

Using vocoders to create training data for speech spoofing countermeasures

Xin Wang, Junichi Yamagishi
National Institute of Informatics

Slides by Xin Wang
National Institute of Informatics

© 2023, Xin Wang. All rights reserved.

This work is licensed under the Creative Commons
Attribution 3.0 license.

See <http://creativecommons.org/> for details.



Contents

<https://arxiv.org/abs/2210.10570>

□ Introduction

- Text-to-speech synthesis
- Speech anti-spoofing

□ Method

- Copy-synthesized data as spoofed data
- Contrastive feature loss over bona fide and copy-synthesized data pair

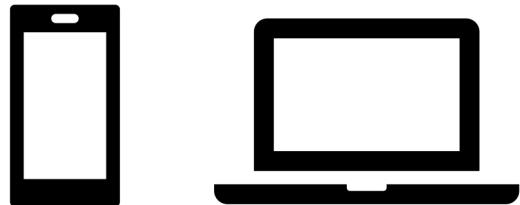
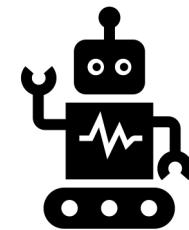
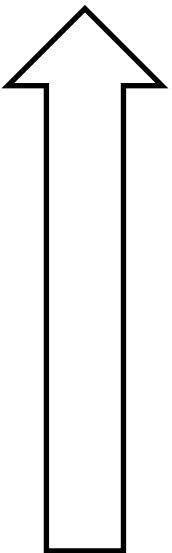
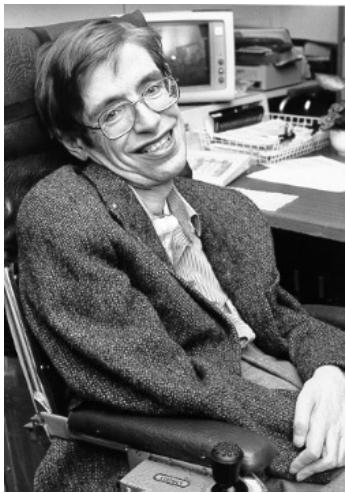
□ Experiment

□ Summary

- We use external training data. Please be careful when interpreting the results.
- We take a data-driven approach. Apologize that we cannot precisely explain the model behavior.

Introduction

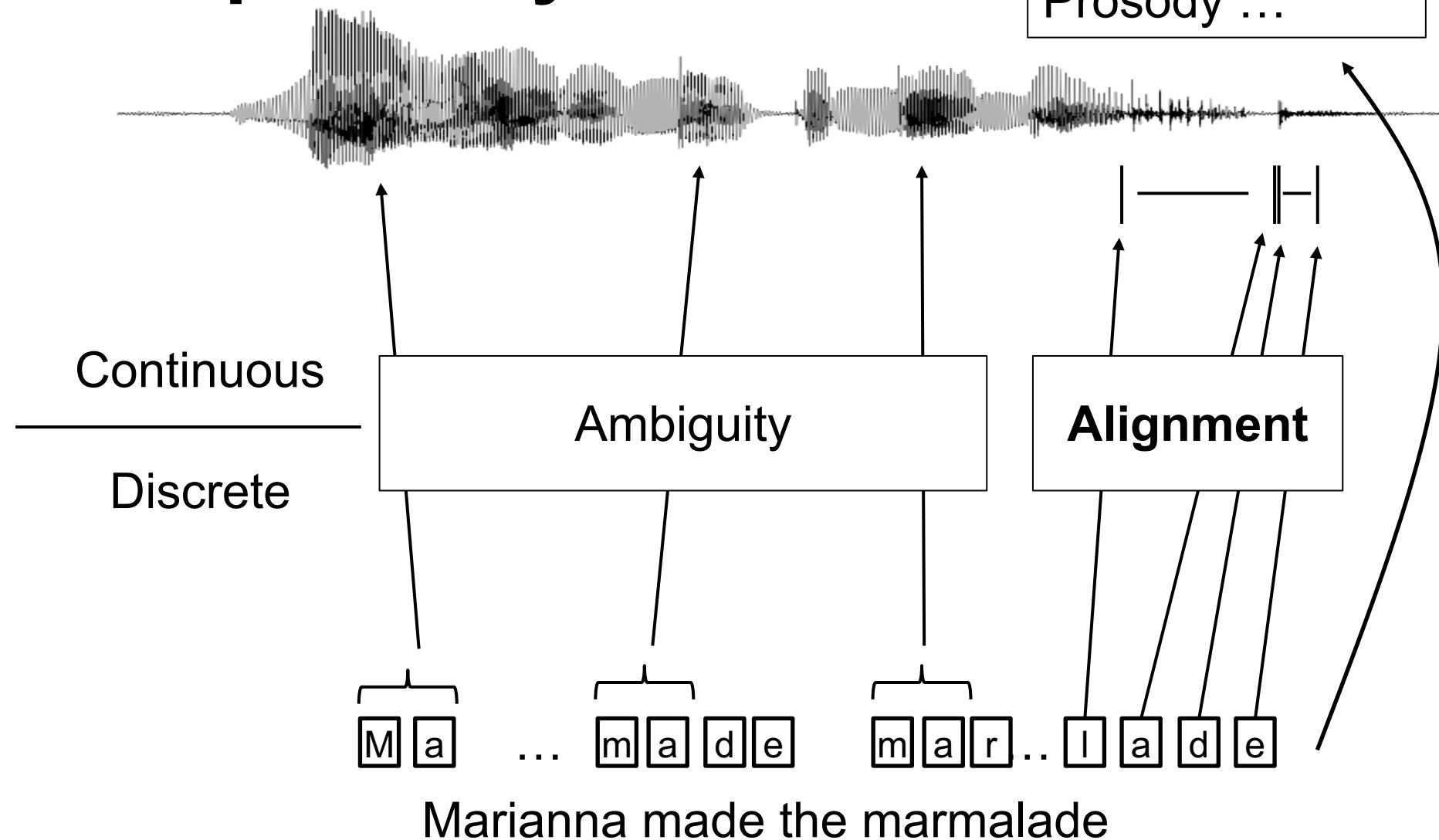
Text-to-speech synthesis



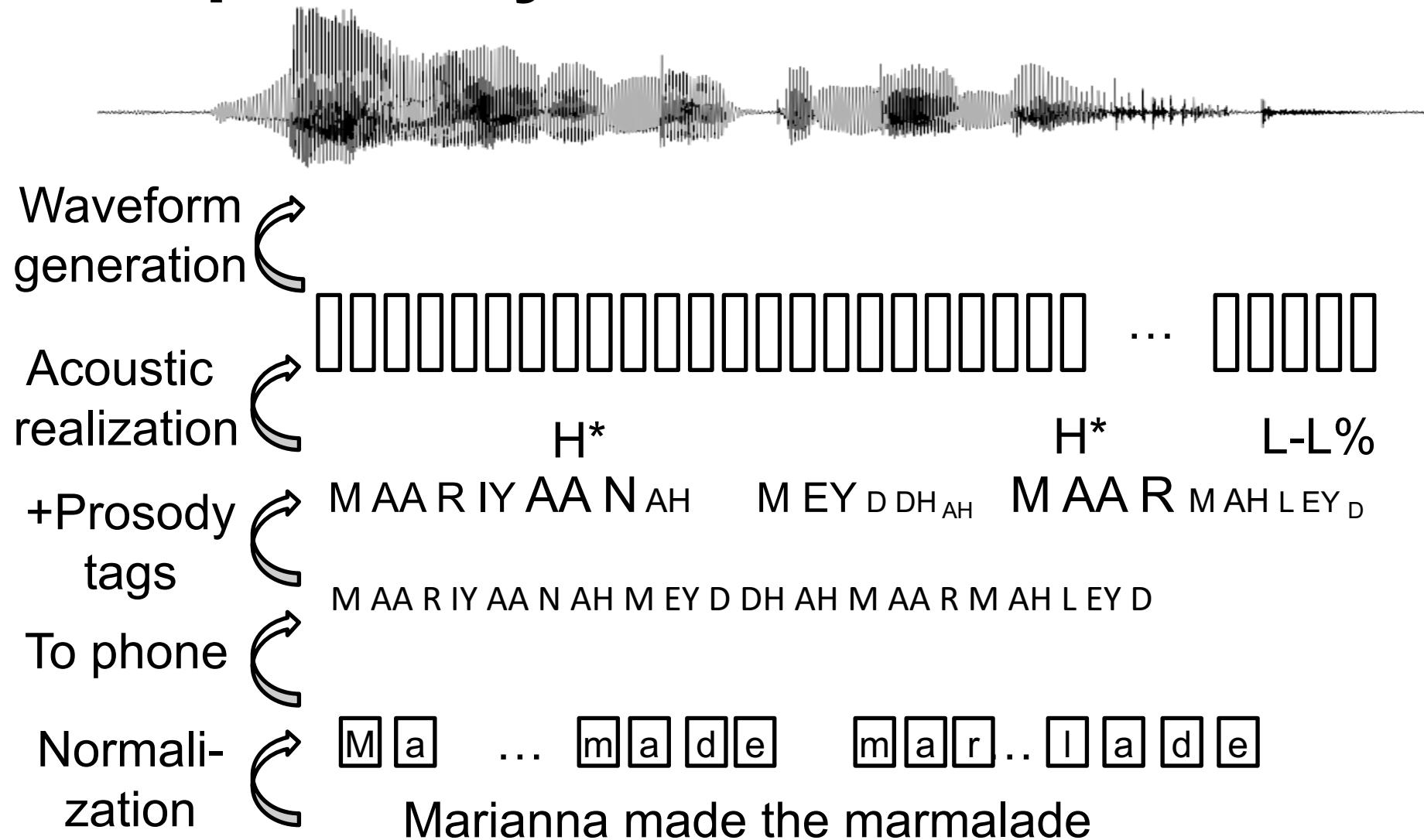
Input text

Text-to-speech synthesis

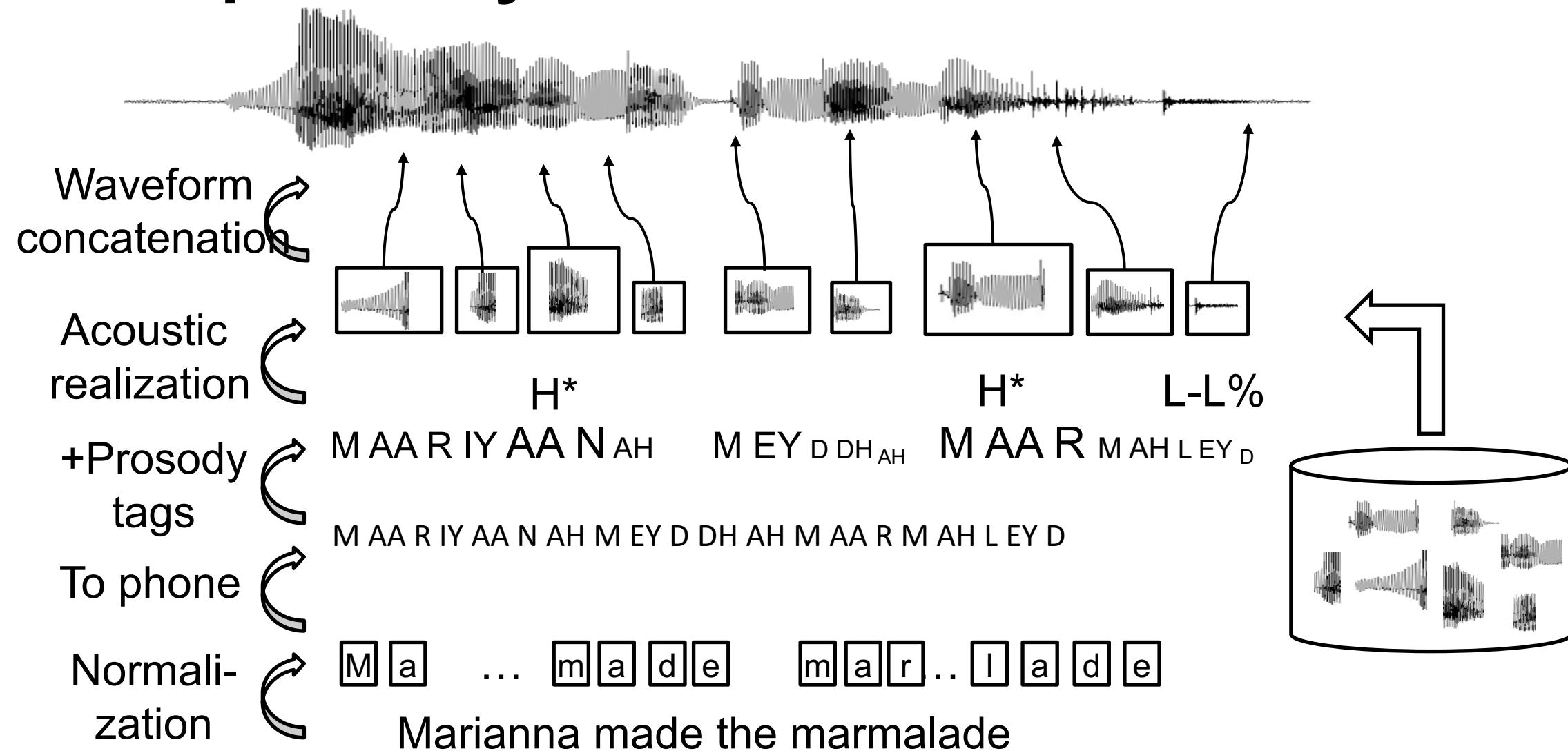
Speaker identity,
Prosody ...



