

# Speaker independence of neural vocoders and their effect on parametric resynthesis speech enhancement

**Soumi Maiti<sup>1</sup>, Michael I Mandel<sup>1,2</sup>**

<sup>1</sup>The Graduate Center, City University of New York

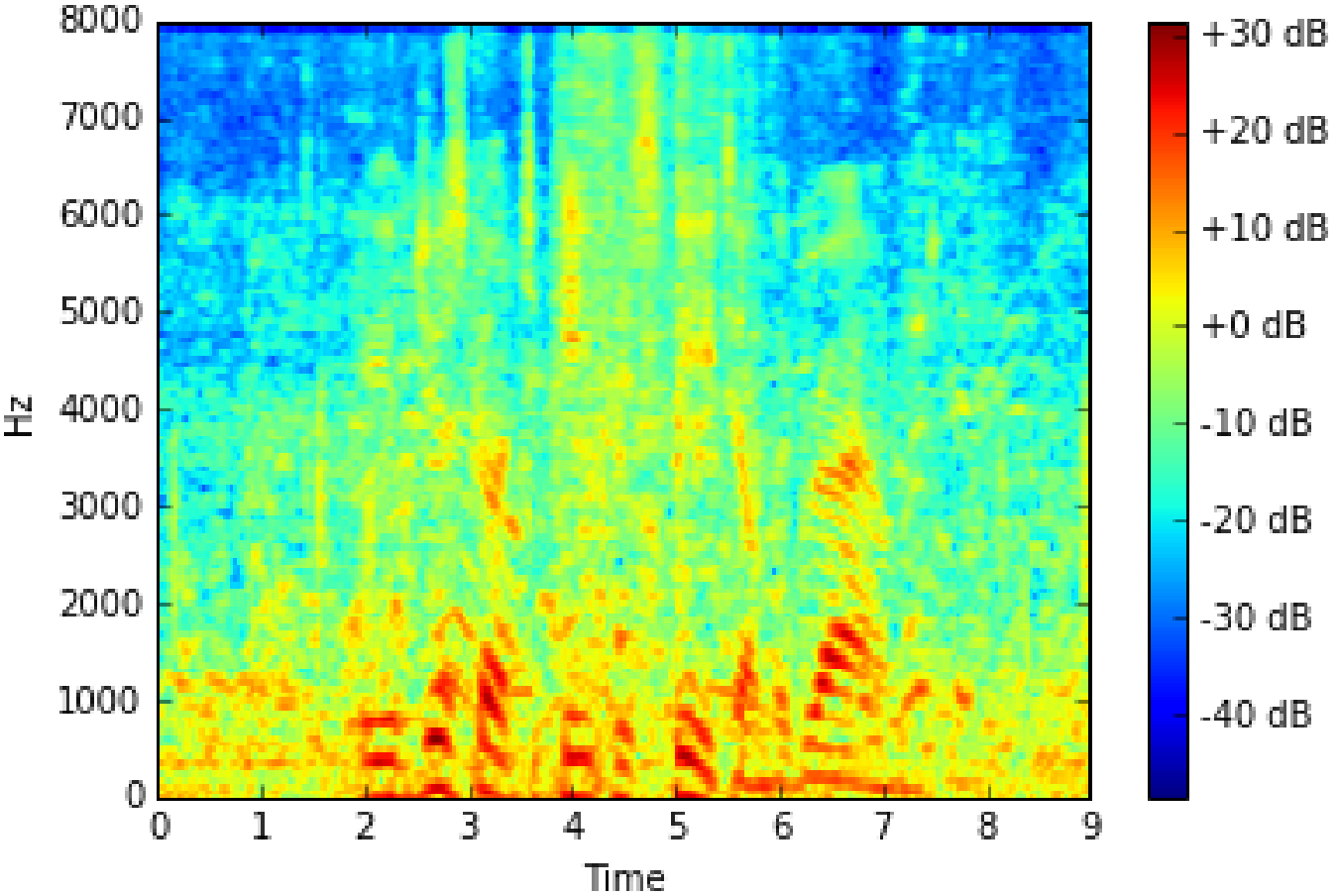
<sup>2</sup>Brooklyn College, City University of New York

