



Inter-University Research Institute Corporation /  
Research Organization of Information and Systems

**National Institute of Informatics**

ICASSP 2024  
SLP.L20.2

# Can Large-scale Vocoded Spoofed Data Improve Speech Spoofing Countermeasure with a Self-supervised Front End?

/wʌn/ /ʃɪn/

Wang Xin, Junichi Yamagishi

National Institute of Informatics

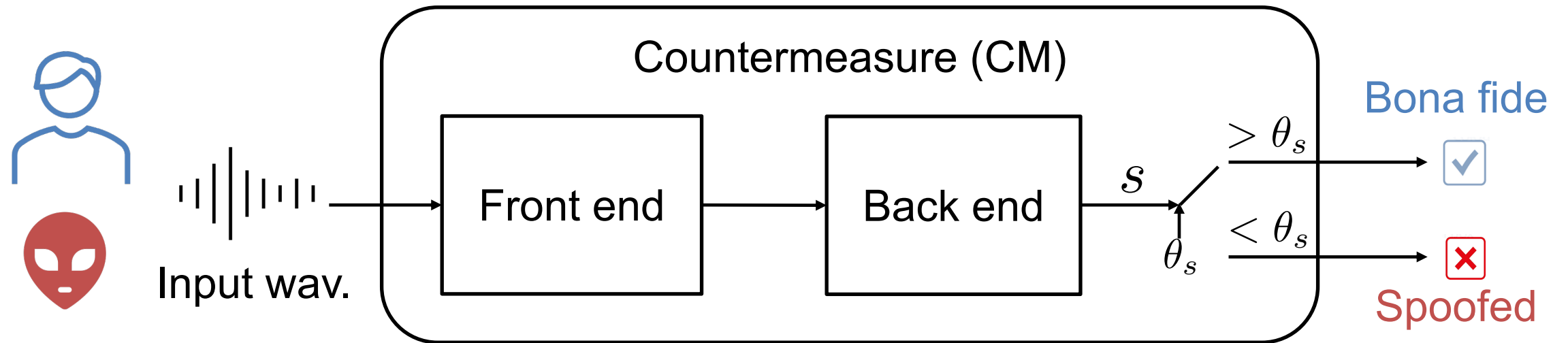
# Summary in one slide

- ❑ Extension of previous work (Wang 2023)
  - Not use any spoofed training data from text-to-speech or voice conversion
- ❑ Method
  - Upstream SSL training, using **vocoded VoxCeleb2**
  - Downstream SSL fine-tuning, using vocoded ASVspoof19
- ❑ Our best overall results

# Introduction

## □ A binary classification task

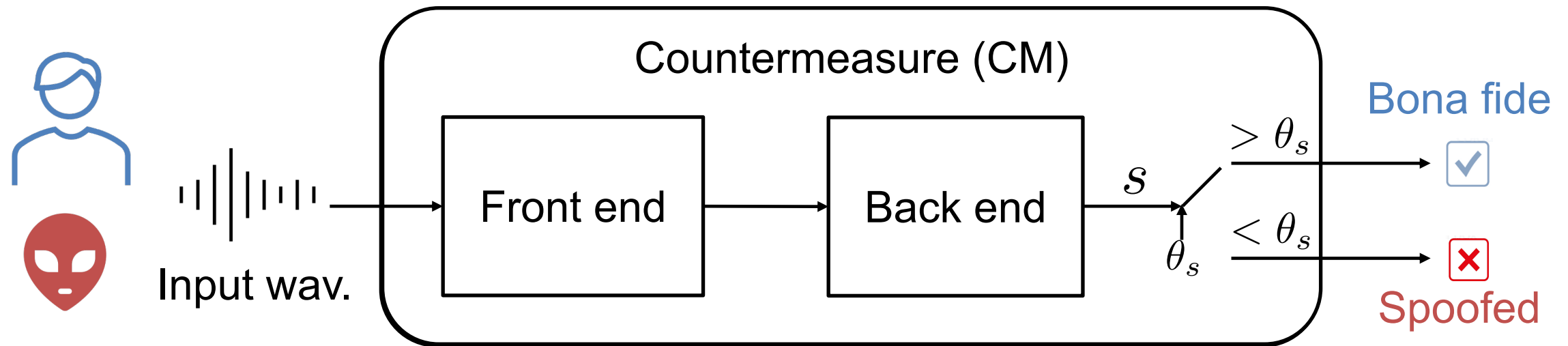
- Bona fide: human voice
- Spoofed: text-to-speech (TTS) or voice conversion (VC) voice
  
- Metric: equal error rate (EER)



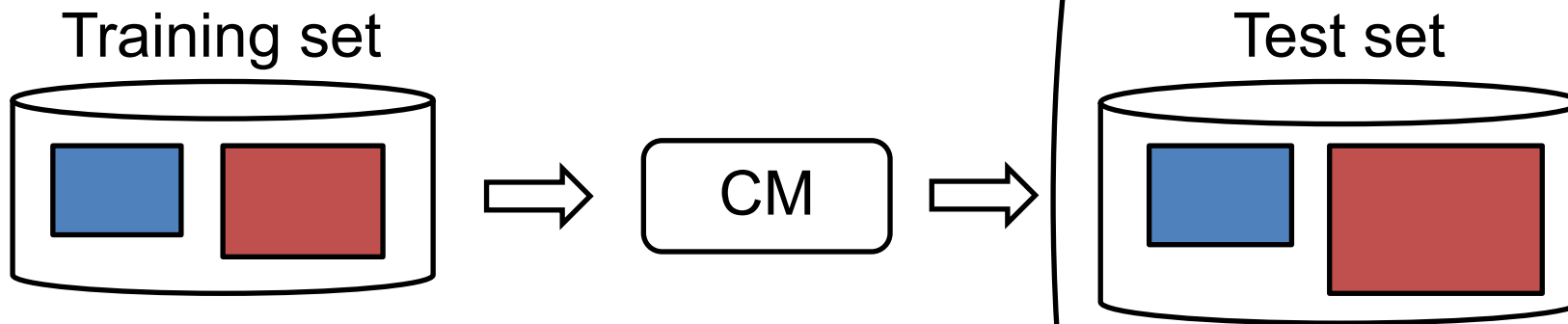
# Introduction

## □ CM architecture in our studies (Wang 2022, Tak 2022)

- **Front end:** self-supervised learning (SSL) model
  - wav2vec 2.0 **XLSR-53** (Conneau 2021)
- **Back end:** global average pooling + 4-layer neural network



# Generalization is desired



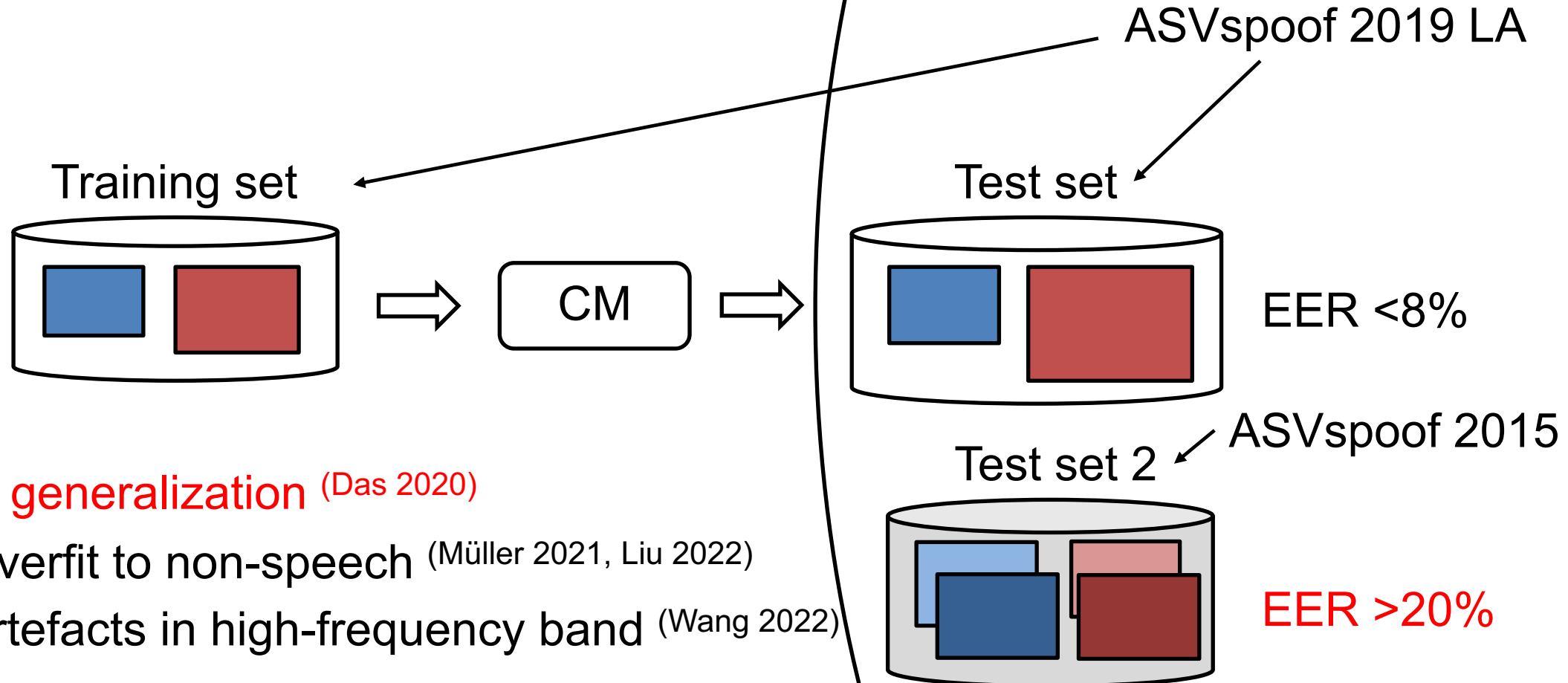
*Space of all possible bona fide and spoofed data*

En, Fr, Ch, Jp, ...

Wav, mp3, m4a ...

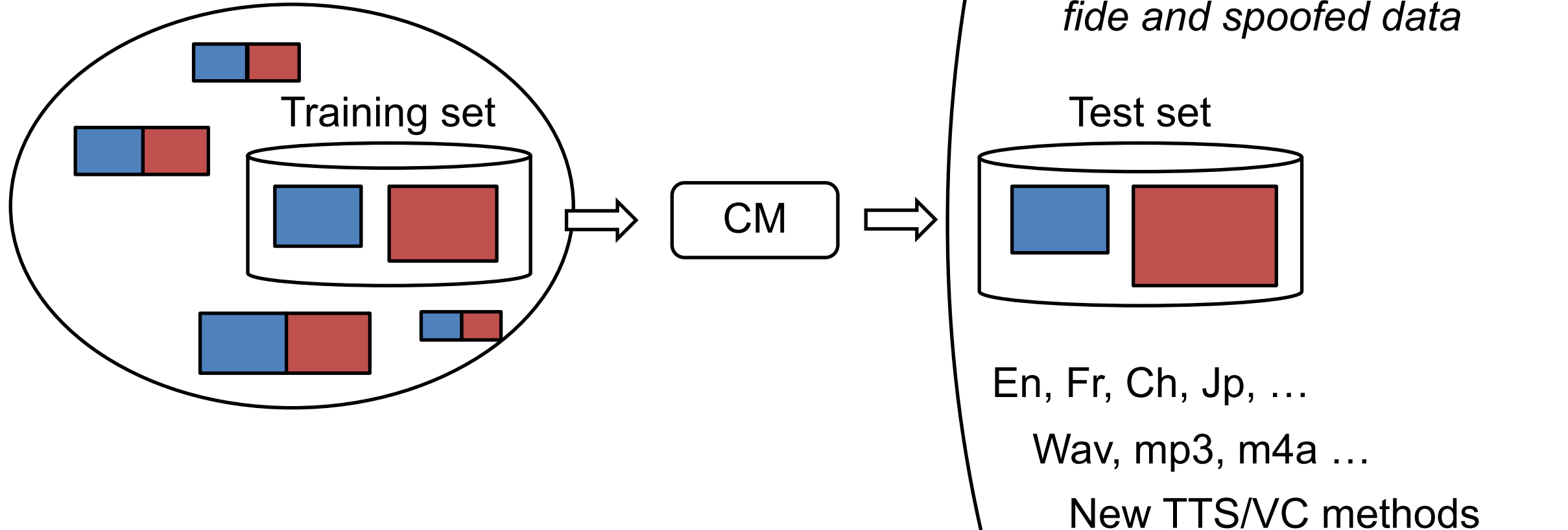
New TTS/VC methods

# Generalization is challenging

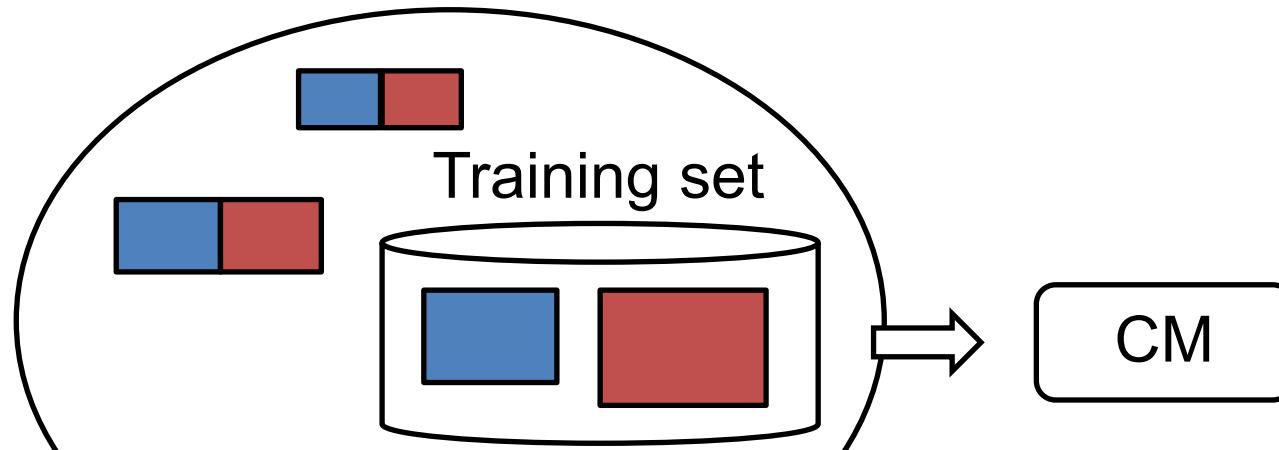


- **Poor generalization** (Das 2020)
  - Overfit to non-speech (Müller 2021, Liu 2022)
  - Artefacts in high-frequency band (Wang 2022)
  - ...