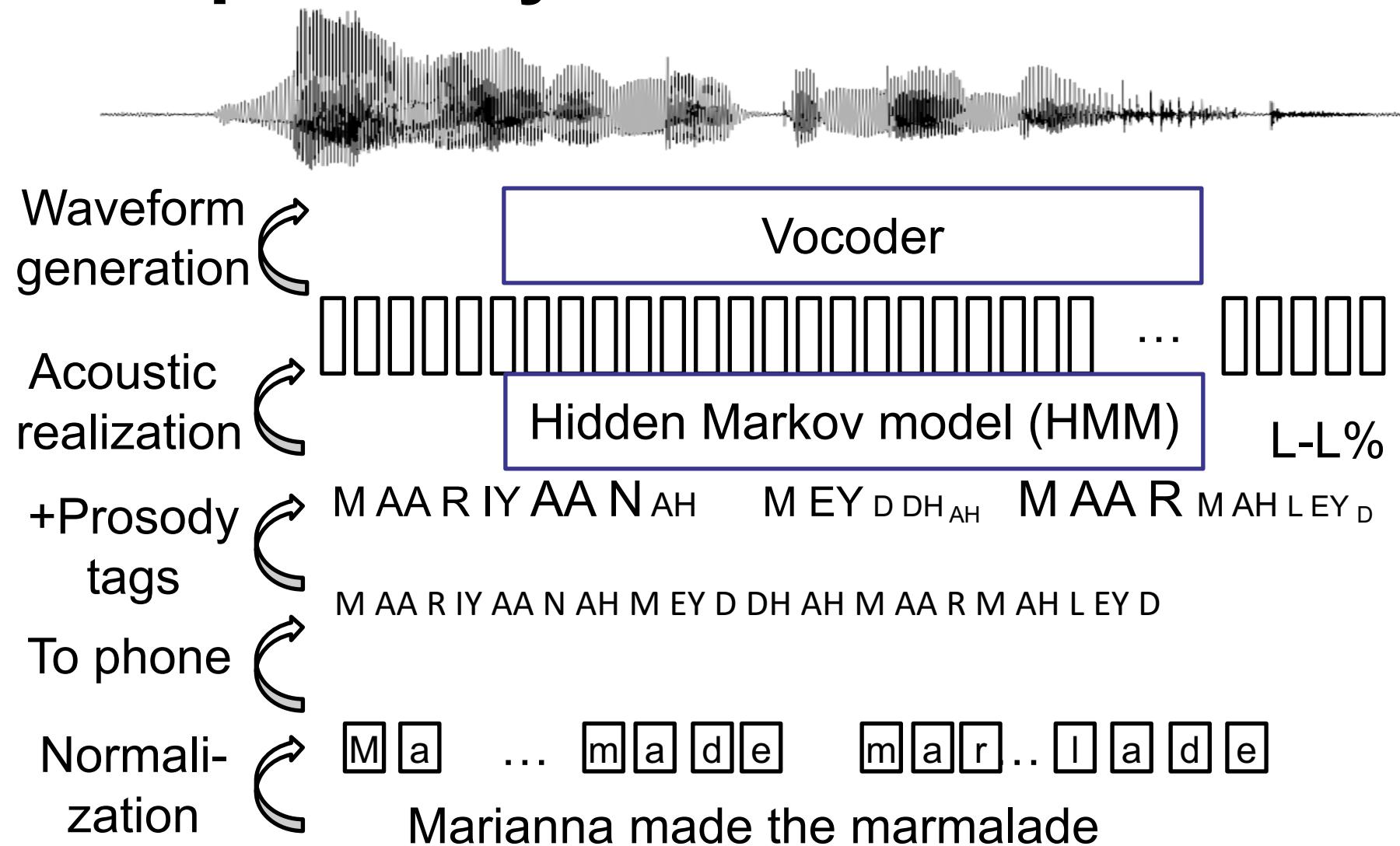
Text-to-speech synthesis



Text-to-speech synthesis – Unit selection



Text-to-speech synthesis – HMM



Text-to-speech synthesis

Waveform
generation

Acoustic
realization

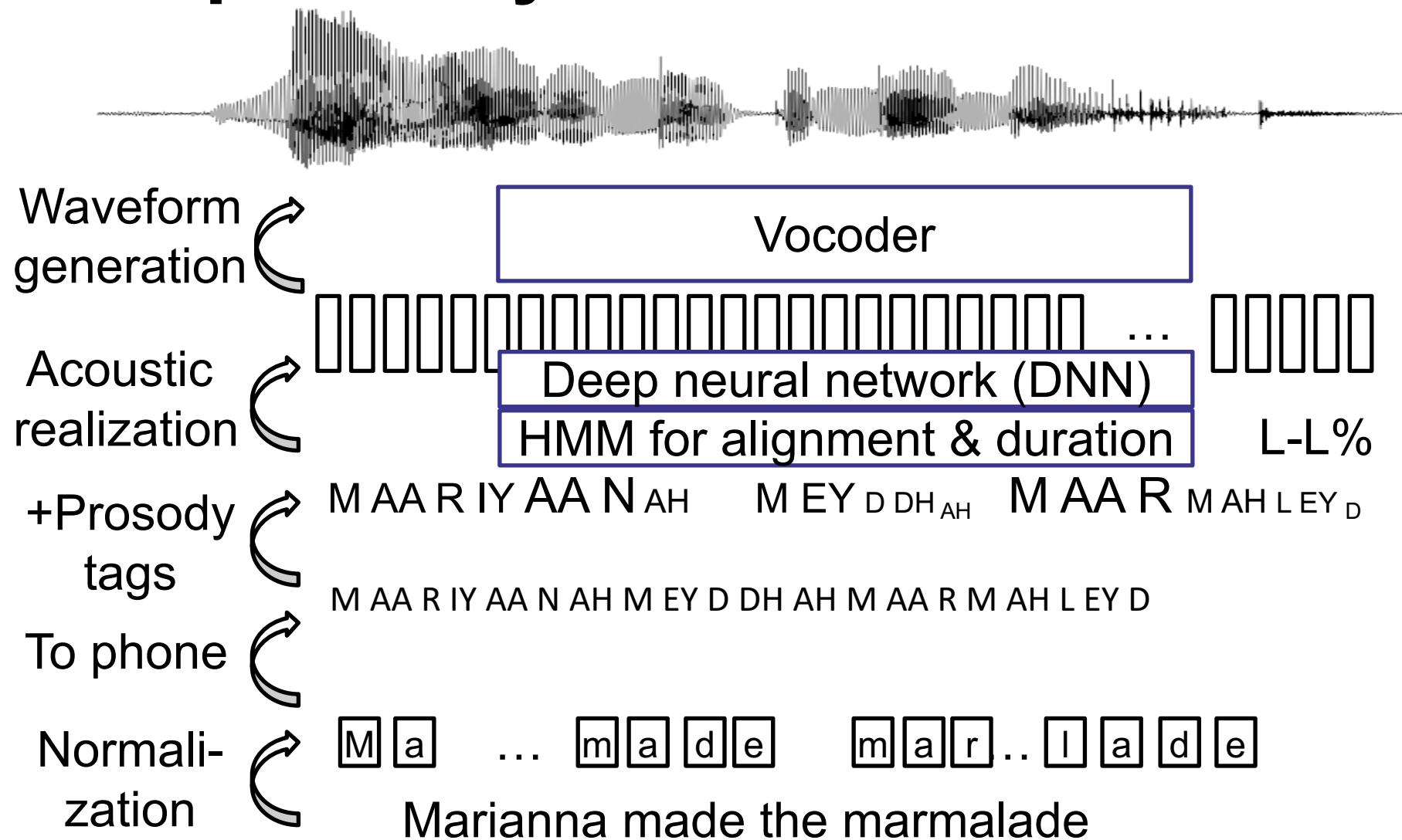
+Prosody
tags

To phonetic
trees

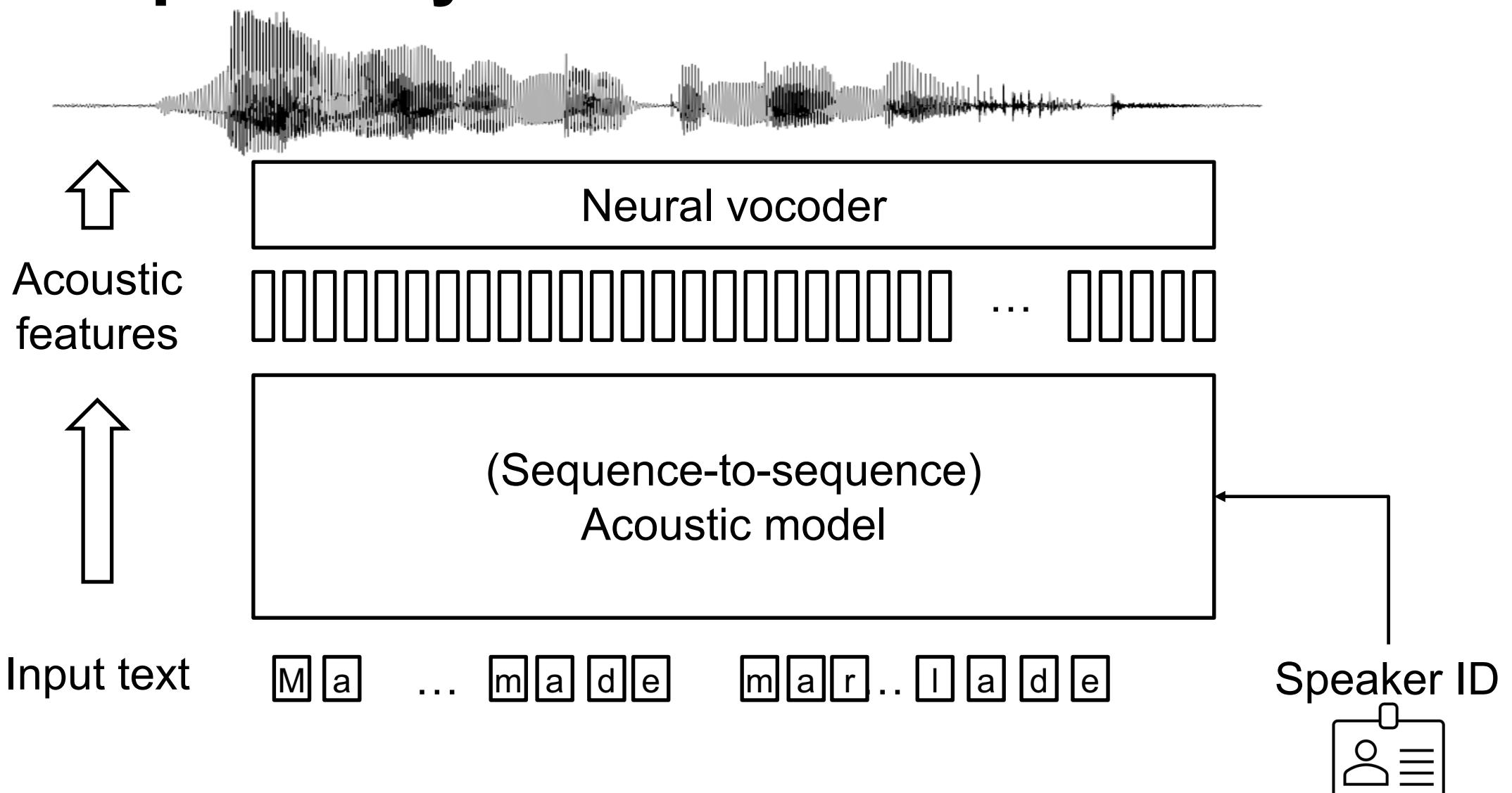
Normaliza-
tion



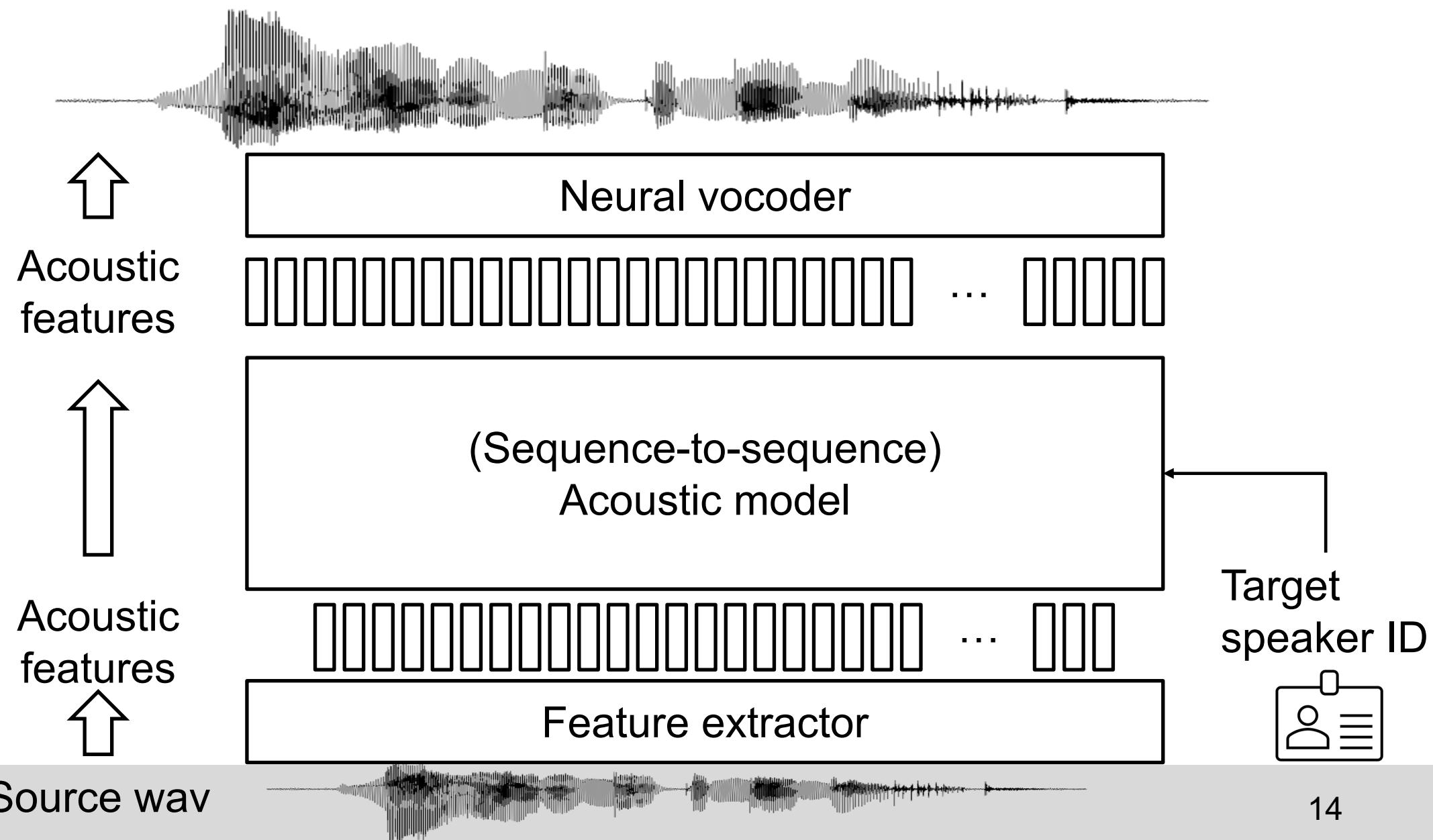
Text-to-speech synthesis – HMM&DNN



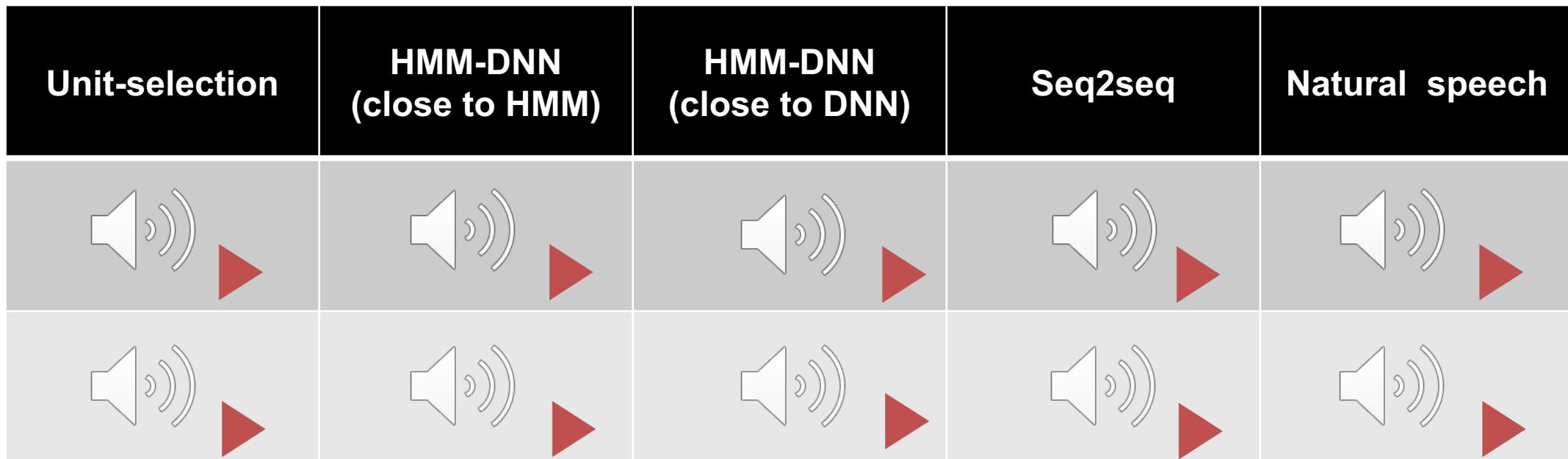
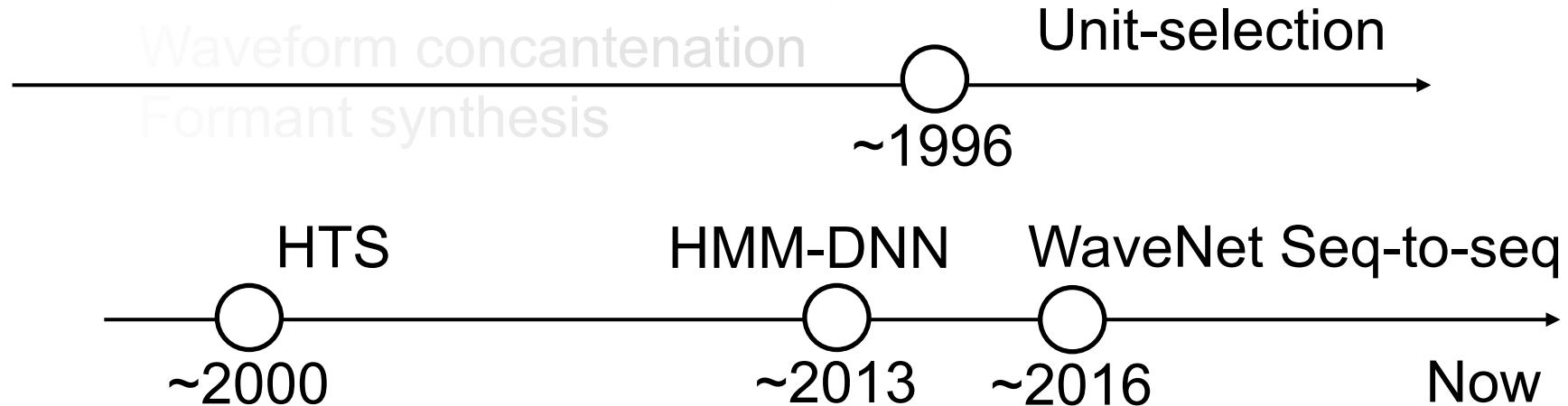
Text-to-speech synthesis – Recent methods