ICASSP 2020



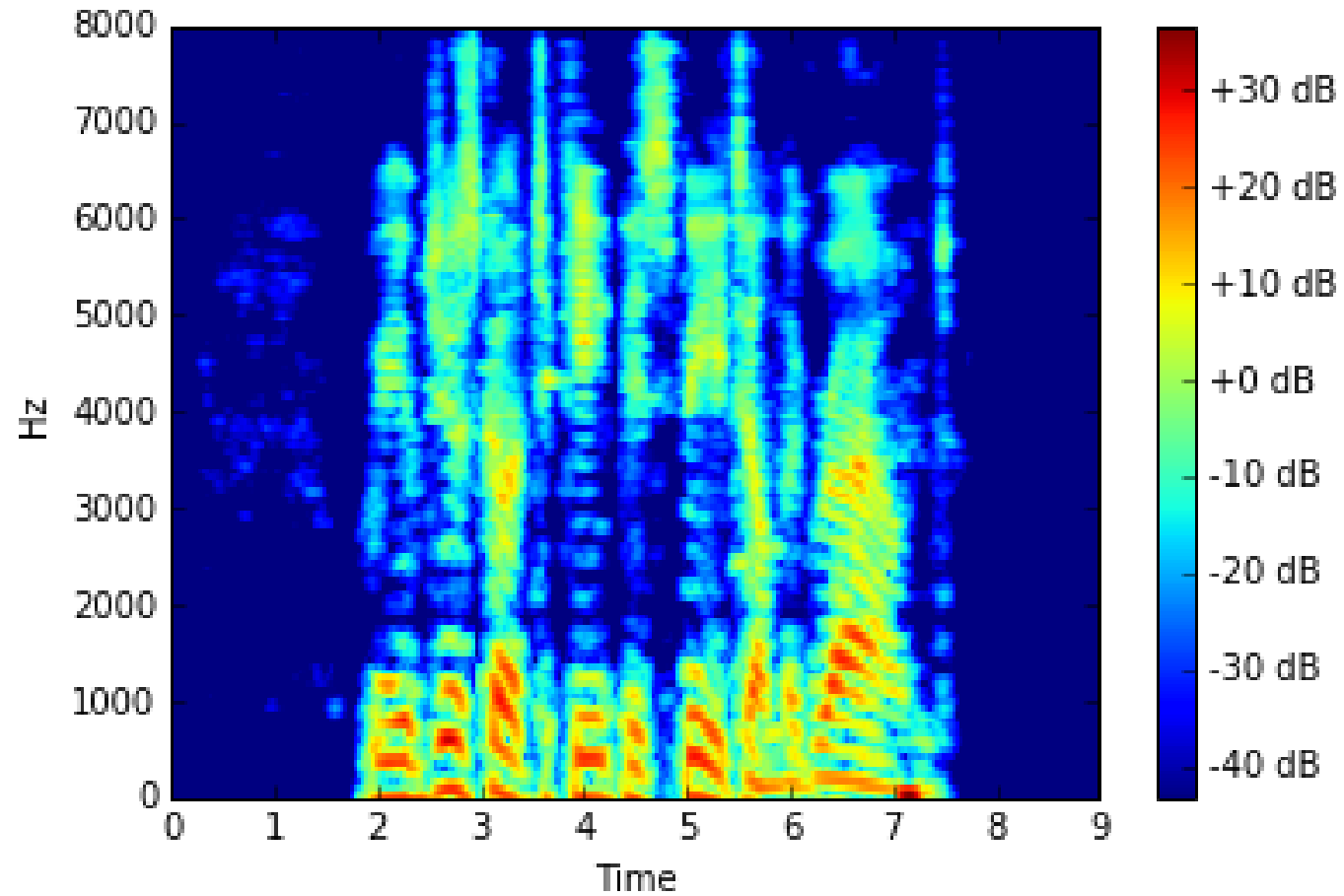
# Example: Noisy file

SNR: 2.5 dB  
Café noise



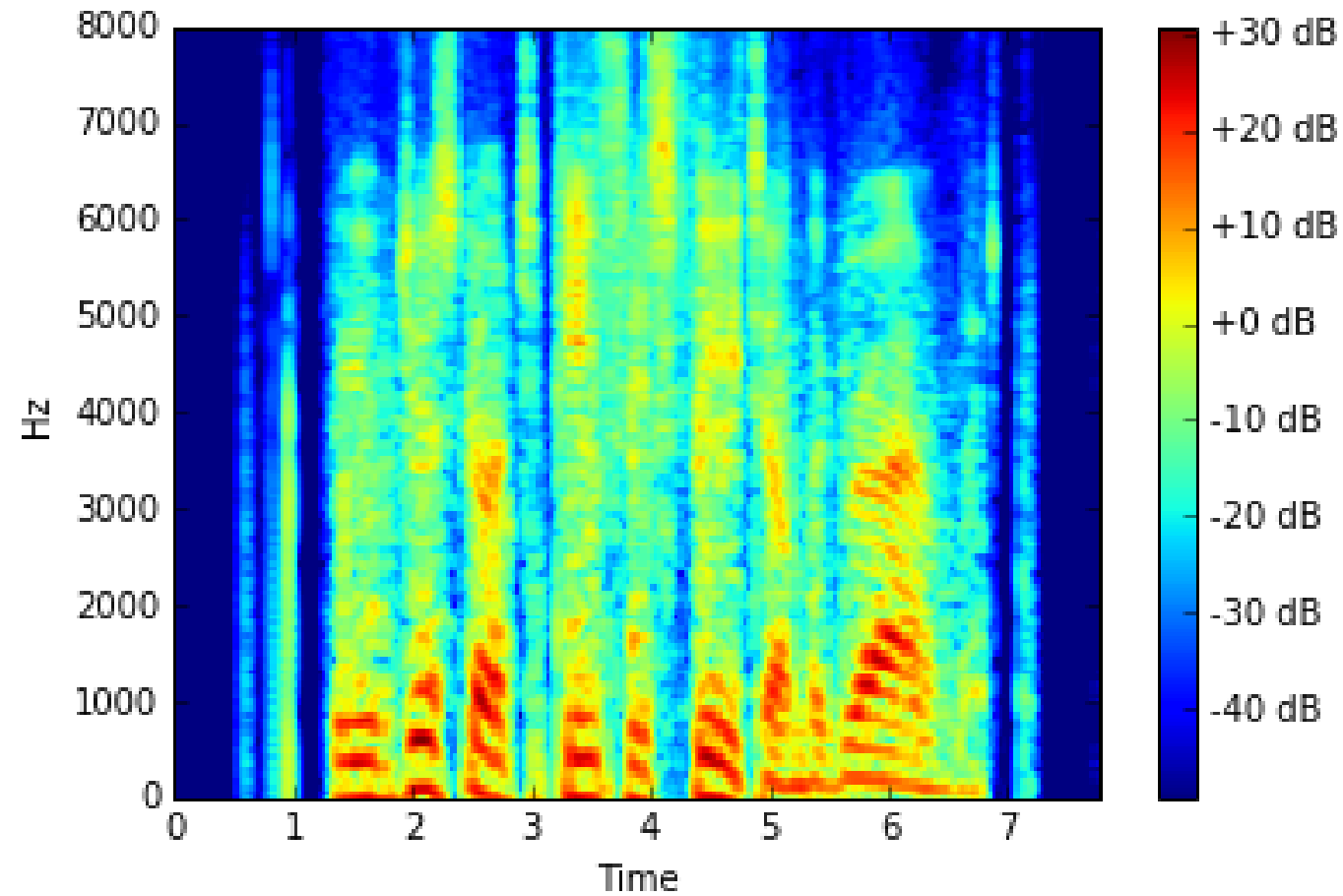
# Example: Oracle mask system

SNR: 2.5 dB  
Café noise



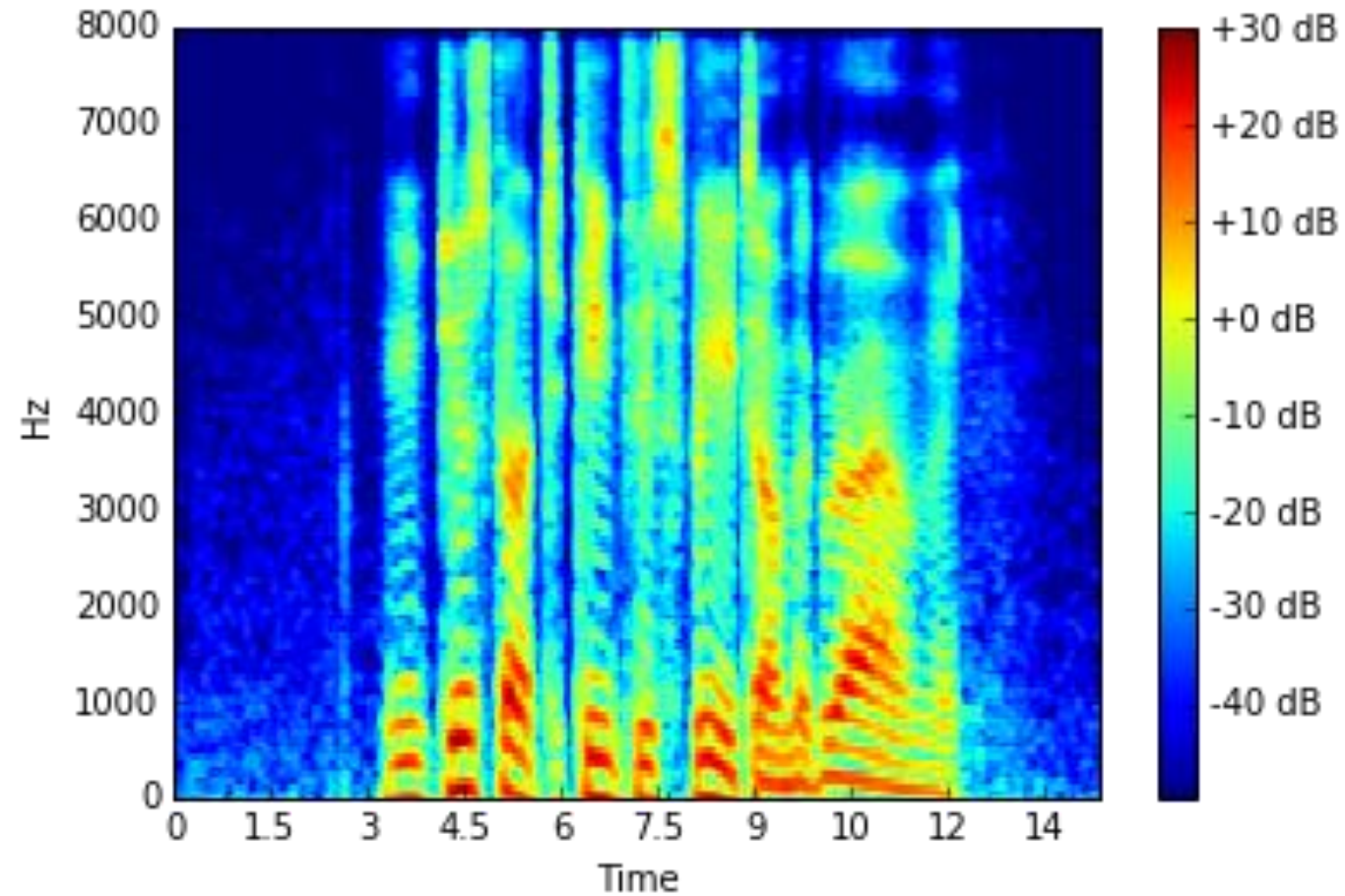