# Generalization may need more data



# Generalization may need more data



ASVspoof 2019: A large-scale public database of synthesized, converted and replayed speech

Xin Wang<sup>a,\*</sup>, Junichi Yamagishi<sup>a,b,\*\*</sup>, Massimiliano Todisco<sup>c,\*\*</sup>, Héctor Delgado<sup>c,\*\*</sup>, Andreas Nautsch<sup>c,\*\*</sup>, Nicholas Evans<sup>c,\*\*</sup>, Md Sahidullah<sup>d,\*\*</sup>, Ville Vestman<sup>c,\*\*</sup>, Tomi Kinnunen<sup>c,\*\*</sup>, Kong Aik Lee<sup>e,\*\*</sup>, Lauri Juvela<sup>g</sup>, Paavo Alku<sup>g</sup>, Yu-Huai Peng<sup>h</sup>, Hsin-Te Hwang<sup>h</sup>, Yu Tsao<sup>h</sup>, Hsin-Min Wang<sup>h</sup>, Sébastien Le Maguer<sup>i</sup>, Markus Becker<sup>i</sup>, Ferguson Henderson<sup>j</sup>, Rob Clark<sup>j</sup>, Yu Zhang<sup>j</sup>, Quan Wang<sup>j</sup>, Ye Jia<sup>k</sup>, Kai Onuma<sup>k</sup>, Koji Mushika<sup>k</sup>, Takashi Kaneda<sup>k</sup>, Yuan Jiang<sup>l</sup>, Li-Juan Liu<sup>l</sup>, Yi-Chiao Wu<sup>m</sup>, Wen-Chin Huang<sup>m</sup>, Tomoki Toda<sup>n</sup>, Kou Tanaka<sup>n</sup>, Hirokazu Kameoka<sup>n</sup>, Ingmar Steiner<sup>o</sup>, Driss Matrouf<sup>p</sup>, Jean-François Bonastre<sup>p</sup>, Avashna Govender<sup>q</sup>, Srikanth Ronanki<sup>q</sup>, Jing-Xuan Zhang<sup>r</sup>, Zhen-Hua Ling<sup>r</sup>

<sup>a</sup>National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo, Japan

<sup>b</sup>Centre for Speech Technology Research, University of Edinburgh, UK

<sup>c</sup>EURECOM, Campus SophiaTech, 450 Route des Chappes, 06410 Biot, France

<sup>d</sup>Université de Lorraine, CNRS, Inria, LORIA, F-54000, Nancy, France

<sup>e</sup>University of Eastern Finland, Joensuu campus, Länssikatu 15, FI-80110 Joensuu, Finland

<sup>f</sup>NEC Corp., 7-1, Shiba 5-chome Minato-ku, Tokyo 108-8001, Japan

<sup>g</sup>Aalto University, Rakentajanaukio 2 C, 00076 Aalto, Finland

<sup>h</sup>Academia Sinica, 128, Sec. 2, Academia Road, Nankang, Taipei, Taiwan

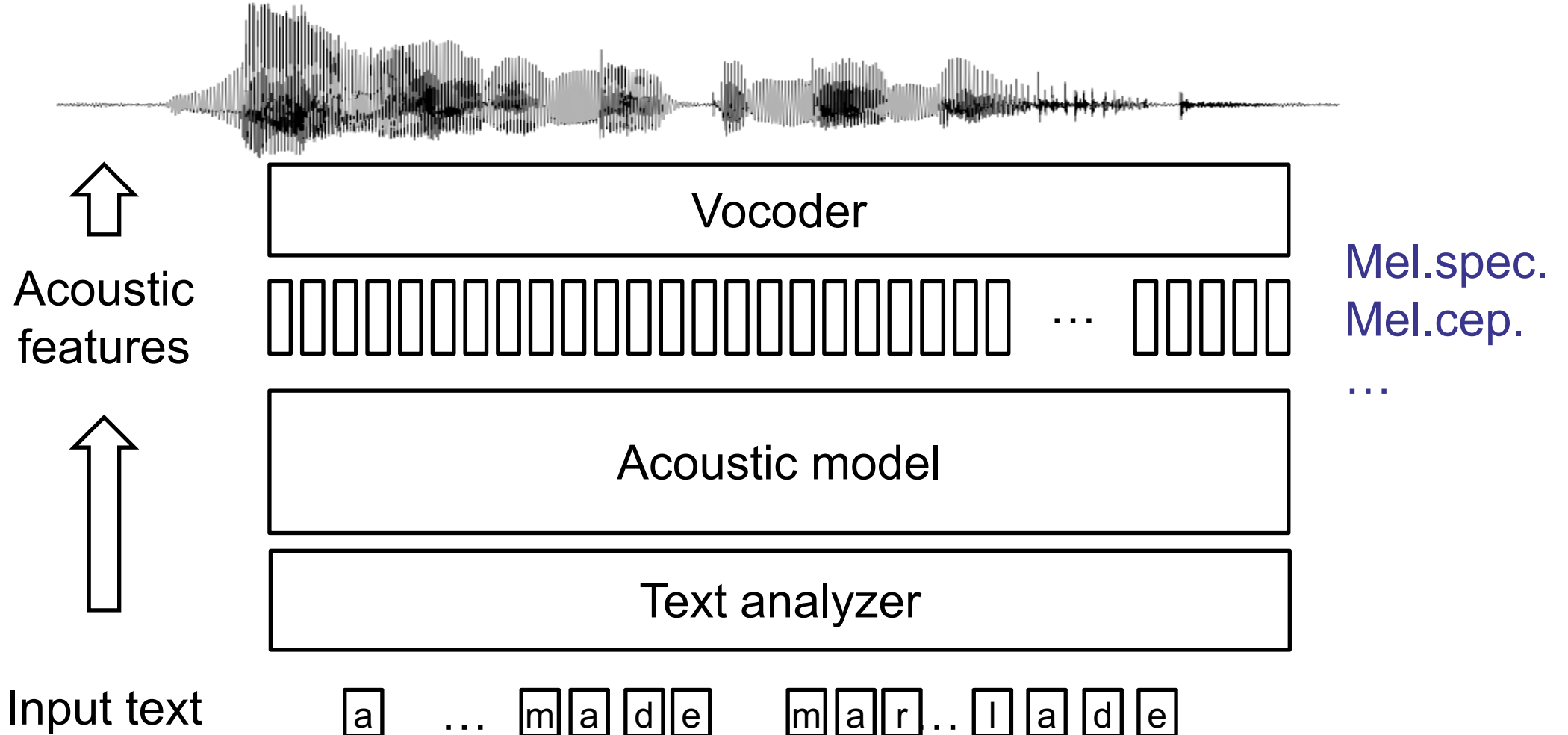
*Easier way to create useful spoofed training data?*

