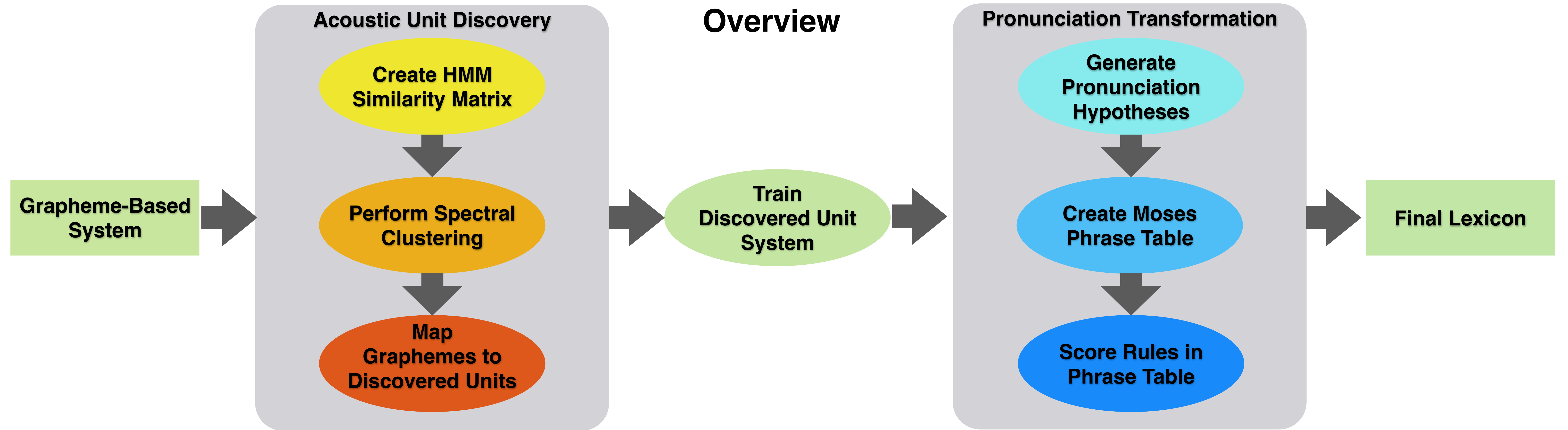# Efficient Rule Scoring for Improved Grapheme-Based Lexicons

## William Hartmann, Lori Lamel, and Jean-Luc Gauvain

Spoken Language Processing Group, LIMSI-CNRS `{hartmann, lamel, gauvain}@limsi.fr`

## Introduction

- Unlike the other main components of an ASR system, the pronunciation lexicon is largely handmade.
- Low-resource languages may not have expert-defined lexicons.
- We propose a two-stage approach to learning both the lexicon and the underlying acoustic units.
- Our approach relies on an initial baseline grapheme-based system.
- Acoustic units are learned by clustering the context-dependent grapheme-based models.
- Pronunciations are generated by transforming the original lexicon with an SMT-based approach.
- Each individual stage produces a significant improvement over the baseline system.
- Combined, the approach reduces the relative word error rate by 16%.

## Overview

**Acoustic Unit Discovery**

- Create HMM Similarity Matrix
- Perform Spectral Clustering
- Map Graphemes to Discovered Units

Grapheme-Based System → Acoustic Unit Discovery → Train Discovered Unit System → Pronunciation Transformation → Final Lexicon

**Pronunciation Transformation**

- Generate Pronunciation Hypotheses
- Create Moses Phrase Table
- Score Rules in Phrase Table

