

# **Unsupervised acoustic unit discovery for speech synthesis using discrete latent-variable neural networks**

Interspeech 2019, Graz, Austria

Ryan Eloff, André Nortje, Benjamin van Niekerk, Avashna Govender,  
Leanne Nortje, Arnu Pretorius, Elan van Biljon, Ewald van der Westhuizen,  
Lisa van Staden, Herman Kamper

Stellenbosch University, South Africa & University of Edinburgh, UK

<https://github.com/kamperh/suzerospeech2019>

# Advances in speech recognition



# Advances in speech recognition



- **Addiction to text:** 2000 hours transcribed speech audio;  
~350M/560M words text [Xiong et al., TASLP'17]

# Advances in speech recognition



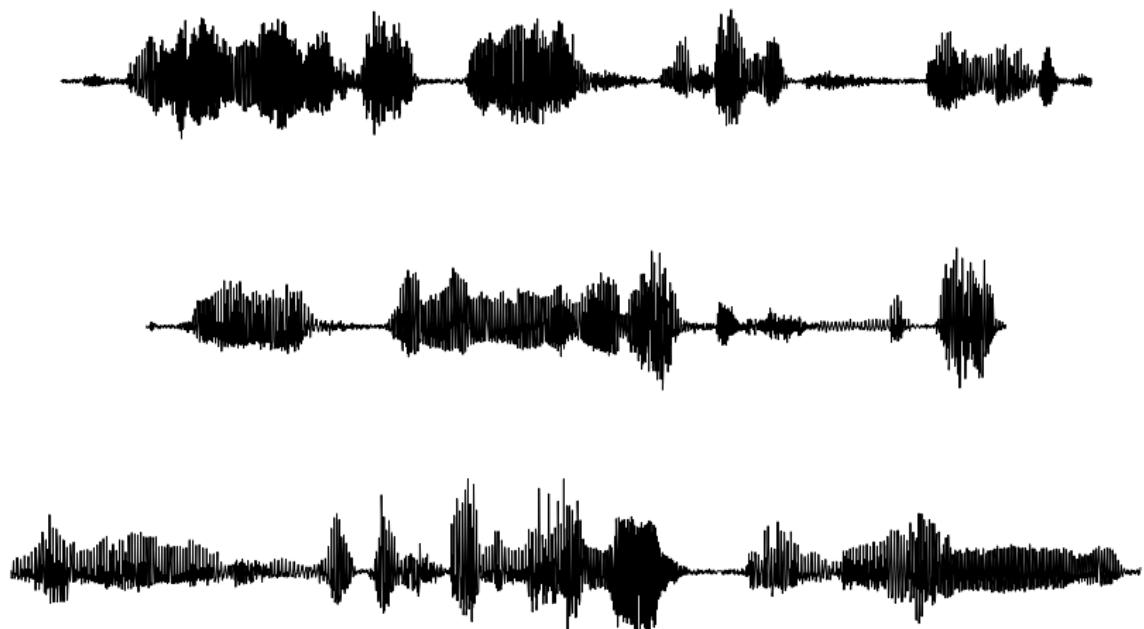
- **Addiction to text:** 2000 hours transcribed speech audio;  
~350M/560M words text [Xiong et al., TASLP'17]
- Sometimes not possible, e.g., for unwritten languages

# Advances in speech recognition

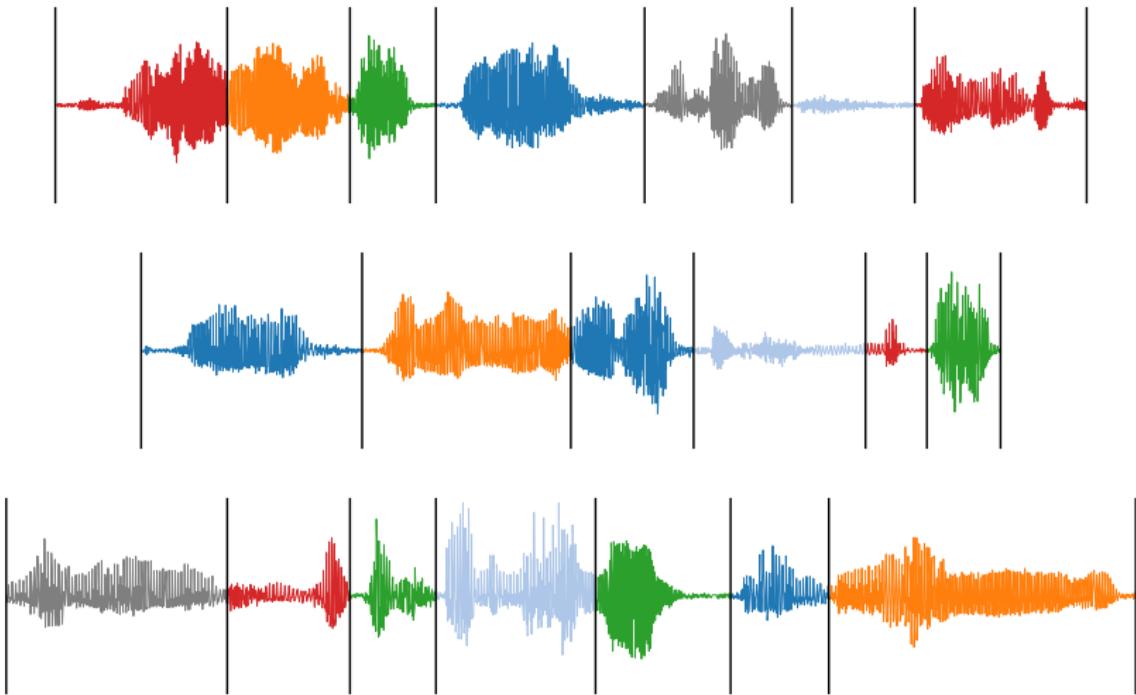


- **Addiction to text:** 2000 hours transcribed speech audio;  
~350M/560M words text [Xiong et al., TASLP'17]
- Sometimes not possible, e.g., for unwritten languages
- Very different from the way human infants learn language