*~6 months of work*

- However, building diverse TTS and VC systems is **time consuming**

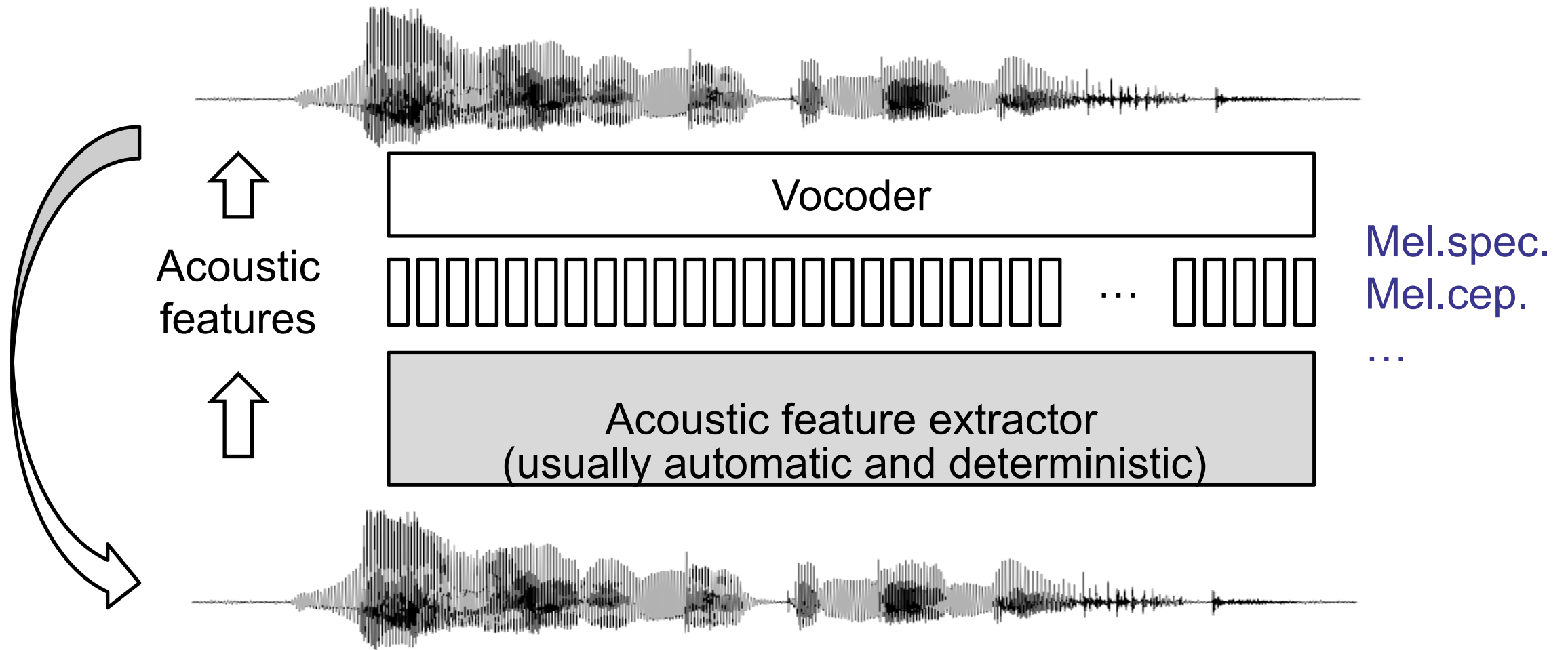


# Idea: creating spoofed training data by vocoding



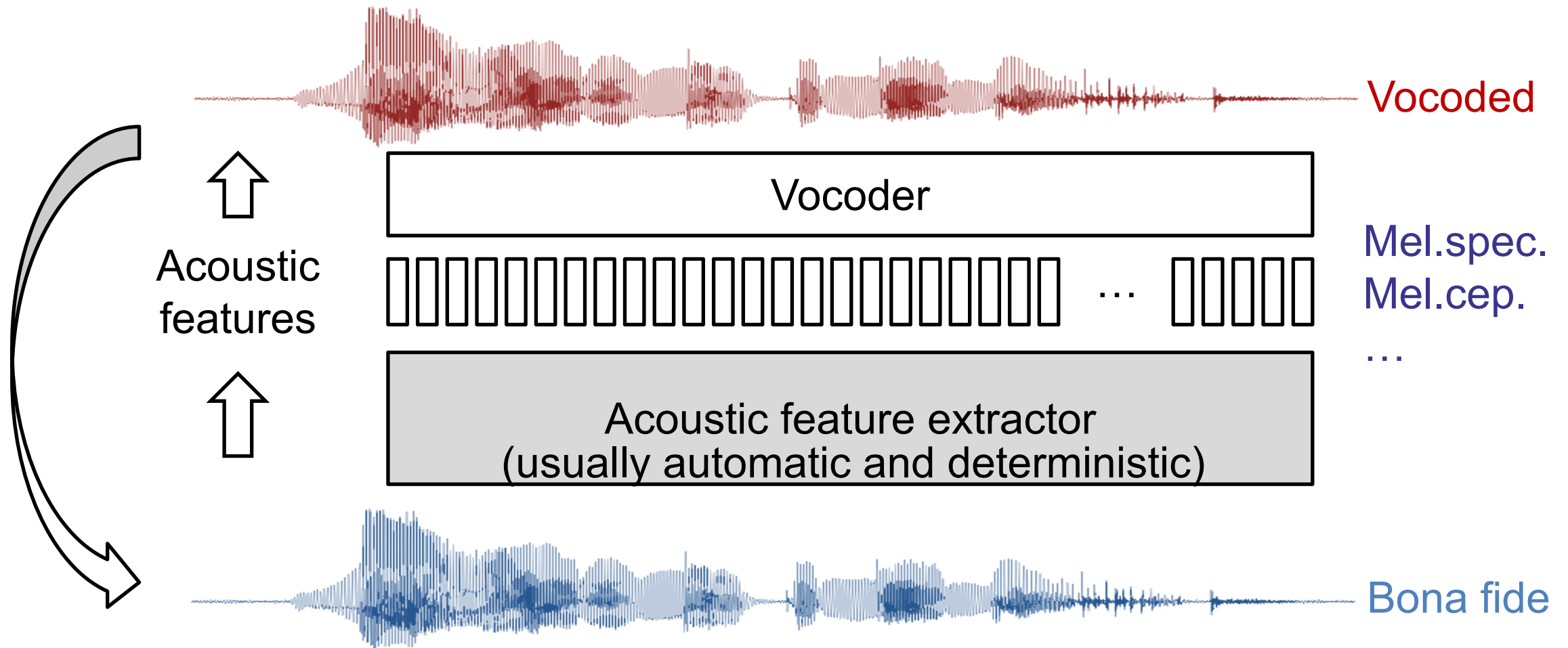
Example of full-fledged TTS

# Idea: creating spoofed training data by vocoding



Vocoding is TTS using a perfect acoustic model

# Idea: creating spoofed training data by vocoding



Vocoding is TTS using a perfect acoustic model

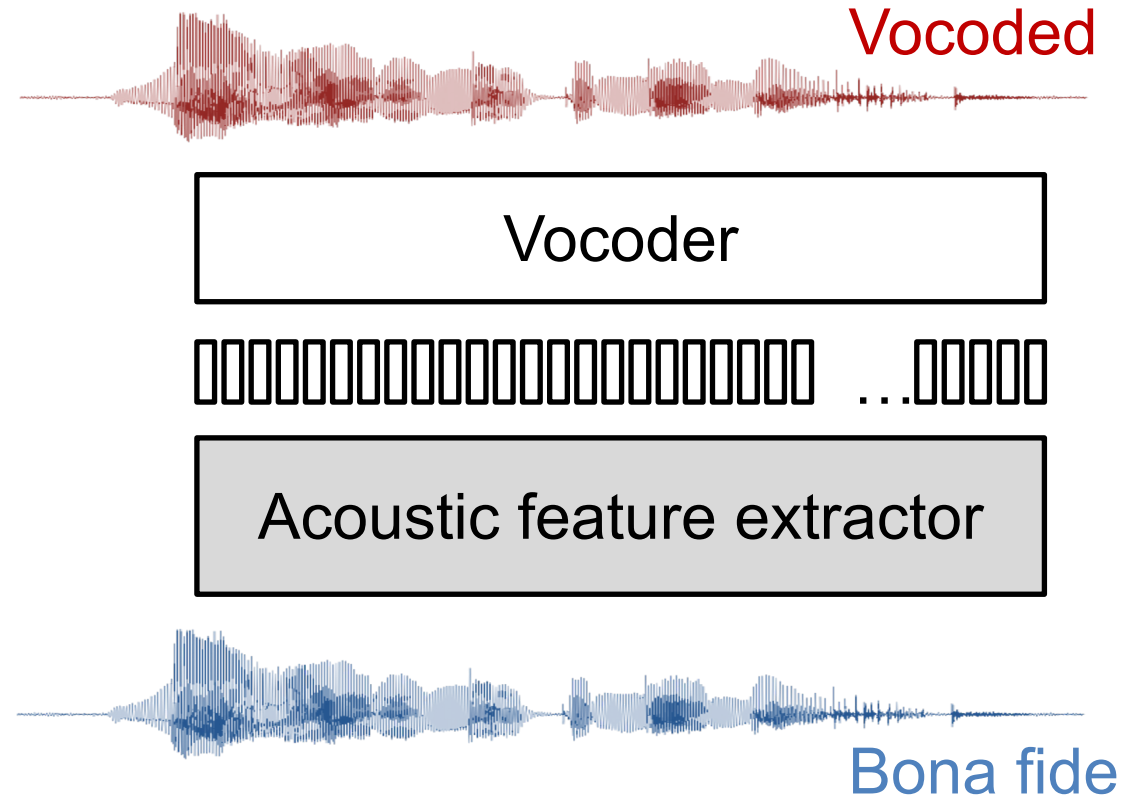
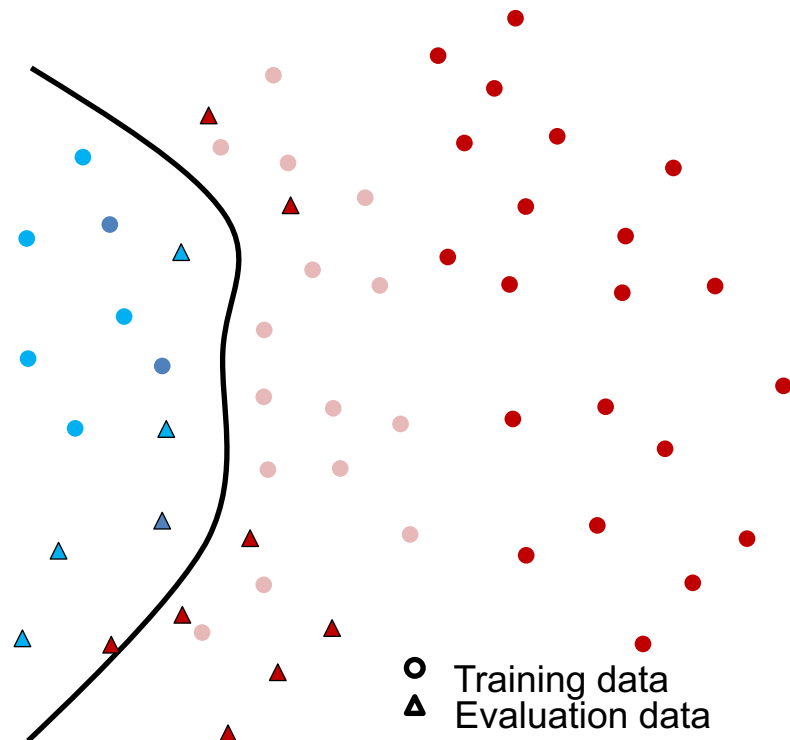
# Idea: creating spoofed training data by vocoding

## □ Assumption (ideally)

Bona fide

Vocoded

TTS/VC

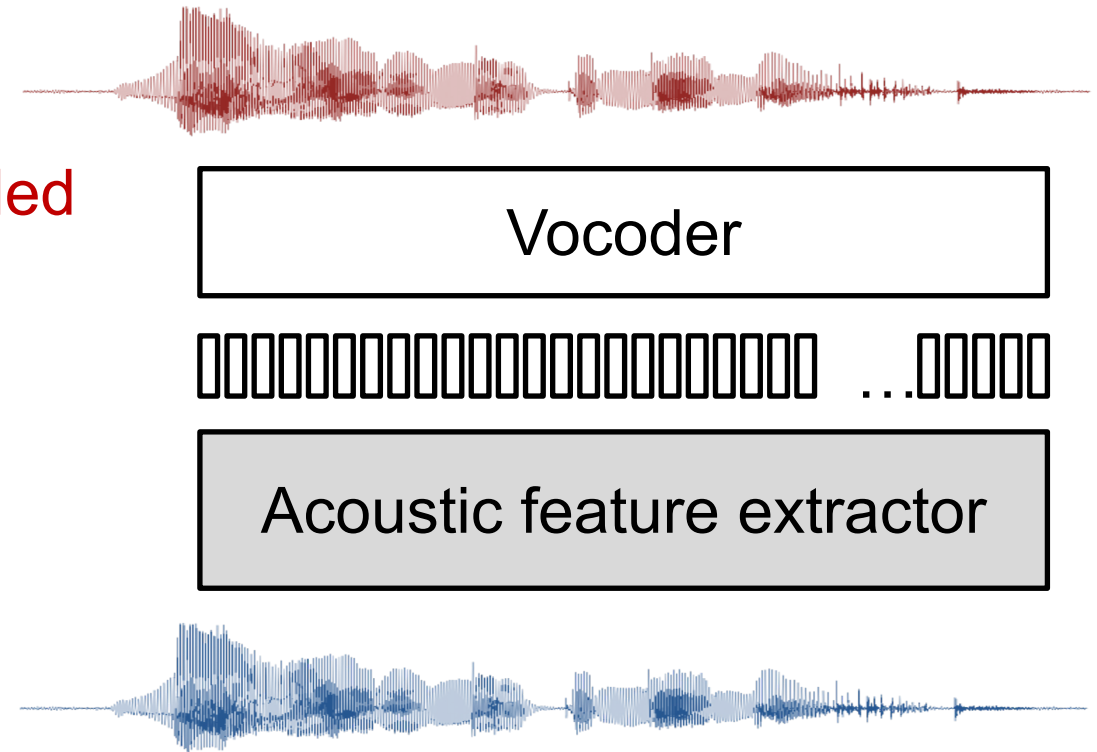
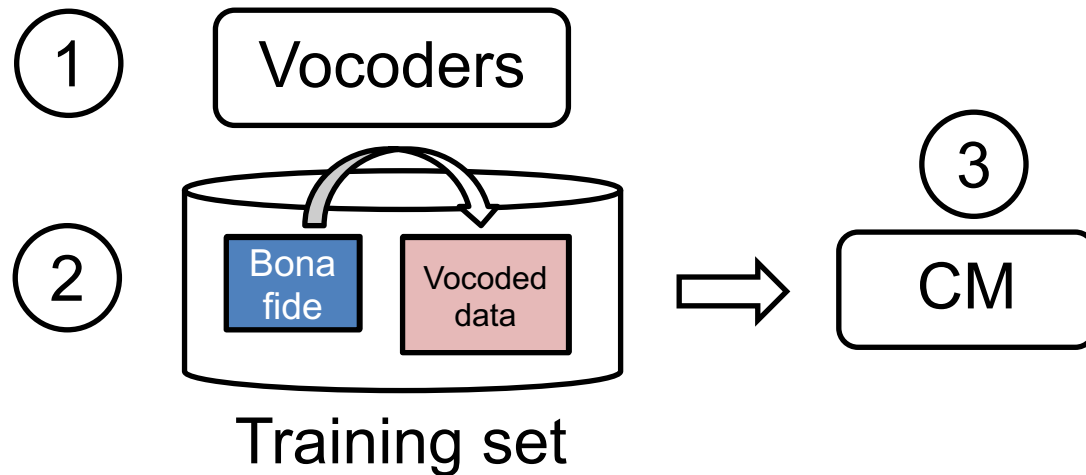