# Example: End to end speech denoise

SNR: 2.5 dB  
Café noise

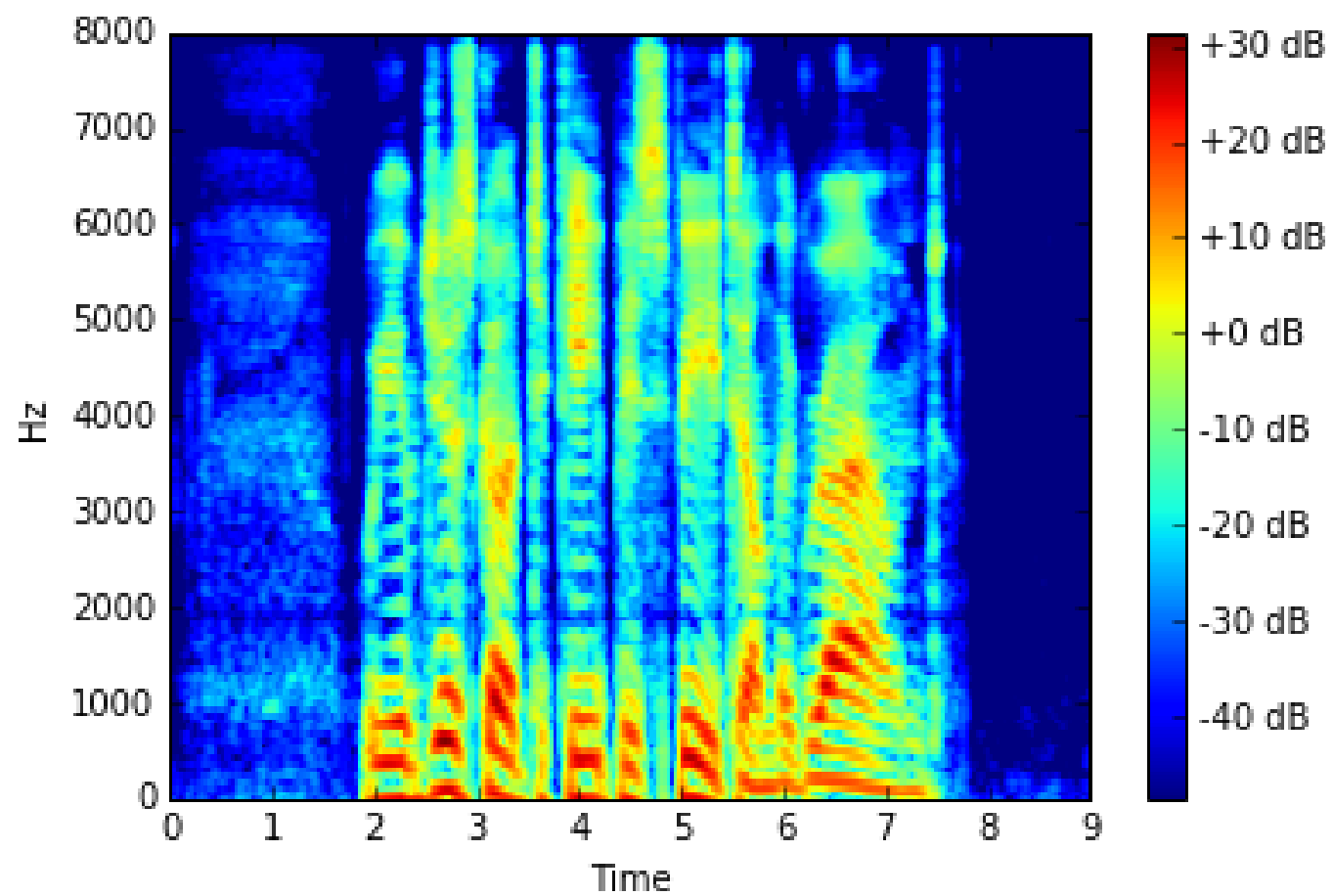


# Example: Parametric resynthesis

SNR: 2.5 dB  
Café noise



# Reference: Clean speech



# Motivation

- Standard speech enhancement → modify noisy recordings
  - Introduce distortions in speech
- Resynthesize clean speech from noisy mixture
- Use speech synthesis for speech enhancement
  - Vocoders to synthesize from acoustic features
  - Easier task than synthesis!<sup>1</sup>

1. Maiti, S. and Mandel, M.I., 2019, May. "Speech denoising by parametric resynthesis", ICASSP, 2019

# Parametric Resynthesis (PR)

- Speech enhancement using vocoders
  - Predict acoustic features from noisy speech
  - Synthesize speech from acoustic features
- High quality speech generation → neural vocoders<sup>2</sup>