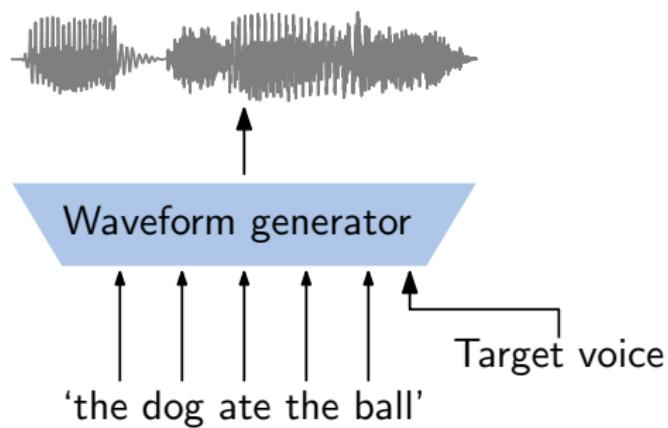
# Zero-Resource Speech Challenges (ZRSC)



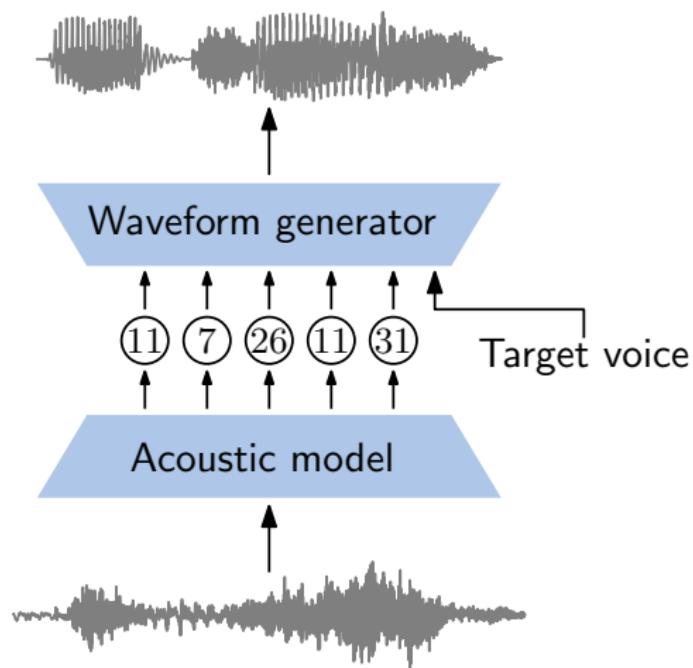
# Zero-Resource Speech Challenges (ZRSC)



# ZRSC 2019: Text-to-speech without text



# ZRSC 2019: Text-to-speech without text



# What do we get for training?

# What do we get for training?

No labels

# What do we get for training?

No labels :)

# What do we get for training?

No labels :)