Vocoding is TTS using a perfect acoustic model

# Idea: creating spoofed training data by vocoding

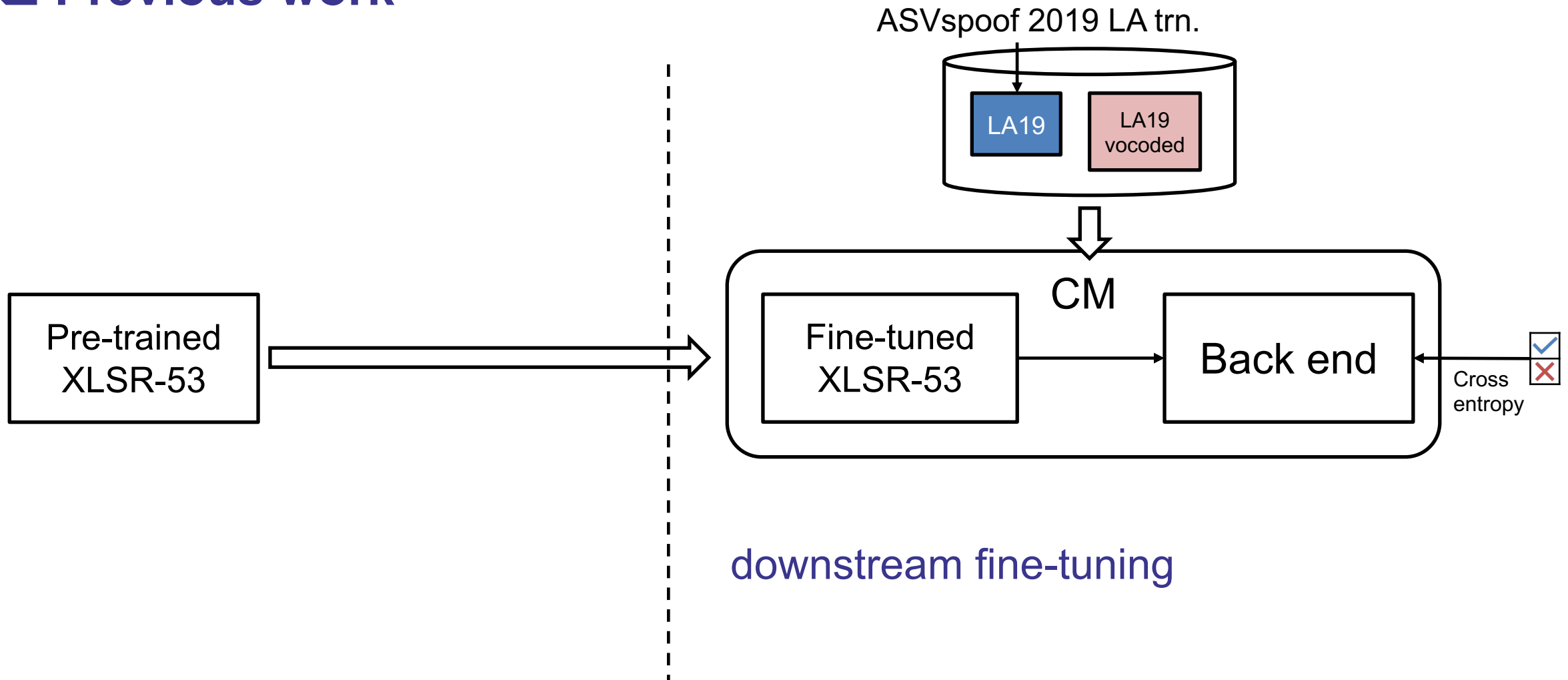
## □ Procedure

1. Prepare (or download) vocoders
2. Vocode the bona fide data
3. Train CM using **bona fide** and **vocoded** data



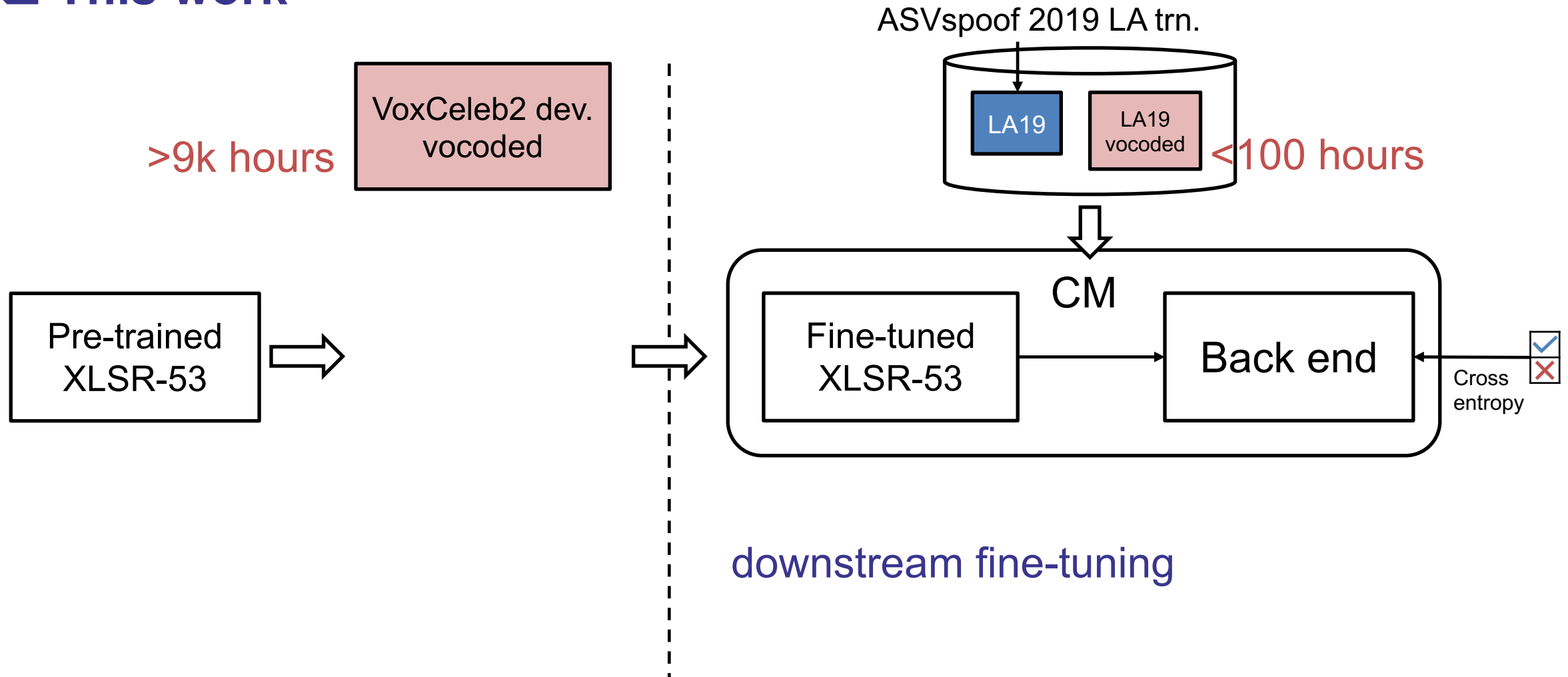
# Method: CM training using bona fide & vocoded data

## □ Previous work (Wang 2023)



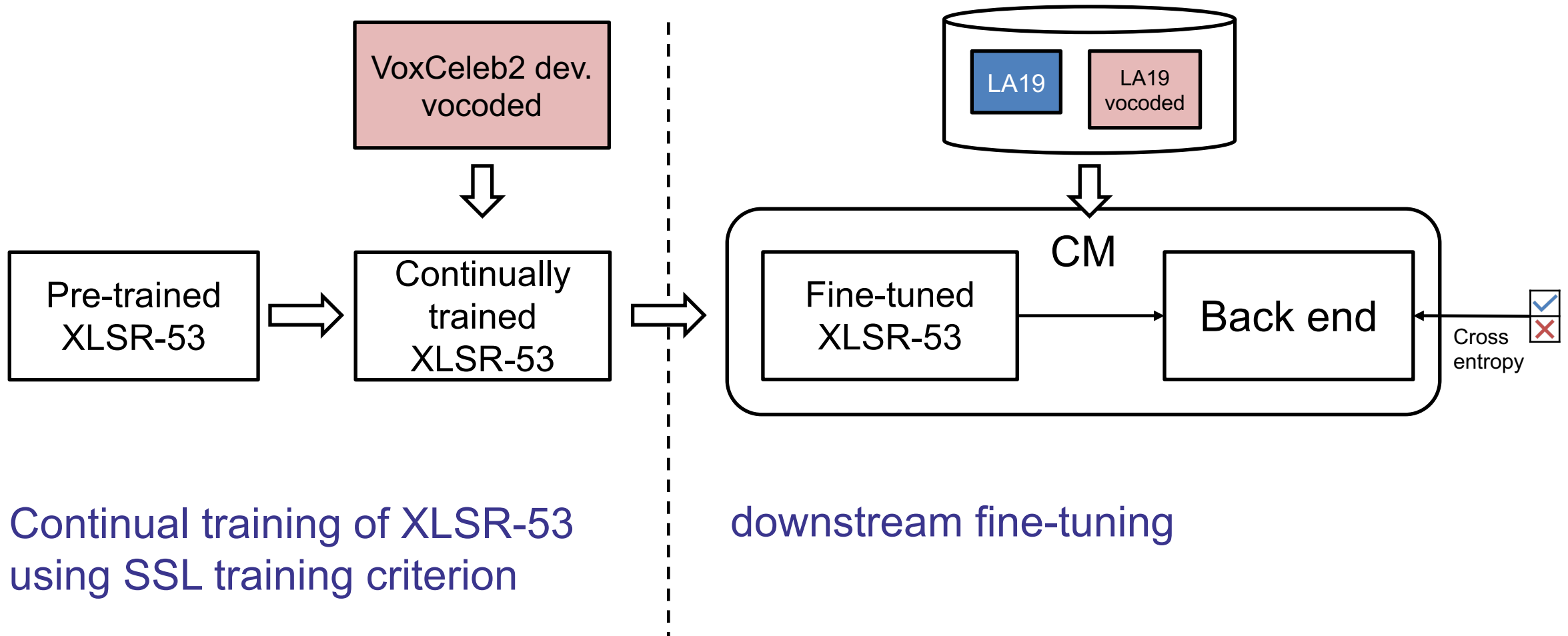
# Method: CM training using bona fide & vocoded data

## □ This work



# Method: CM training using bona fide & vocoded data

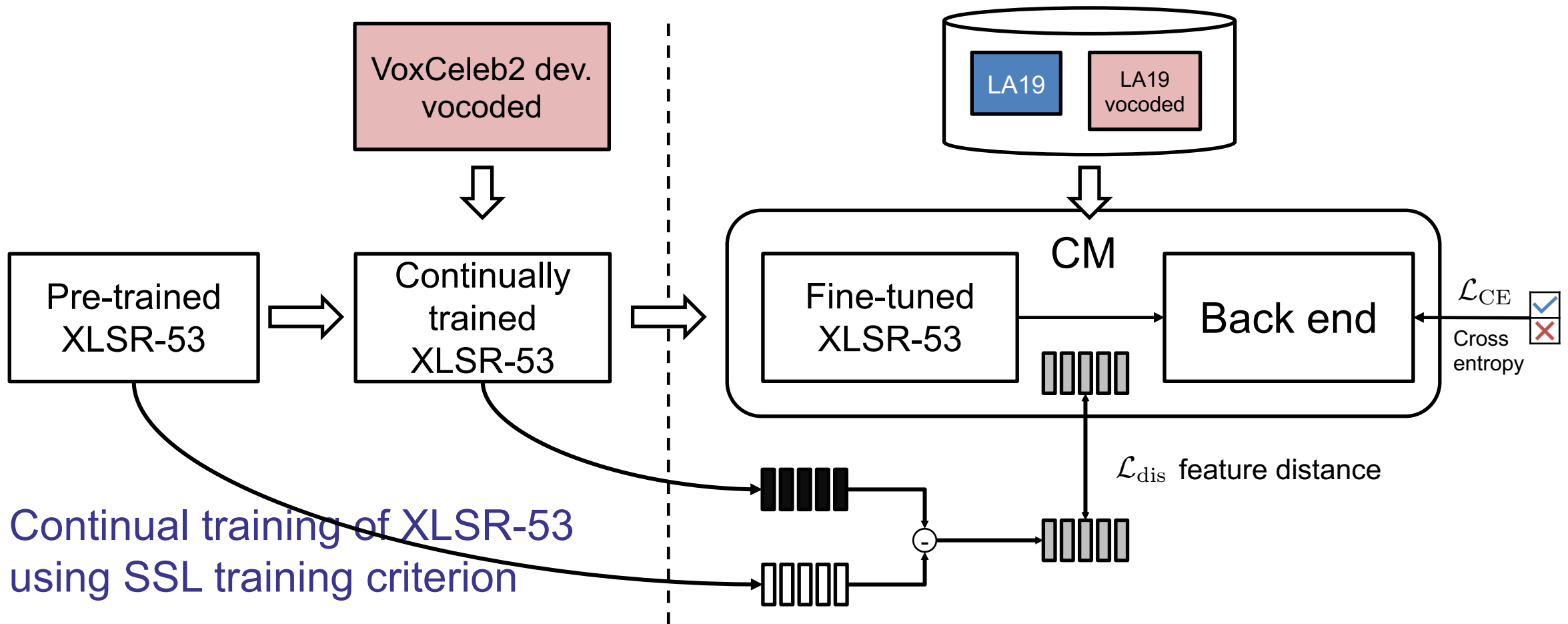
## □ This work – method 1





# Method: CM training using bona fide & vocoded data

## □ This work – method 2



# Method: CM training using bona fide & vocoded data

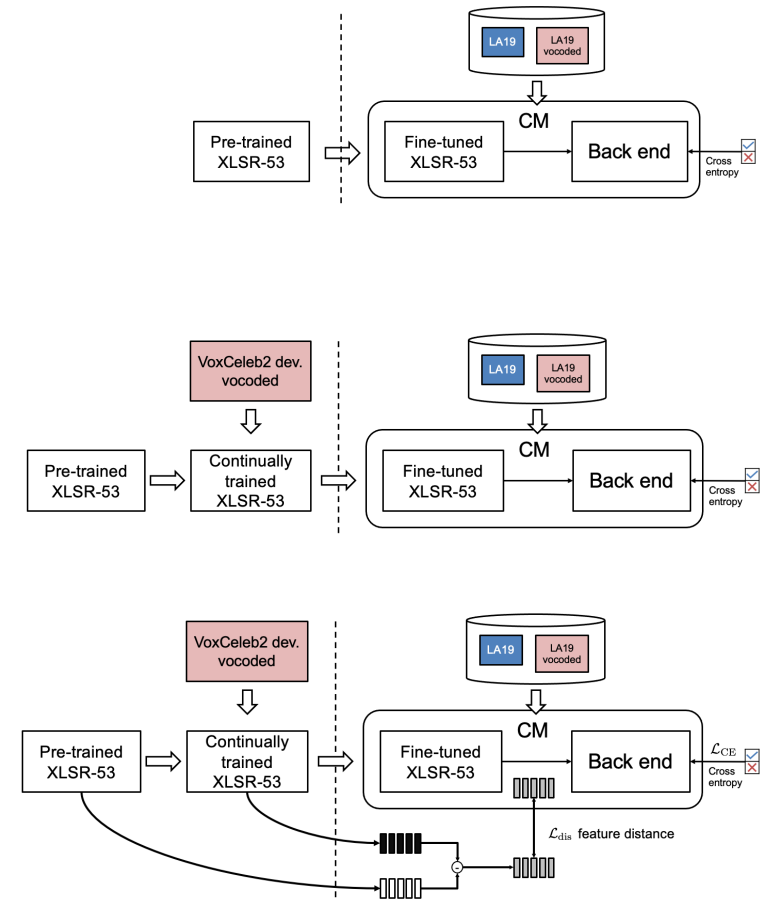
## □ Previous work (Wang 2023)

- Vocoded ASVspoof 2019 LA trn.
- downstream fine-tuning of SSL model

## □ This work

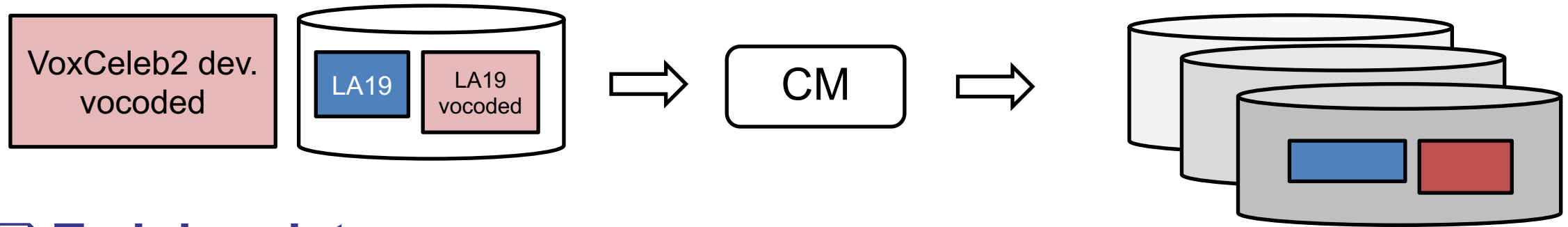
- Vocoded VoxCeleb2 dev.
- Upstream training of SSL
- Downstream training + dis\*\*\*

## □ Related studies using DSP-based vocoders (Wu 2013, Khoury 2014, Sizov 2015, Saratxaga 2016, Pal 2018)



# Experiment

*Training data*

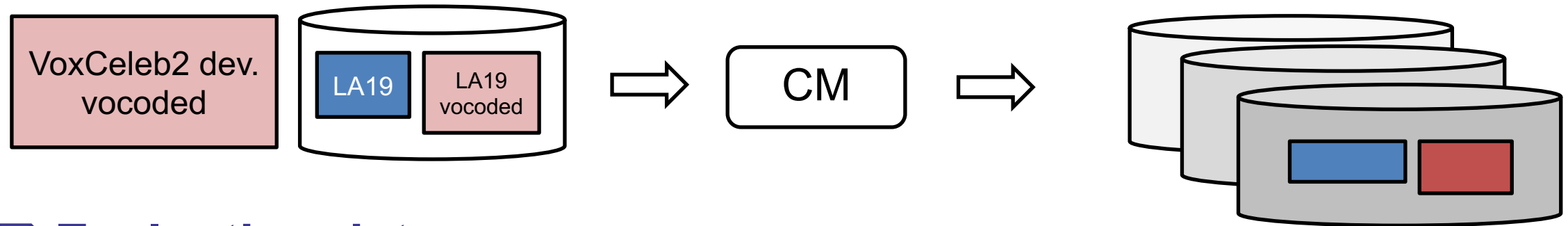


## □ Training data

- SSL upstream training: vocoded VoxCeleb2 dev.
- Downstream fine-tuning: bonafide + vocoded ASVspoof 2019 LA trn.
- Vcoders: HiFi-GAN (Kong 2020), NSF (Wang 2019), NSF-GAN, WaveGlow (Prenger 2019)