**Neural vocoders work on unseen speakers?**



# Neural Vocoders

- Models speech in time domain directly

WaveNet <sup>3</sup>	WaveGlow <sup>4</sup>	LPCNet <sup>5</sup>
<ul style="list-style-type: none"><li>• Autoregressive</li><li>• High Quality</li><li>• Slower synthesis</li></ul>	<ul style="list-style-type: none"><li>• Generates samples in parallel</li><li>• Glow based model</li><li>• Fast synthesis</li></ul>	<ul style="list-style-type: none"><li>• Autoregressive</li><li>• Faster synthesis → written in C</li><li>• Hybrid model<ul style="list-style-type: none"><li>○ Models vocal response with LPC coefficients</li><li>○ Predicts excitation → simpler task</li></ul></li></ul>

WaveNet: we use GPU accelerated nv-WaveNet – for faster synthesis

3. Oord, Aaron van den, et al. "Wavenet: A generative model for raw audio." , SSW, 2016.

4. Prenger, R, Rafael V, and Bryan C. "Waveglow: A flow-based generative network for speech synthesis.", ICASSP 2019

5. Valin, J. M., & Skoglund, J. "LPCNet: Improving neural speech synthesis through linear prediction." , ICASSP 2019

# Train neural vocoders with large number of speakers

- **Train<sup>6</sup>: 56 speakers**
  - Voices from VCTK dataset
  - 28 male / 28 female
  - Accent: US and Scotland
- **Test<sup>6</sup>: 6 unseen speakers**
  - 3 male / 3 female
  - Accent: England
- Sampling rate 16 kHz
- **Objective quality metrics**
  - **CSIG, CBAK, COVL<sup>7</sup>**
  - 0-5 → higher is better

6. Valentini-Botinhao, C., 2017. Noisy speech database for training speech enhancement algorithms and TTS models.

7. Loizou, P.C., 2013. *Speech enhancement: theory and practice*. CRC press.

# Speaker independence of neural vocoders

Scores averaged over 10 files/speaker

**nv-Wavenet → low generation quality**

Neural vocoders generalizes to unseen speakers



	CSIG	CBAK	COVL
<b>1 Seen speaker</b>			
WaveGlow	4.7	3.0	4.0
LPCNet	3.8	2.2	2.9
nv-WaveNet	3.3	2.1	2.5
<b>6 Unseen speakers</b>			
WaveGlow	4.6	2.8	3.9
LPCNet	4.0	2.4	3.1
nv-WaveNet	3.2	2.1	2.5

# Effect of gender on neural vocoders

Scores averaged over 10 files/speaker

WaveGlow models female speakers slightly better

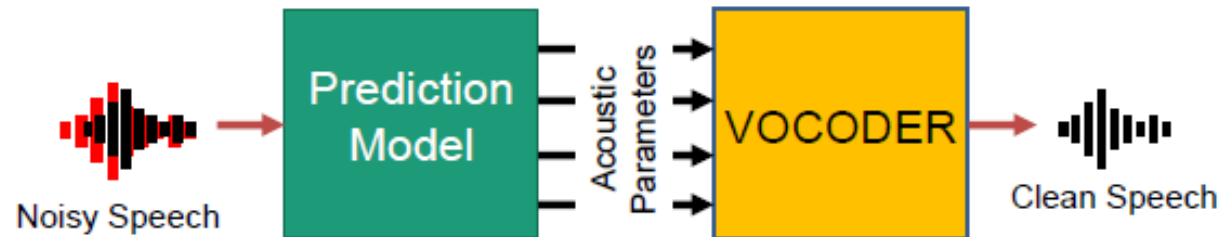
LPCNet, WaveNet no such difference

		CSIG	CBAK	COVL
Male				
	WaveGlow	4.5	2.8	3.8
	LPCNet	4.0	2.3	3.1
	nv-WaveNet	3.2	2.1	2.5
Female				
	WaveGlow	4.6	2.8	3.9
	LPCNet	4.0	2.4	3.2
	nv-WaveNet	3.3	2.0	2.5

# Parametric Resynthesis (PR)

1. Predict “clean” acoustic features from noisy speech
2. Reconstruct speech from acoustic features
  - Acoustic features are different for different vocoders