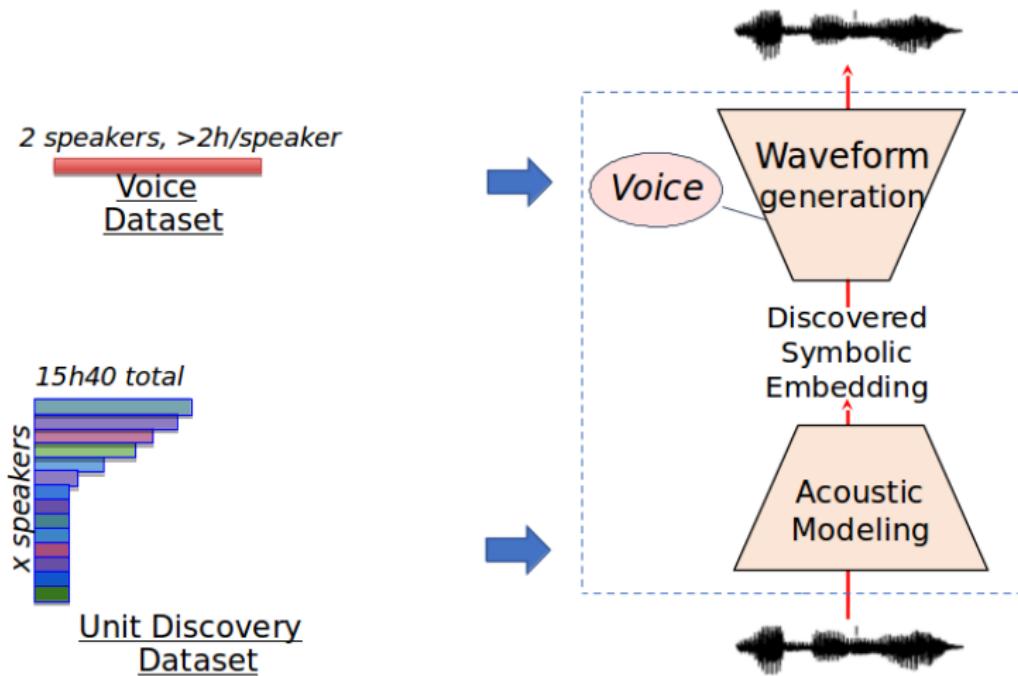
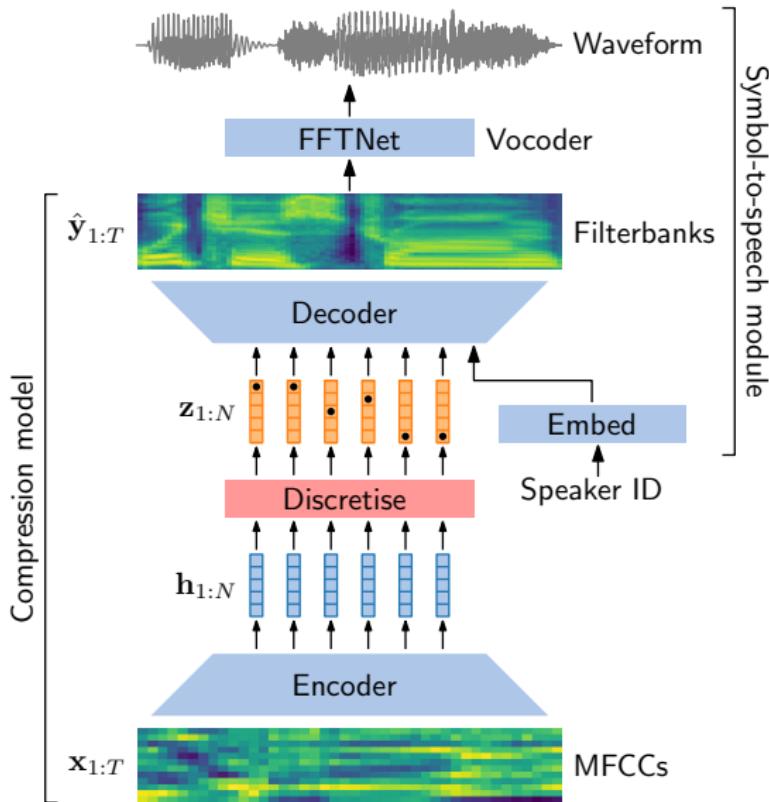
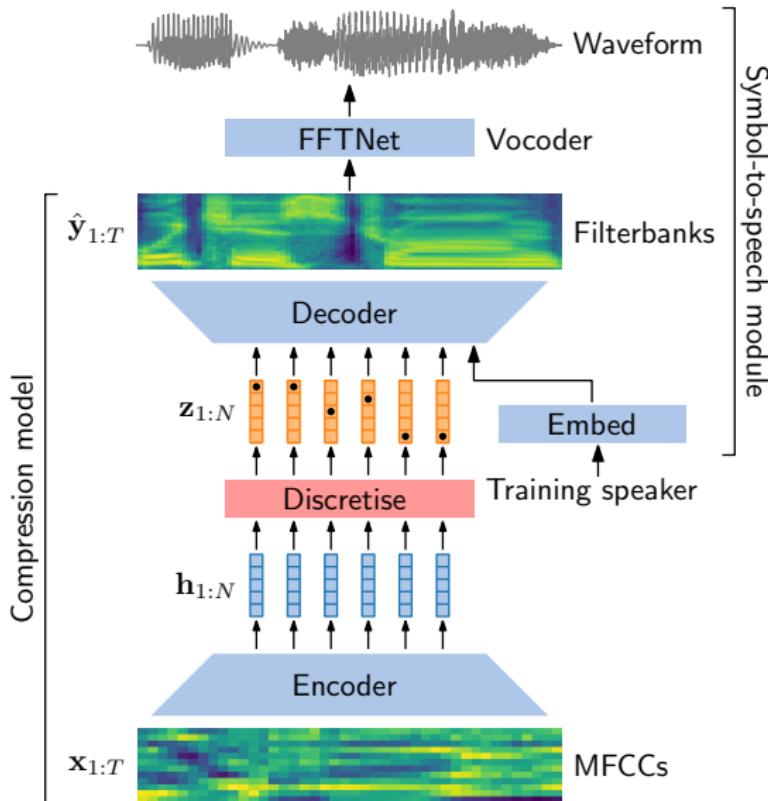


Figure adapted from: <http://zerospeech.com/2019>

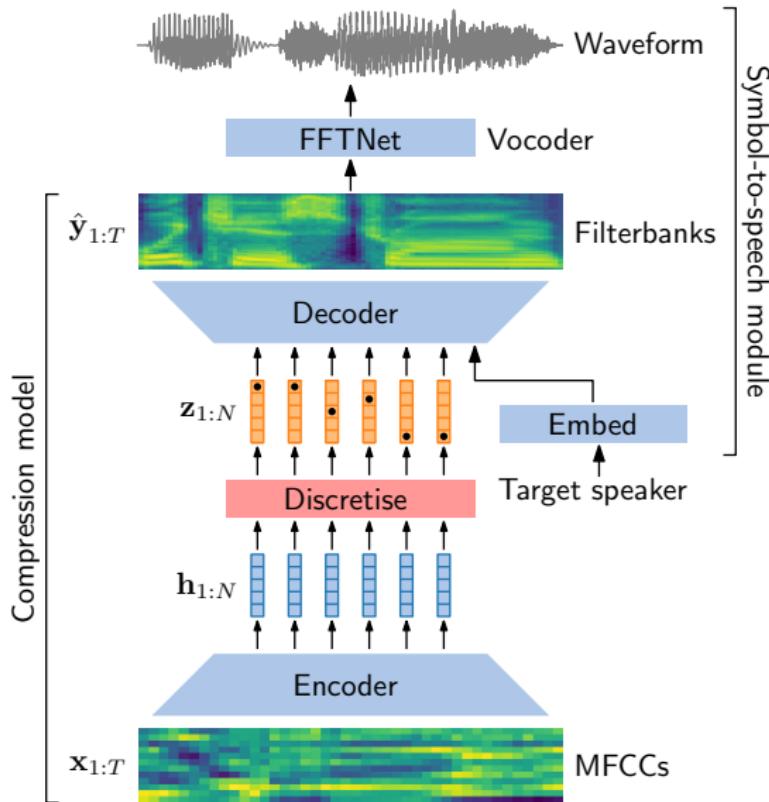
# Approach: Compress, decode and synthesise