## Acoustic Unit Discovery

- Acoustic units are discovered by clustering context-dependent grapheme-based HMMs.
- Requires defining a similarity measure between individual HMMs (Equation 1).
- CSD is the Cauchy-Schwarz Divergence measure (Equation 2).
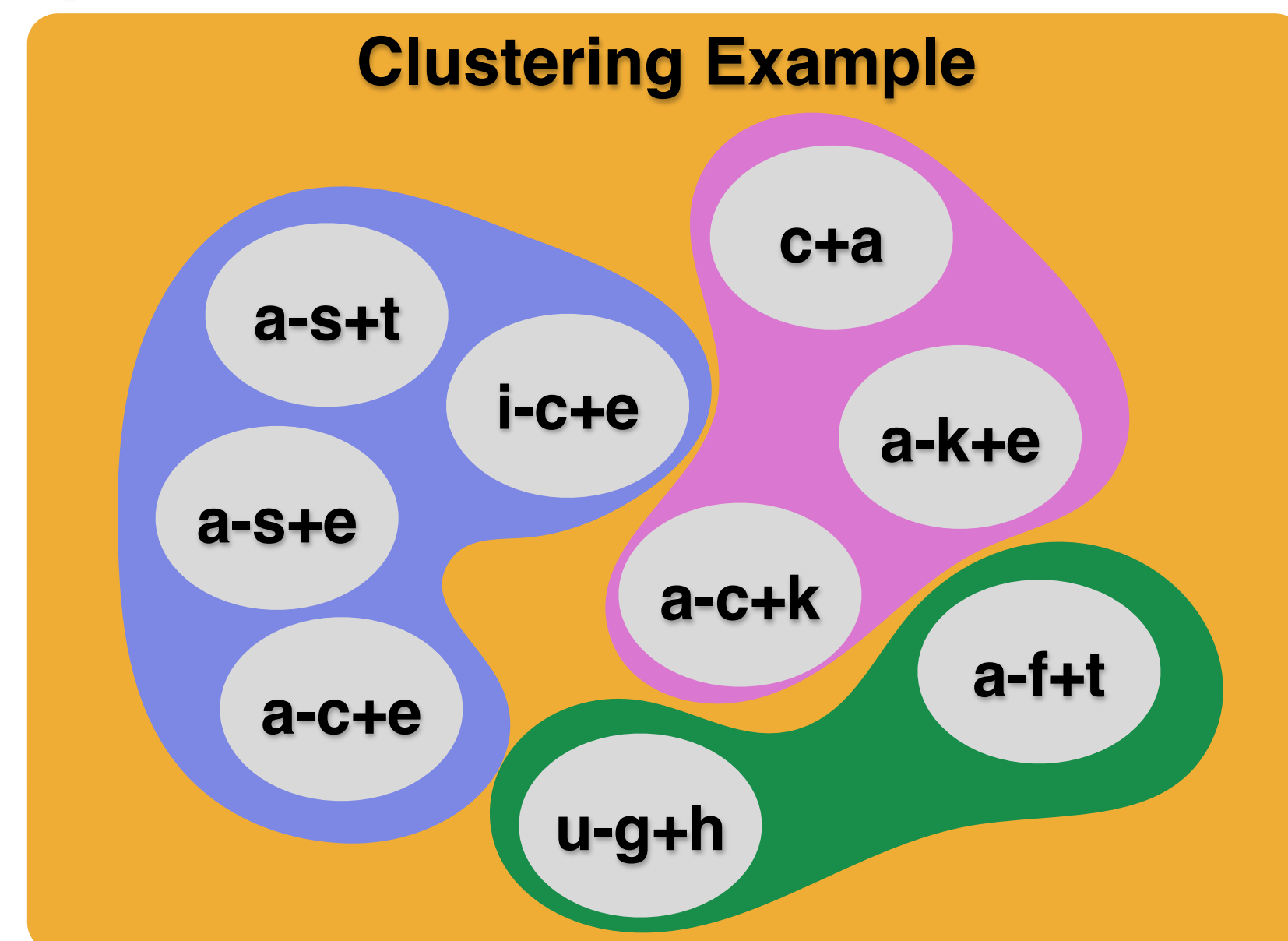- We use the CSD because a closed form solution for a Mixture of Gaussians exists.