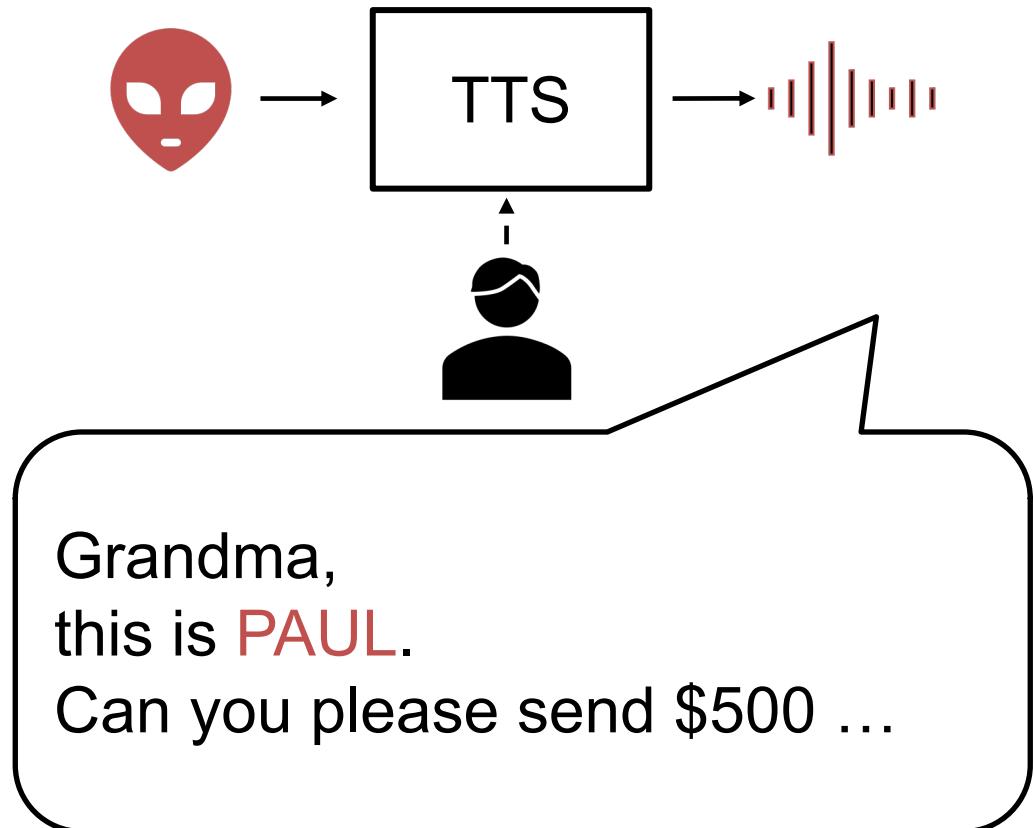
Voice conversion – Recent methods



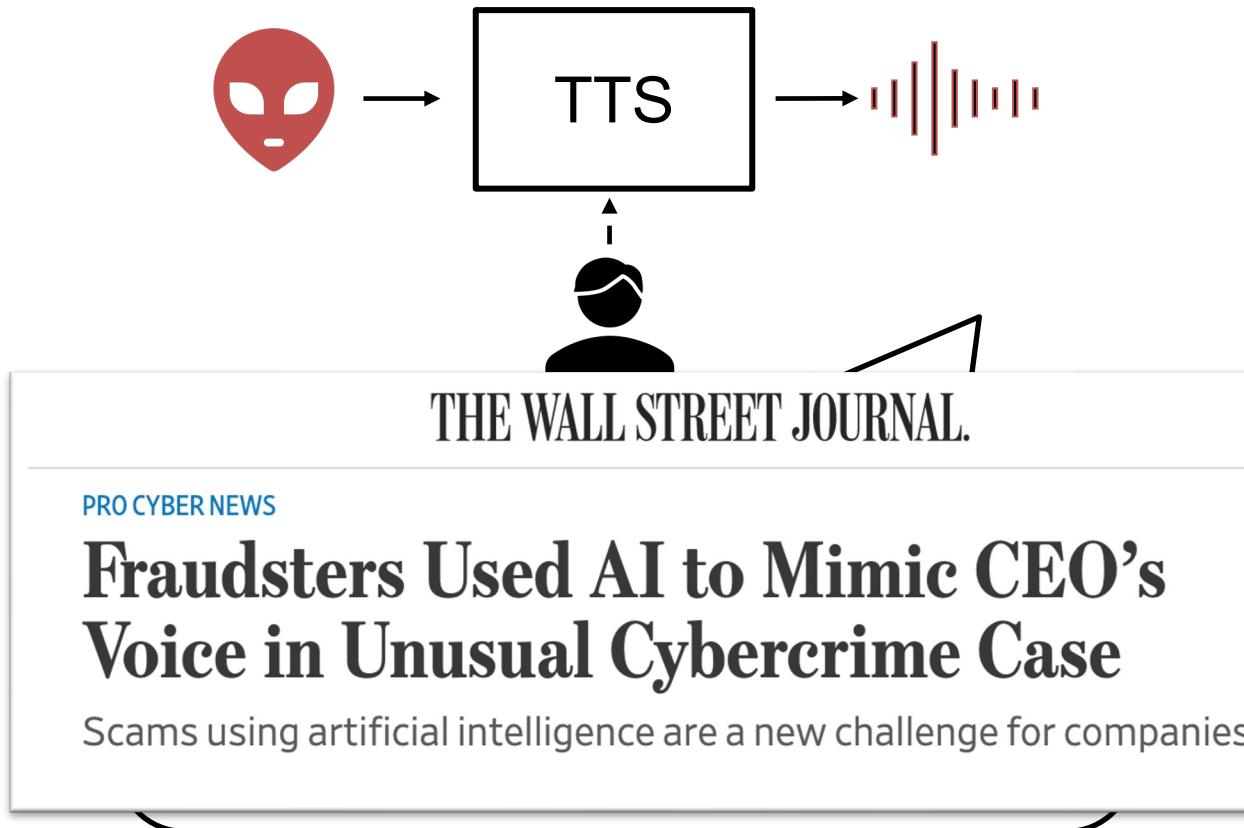
Rapid progress of TTS



TTS/VC may be misused



TTS/VC may be misused



Forbes

Sep 3, 2019, 04:42pm EDT | 51,807 views

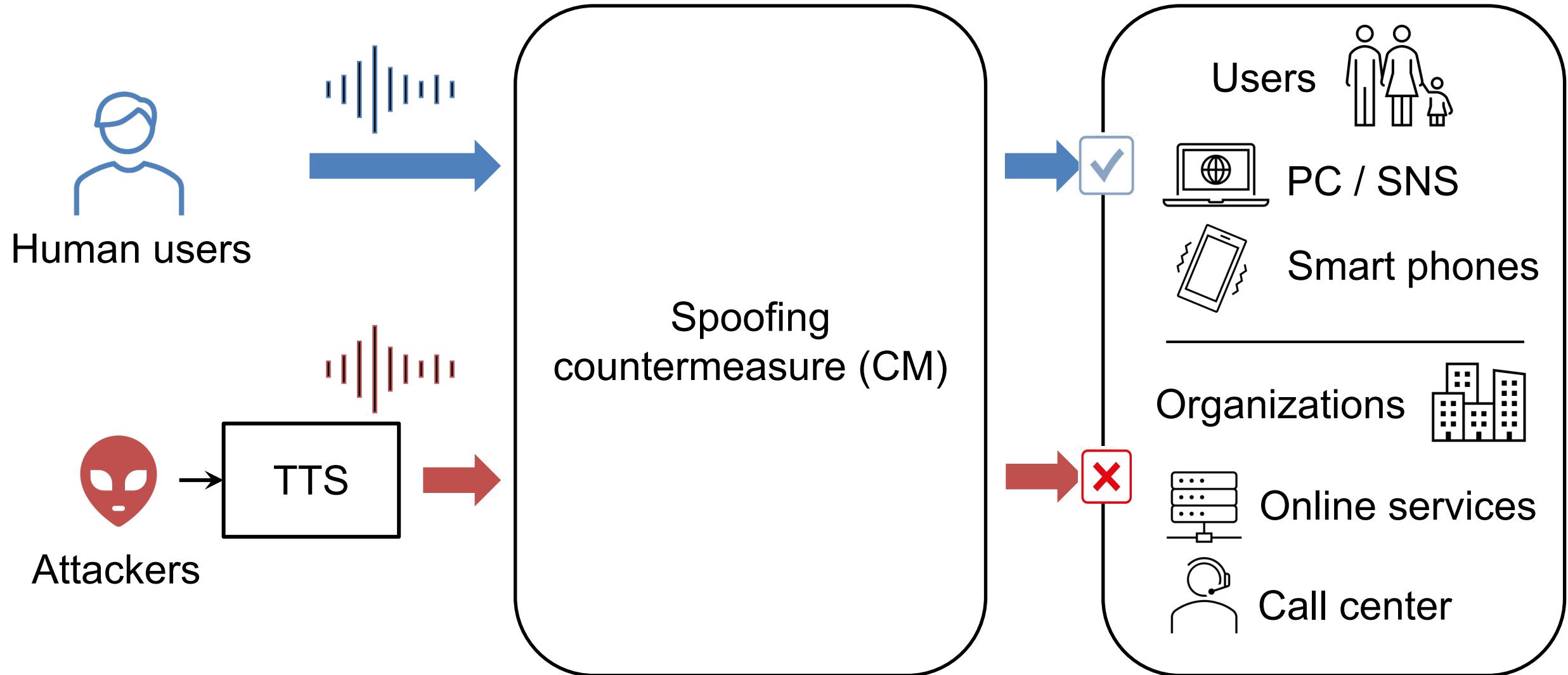
A Voice Deepfake Was Used To Scam A CEO Out Of \$243,000

Jesse Damiani Contributor ⓘ
Consumer Tech
I cover the human side of VR/AR, Blockchain, AI, Startups, & Media.

Follow

Listen to article 3 minutes

Spoofing countermeasure



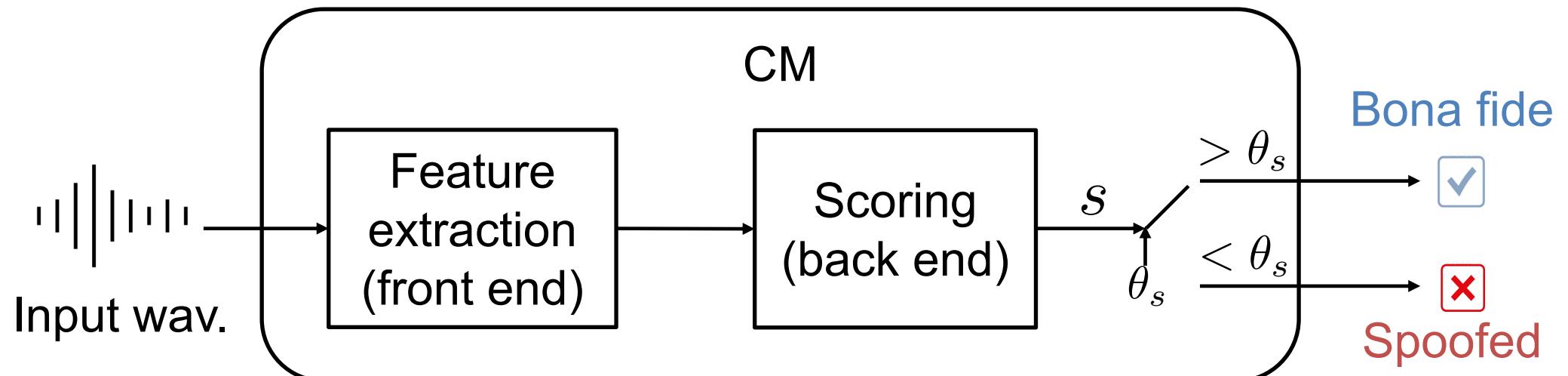
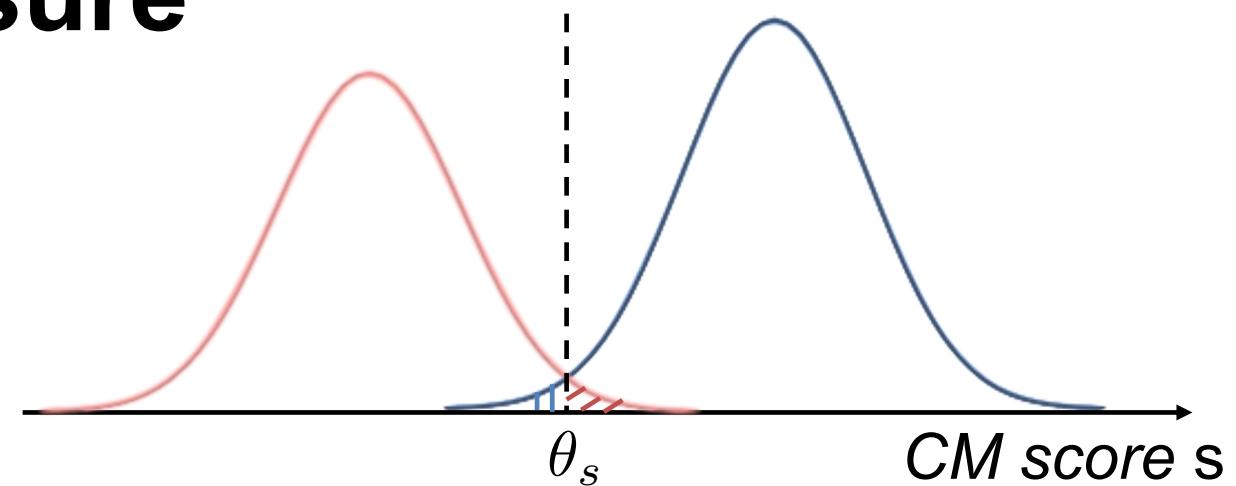
Spoofing countermeasure

□ Model construction

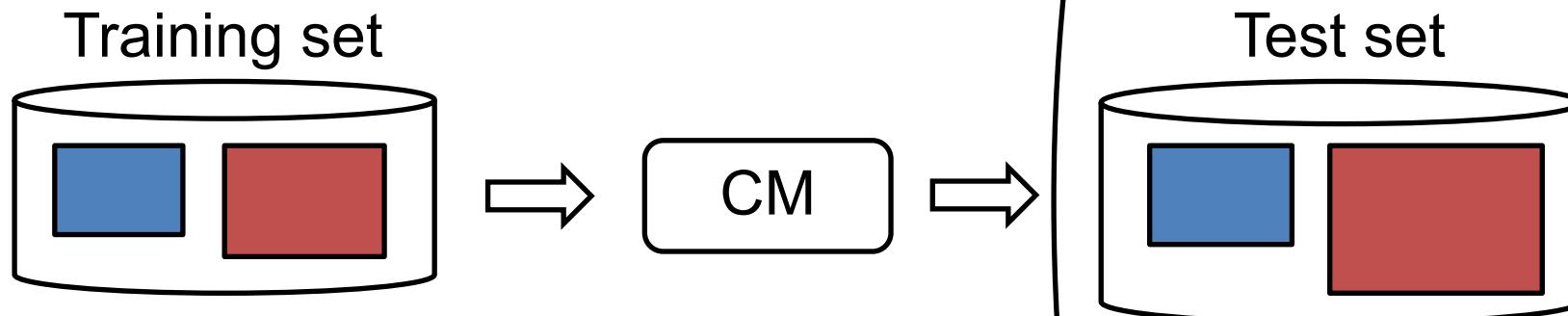
- a binary classification model

□ Model evaluation

- accuracy
- **equal error rate (EER)**
- **tandem detection cost function (t-DCF)** (Kinnunen 2020)

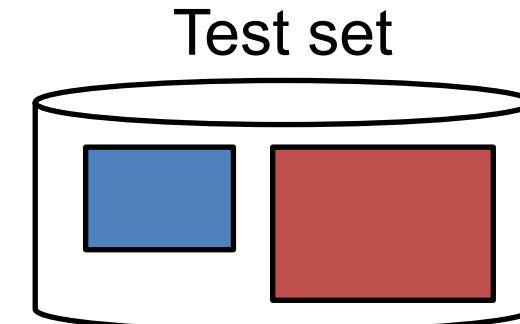


Speech anti-spoofing



- training and test sets from a database
 - ASVspoof www.asvspoof.org
 - BTAS2016
 - FMFCC-A
 - ...

