

UNIVERSITEIT iYUNIVESITH1 STELLENBOSCH **UNIVERSITY** 



## Background • Current speech recognition methods require large labelled data sets. • Zero-resource speech processing aims to develop methods that can discover linguistic structure directly from unlabelled speech. • **Problem:** Need to compare speech segments of variable duration. • Dynamic time warping (DTW) is one option: Can be slow. • Acoustic word embedding methods map variable-length segments into fixed-dimensional space to enable efficient comparisons: Acoustic word embeddings $\boldsymbol{z} \in \mathbb{R}^{M}$ $X^{(1)}$ $oldsymbol{z}^{(1)}$ $X^{(2)}$

### **Example application:** Query-by-example

Find utterances in a speech collection containing a given spoken query:





IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Brighton, UK, 2019

# Truly unsupervised acoustic word embeddings using weak top-down constraints in encoder-decoder models

Herman Kamper E&E Engineering, Stellenbosch University, South Africa

![](_page_0_Figure_11.jpeg)

![](_page_0_Figure_12.jpeg)

![](_page_0_Figure_14.jpeg)

![](_page_0_Picture_24.jpeg)