$$\mathrm{HMM}_{\mathrm{sim}}(\mathbf{h}, \mathbf{h}') = \sum_{a=1}^{A} \sum_{b=1}^{B} \frac{\alpha_{a,b}}{\mathrm{CSD}(h_a, h_b') + 1}$$
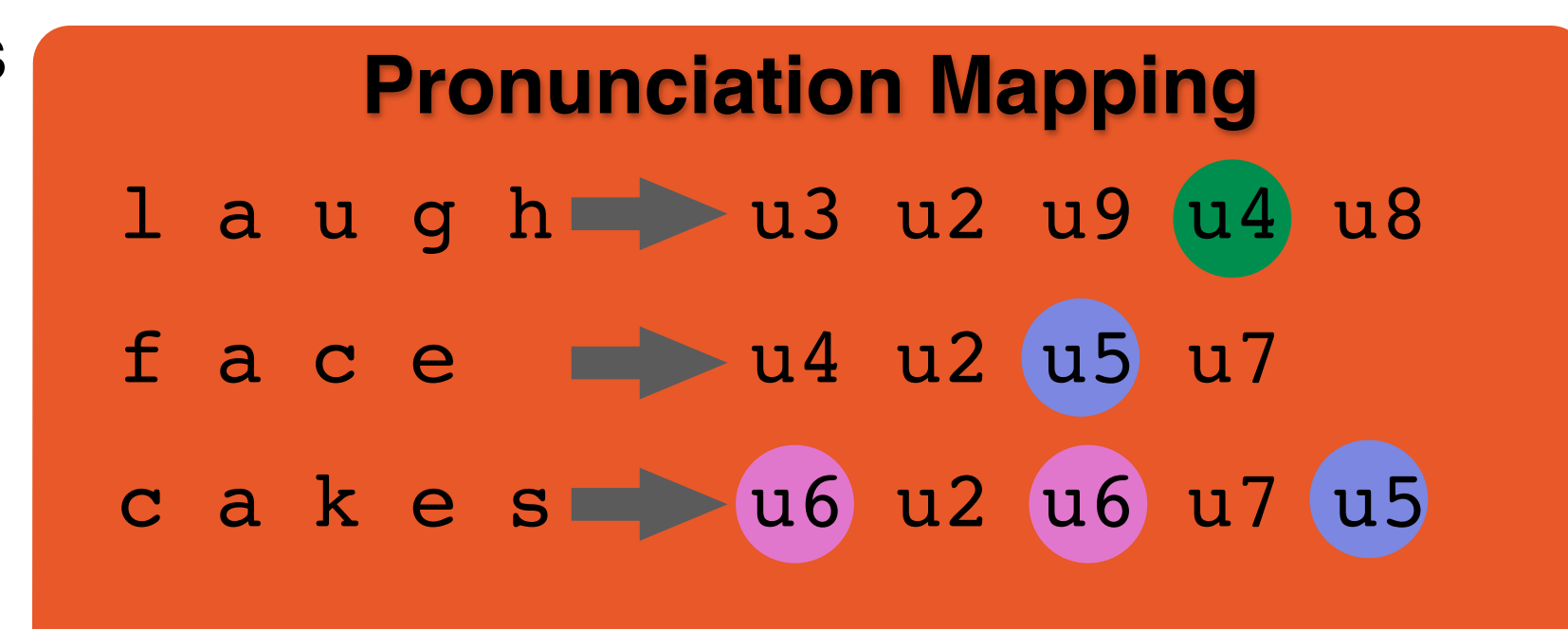
**Equation 1**

$$\mathrm{CSD}(\mathbf{p}, \mathbf{q}) = -\log \frac{\sum_i p_i q_i}{\sqrt{\sum_i p_i^2 \sum_i q_i^2}}$$

**Equation 2**

- Clustering is performed using spectral clustering.
- We achieved better performance with a k-nearest neighbor similarity graph rather than a similarity matrix.
- Since the optimal number of acoustic units is not known a priori, we tried various numbers of clusters.
- The final clusters group acoustically similar context-dependent HMMs into a single acoustic unit.

**Clustering Example**

a-s+t, i-c+e, c+a, a-k+e, a-s+e, a-c+k, a-c+e, a-f+t, u-g+h

- Based on the clustering, pronunciations are mapped to the new acoustic units.
- Each pronunciation will have the same number of units as in the baseline grapheme-based lexicon.
- The new acoustic units are labeled as numbers since no other label exists.

**Pronunciation Mapping**

```
l a u g h  →  u3 u2 u9 u4 u8
f a c e    →  u4 u2 u5 u7
c a k e s  →  u6 u2 u6 u7 u5
```

## Pronunciation Transformation

- Context-dependent acoustic models are trained.
- The training data is decoded in terms of the acoustic units.
- Based on the time-aligned results, each word in the training set has one or more pronunciation hypotheses.
- The example shown uses grapheme units for clarity.

```
lack       l a k
lack       l e k
lochs      l o c x s
necessary  n e s e s r y
ford       f r d
ford       f r n
caught     k o t
```

- Using Moses, a phrase table is learned from the pronunciation hypotheses.
- Each entry in the table represents the translation of a sequence of acoustic units into an alternate sequence.
- Moses allows for reordering of units, but we disable this option for our work.
- The phrase table can transform the pronunciations in the original lexicon.

| Source | Target | $p(t \mid s)$ |
|---|---|---|
| a c k | a k | 0.19 |
| c h s | c x s | 0.13 |
| c e s s a | s e s e | 0.36 |
| f o r d | f r d | 0.17 |
| a u g h t | o t | 0.25 |