Space of all possible bona fide and spoofed data



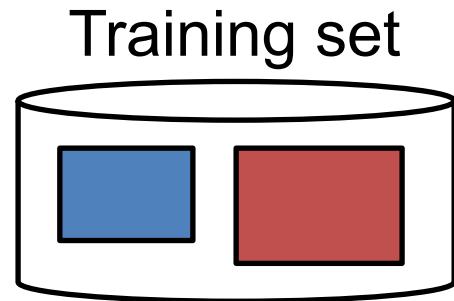
En, Fr, Ch, Jp, ...

Wav, mp3, m4a ...

New TTS/VC methods

Speech anti-spoofing

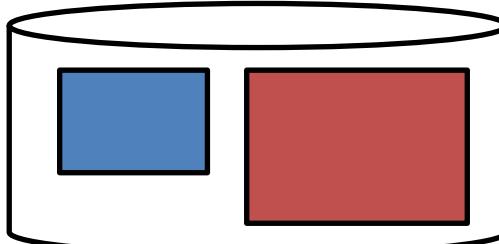
Just 25k utterances from 20 speakers!?



Test set

ASVspoof 2019 LA

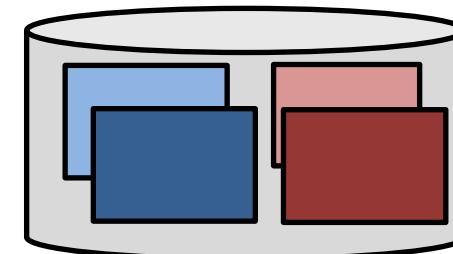
EER <8%



Test set 2

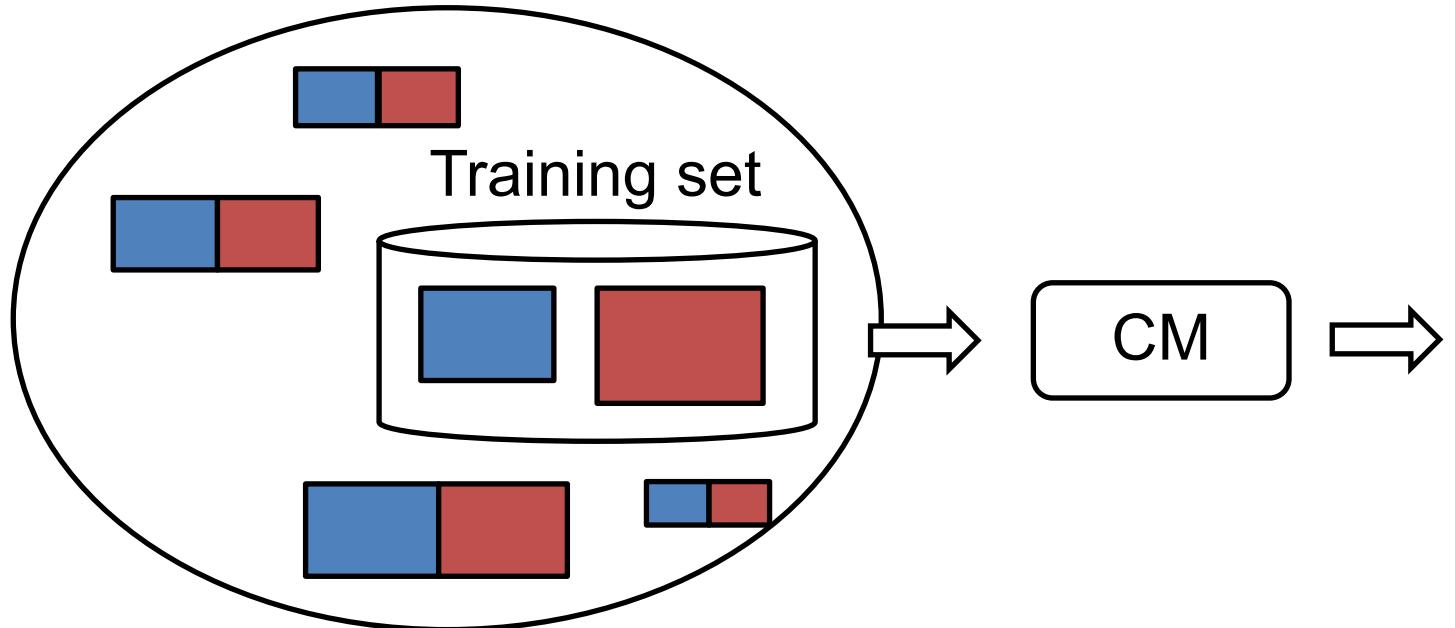
ASVspoof 2015

EER >20%



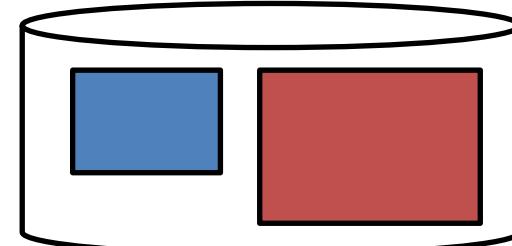
- poor generalization (Das 2020)
- dataset-specific bias?
 - biased dist. of non-speech (Müller 2021, Liu 2022)
 - artefacts in high-frequency band (Wang 2022)
 - ...

Speech anti-spoofing



Space of all possible bona fide and spoofed data

Test set



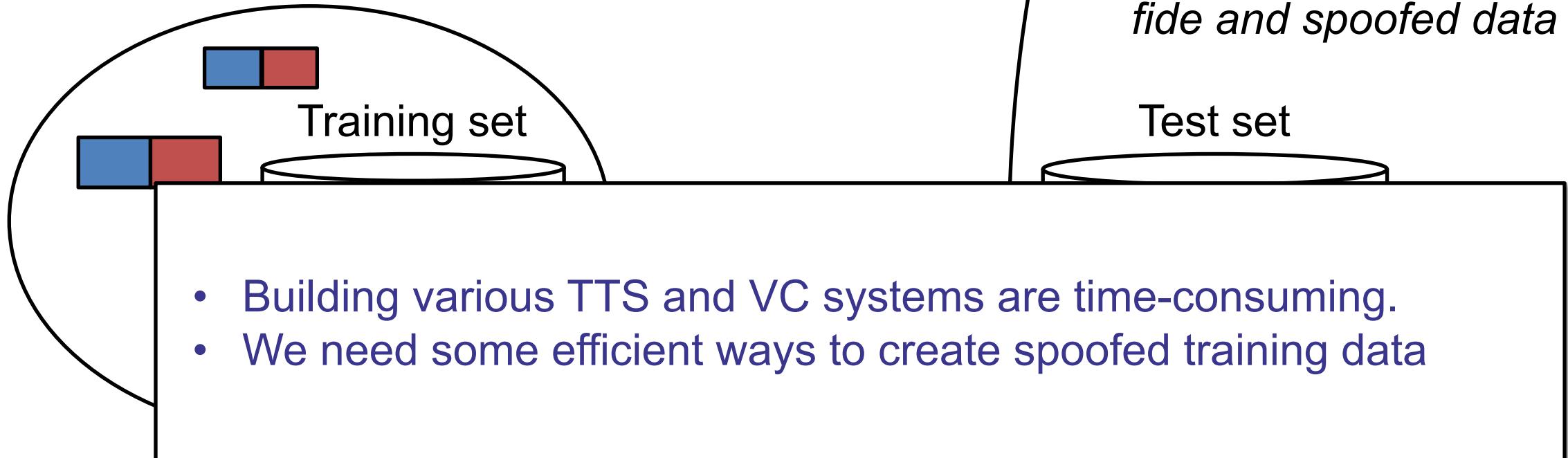
En, Fr, Ch, Jp, ...

Wav, mp3, m4a ...

New TTS/VC methods

- “Dead” data in self-supervised speech model
(Xie 2021, Wang 2022, Tak 2022, Donas 2022)
- Can we add more diverse training data?

Speech anti-spoofing



- Building various TTS and VC systems are time-consuming.
- We need some efficient ways to create spoofed training data

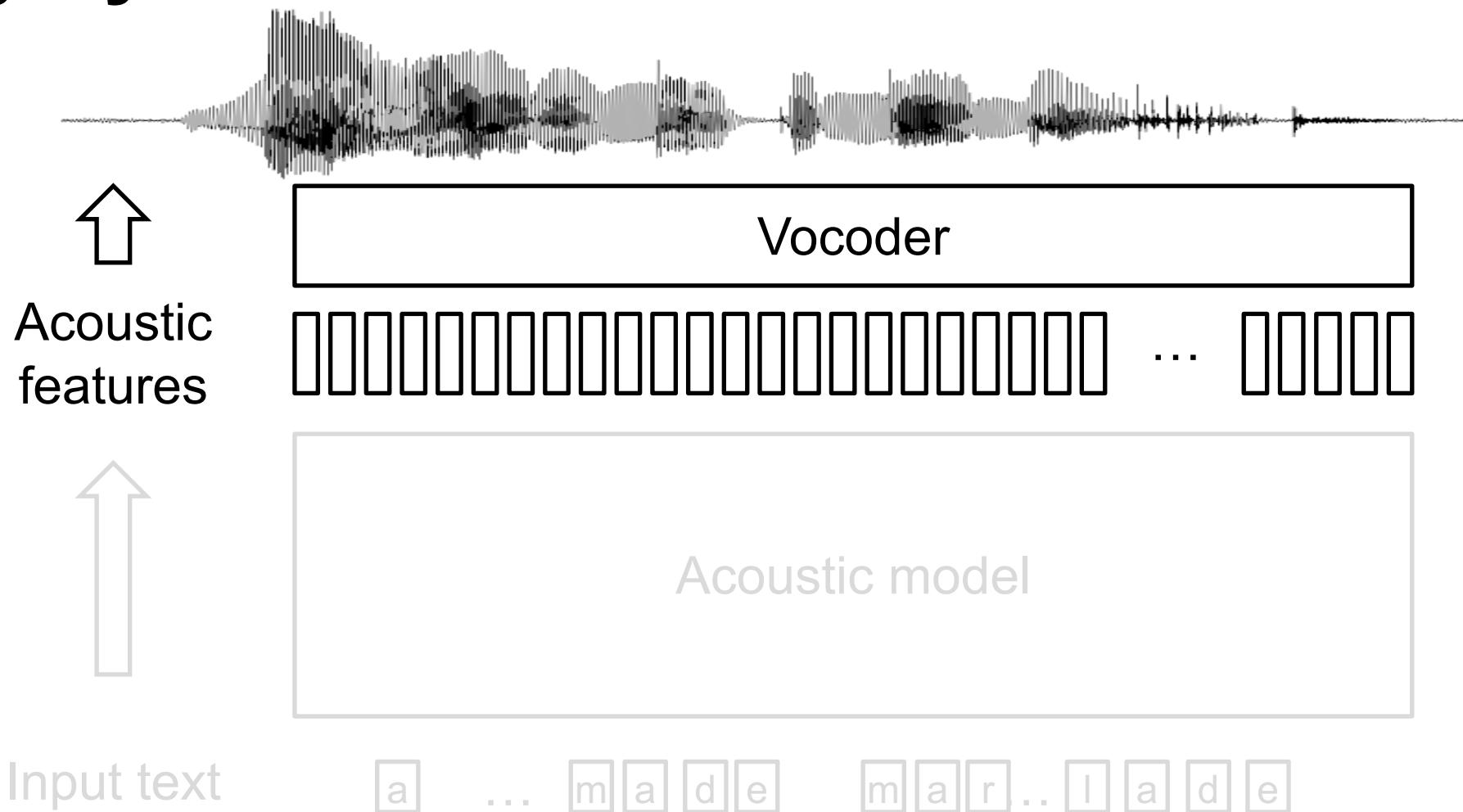
- “Dead” data in self-supervised speech model
(Xie 2021, Wang 2022, Tak 2022, Donas 2022)
- Can we add more diverse training data?

Wav, mp3, m4a ...