# Experiment

*Training data*



## □ Evaluation data

- ASVspoof 2019 LA test set, 2021 LA & DF eval sets
- ASVspoof 2019 LA test set w/o non-speech, 2021 LA & DF hidden track
- WaveFake (Frank 2021) , In-the-Wild (Müller 2022)
- three independent training-evaluation rounds
- averaged EERs

More challenging due to domain mismatch

# Experiment results

😊 Low EER

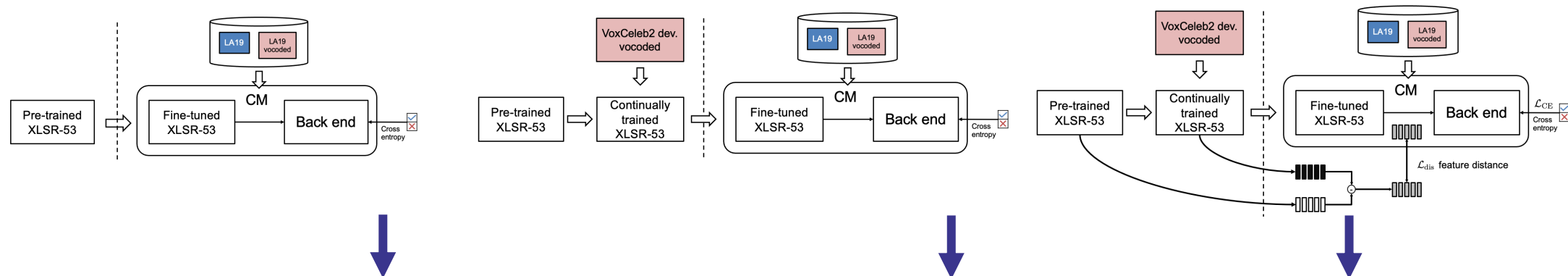


☹️ High EER

## Systems using different training configurations

CM	ID	B1			B2			B3			P1		P2		P3	
	Front end SSL(s)	xlsr	xlsr,	w2v	xlsr,	w2v	xlsr,	w2v	v.vox	xlsr,	v.vox	xlsr,	v.vox	xlsr,	v.vox	
	SSL distilling	-		×		✓		✓	-		×		✓		✓	
Data for fine-tune CM		voc.LA						voc.LA								
EER on each test set	LA19eval	3.45	1.97	1.26	2.09	2.01	1.91									
	LA21eval	17.59	13.94	21.09	16.88	14.94	15.92									
	DF21eval	6.53	4.04	14.72	4.34	5.28	5.67									
	LA19etrim	2.69	2.80	3.74	3.33	2.79	3.28									
	LA21hid	13.93	14.05	20.03	16.02	13.95	14.97									
	DF21hid	8.89	9.10	15.27	7.71	8.40	8.84									
	WaveFake	7.33	1.48	5.88	1.94	0.89	1.30									
	InWild	6.78	4.25	13.20	5.84	4.07	6.10									
	Pooled EER single threshold	Pooled	11.13	12.95	14.06	10.54	9.07	9.98								

# Experiment results

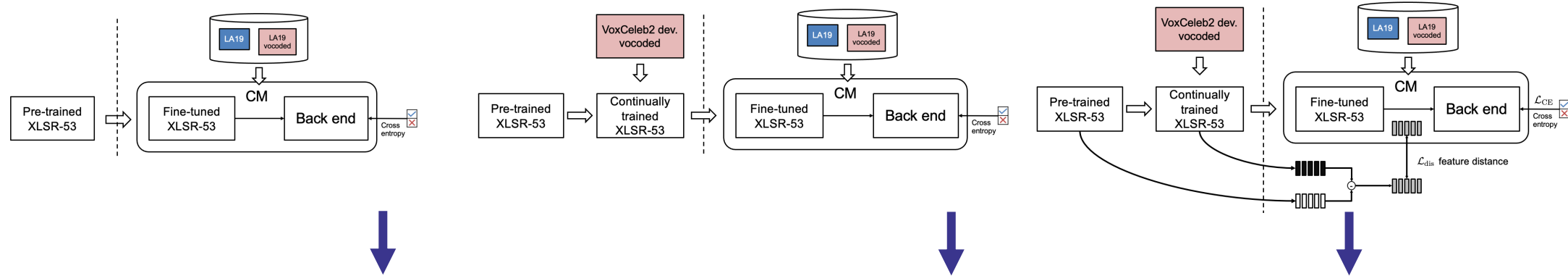


I	ID	B1	B2	B3	P1	P2	P3
	Test sets	LA19eval	3.45	1.97	1.26	2.09	2.01
LA21eval		17.59	13.94	21.09	16.88	14.94	15.92
DF21eval		6.53	4.04	14.72	4.34	5.28	5.67
LA19etrim		2.69	2.50	3.74	3.33	2.79	3.28
LA21hid		13.93	10.16	20.03	16.02	13.95	14.97
DF21hid		8.89	8.16	11.97	7.71	8.40	8.84
WaveFake		7.33	1.48	5.88	1.94	0.89	1.30
InWild		6.78	4.25	13.20	5.84	4.07	6.10
Pooled		11.13	12.55	11.06	10.54	9.07	9.98

**B1 vs P1:  
continually  
trained SSL is  
not useless**



# Experiment results



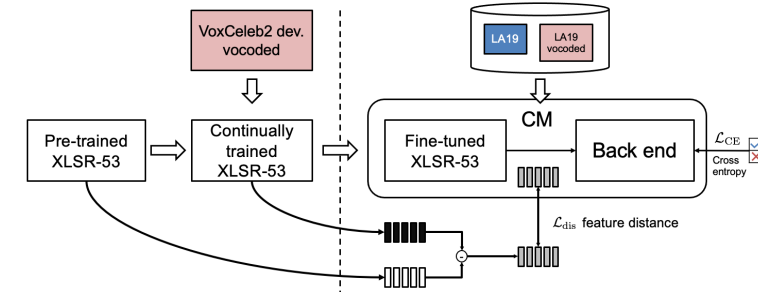
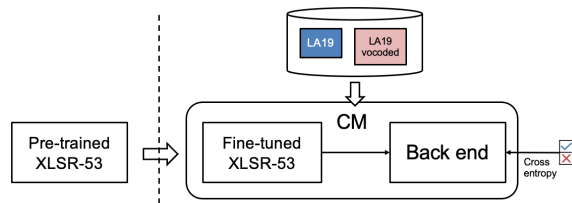
ID	B1	B2	B3	P1	P2	P3
	LA19eval	3.45	1.97	1.26	2.09	2.01
LA21eval	17.59	13.94	21.09	16.88	14.94	15.92
DF21eval	6.53	4.04	14.72	4.34	5.28	5.67
LA19etrim	2.69	2.50	3.74	3.33	3.74	3.28
LA21hid	13.93	14.16	20.03	16.02	13.95	14.97
DF21hid	8.89	9.16	11.97	7.71	8.19	8.84
WaveFake	7.33	1.48	5.88	1.94	0.89	1.30
InWild	6.78	4.25	13.20	5.84	4.07	6.10
Pooled	11.13	11.06	12.06	10.54	10.54	9.98

**B1 vs P1:  
continually  
trained SSL is  
not useless**

**Merging  
two  
SSLs is  
helpful**



# Experiment results



ID		B1	B2	B3	P1	P2	P3
	LA19eval	3.45	1.07	1.26	2.09	2.01	1.91
	LA21eval	17.59	13.94	21.09	16.88	14.94	15.92
	DF21eval	6.53	1.01	11.72	4.34	5.28	5.67
Test sets	LA19etrim	2.69	2.80	3.74	3.33	2.79	3.28
	LA21hid	13.93	11.93	17.03	6.03	13.95	14.97
	DF21hid	8.89	1.01	5.27	7.71	8.40	8.84
	WaveFake	7.33	1.48	5.88	1.94	0.89	1.30
	InWild	6.78	4.25	13.20	5.84	4.07	6.10
Pooled		11.13	11.83	11.83	10.94	9.07	9.98

B1 is the best model in our previous work (Wang 2023)

It is better than using TTS/VC spoofed data

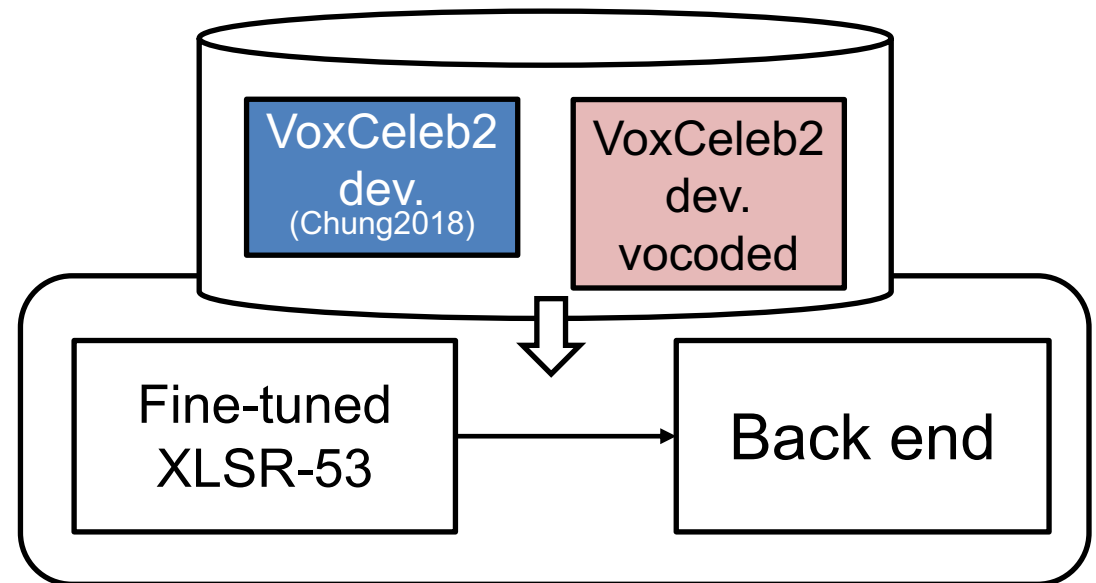
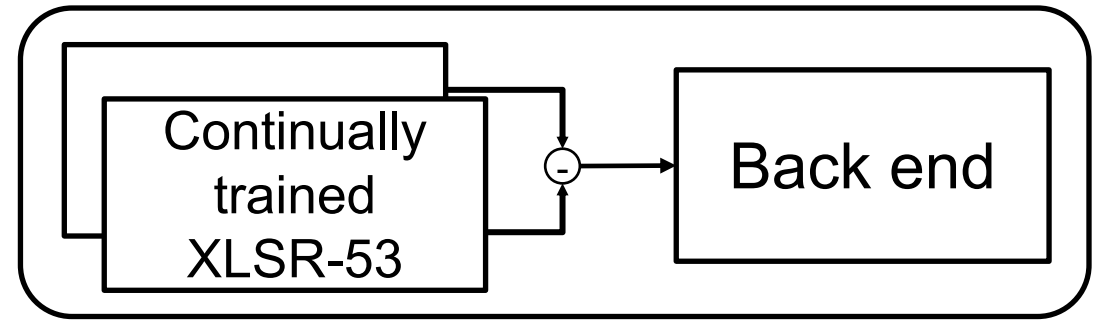




# Experiment results

## ❑ Other results in the paper

- Using two SSLs without distillization?
- Downstream fine-tuning using vocoded voxceleb2?
- ...



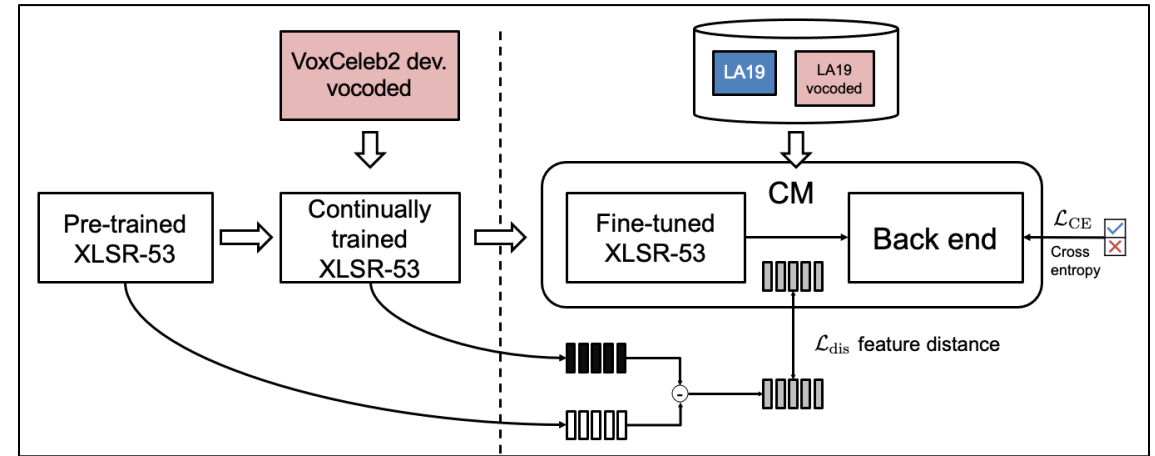
# Summary

## □ Method

- Large scale vocoded VoxCeleb2
- Upstream SSL training
- Downstream fine-tuning + distilling

## □ Results

- Slightly outperformed previous work (pooled EER)
- Limitation: only 4 types of vocoders



Thank you



project/[10-asvspoof-vocoded-trn-ssl](https://github.com/nii-yamagishilab/project-NN-Pytorch-scripts/tree/master/project/10-asvspoof-vocoded-trn-ssl)

