# Multilingual acoustic word embedding models for processing zero-resource languages

ICASSP 2020

Herman Kamper[1], Yevgen Matusevych[2], Sharon Goldwater[2]

[1]Stellenbosch University, South Africa
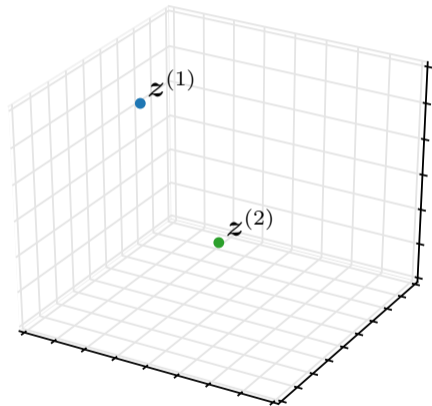[2]University of Edinburgh, UK

http://www.kamperh.com/

# Background: Why acoustic word embeddings?

- Current speech recognition methods require large labelled data sets

- *Zero-resource speech processing* aims to develop methods that can discover linguistic structure from unlabelled speech [Dunbar et al., ASRU'17]

- Example applications: Unsupervised term discovery, query-by-example

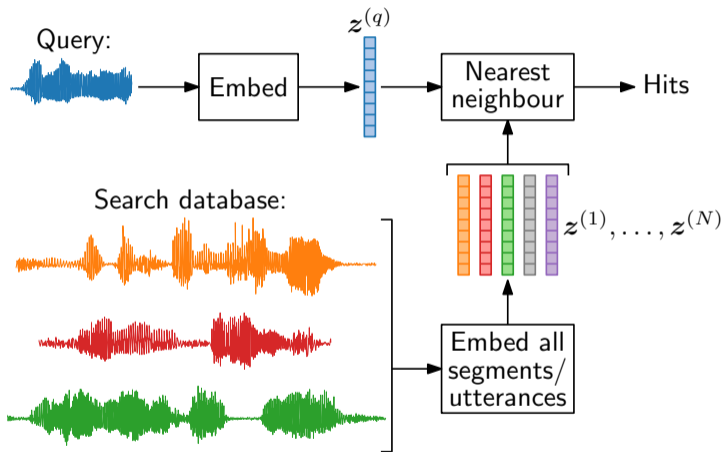- **Problem:** Need to compare speech segments of variable duration

# Acoustic word embeddings

# Example application: Query-by-example search



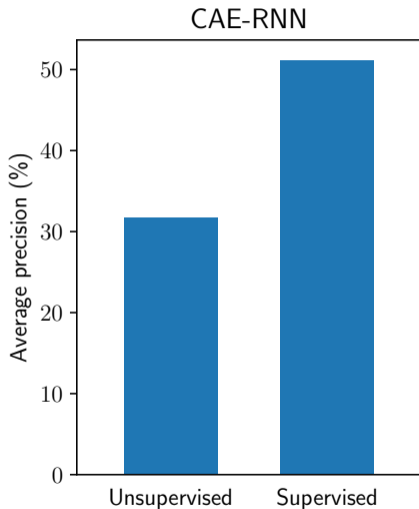[Levin et al., ICASSP'15]

# Supervised and unsupervised acoustic embeddings

- Growing body of work on acoustic word embeddings

- Supervised and unsupervised methods

- Unsupervised methods can be applied in zero-resource settings

- But there is still a large performance gap

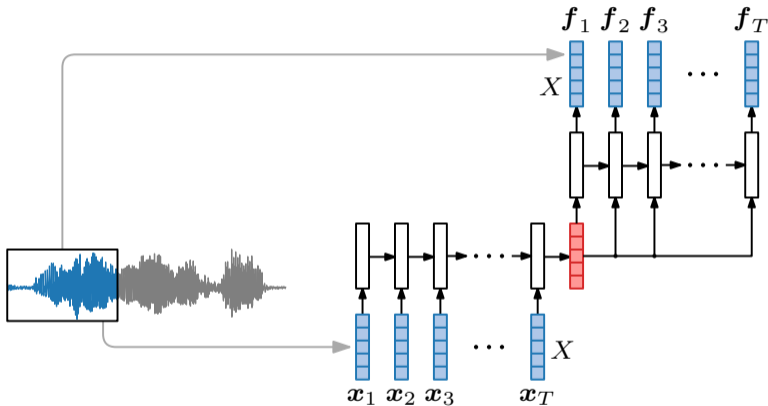# Supervised and unsupervised acoustic embeddings