Resynthesize clean speech by predicting acoustic parameters



# Vocoder acoustic features

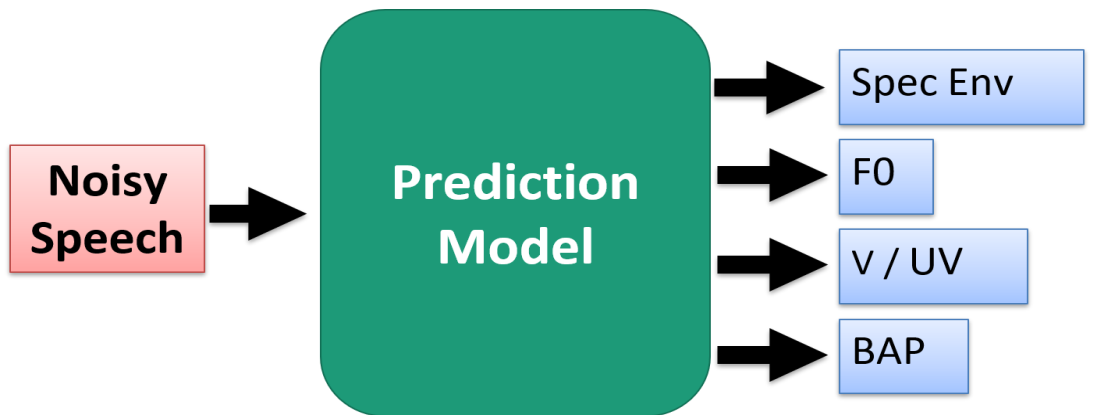
Vocoder	Acoustic features	Dim
WaveNet	Mel-spectrogram	80
WaveGlow	Mel-spectrogram	80
LPCNet	BFCC, F0 period, F0 correlation	20
WORLD	Spectral envelope, aperiodicity, F0, v/uv	63

Neural vocoders →

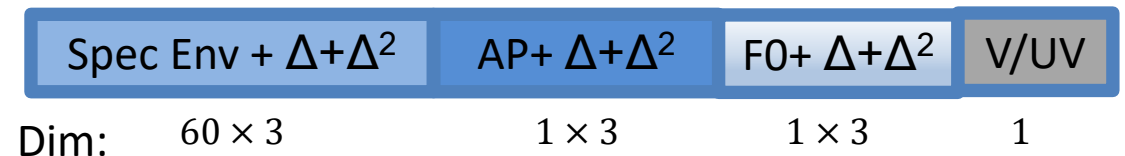
WaveNet, WaveGlow, LPCNet

# Prediction model

- Predicts acoustic features ( $X$ ) at a fixed frame rate
  - **Input:** noisy mel-spectrogram ( $Y$ )
  - **Loss:**  $\text{MSE} := |X - X'|^2$
  - $X' \rightarrow$  predicted features



PR-WORLD



# Training of Parametric resynthesis

- **Noisy trainset:** 56 speaker set
  - 8 noises from DEMAND
  - 2 artificial noises
  - SNR range: 15 – 0 dB
- **Test:** 2 unseen speakers
  - 8 unseen noises from DEMAND
  - 824 files
  - 4 SNR level: 17.5 dB, 12.5 dB, 7.5 dB, 2.5 dB
- Comparison models:
  - SEGAN<sup>8</sup>, Wave-U-Net<sup>9</sup>, Wavenet-denoise<sup>10</sup>
  - Oracle Wiener mask
    - Has access to clean speech