# Appendix

CM	ID	B1	B2	B3	P1	P2	P3	B1-b	P3-b	B1-c	P3-c
	Front end SSL(s)	xlsr	xlsr, w2v	xlsr, w2v	v.vox	xlsr, v.vox	xlsr, v.vox	xlsr	xlsr, v.vox	xlsr	xlsr, v.vox
SSL distilling	-	×	✓	-	×	✓	-	✓	-	✓	
Data for fine-tune CM	voc.LA			voc.LA			LA19trn		voc.VoxCel		
Test sets	LA19eval	3.45	1.97	1.26	2.09	2.01	1.91	0.22	0.13	3.59	3.71
	LA21eval	17.59	13.94	21.09	16.88	14.94	15.92	2.69	3.29	15.22	12.37
	DF21eval	6.53	4.04	14.72	4.34	5.28	5.67	4.27	3.45	5.99	3.31
	LA19etrim	2.69	2.80	3.74	3.33	2.79	3.28	7.37	7.37	2.74	3.63
	LA21hid	13.93	14.05	20.03	16.02	13.95	14.97	15.56	24.23	10.14	9.53
	DF21hid	8.89	9.10	15.27	7.71	8.40	8.84	9.16	13.95	9.03	7.77
	WaveFake	7.33	1.48	5.88	1.94	0.89	1.30	23.75	15.44	13.41	24.17
	InWild	6.78	4.25	13.20	5.84	4.07	6.10	13.52	12.32	6.90	7.00
Pooled	11.13	12.95	14.06	10.54	9.07	9.98	12.76	12.50	10.92	12.26	

# Reference

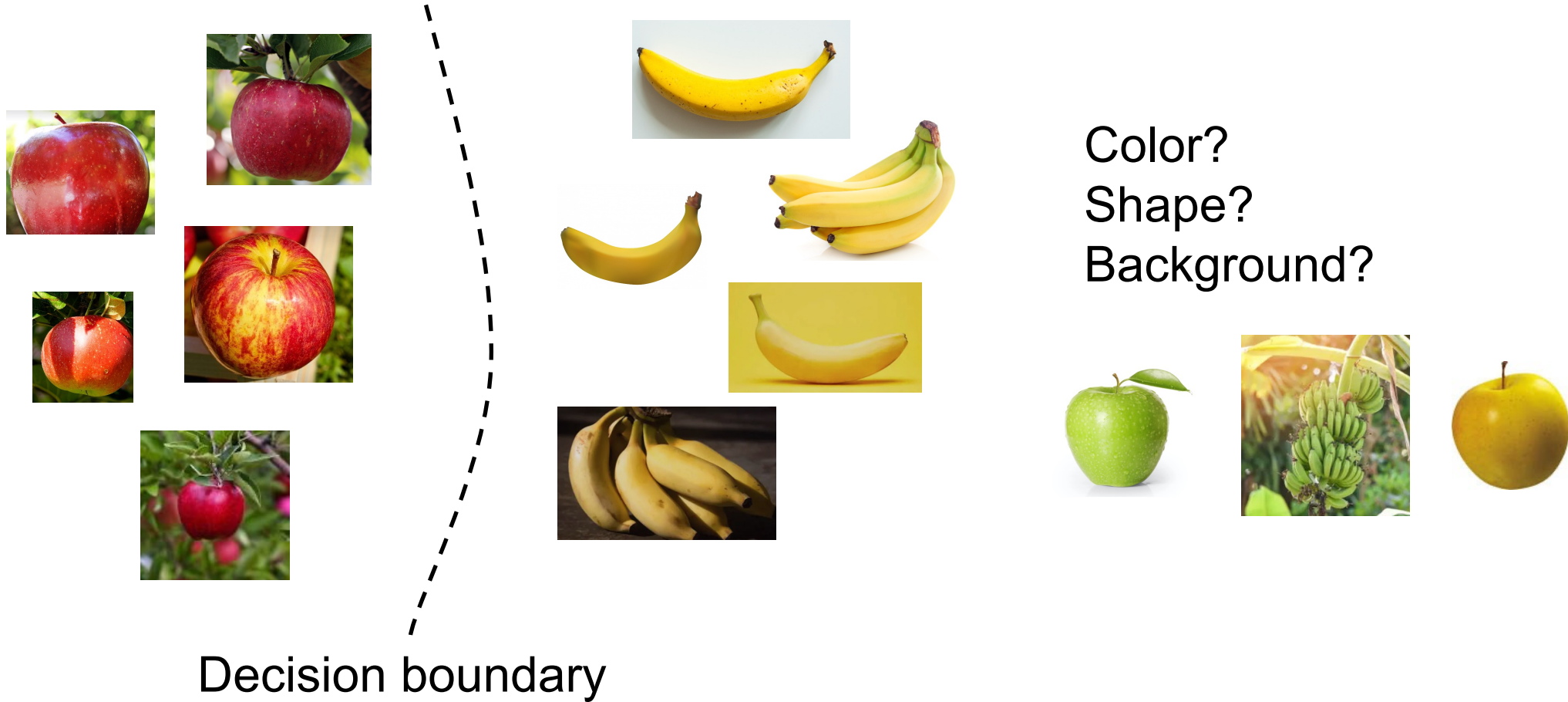
## DSP-based vocoders

- 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.

## neural vocoders

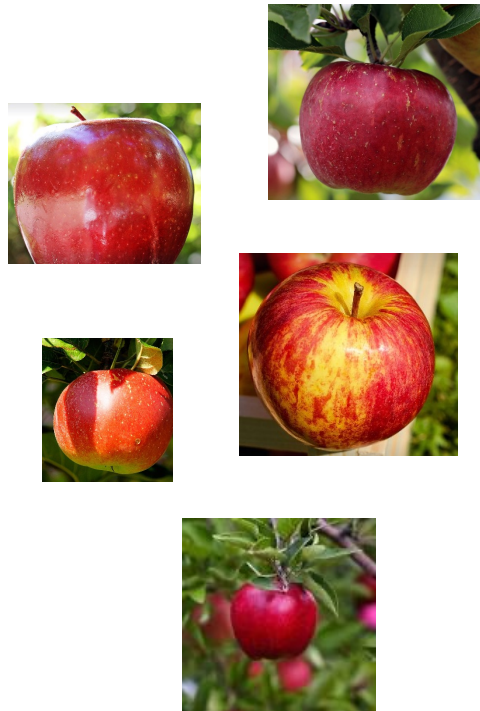
- Joel Frank, and Lea Schönherr. WaveFake: A Data Set to Facilitate Audio DeepFake Detection. In *Proc. NeurIPS Datasets and Benchmarks 2021*. 2021.
- Chengzhe Sun, Shan Jia, Shuwei Hou, Ehab AlBadawy, and Siwei Lyu. Exposing AI-Synthesized Human Voices Using Neural Vocoder Artifacts. ArXiv Preprint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

# Introduction – one challenge in my opinion

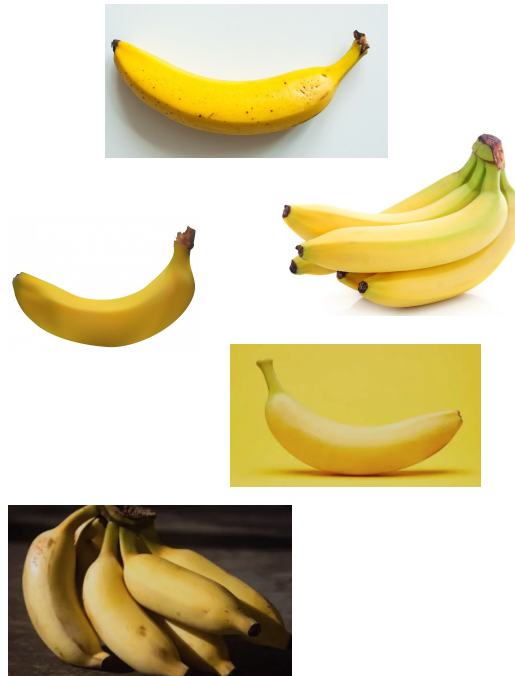


# Introduction – one challenge in my opinion

Spooferd



Bona fide

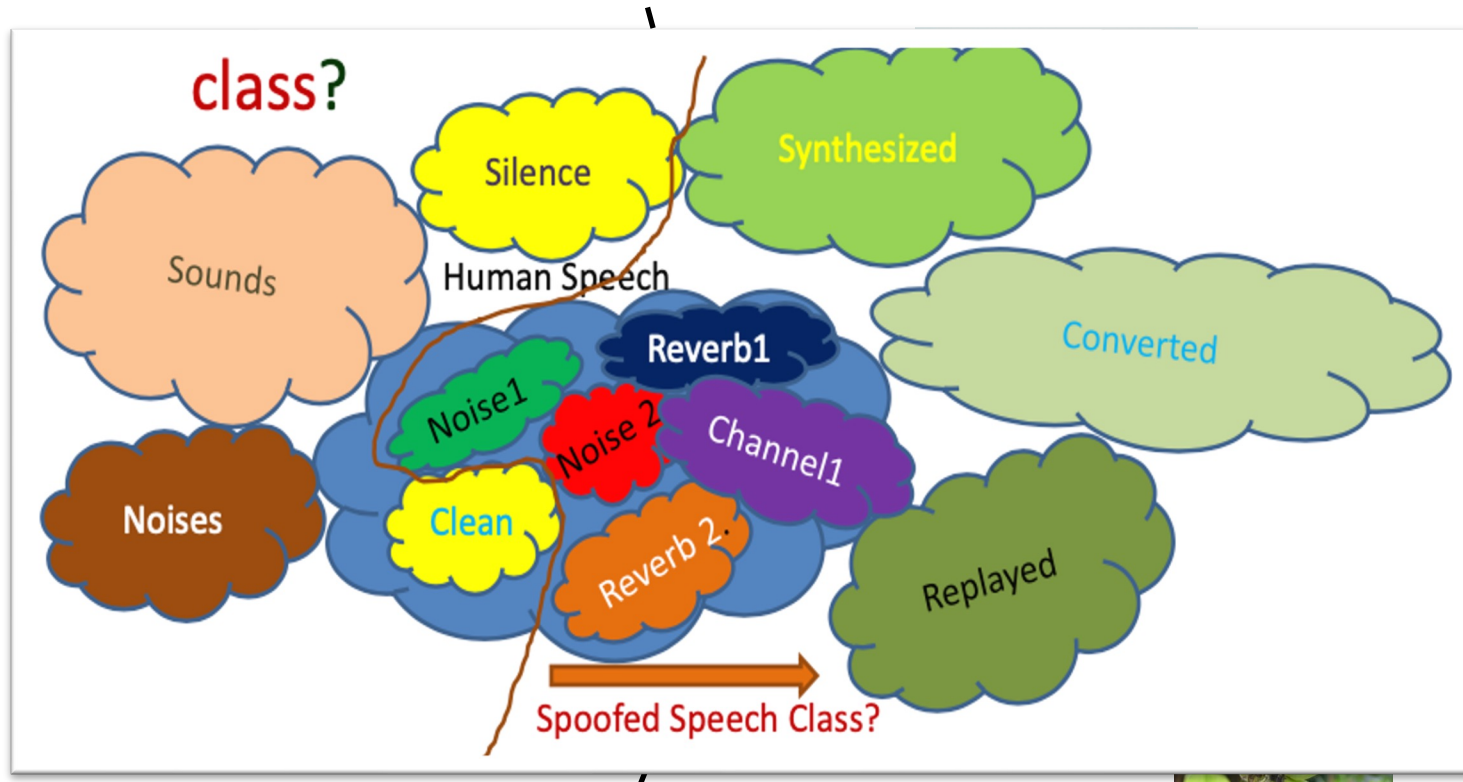


Non-speech (Muller 2021) ?  
Mp3 & additive noise (Shim 2023)?  
Single source data domain?



Decision boundary

# Introduction – one challenge in my opinion



Speech is more complicated

En, Fr, Ch, Jp, ...

LJ-speech, Librispeech ...

Wav, mp3, m4a ...

New speech synthesis methods

*For spoofed trials, how is it possible to “well define” something which is unknown?*

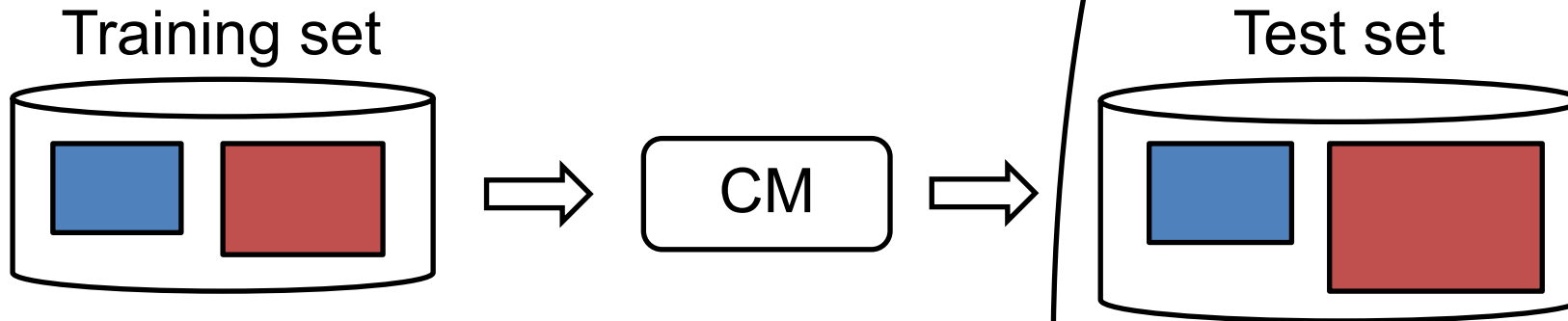
From Jean-Francois Bonastre's talk



# Mini-tutorial on TTS

Building diverse TTS and VC systems is not that easy

# More spoofing data?



***Space of all possible bona fide and spoofed data***

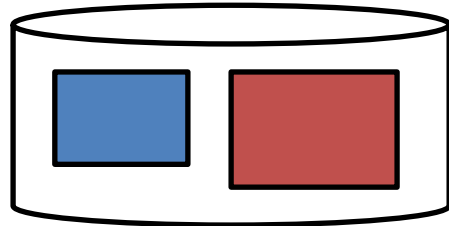
En, Fr, Ch, Jp, ...

wav, mp3, m4a ...