- Growing body of work on acoustic word embeddings

- Supervised and unsupervised methods

- Unsupervised methods can be applied in zero-resource settings

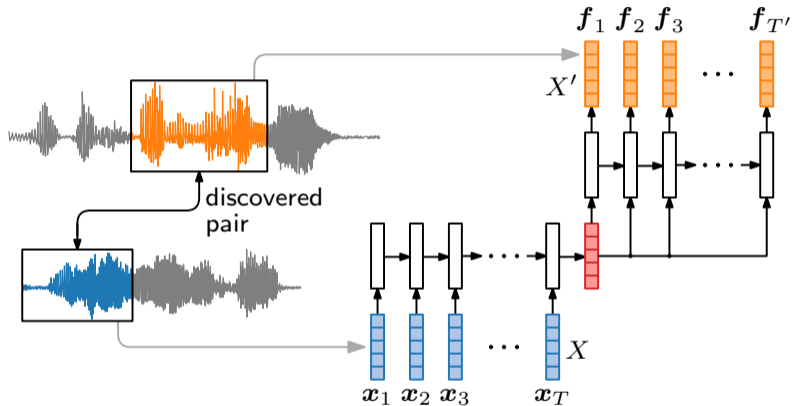- But there is still a large performance gap



CAE-RNN

[Kamper, ICASSP'19]

# Unsupervised monolingual acoustic word embeddings



[Chung et al., Interspeech'16; Kamper, ICASSP'19]

# Unsupervised monolingual acoustic word embeddings
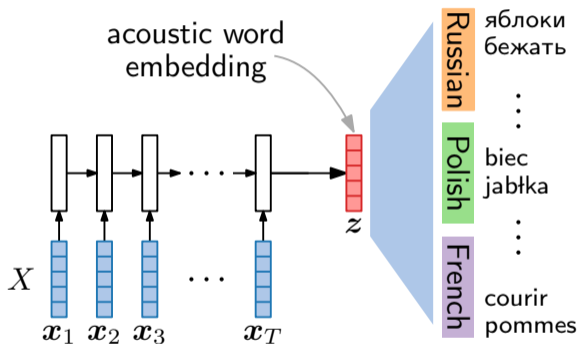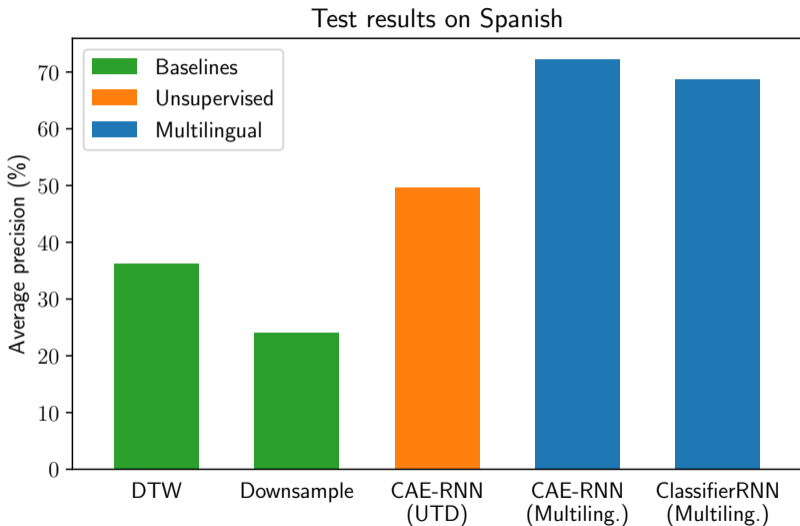


[Chung et al., Interspeech'16; Kamper, ICASSP'19]