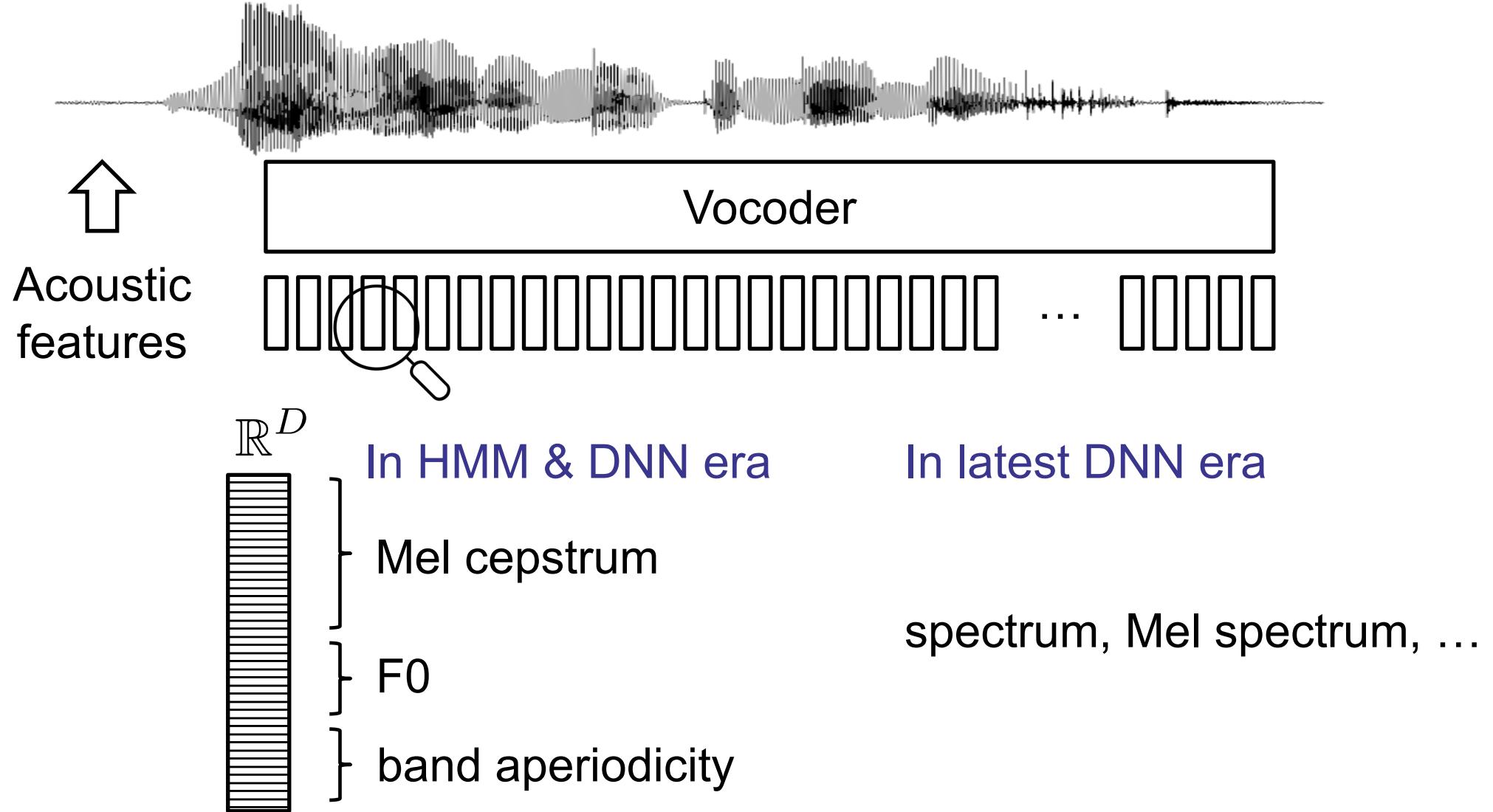
New TTS/VC methods

Method

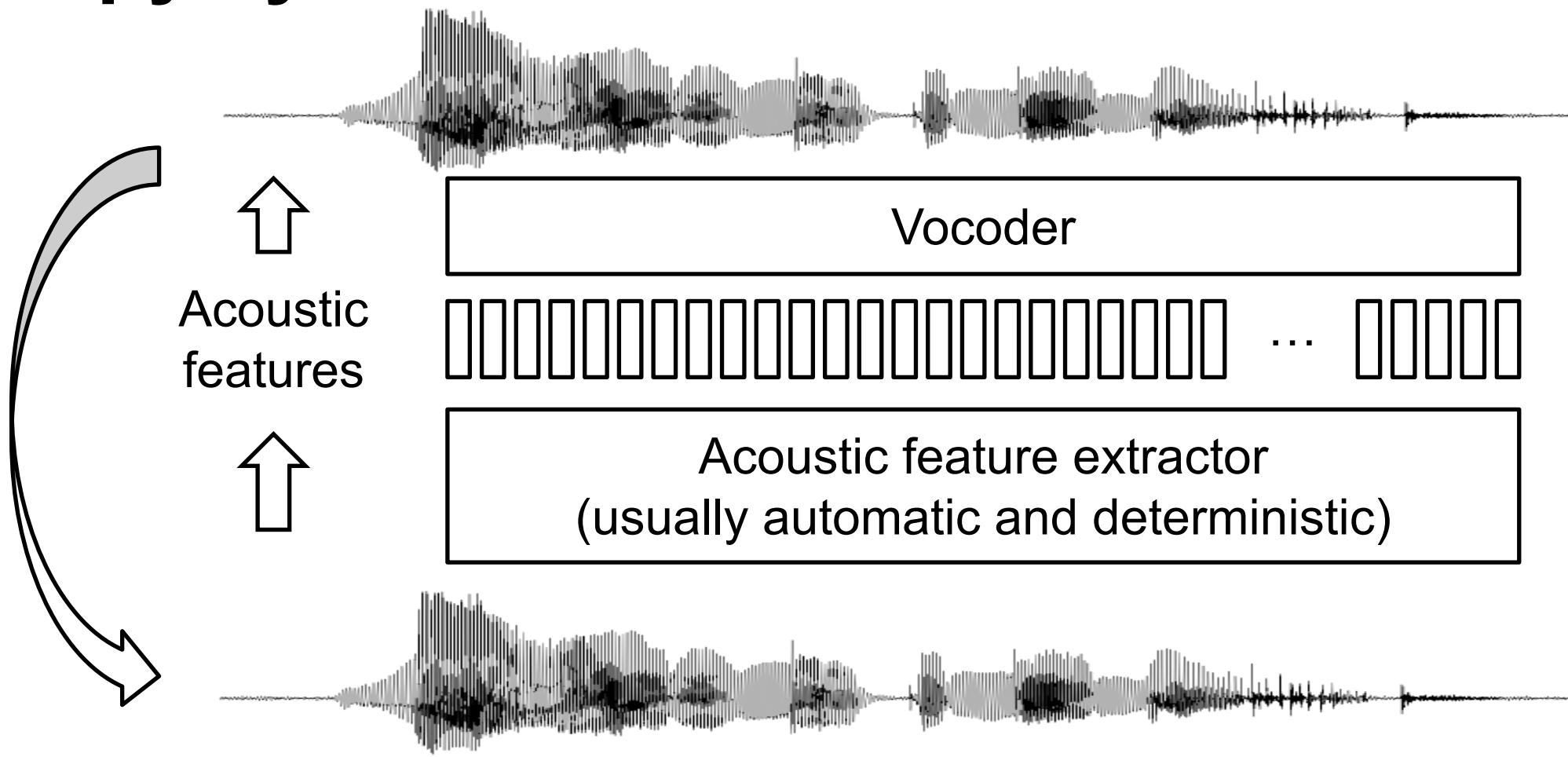
Copy synthesis



Copy synthesis

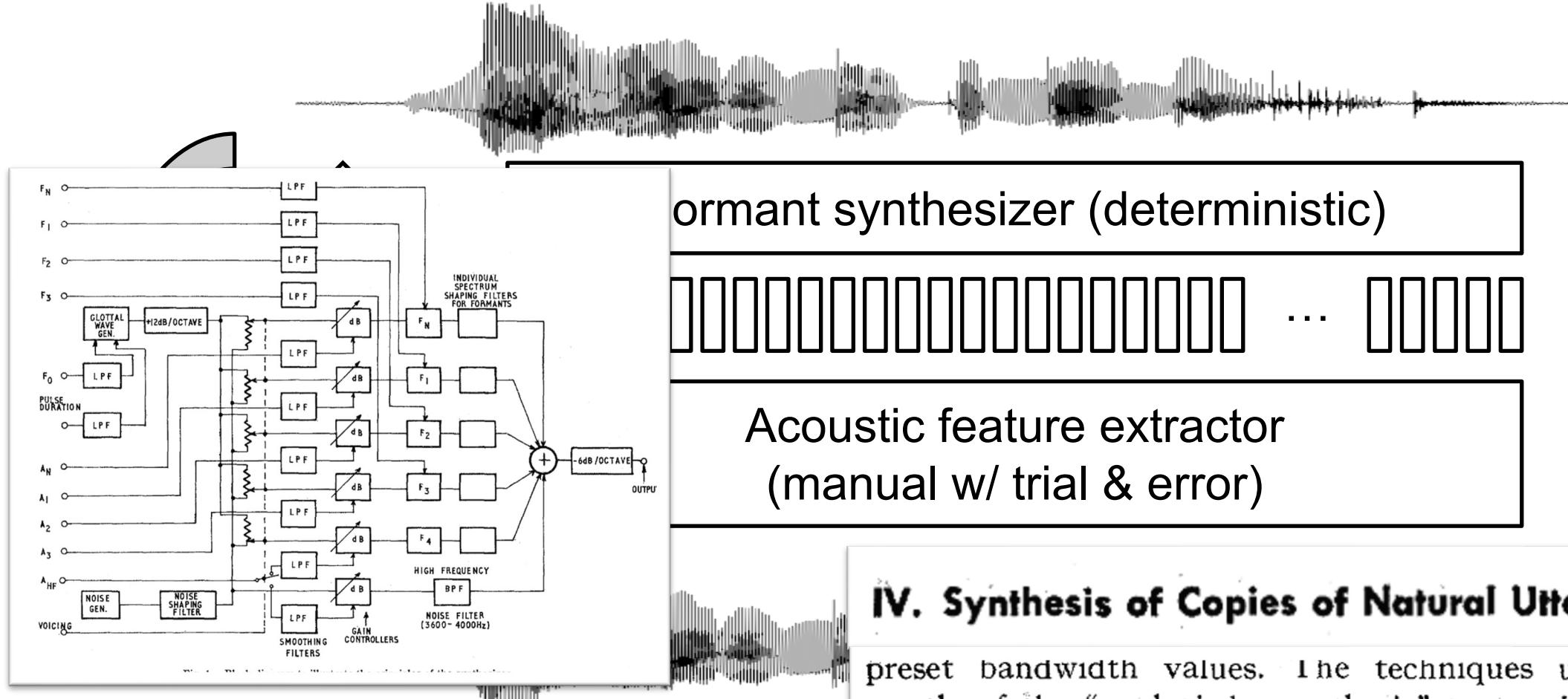


Copy synthesis



Copy-synthesis, analysis-by-synthesis, copy-resynthesis, vocoding ...

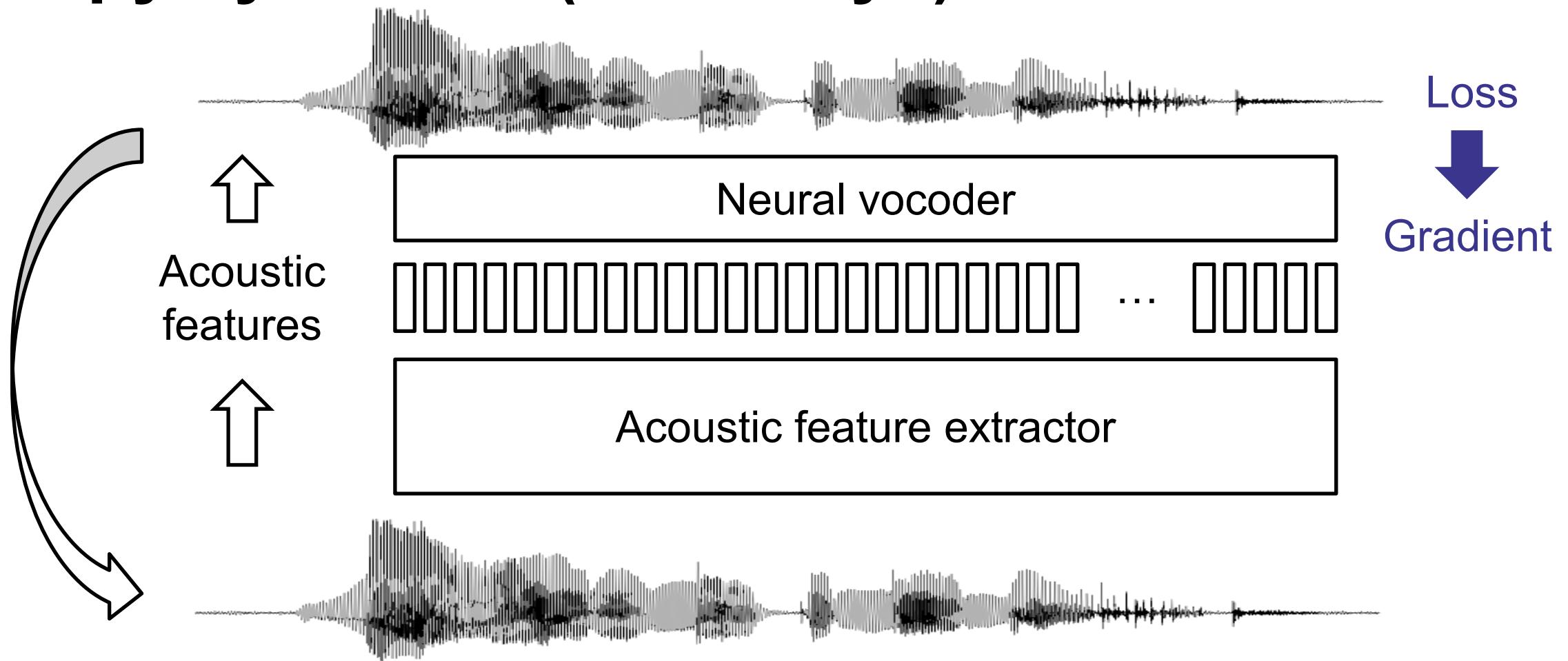
Copy synthesis (in history)



Copy-synthesis, analysis-by-syn

preset bandwidth values. The techniques used are mostly of the “analysis-by-synthesis” type, with a human interpreter of differences between natural and synthetic speech in the feedback loop.

Copy synthesis (nowadays)

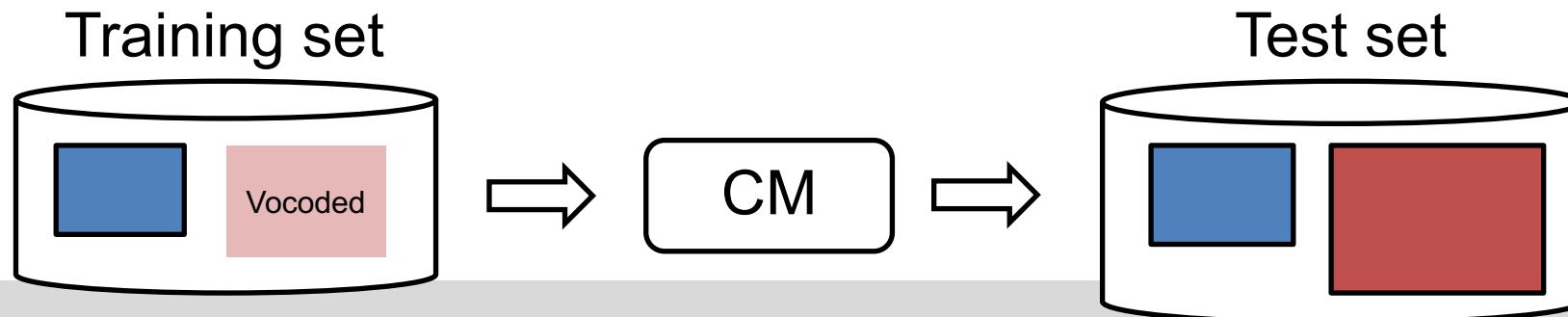
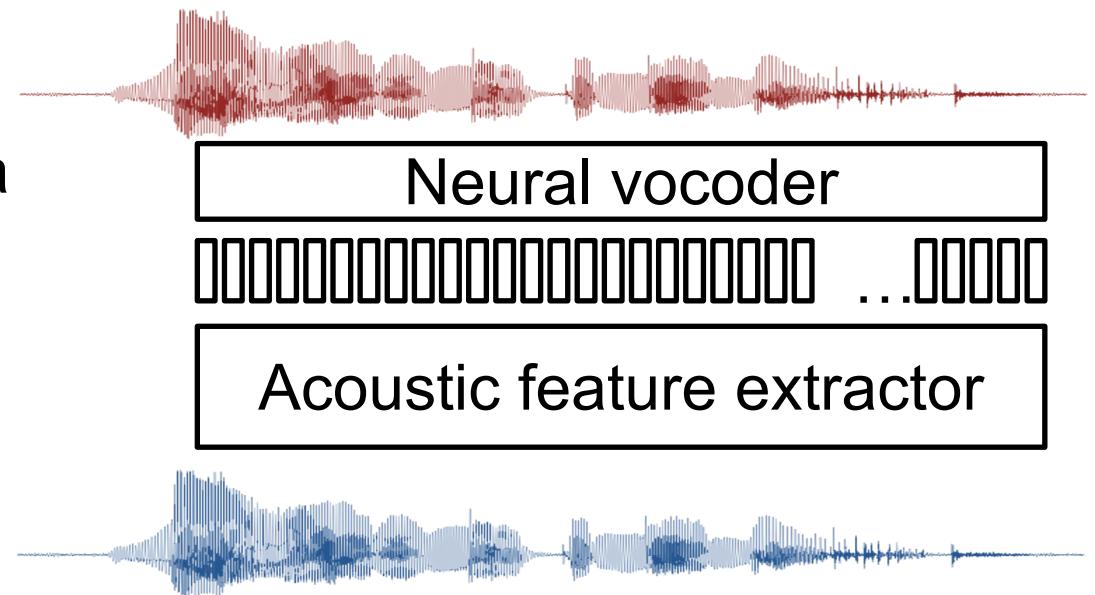