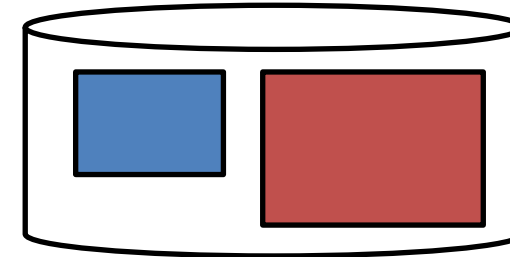
New spoofing algorithms

# More spoofing data?

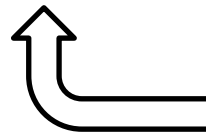
ASVspoof 2019 LA  
Training set



ASVspoof 2019 LA  
Test set



6 TTS/VC



ASVspoof 2019: A large-scale public database of synthesized, converted and replayed speech

Xin Wang<sup>a\*</sup>, Junichi Yamagishi<sup>ab\*\*</sup>, Massimiliano Todisco<sup>c\*\*</sup>, Héctor Delgado<sup>d\*\*</sup>, Andreas Nautsch<sup>e\*\*</sup>, Nicholas Evans<sup>f\*\*</sup>, Md Sahidullah<sup>g\*\*</sup>, Ville Vestman<sup>h\*\*</sup>, Tomi Kinnunen<sup>i\*\*</sup>, Kong Aik Lee<sup>j\*\*</sup>, Lauri Juvela<sup>k</sup>, Paavo Alku<sup>l</sup>, Yu-Huai Peng<sup>h</sup>, Hsin-Te Hwang<sup>h</sup>, Yu Tsao<sup>h</sup>, Hsin-Min Wang<sup>h</sup>, Sébastien Le Maguer<sup>l</sup>, Markus Becker<sup>l</sup>, Fergus Henderson<sup>l</sup>, Rob Clark<sup>l</sup>, Yu Zhang<sup>l</sup>, Quan Wang<sup>l</sup>, Ye Jia<sup>l</sup>, Kai Onuma<sup>k</sup>, Koji Mushioka<sup>k</sup>, Takashi Kaneda<sup>k</sup>, Yuan Jiang<sup>l</sup>, Li-Juan Liu<sup>l</sup>, Yi-Chiao Wu<sup>m</sup>, Wen-Chin Huang<sup>m</sup>, Tomoki Toda<sup>n</sup>, Kou Tanaka<sup>o</sup>, Hirokazu Kameoka<sup>o</sup>, Ingmar Steiner<sup>o</sup>, Driss Matrouf<sup>o</sup>, Jean-François Bonastre<sup>o</sup>, Avashna Govender<sup>o</sup>, Srikanth Ronanki<sup>o</sup>, Jing-Xuan Zhang<sup>o</sup>, Zhen-Hua Ling<sup>o</sup>

<sup>a</sup>National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo, Japan

<sup>b</sup>Centre for Speech Technology Research, University of Edinburgh, UK

<sup>c</sup>EURECOM, Campus SophiaTech, 450 Route des Chappes, 06410 Biot, France

<sup>d</sup>Université de Lorraine, CNRS, Inria, LORIA, F-54000, Nancy, France

<sup>e</sup>University of Eastern Finland, Joensuu campus, Länsskatu 15, FI-80110 Joensuu, Finland

<sup>f</sup>NEC Corp., 7-1, Shiba 5-chome Minato-ku, Tokyo 108-8001, Japan

<sup>g</sup>Aalto University, Raketajamaukio 2 C, 00076 Aalto, Finland

<sup>h</sup>Academia Sinica, 128, Sec. 2, Academia Road, Nankang, Taipei, Taiwan

<sup>i</sup>Sigmedia, ADAPT Centre, School of Engineering, Trinity College Dublin, Ireland

<sup>j</sup>Google Inc., 1600 Amphitheatre Parkway, Mountain View, CA 94043, USA

<sup>k</sup>HOYA, Shinjuku Park Tower 33F, 3-7-1 Nishi-Shinjuku, Shinjuku-ku, Tokyo 163-1035 Japan

<sup>l</sup>Flytek Research, High-tech Development Zone, No. 666 Wangjiang West Road, Hefei, 230088, China

<sup>m</sup>Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8601, Japan

<sup>n</sup>NTT Communication Science Laboratories, 3-1, Morinosato Wakamiya Atsugi-shi, Kanagawa, 243-0198 Japan

<sup>o</sup>audEERING GmbH, Friedrichshafener Str. 1 82205 Gilching, Germany

<sup>p</sup>Avignon University, LIA, 339 Chemin des Meinajariés, 84911 Avignon, France

<sup>q</sup>Centre for Speech Technology Research, University of Edinburgh, UK (Currently with Amazon)

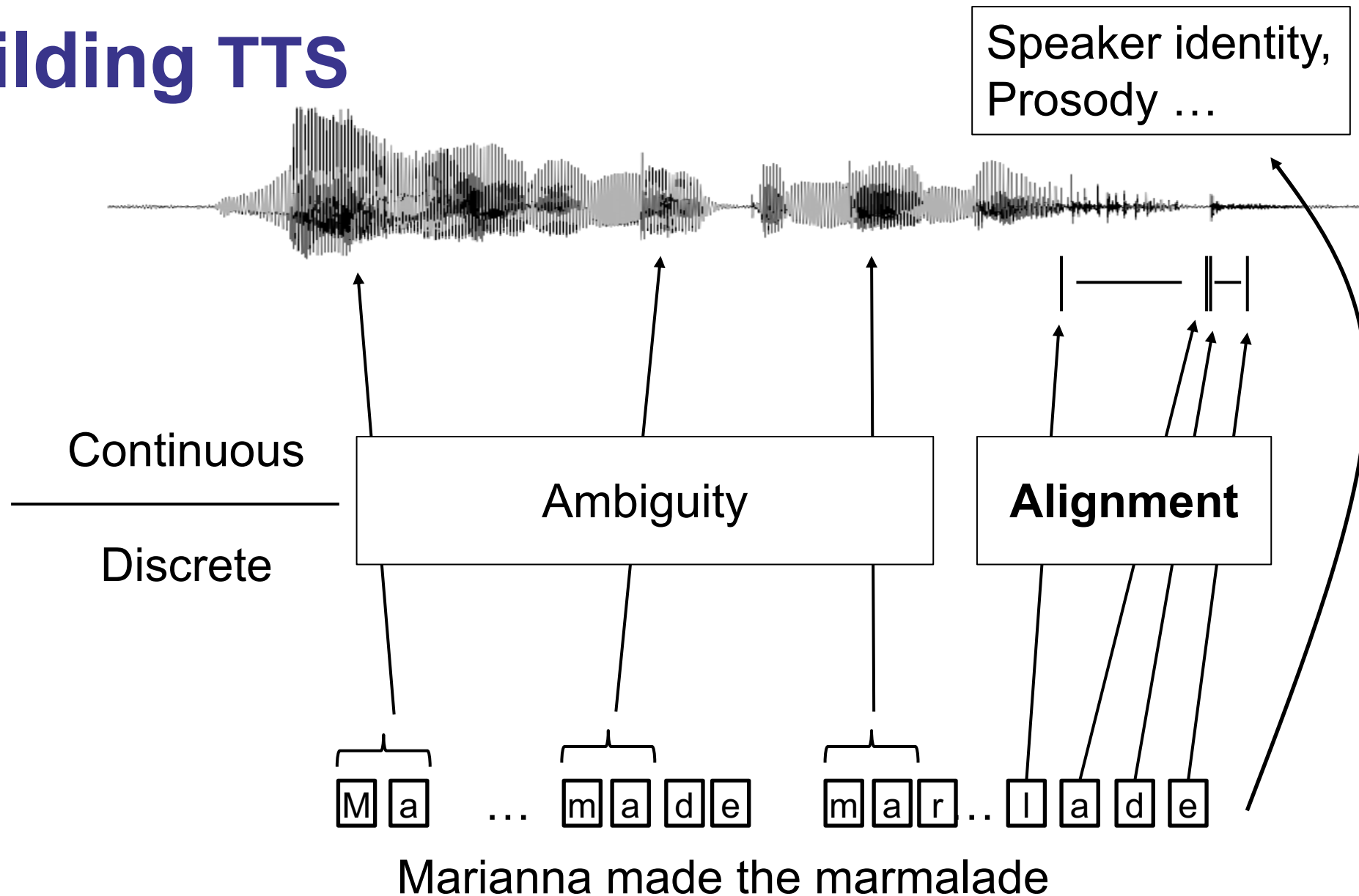
<sup>r</sup>University of Science and Technology of China, No.96, JinZhai Road Baohe District, Hefei, Anhui, 230026, China

11 + 2 TTS/VC

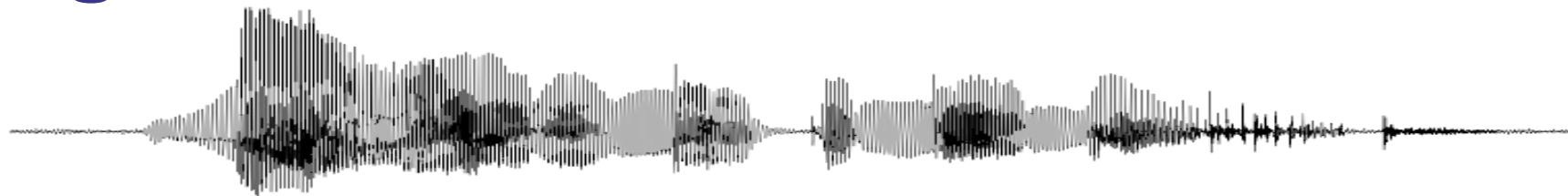


~6 months of work

# Building TTS



# Building TTS



Waveform generation

Acoustic realization

+Prosody tags

To phone

Normalization



H\* H\* L-L%

M A A R I Y A A N A H M E Y D D H A H M A A R M A H L E Y D

M A A R I Y A A N A H M E Y D D H A H M A A R M A H L E Y D

M a ... m a d e m a r . . l a d e

Marianna made the marmalade

# Building TTS

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

DWIGHT BOLINGER

Harvard University




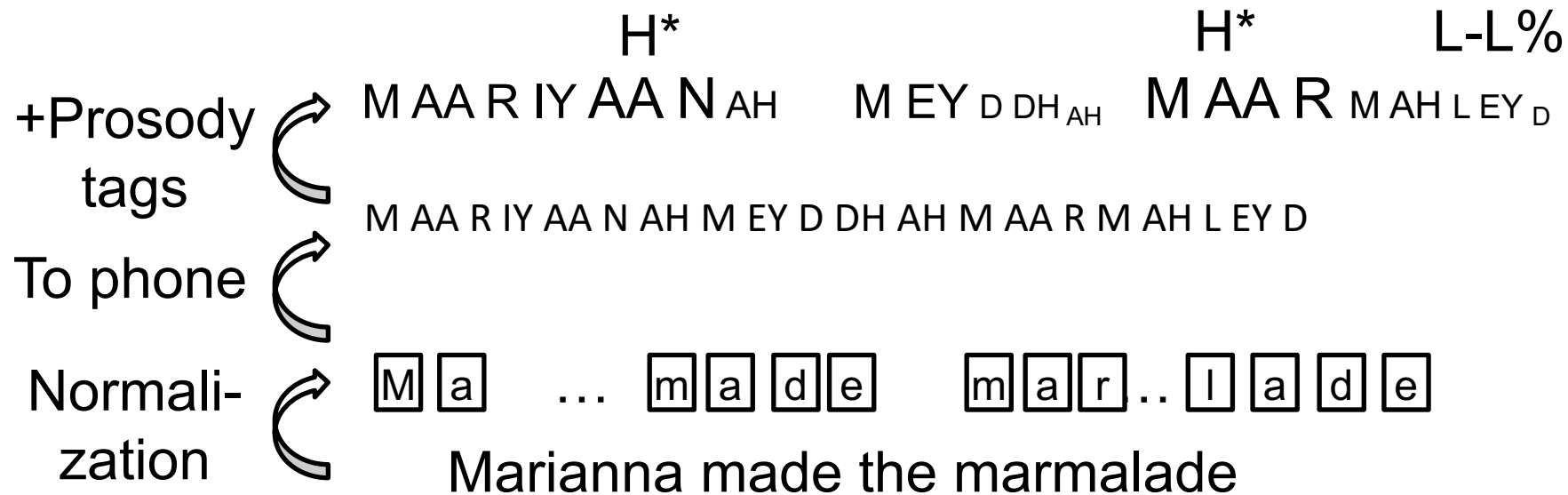
Speaker A: Who made the marmalade.

Speaker A: Bob made the marmalade.

Speaker B: (No,) Mari<sup>an</sup>na made the marmalade. 

