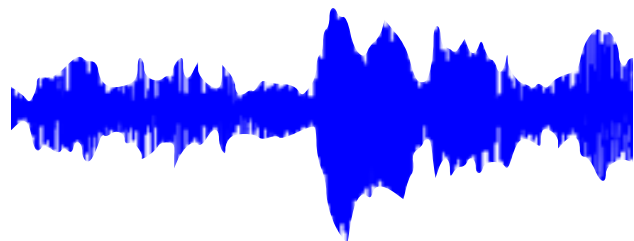
DTW Barycenter Averaging (DBA)

# M-step:

Prototype estimation with DTW Barycenter Averaging

# M-step:

Prototype estimation with DTW Barycenter Averaging

# M-step:

Prototype estimation with DTW Barycenter Averaging