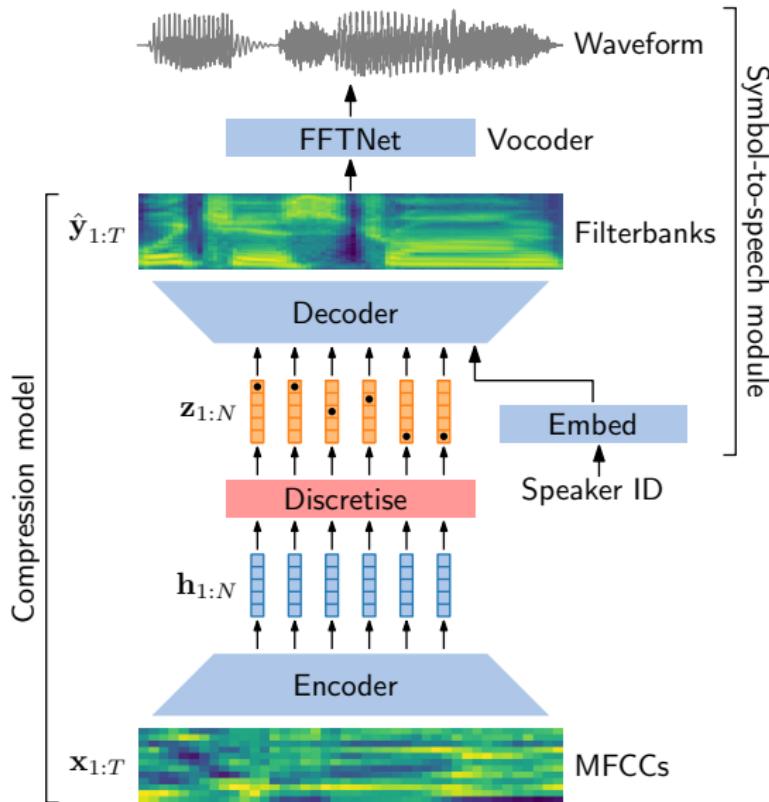
# Approach: Compress, decode and synthesise



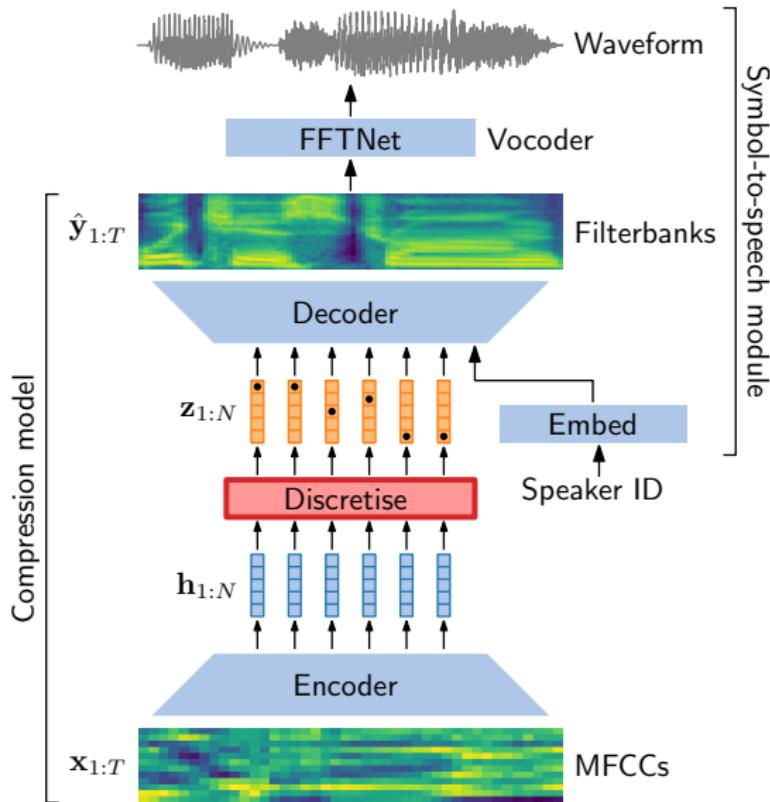
# Approach: Compress, decode and synthesise



# Approach: Compress, decode and synthesise



# Approach: Compress, decode and synthesise



# Discretisation methods

- Straight-through estimation (STE) binarisation:

$$\begin{bmatrix} \mathbf{h} \\ 0.9 \\ -0.1 \\ 0.3 \\ 0.7 \\ -0.8 \end{bmatrix} \xrightarrow{\text{threshold}} \begin{bmatrix} \mathbf{z} \\ 1 \\ -1 \\ 1 \\ 1 \\ -1 \end{bmatrix}$$