# Supervised multilingual acoustic word embeddings
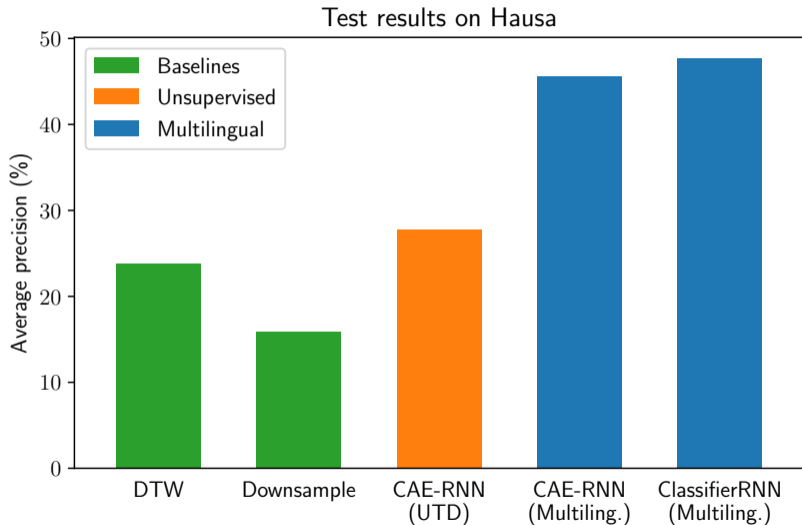


acoustic word embedding

# Experimental setup

- **Training data:** Six well-resourced languages
  Czech (CS), French (FR), Polish (PL), Portuguese (PT), Russian (RU), Thai (TH)

- **Test data:** Six languages treated as zero-resource
  Spanish (ES), Hausa (HA), Croatian (HR), Swedish (SV), Turkish (TR), Mandarin (ZH)

- **Evaluation:** Same-different isolated word discrimination

- **Embeddings:** $M = 130$ for all models

- **Baselines:**
  — Downsampling: 10 equally spaced MFCCs flattened
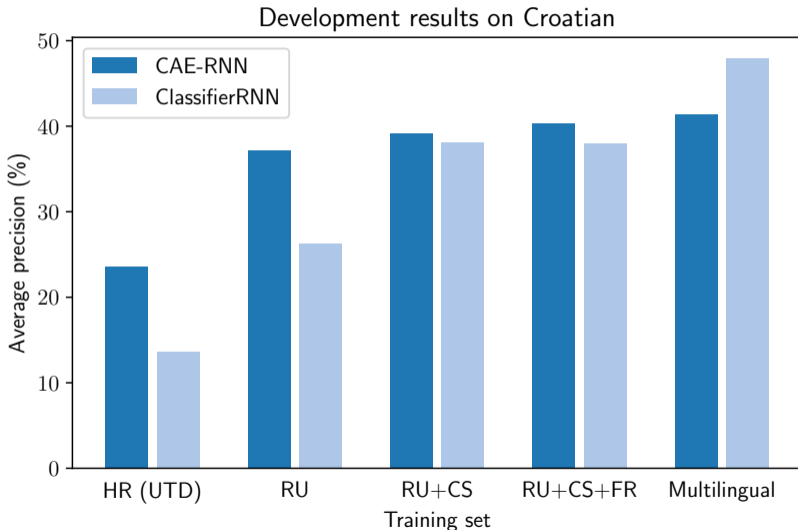  — Dynamic time warping (DTW) alignment cost between test segments

# 1. Is multilingual supervised > monolingual unsupervised?



Test results on Spanish

# 1. Is multilingual supervised $>$ monolingual unsupervised?



Test results on Hausa

# 2. Does training on more languages help?



Development results on Croatian

# 3. Is the choice of training language important?

# Conclusions and future work

Conclusions:

- Proposed to train a supervised multilingual acoustic word embedding model on well-resourced languages and then apply to zero-resource languages

- Multilingual CAE-RNN and ClassifierRNN consistently outperform unsupervised models trained on zero-resource languages

# Conclusions and future work

Conclusions:

- Proposed to train a supervised multilingual acoustic word embedding model on well-resourced languages and then apply to zero-resource languages

- Multilingual CAE-RNN and ClassifierRNN consistently outperform unsupervised models trained on zero-resource languages

Future work:

- Different models both for multilingual and unsupervised training

- Analysis to understand the difference between CAE-RNN and ClassifierRNN

- Does language conditioning help during decoding?

https://arxiv.org/abs/2002.02109

https://github.com/kamperh/globalphone_awe