We do copy-synthesis when training the neural vocoders

Creating copy-synthesized spoofed data

□ Steps

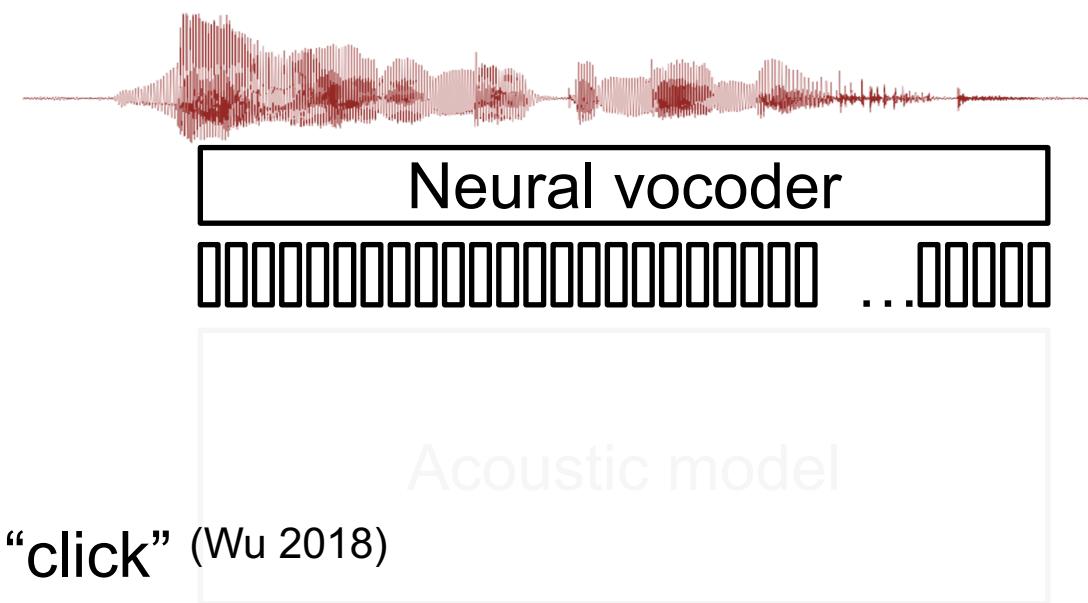
- prepare (or training) vocoders
 - not necessarily neural vocoders
- Do copy synthesis on **bona fide** data
- use output as **copy-synthesized (or vocoded) spoofed** data
- train a CM using {**bona fide, copy-synthesized spoofed**}



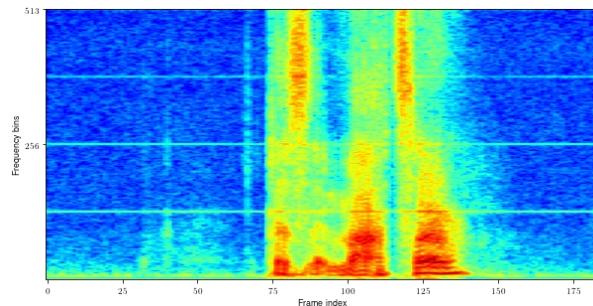
Creating copy-synthesized spoofed data

□ Hypothesis

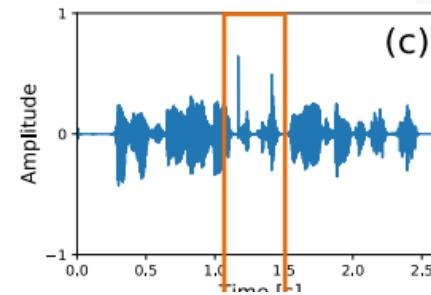
- *Copy-synthesis is TTS/VC with a perfect acoustic model*
 - ✗ artefacts by the acoustic model
 - ✓ artefacts by the vocoder
- TTS/VC spoofed data may contain artefacts by the vocoder



WaveGlow “bar” (Prenger 2019)



WaveNet “click” (Wu 2018)



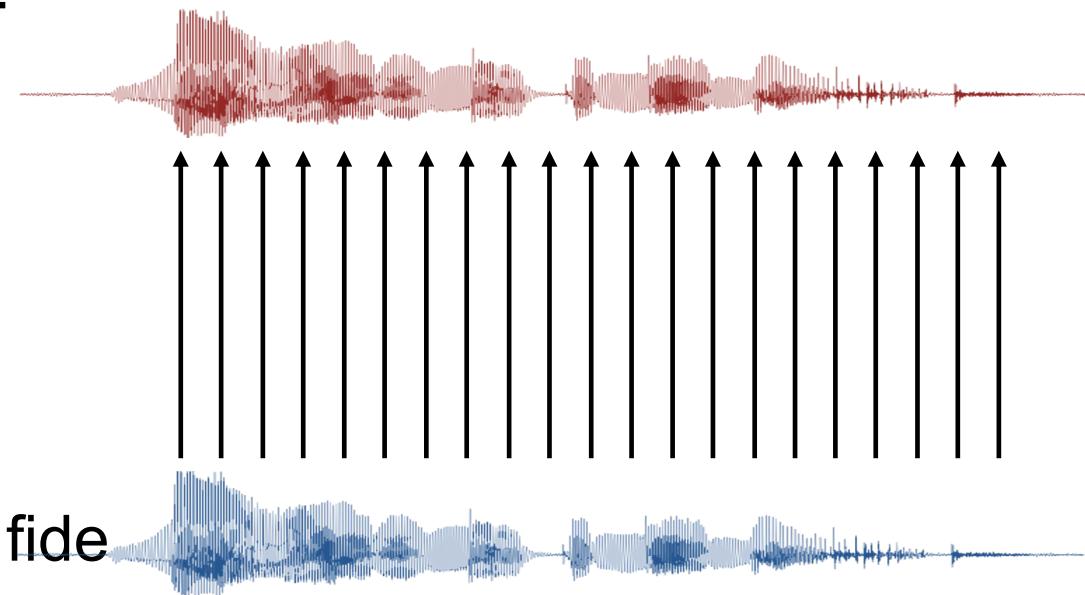
More artefacts not perceptible to ears

Creating copy-synthesized spoofed data

□ Questions

- How to prepare or train the vocoder?
 - pre-trained vocoder(s)?
 - Fine tuning?

Experiment I



- Can we better use the aligned bona fide and spoofed data pairs?

Experiment II

Creating copy-synthesized spoofed data

□ Related studies using DSP-based vocoders

- Xingming Wang, Xiaoyi Qin, Tinglong Zhu, Chao Wang, Shilei Zhang, and Ming Li. The DKU-CMRI System for the ASVspoof 2021 Challenge: Vocoder Based Replay Channel Response Estimation. In *Proc. ASVspoof challenge workshop*, 16–21. 2021.
- Monisankha Pal, Dipjyoti Paul, and Goutam Saha. Synthetic Speech Detection Using Fundamental Frequency Variation and Spectral Features. *Computer Speech & Language* 48. Elsevier: 31–50. 2018.
- Ibon Saratxaga, Jon Sanchez, Zhizheng Wu, Inma Hernaez, and Eva Navas. Synthetic Speech Detection Using Phase Information. *Speech Communication* 81 (July): 30–41. doi:10.1016/j.specom.2016.04.001. 2016.
- Aleksandr Sizov, Elie Khoury, Tomi Kinnunen, Zhizheng Wu, and Sébastien Marcel. Joint Speaker Verification and Antispoofing in the I-Vector Space. *IEEE Transactions on Information Forensics and Security* 10 (4). IEEE: 821–832. doi:10.1109/TIFS.2015.2407362. 2015.
- Elie Khoury, Tomi Kinnunen, Aleksandr Sizov, Zhizheng Wu, and Sébastien Marcel. Introducing I-Vectors for Joint Anti-Spoofing and Speaker Verification. In *Proc. Interspeech*, 61–65. 2014.
- Jon Sanchez, Ibon Saratxaga, Inma Hernaez, Eva Navas, and Daniel Erro. A Cross-Vocoder Study of Speaker Independent Synthetic Speech Detection Using Phase Information. In *Proc. Interspeech*. 2014.
- Zhizheng Wu, Xiong Xiao, Eng Siong Chng, and Haizhou Li. Synthetic Speech Detection Using Temporal Modulation Feature. In *Proc. ICASSP*, 7234–7238. 2013.

□ Related studies using neural vocoders

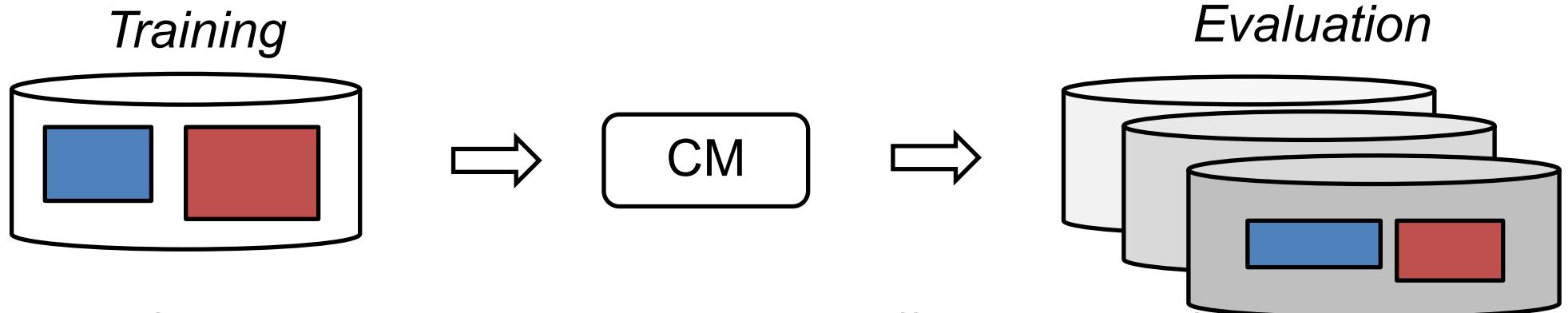
- Joel Frank, and Lea Schönherr. **WaveFake**: A Data Set to Facilitate Audio DeepFake Detection. In *Proc. NeurIPS Datasets and Benchmarks 2021*. 2021.
- Chengzhe Sun, Shan Jia, Shuwei Hou, Ehab AlBadawy, and Siwei Lyu. Exposing AI-Synthesized Human Voices Using Neural Vocoder Artifacts. ArXiv Preprint ArXiv:2302.09198. 2023.

Experiment I

How to prepare or train the vocoder?

Experiment I

□ Design

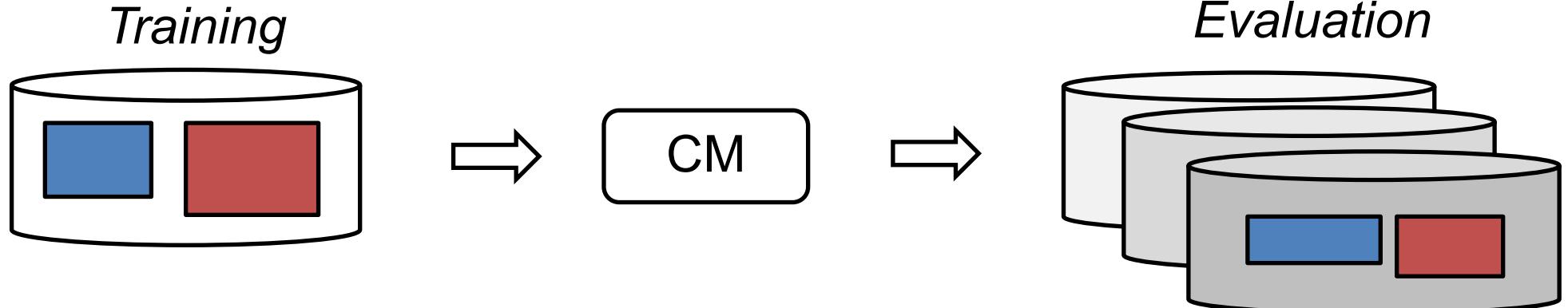


- changed factor: training data created by different sets of vocoders
- unchanged: CM (Wang 2022) using a Wav2vec2.0-based front end (Baevski 2020)
- unchanged: multiple test sets

- three independent training & evaluation rounds
- averaged EER

Experiment I

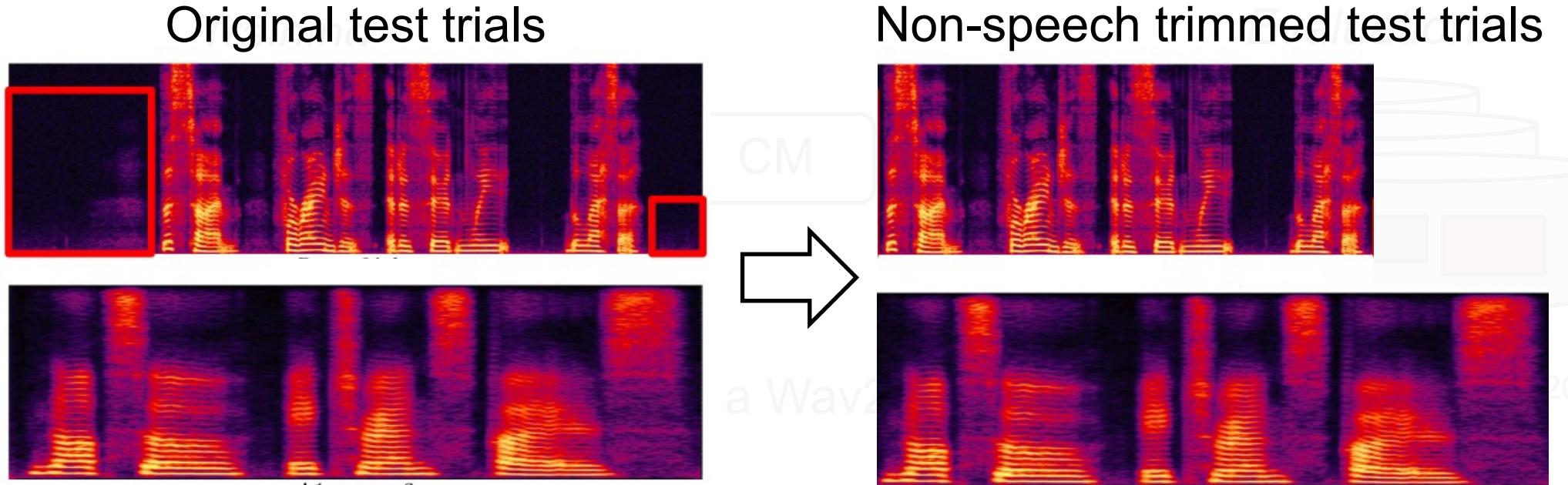
□ Design



- changed factor: training data created by different sets of vocoders
- unchanged: CM (Wang 2022) using a Wav2vec2.0-based front end (Baevski 2020)
- unchanged: multiple test sets
 - ASVspoof 2019 LA test set, 2021 LA&DF eval track, 2015 test set
 - ASVspoof 2019 LA test set w/o non-speech, 2021 LA & DF hidden track
 - WaveFake^(Frank 2021), In-the-Wild (Müller 2022)

Experiment I

□ Design



- ASVspoof 2019 LA test set, 2021 LA&DF eval track, 2015 test set
- ASVspoof 2019 LA test set w/o non-speech, 2021 LA & DF hidden track