- Categorical variational autoencoder (CatVAE):

$$\begin{bmatrix} \mathbf{h} \\ 0.9 \\ -0.1 \\ 0.3 \\ 0.7 \\ -0.8 \end{bmatrix} \xrightarrow{\frac{e^{(h_k+g_k)/\tau}}{\sum_{k=1}^K e^{(h_k+g_k)/\tau}}} \begin{bmatrix} \mathbf{z} \\ 0.86 \\ 0.01 \\ 0.02 \\ 0.11 \\ 0.00 \end{bmatrix}$$

- Vector-quantised variational autoencoder (VQ-VAE):

$$\begin{bmatrix} \mathbf{h} \\ 0.9 \\ -0.1 \\ 0.3 \\ 0.7 \\ -0.8 \end{bmatrix} \xrightarrow{\substack{\text{Choose closest} \\ \text{embedding } \mathbf{e}}} \begin{bmatrix} \mathbf{z} \\ 0.8 \\ -0.2 \\ 0.3 \\ 0.5 \\ -0.6 \end{bmatrix}$$

# Neural network architectures

- Encoder: Convolutional layers, each layer with a stride of 2
- Decoder: Transposed convolutions mirroring encoder
- Waveform generation: FFTNet autoregressive vocoder
- Also experimented with WaveNet: Sometimes gave noisy output
- Bitrate: Set by number of symbols  $K$  and number of striding layers

# Evaluation

Human evaluation metrics:

- Mean opinion score (MOS)
- Character error rate (CER)
- Similarity to the target speaker's voice

# Evaluation

Human evaluation metrics:

- Mean opinion score (MOS)
- Character error rate (CER)
- Similarity to the target speaker's voice

Objective evaluation metrics:

- ABX discrimination
- Bitrate

# Evaluation

Human evaluation metrics:

- Mean opinion score (MOS)
- Character error rate (CER)
- Similarity to the target speaker's voice

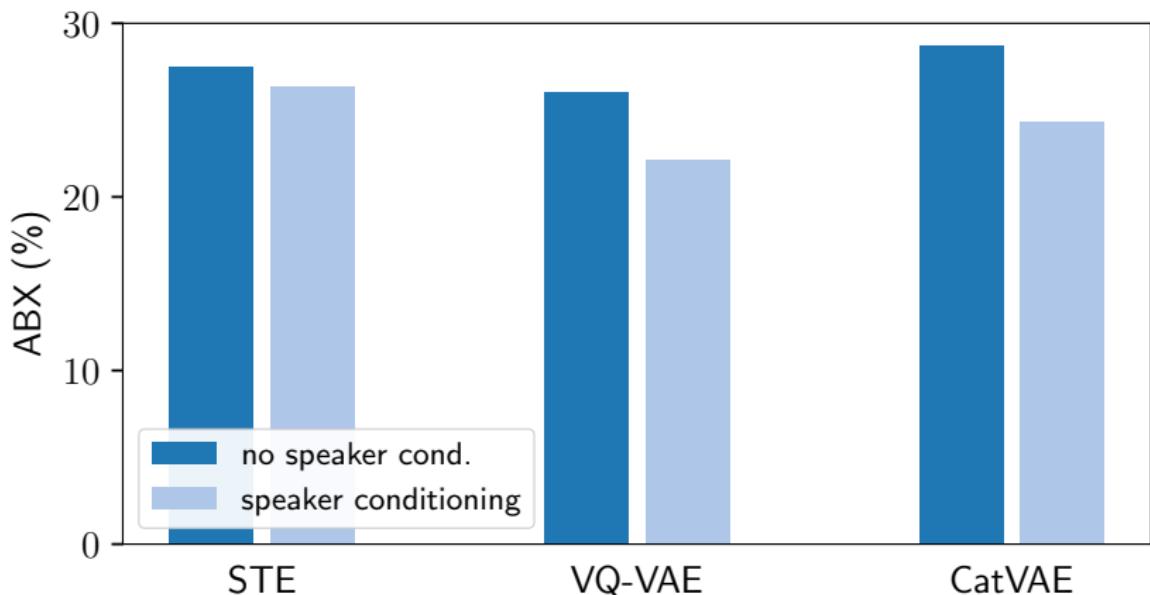
Objective evaluation metrics:

- ABX discrimination
- Bitrate

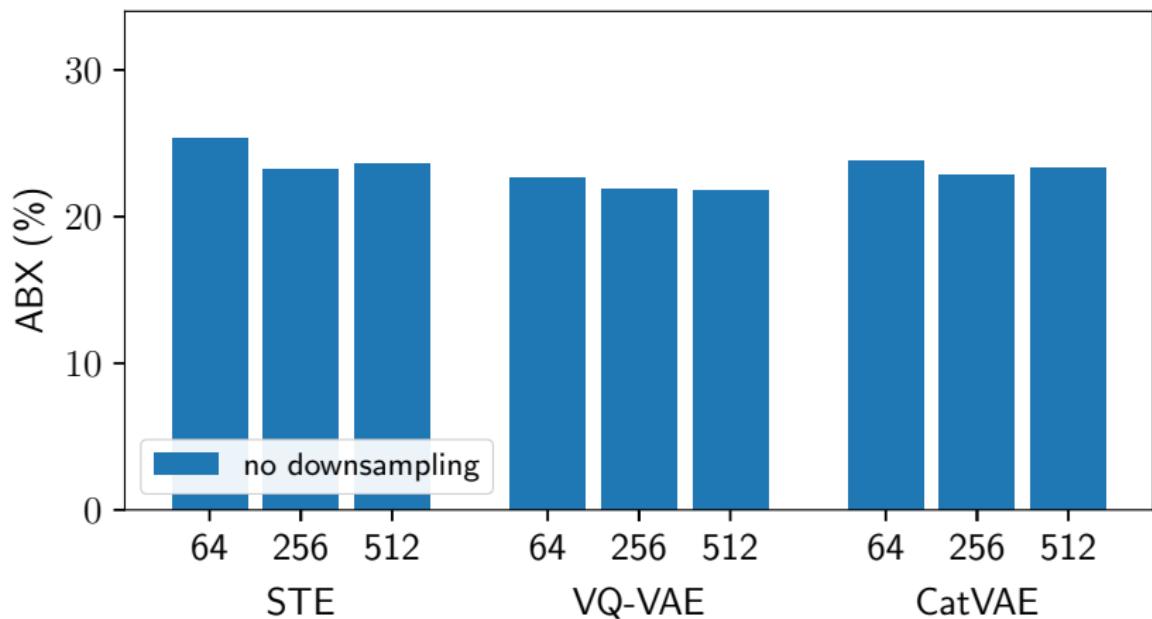
Two evaluation languages:

- English: Used for development
- Indonesian: Held out “surprise language”

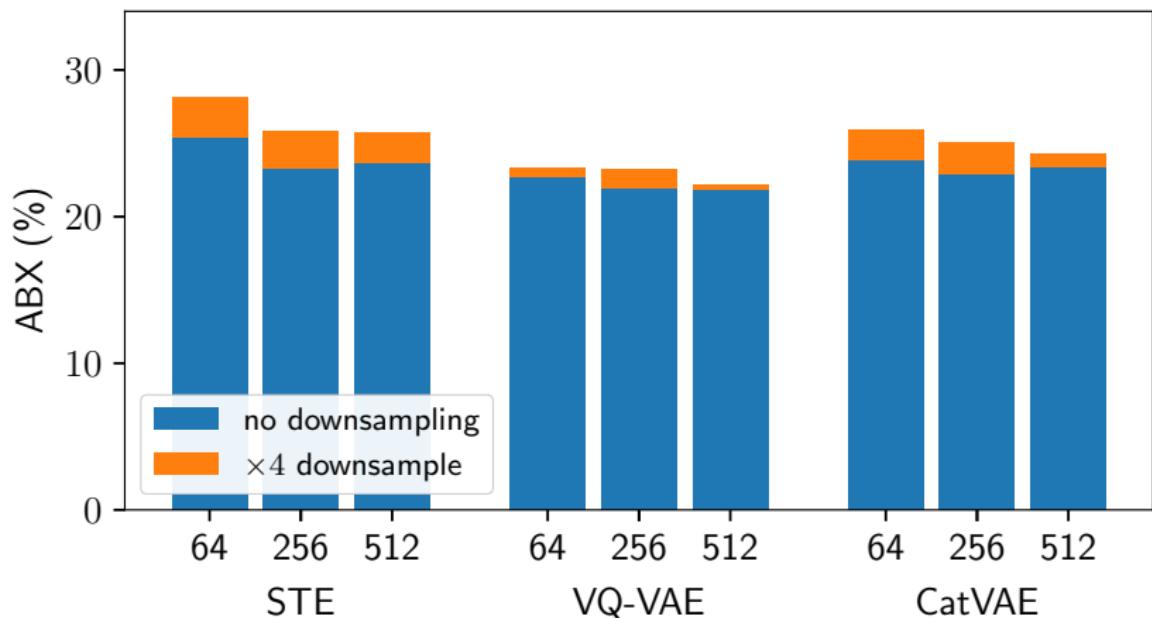
# ABX on English with speaker conditioning



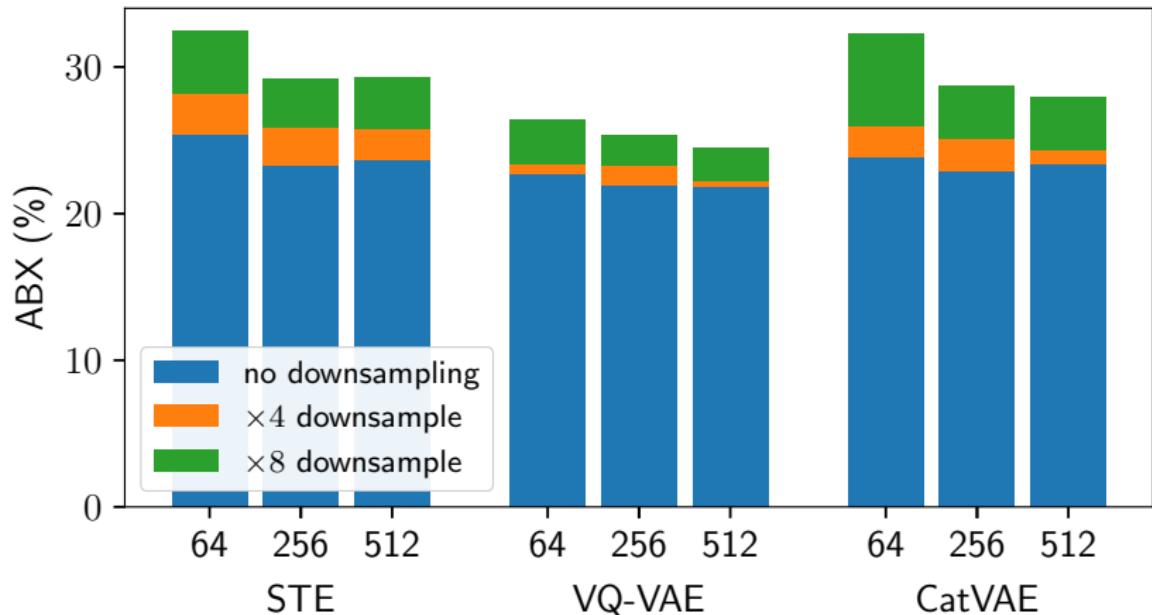
# ABX on English for different compression rates