- The initial translation table decreases performance, so we prune the table.
- Rules are rescored with a single pass.
- The score is based on how often a rule is used for the best pronunciation for a word during forced alignment.
- Only rules that surpass a certain threshold are kept.
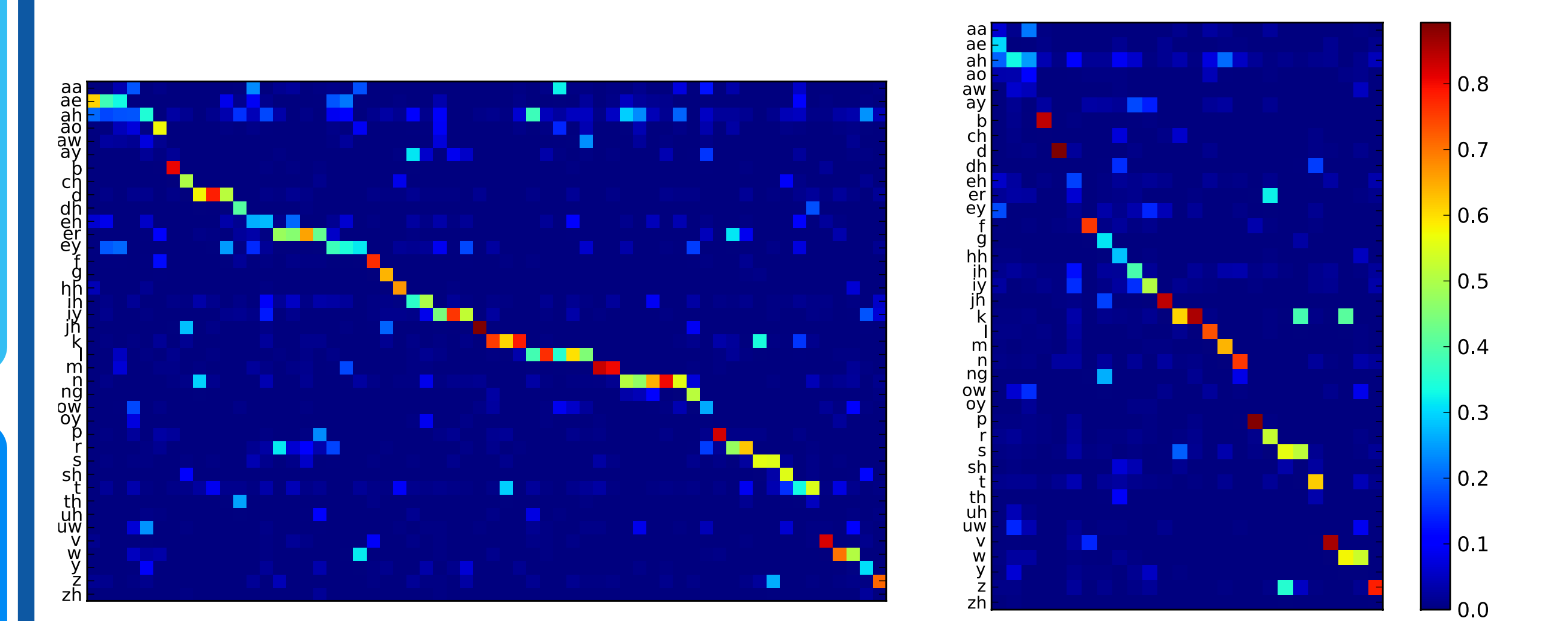- The final transformation works for words unseen during training.

| Source | Target | Score |
|---|---|---|
| a c k | a k | 0.65 |
| c h s | c x s | 0.48 |
| c e s s a | s e s e | 0.51 |
| f o r d | f r d | 0.03 |
| a u g h t | o t | 0.08 |

## Results

- WER results are presented on the WSJ0 5k word task.
- Acoustic models are trained with HTK.
- All recognizers use context-dependent models with 2000 tied states.

| Unit Type | # Units | Direct | Trans. |
|---|---|---|---|
| Grapheme | 26 | 15.8 | 14.2 |
| Discovered | 39 | 15.0 | 13.3 |
| Discovered | 50 | 15.2 | 14.1 |
| Discovered | 60 | **14.4** | **13.2** |

- Decoding is performed with a bigram LM.
- Direct refers to the lexicon after mapping, and Trans. refers to the lexicon after pronunciation transformation.
- The figures below demonstrate the correlation between the discovered units and phones, and the grapheme units and phones.



**Discovered 60**



**Grapheme**

- We have proposed a two-stage approach for acoustic unit discovery and pronunciation generation that reduces relative WER by 16% compared to a baseline grapheme-based system.
- We are currently working to apply these techniques to other languages.