We recommend testing on both versions of ASVspoof test sets

Experiment I

☐ Training data

ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data	Vocoder SR
LA19trn	20	2,580	22,800	-	-	-	16 kHz
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -	24 kHz
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -	24 kHz
Voc.v2	20 same as	2,580	10,320	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / -	16 kHz
Voc.v3	LA19trn			HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LA19trn bona. / -	16 kHz
Voc.v4				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / LA19trn bona.	16 kHz

- LA19trn: ASVspoof 2019 LA training set
- WFtrn: WaveFake English subset, down-sampled to 16kHz
- Voc.v*: ASVspoof 2019 LA training set bona fide data + its vocoded data

Experiment I

😊 Low EER

😢 High EER

☐ Results in EER (%, mean of three runs)

	LA19 trn	WF trn	Training set			
	Voc. v1	Voc. v2	Voc. v3	Voc. v4		
LA19eval	2.98	44.48	5.78	5.32	8.74	4.36
LA21eval	7.53	41.57	26.30	17.98	19.29	24.39
DF21eval	6.67	24.26	11.95	11.54	9.71	13.31
Test sets	LA19etrim	15.56	31.62	23.29	16.16	14.99
	LA21hid	28.80	27.60	28.30	19.49	17.62
	DF21hid	23.62	26.18	22.01	13.92	13.50
	WaveFake	15.76	-	39.27	34.05	17.10
	InWild	26.65	19.98	41.06	36.46	22.26
Pooled	14.24	-	36.57	39.95	19.39	16.35

↑ Single threshold

Experiment I

□ Results in EER (%, mean of three runs)

	LA19 trn	WF trn	Training set			
			Voc. v1	Voc. v2	Voc. v3	Voc. v4
Test sets	LA19eval	2.98	14.48	5.78	5.32	8.74
	LA21eval	7.53	11.30	10.99	10.90	4.36
	DF21eval	6.67	24.26	11.95	11.54	9.71
	LA19etrim	15.56	31.62	23.29	16.16	14.99
	LA21hid	28.80	27.60	28.30	19.49	17.62
	DF21hid	23.62	26.18	22.01	13.92	13.50
	WaveFake	15.76	39.37	34.05	17.10	10.89
	InWild	26.65	19.98	36.46	22.26	19.45
	Pooled	14.24	-	36.57	39.95	19.39
Pooled: 14.24						
Similar results to our previous work (Wang 2022)						
Same trend as previous studies (Müller 2021, Liu 2022)						
Not 50% :)						

ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data	Vocoder SR
LA19trn	20	2,580	22,800	-	-	-	16 kHz
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -	24 kHz
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -	24 kHz
Voc.v2	²⁰ same as Voc.v3 LA19trn	2,580	10,320	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / -	16 kHz
Voc.v4				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LA19trn bona. / -	16 kHz
				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / LA19trn bona.	16 kHz

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
	LA19eval	2.98	44.48	5.78	5.32	8.74	4.36
	LA21eval	7.53	41.57	26.30	17.98	19.29	24.39
	DF21eval	6.67	24.26	11.95	11.54	9.71	13.31
Test sets	LA19etrim	15.56	31.62	23.29	Vocoders pre-trained by ESPNet (Hayashi 2020)		
	LA21hid	28.80	27.60	28.30			
	DF21hid	23.62	26.18	22.01			
	WaveFake	15.76	-	39.27			
	InWild	26.65	19.98	41.06			
	Pooled	14.24	-	36.57	39.95	19.39	16.35

ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data	Vocoder SR
LA19trn	20	2,580	22,800	-	-	-	16 kHz
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -	24 kHz
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -	24 kHz
Voc.v2	²⁰ same as Voc.v3 LA19trn	2,580	10,320	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / -	16 kHz
Voc.v3				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LA19trn bona. / -	16 kHz
Voc.v4				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / LA19trn bona.	16 kHz

	LA19 trn	WF trn	Voc. v1	Training set		
				Voc. v2	Voc. v3	Voc. v4
LA19eval	2.98	44.48	5.78	5.32	8.74	4.36
LA21eval	7.53	41.57	26.30	17.98	19.29	24.39
DF21eval	6.67	24.26	11.95	11.54	9.71	13.31
Test sets	LA19etrim	15.56	31.62	23.29	16.16	14.99
	LA21hid	28.80	27.60	28.30	19.49	17.62
	DF21hid	23.62	26.18	22.01	13.92	13.50
	WaveFake	15.76	-	39.27	34.05	17.10
	InWild	26.65	19.98	41.06	36.46	22.26
Pooled		14.24	-	36.57	39.95	16.35

Was vocoder trained on
the bona fide data?

No

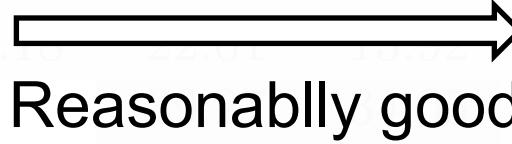
Yes

ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data	Vocoder SR
LA19trn	20	2,580	22,800	-	-	-	16 kHz
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -	24 kHz
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -	24 kHz
Voc.v2	²⁰ same as Voc.v3 LA19trn	2,580	10,320	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / -	16 kHz
Voc.v4				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LA19trn bona. / -	16 kHz
				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / LA19trn bona.	16 kHz

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
Test sets	LA19eval	2.98					4.36
	LA21eval	7.53					24.39
	DF21eval	6.67					13.31
	LA19etrim	15.56					9.52
	LA21hid	28.80					21.43
	DF21hid	23.62					16.99
	WaveFake	15.76					10.89
	InWild	26.65					19.45
Pooled		14.24					16.35



We cannot exploit non-speech length



Reasonably good

Experiment I

□ How to prepare or train the vocoder?

- Pre-trained vocoders may not work --- WFtrn, Voc.v1
- It is better to fine tune vocoders on the bona fide to be copy-synthesized

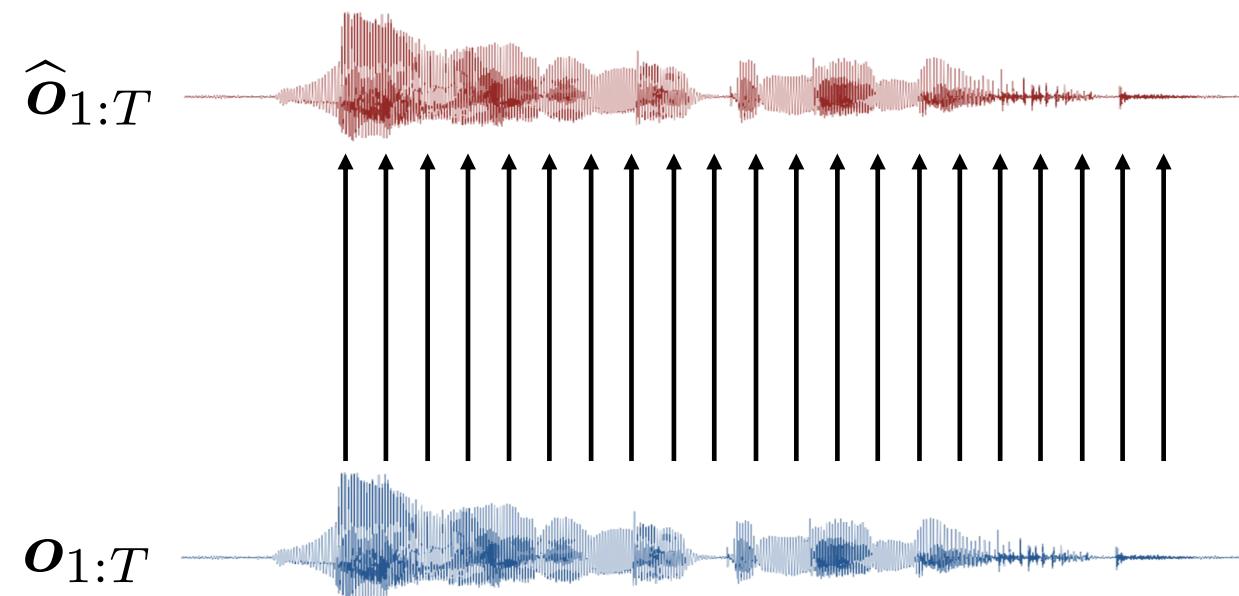
ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data	Vocoder SR
LA19trn	20	2,580	22,800	-	-	-	16 kHz
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -	24 kHz
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -	24 kHz
Voc.v2	same as LA19trn	20	2,580	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / -	16 kHz
Voc.v3	LA19trn	2,580	10,320	HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LA19trn bona. / -	16 kHz
Voc.v4				HiFiGAN, NSF-HiFiGAN, Hn-NSF, WaveGlow	in-house	LibriTTS / LA19trn bona.	16 kHz

□ CAUTION: other CMs (e.g., LCNN and AASIST) did not work well on vocoded data

- See results in Appendix, <https://arxiv.org/abs/2210.10570>

Experiment II

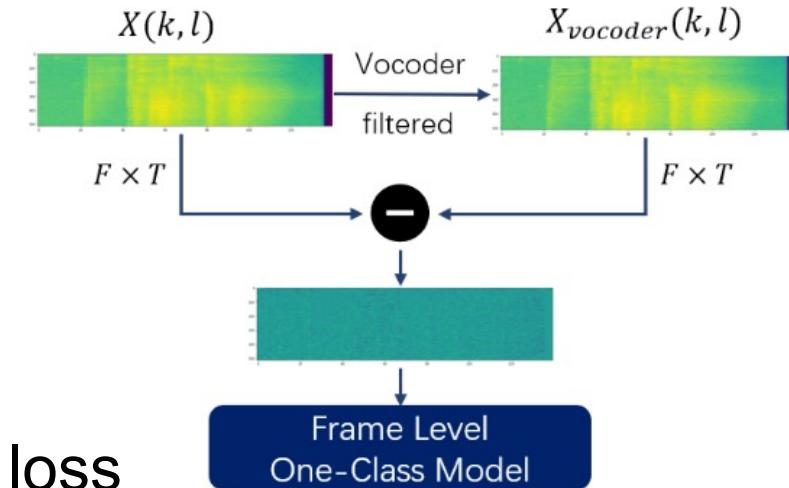
How to make good use of the aligned bona fide and vocoded data pair?



Experiment II

☐ Method: make use of the aligned bona fide and vocoded pair

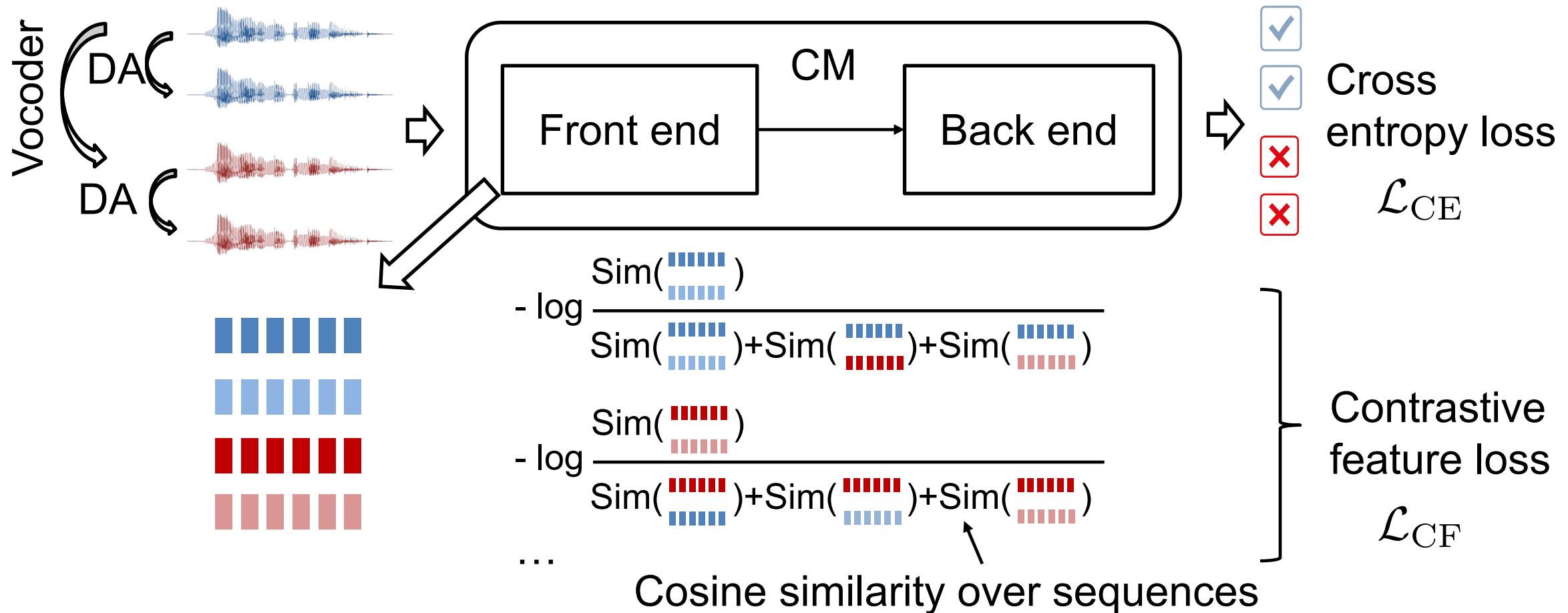
- Existing method: use their differences in frequency domain (Wang 2021)
 - Vocoder are needed during inference
 - Too slow to score the test sets



- We proposed an auxiliary contrastive feature loss
 - It is based on supervised contrastive loss (Khosla 2020)
 - It needs data augmentation (DA)
 - No need to run vocoders during inference

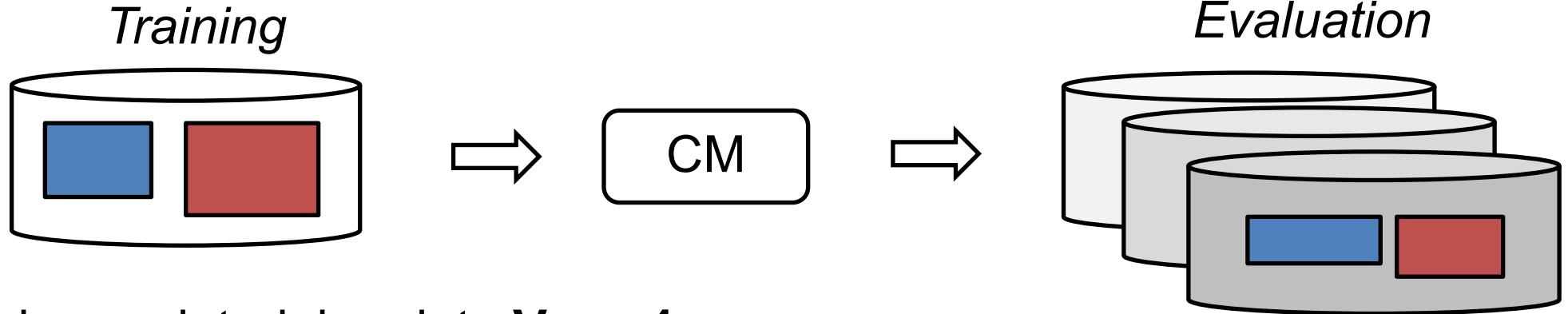
Experiment II

Method: an auxiliary contrastive feature loss



Experiment II

□ Design



- unchanged: training data **Voc.v4**
- unchanged: same CM architecture as Experiment I
- unchanged: multiple test sets

- changed: training criterion and method
 - data augmentation based on RawBoost^(Tak 2022)

Experiment II

□ Results

	from Experiment I		Control groups		Best	
Training criterion	\mathcal{L}_{CE}				$\mathcal{L}_{CE} + \mathcal{L}_{CF}$	
Data augmentation	×	RawBoost		RawBoost		
Training set	LA19 trn	Voc. v4	LA19 trn	Voc. v4	LA19 trn	Voc. v4
Bona-spoof paired	×	×	×	×	×	✓
ID	①	②	③	④	⑤	⑦
LA19eval	2.98	4.36	0.22	3.46	0.21	2.63
LA21eval	7.53	24.39	3.63	16.55	3.30	16.67
DF21eval	6.67	13.31	3.65	9.60	4.12	6.92
Test sets	LA19etrim	15.56	9.52	9.16	6.09	9.00
	LA21hid	28.80	21.43	21.18	19.37	26.98
	DF21hid	23.62	16.99	13.64	14.29	16.85
	WaveFake	15.76	10.89	26.37	6.87	24.62
	InWild	26.65	19.45	16.17	12.08	17.07
	Pooled	14.24	16.35	13.12	13.13	13.68
						11.27

② vs ④
RawBoost is useful

④ vs ⑦
Contrastive feature loss is useful

Experiment II

□ Results

	from Experiment I		Control groups		Best		
	\mathcal{L}_{CE}		$\mathcal{L}_{CE} + \mathcal{L}_{CF}$				
Training criterion	LA19	Voc. v4	LA19	Voc. v4	LA19	Voc. v4	Voc. v4
Data augmentation	×		RawBoost		RawBoost		
Training set							
Bona-spoof paired	×	×	×	×	×	×	✓
ID	①	②	③	④	⑤	⑥	⑦
LA19eval	2.98	4.36	0.22	3.46	0.21	2.63	2.21
LA21eval	7.53	24.39	3.63	16.55	3.30	16.67	17.90
DF21eval	6.67	13.31	3.65	9.60	4.12	6.92	5.04
Test sets							
LA19etrim	15.56	9.52	9.16	6.09	9.00	4.48	3.79
LA21hid	28.80	21.43	21.18	19.37	26.98	15.05	14.57
DF21hid	23.62	16.99	13.64	14.29	16.85	8.17	7.78
WaveFake	15.76	10.89	26.37	6.87	24.62	4.03	2.50
InWild	26.65	19.45	16.17	12.08	17.07	9.37	7.55
Pooled	14.24	16.35	13.12	13.13	13.68	13.15	11.27

- ② vs ④
RawBoost is useful
- ④ vs ⑦
Contrastive feature loss is useful
- ⑥ vs ⑦
Use aligned {bona, vocoded} pairs!