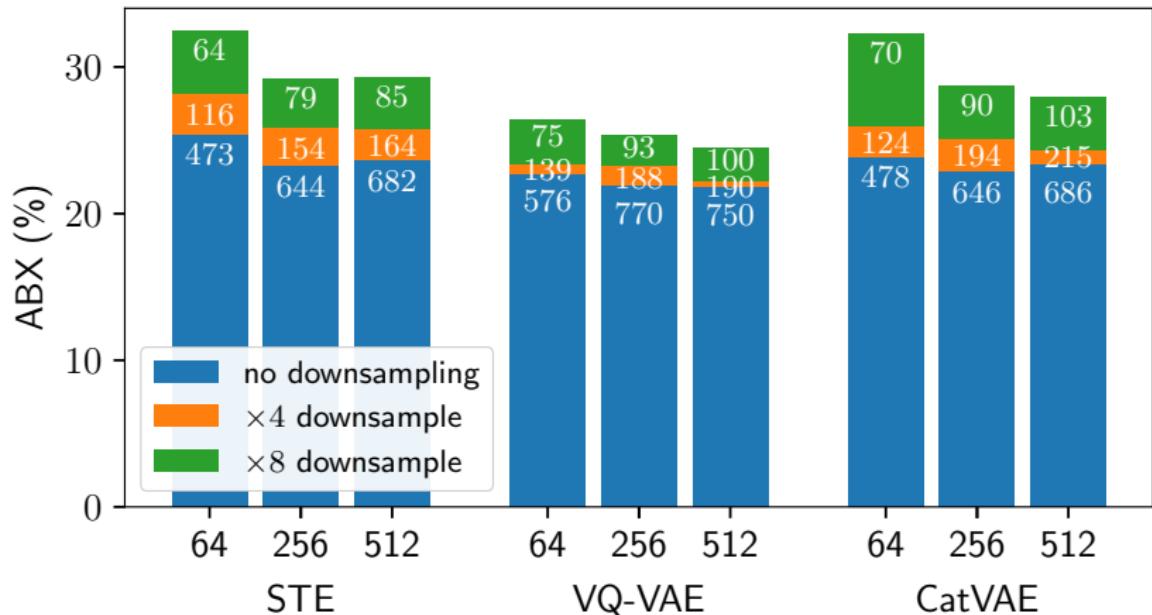
# ABX on English for different compression rates



# ABX on English for different compression rates



# ABX on English for different compression rates



# Official evaluation results

Model	CER (%)	MOS [1, 5]	Similarity [1, 5]	ABX (%)	Bitrate
<i>English:</i>					
DPGMM-Merlin	75	<b>2.50</b>	<b>2.97</b>	35.6	<b>72</b>
VQ-VAE-x8	75	2.31	2.49	25.1	88
VQ-VAE-x4	<b>67</b>	2.18	2.51	<b>23.0</b>	173
Supervised	44	2.77	2.99	29.9	38
<i>Indonesian:</i>					
DPGMM-Merlin	62	<b>2.07</b>	<b>3.41</b>	27.5	75
VQ-VAE-x8	<b>58</b>	1.94	1.95	17.6	<b>69</b>
VQ-VAE-x4	60	1.96	1.76	<b>14.5</b>	140
Supervised	28	3.92	3.95	16.1	35

# Synthesised examples

Model	Input	Synthesised output	Target speaker
<i>English:</i>			
VQ-VAE-x4	<a href="#">Play</a>	<a href="#">Play</a>	<a href="#">Play</a>
VQ-VAE-x4-new		<a href="#">Play</a>	
VQ-VAE-x4	<a href="#">Play</a>	<a href="#">Play</a>	<a href="#">Play</a>
VQ-VAE-x4-new		<a href="#">Play</a>	
<i>Indonesian:</i>			
VQ-VAE-x4	<a href="#">Play</a>	<a href="#">Play</a>	<a href="#">Play</a>
VQ-VAE-x4-new		<a href="#">Play</a>	
VQ-VAE-x4	<a href="#">Play</a>	<a href="#">Play</a>	<a href="#">Play</a>
VQ-VAE-x4-new		<a href="#">Play</a>	

# Conclusions

- Speaker conditioning consistently improves performance
- Different discretisation methods are similar (VQ-VAE slightly better)
- Different models difficult to compare because of bitrate
- Future: Does discretisation actually benefit feature learning?