8. Pascual, S., Bonafonte, A. and Serra, J., "SEGAN: Speech enhancement generative adversarial network". *arXiv*, 2017.

9. Macartney, C. and Weyde, T., "Improved speech enhancement with the wave-u-net". *arXiv*, 2018.

10. Rethage, D., Pons, J. and Serra, X. "A wavenet for speech denoising", *ICASSP*, 2018.

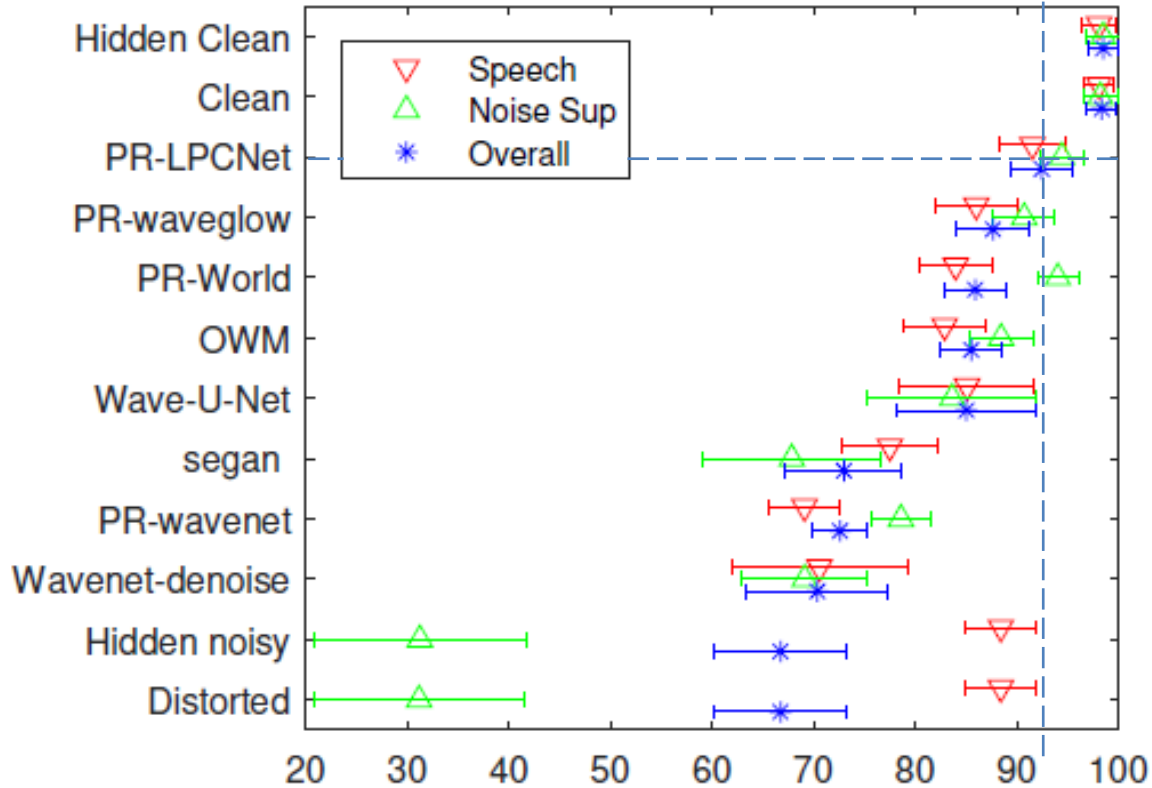


# Objective metrics for speech enhancement

PR-WaveGlow performs best  
in CSIG and CBAK

	CSIG	CBAK	COVL	STOI
Oracle Wiener	4.3	3.8	3.8	0.98
PR-WaveGlow	3.8	2.4	3.1	0.91
PR-LPCNet (noisy F0)	3.5	2.1	2.7	0.88
PR-LPCNet	3.1	1.8	2.2	0.88
Wave-U-Net	3.5	3.2	3.0	
SEGAN	3.5	2.9	2.8	

# Subjective Quality



- MUSHRA listening test
- Number of files: 12
- SNR: 12.5 dB to 2.5 dB

**PR-LPCNet outperforms all systems!**

# Objective metrics on 12 listening test files

	CSIG	CBAK	COVL
OWM	4.3	3.8	3.9
<b>PR-WaveGlow</b>	<b>3.8</b>	<b>2.4</b>	<b>3.0</b>
PR-World	3.1	1.9	2.2
<b>PR-LPCNet</b>	<b>3.0</b>	<b>1.8</b>	<b>2.2</b>

**Subjective quality scores does not match objective scores**  
**LPCNet scores 0.8 lower than WaveGlow!**

# In a nutshell

- Neural vocoders
  - Speaker Independent when trained on large number of speakers
  - All 3 vocoders were able to generalize to unseen speakers
- Speech enhancement
  - PR-LPCNet
    - Outperforms Oracle Wiener mask in subjective quality scores
  - PR-WaveGlow
    - Higher objective metrics than LPCNet

Thank You