Summary

Summary

- ❑ Spooferd training data can be created using neural vocoders
- ❑ Recommendations (from this study)
 - Fine-tune vocoders on the bona fide data to be copy-synthesized
 - Exploit the aligned {bona fide, vocoded} pairs
 - for example, by using contrastive feature loss
 - The best trained CM showed promising generalization performance

Summary

❑ Additional findings (see <https://arxiv.org/abs/2210.10570>)

- Do we generalize to old TTS and VC?
 - Yes, EER on the ASVspoof 2015 test set < 1%
- Is the CM sufficiently generalizable?
 - No, TTS/VC with autoregressive (AR) vocoders are challenging

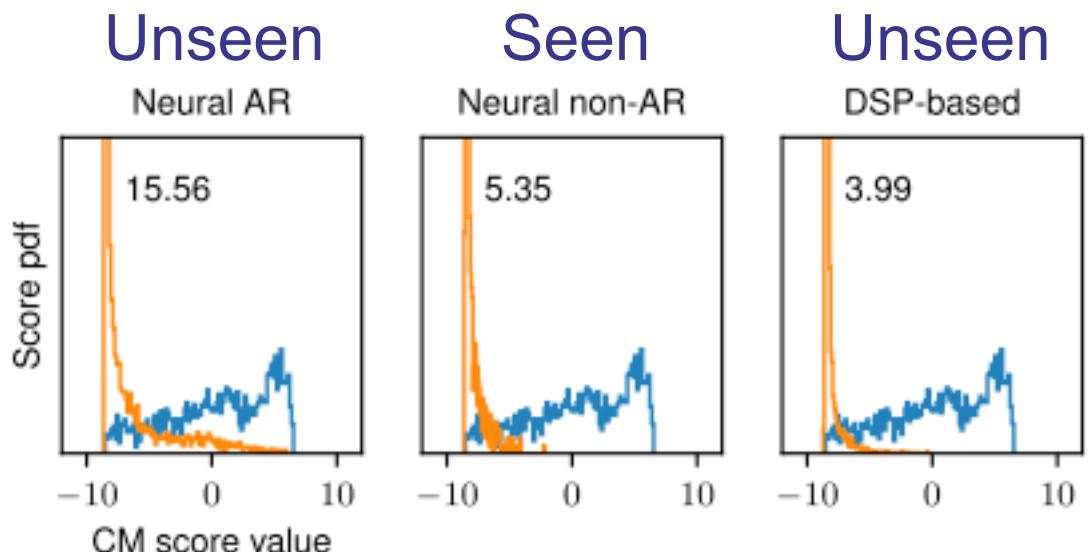


Fig. 1: Score distributions of \mathcal{I} on **bona fide** and **spoofed** trials in DF21hid. Number in each sub-figure is EER (%).

In progress

☐ Scale to large data

- CM: same as Experiment II
- ASVspoof 2019 LA trn bona fide + vocoded data: 2 * 5 hours
- VoxCeleb2 dev set + vocoded data: 2300 * 5 hours

Training steps	ASVspoof 2019 trn			VoxCeleb2 dev			
	N	0.3N	0.6N	1.2N	2.4N	6.0N	12.0N
LA19eval	2.21	8.20	7.04	5.40	6.53	5.40	5.63
LA21eval	17.90	20.19	16.73	14.33	18.10	17.44	17.84
DF21eval	5.04	7.49	5.41	5.39	6.00	5.65	5.64
LA19etrim	3.79	6.17	5.53	5.16	5.25	5.22	5.14
LA21hid	14.57	13.98	12.34	11.47	11.64	11.37	11.43
DF21hid	7.78	11.02	9.71	9.90	10.05	9.99	10.04
WaveFake	2.50	14.94	10.39	8.38	5.52	4.88	4.94
InWild	7.55	16.12	15.63	14.19	13.43	13.32	13.77
Pooled	11.27	13.52	11.79	9.98	9.01	8.36	8.37

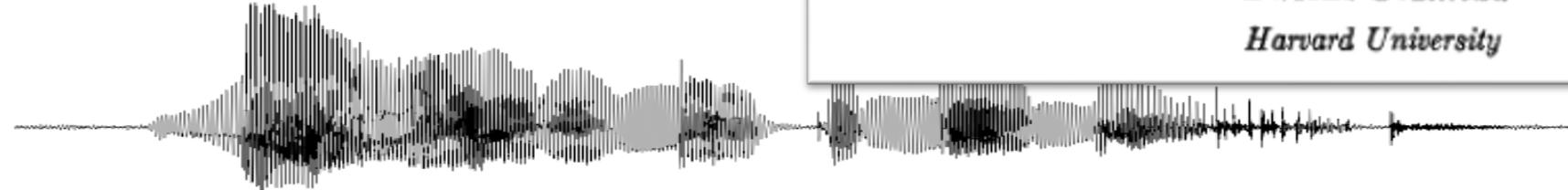
Resources

- **Code, vocoded data, vocoders, and trained CMs** [Git](#)
- **Tutorials on neural vocoders (AR, flow, GAN, DSP ...)** [Git](#)
 - Jupyter notebooks & pre-trained models
- **ASVspoof 2021 hidden track**
 - They are already in the data packages you've downloaded
 - You just need to find them using the official meta labels
 - [See https://github.com/asvspoof-challenge/2021](https://github.com/asvspoof-challenge/2021)

Thank you!

Appendix

Why TTS is difficult?



Speaker A: Who made the marmalade.

Speaker A: Bob made the marmalade.

Speaker B: (No,) Mari ^{anna} made the marmalade.

ACCENT IS PREDICTABLE (IF YOU'RE A MIND-READER)

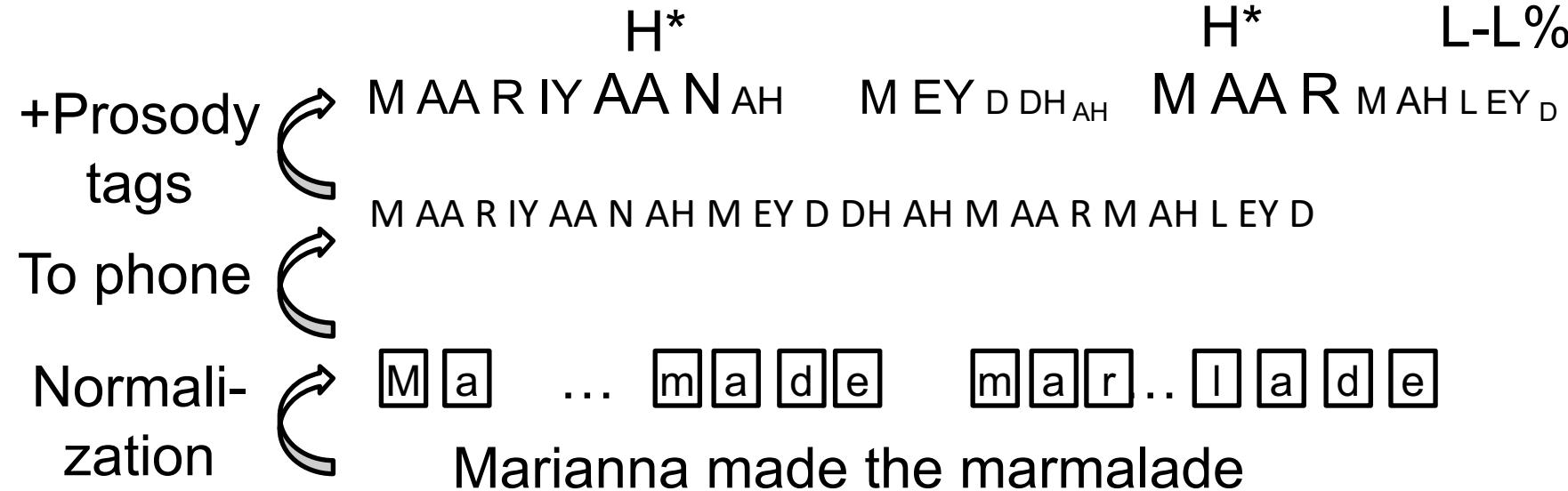
DWIGHT BOLINGER

Harvard University

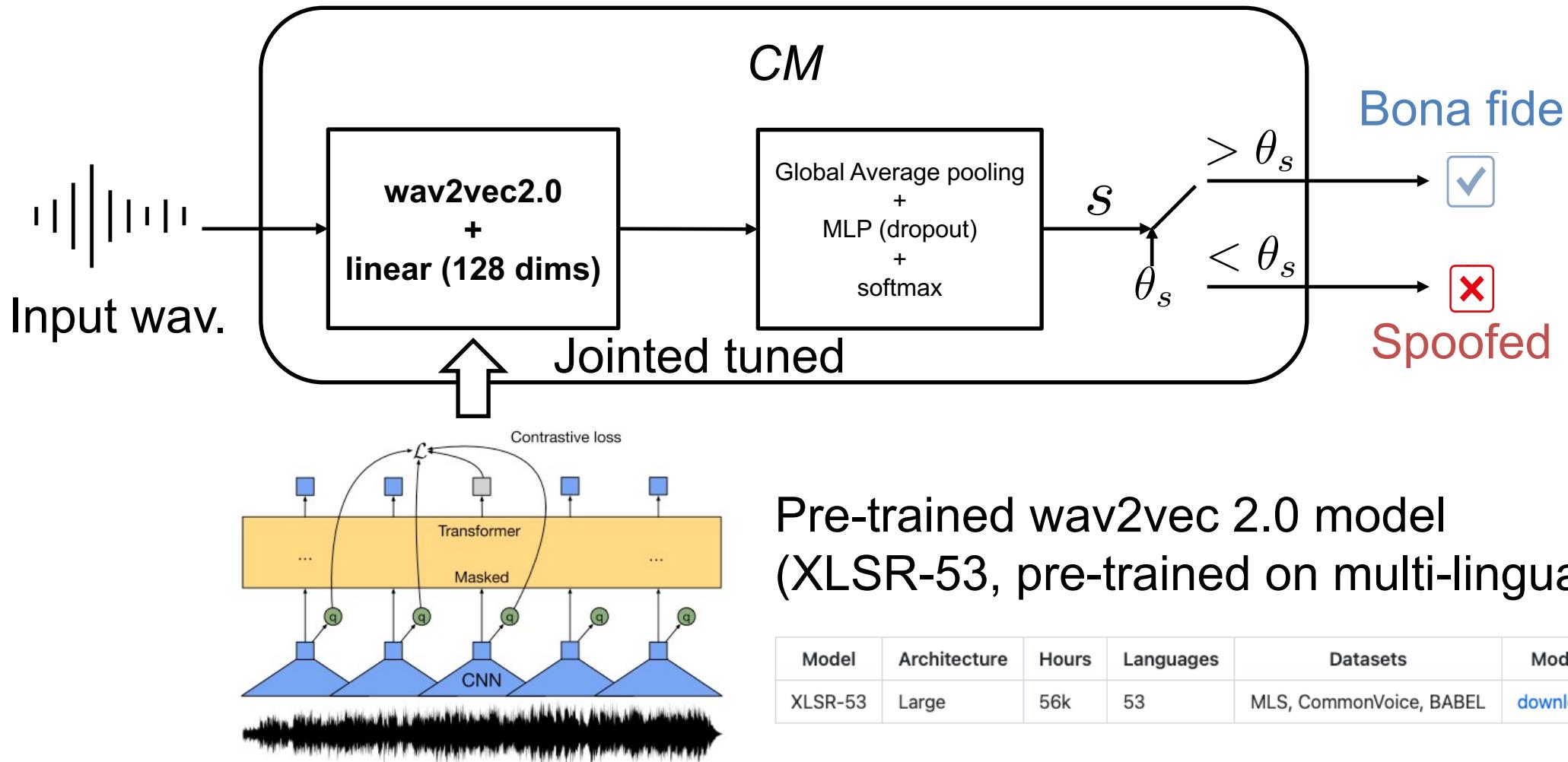


Speaker B: Marianna made the mar malade.

Speaker B: Marianna made the mar malade.



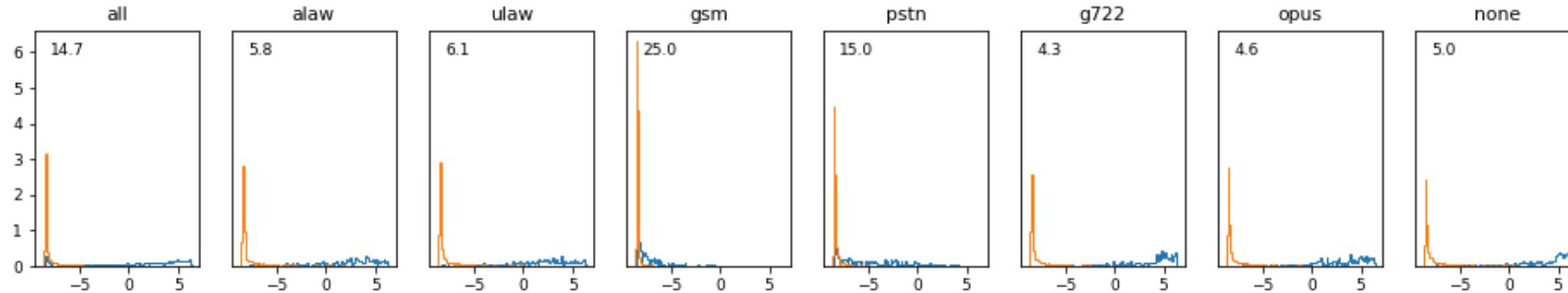
CM architecture



Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. In *Proc. NIPS*, 33:12449–12460. 2020.
<https://github.com/pytorch/fairseq/tree/main/examples/wav2vec>

Analysis

□ System 7 on LA 2021



□ System 7 on DF 2021