# Why do we have ten authors on this paper?



Ryan  
Eloff



André  
Nortje



Benjamin  
van Niekerk



Avashna  
Govender



Leanne  
Nortje



Arnu  
Pretorius



Elan van  
Biljon



Ewald van der  
Westhuizen



Lisa van  
Staden



Herman  
Kamper

<https://github.com/kamperh/suzerospeech2019>

[`https://github.com/kamperh/suzerospeech2019`](https://github.com/kamperh/suzerospeech2019)  
(Update coming soon)

# Straight-through estimation (STE) binarisation

- STE binarisation:

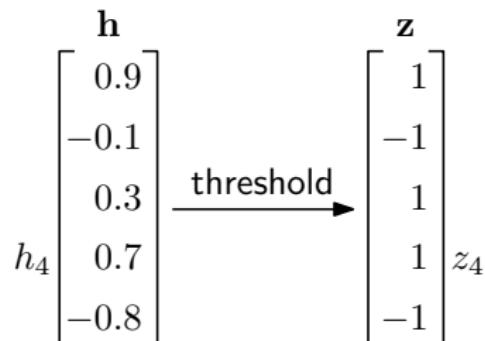
$$z_k = 1 \text{ if } h_k \geq 0 \text{ or } z_k = -1 \text{ otherwise}$$

- For backpropagation we need:  $\frac{\partial J}{\partial \mathbf{h}}$

- For single element:  $\frac{\partial J}{\partial h_k} = \frac{\partial z_k}{\partial h_k} \frac{\partial J}{\partial z_k}$

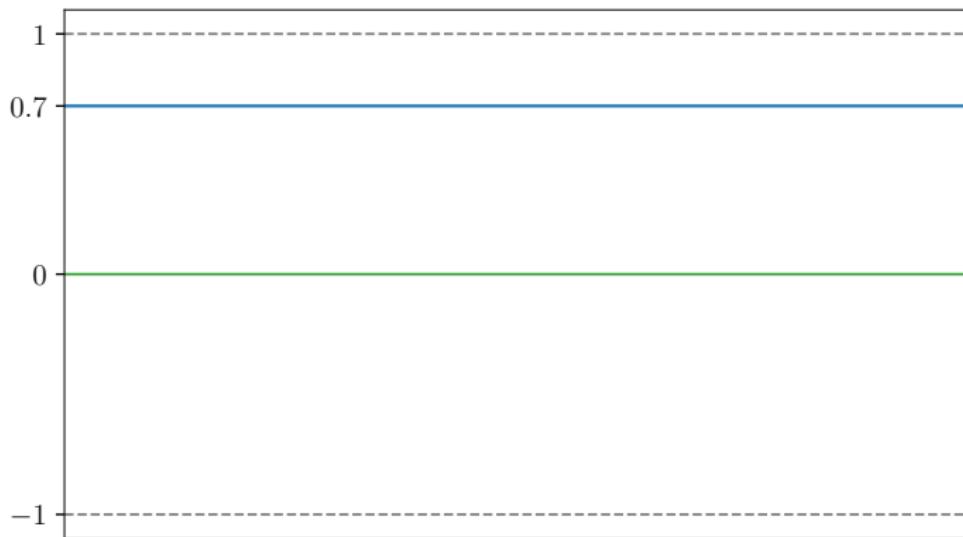
- What is  $\frac{\partial z_k}{\partial h_k}$  with  $z_k = \text{threshold}(h_k)$ ? Cannot solve directly

- Idea: If  $z_k \approx h_k$  then we could use  $\frac{\partial J}{\partial h_k} \approx \frac{\partial J}{\partial z_k}$



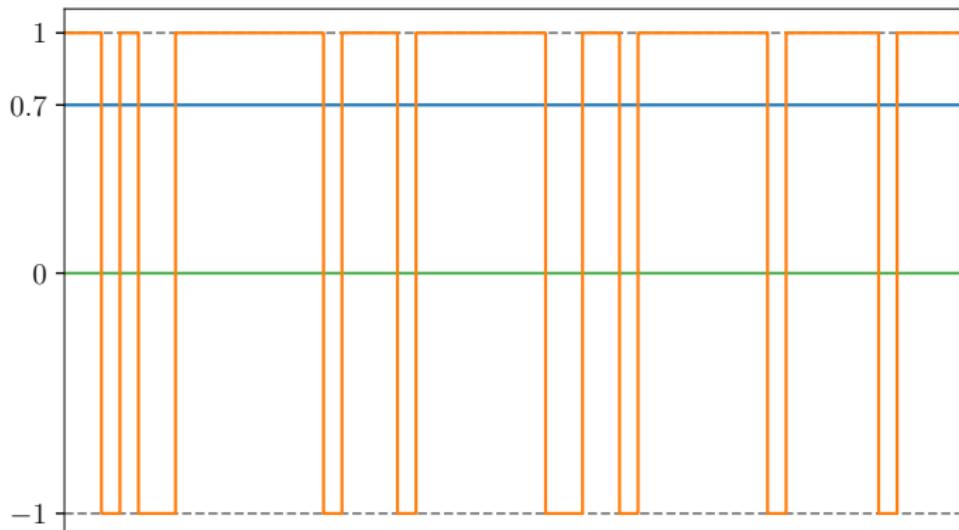
# Straight-through estimation (STE) binarisation

As an example, let us say  $h_k = 0.7$ :



# Straight-through estimation (STE) binarisation

Instead of direct thresholding, let us set  $z_k = 1$  with probability 0.85 and  $z_k = -1$  with probability 0.15:



Estimated mean of  $z_k$  over 500 samples: 0.668

# Straight-through estimation (STE) binarisation

- So, instead of direct thresholding, we set  $z_k = h_k + \epsilon$ , where  $\epsilon$  is sampled noise:

$$\epsilon = \begin{cases} 1 - h_k & \text{with probability } \frac{1+h_k}{2} \\ -h_k - 1 & \text{with probability } \frac{1-h_k}{2} \end{cases}$$

- Since  $\epsilon$  is zero-mean, the derivative of the expected value of  $z_k$  is:  $\frac{\partial \mathbb{E}[z_k]}{\partial h_k} = 1$
- Therefore, gradients are passed unchanged through the thresholding operation:  $\frac{\partial J}{\partial \mathbf{h}} \approx \frac{\partial J}{\partial \mathbf{z}}$