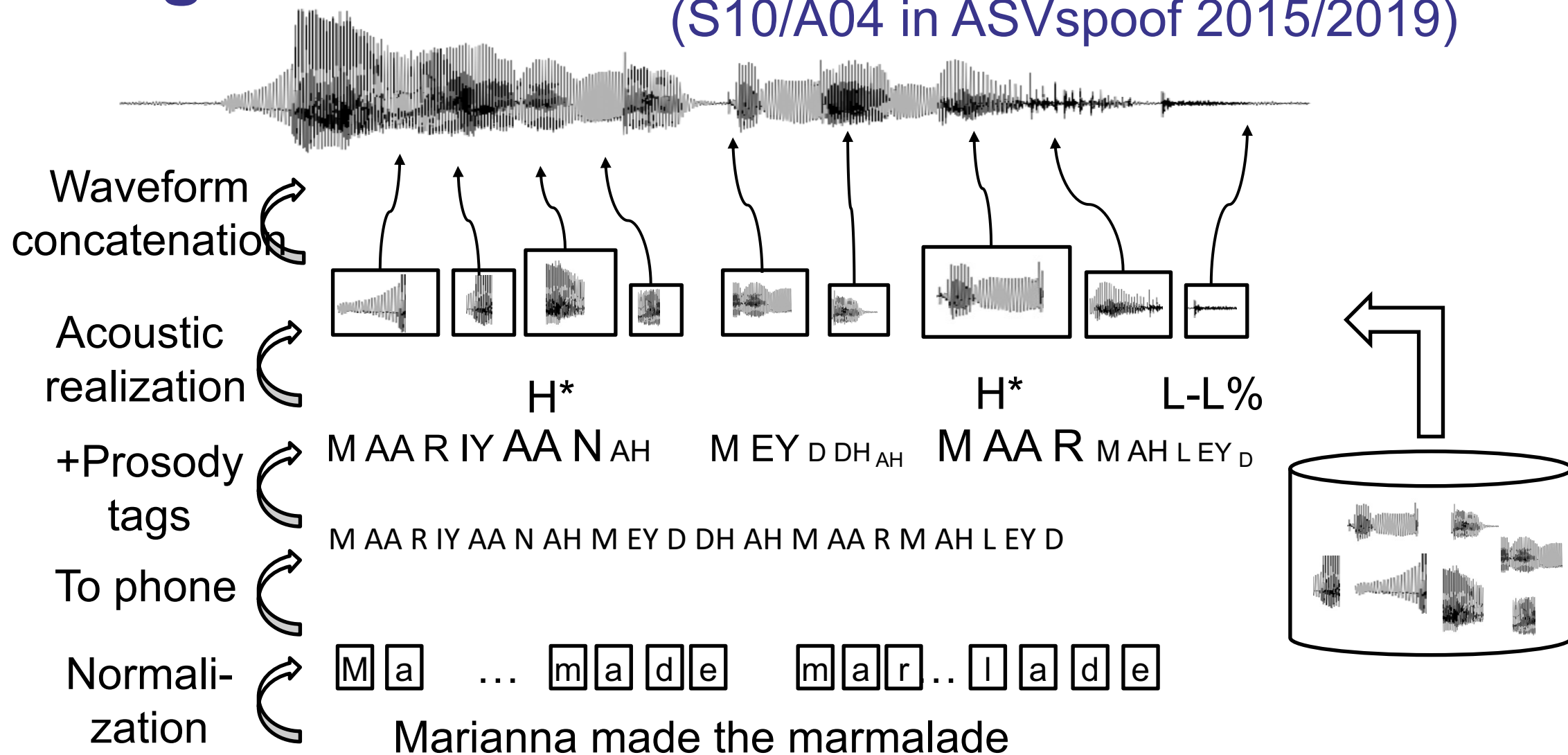
Speaker B: Marianna made the mar<sup>malade</sup>.

Speaker B: Marianna<sup>made</sup> the mar<sup>malade</sup>. 



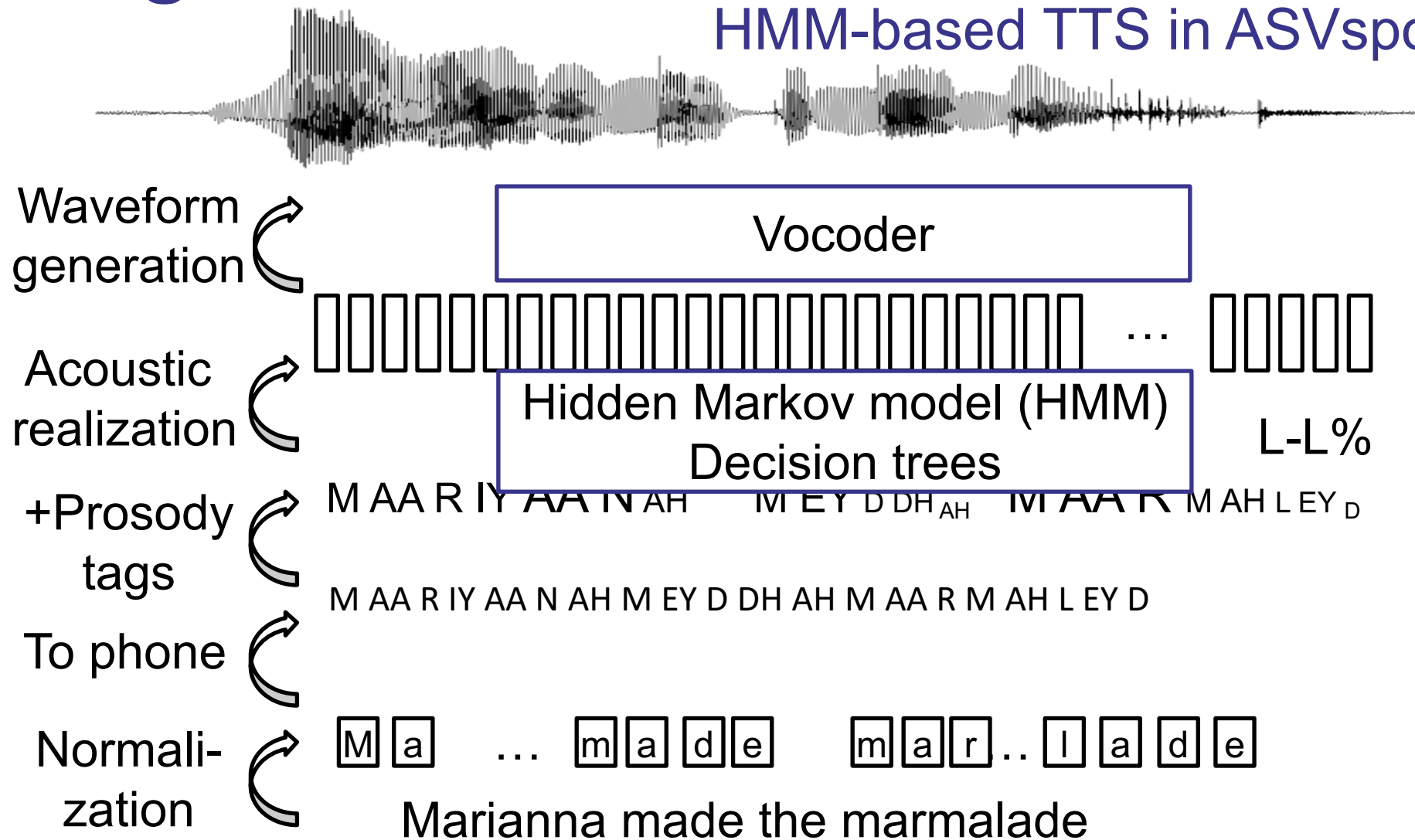
# Building TTS

## Unit-selection (S10/A04 in ASVspoof 2015/2019)



# Building TTS

## HMM-based TTS in ASVspoof 2015



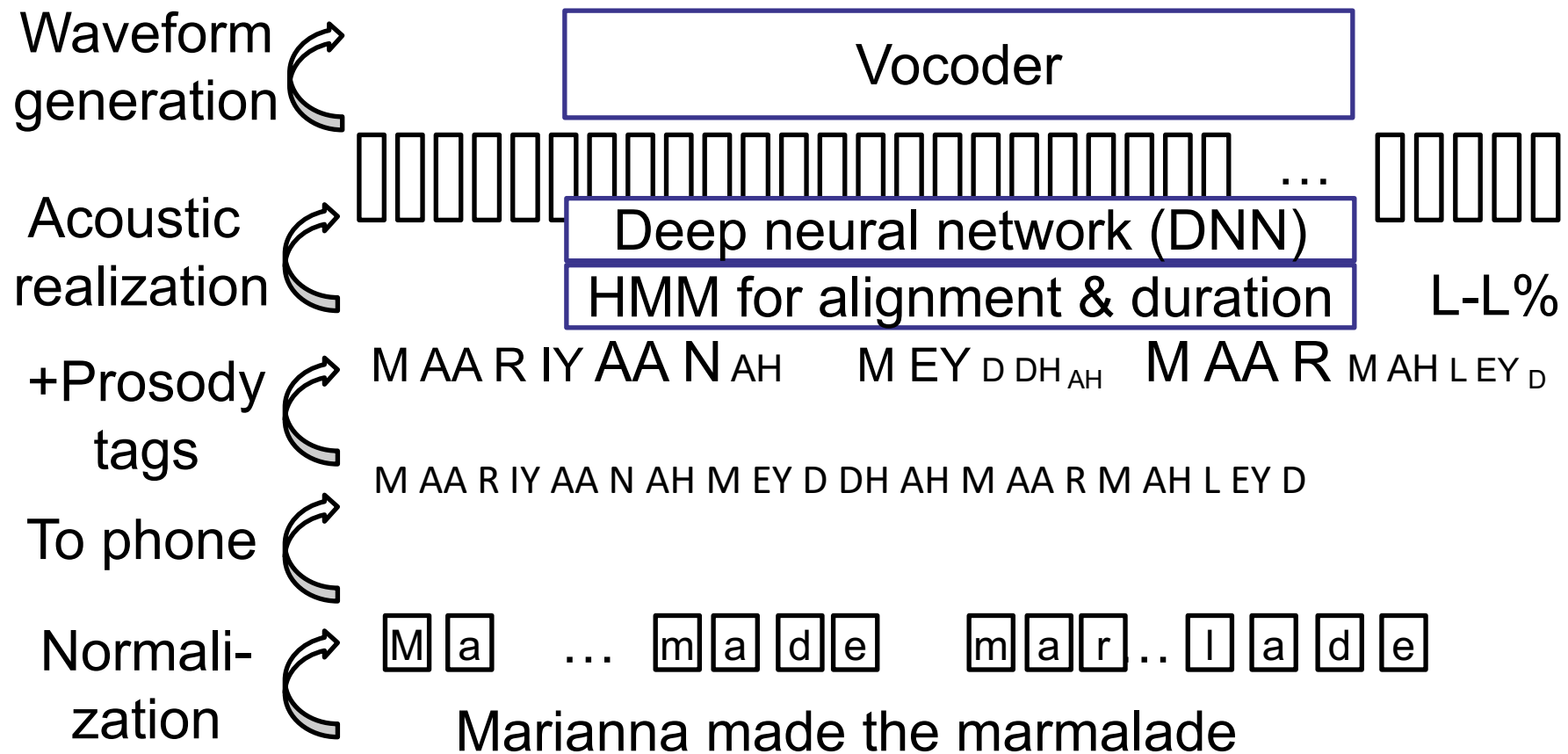


# Building TTS



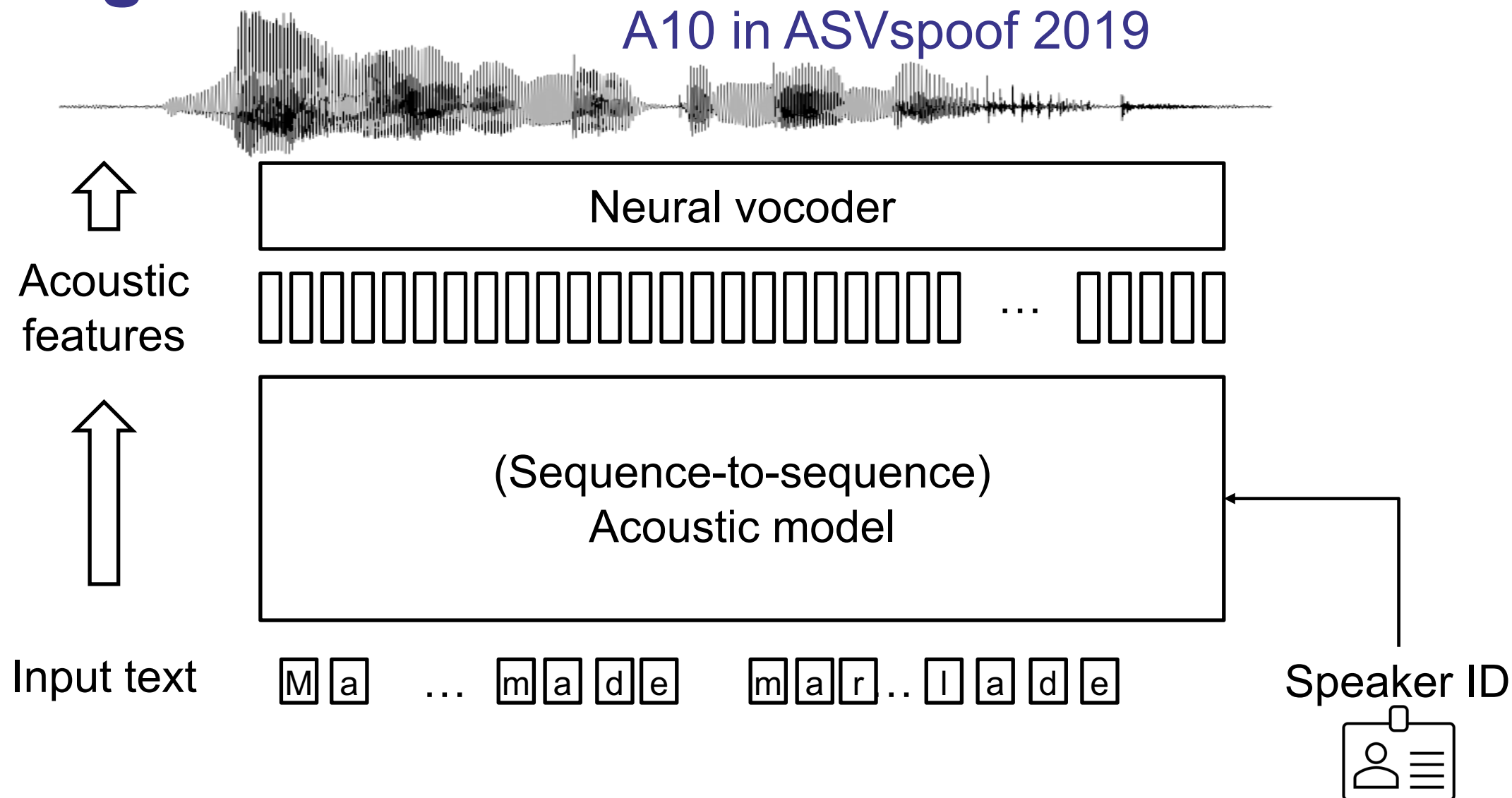
# Building TTS

A02/A03... in ASVspoof 2019

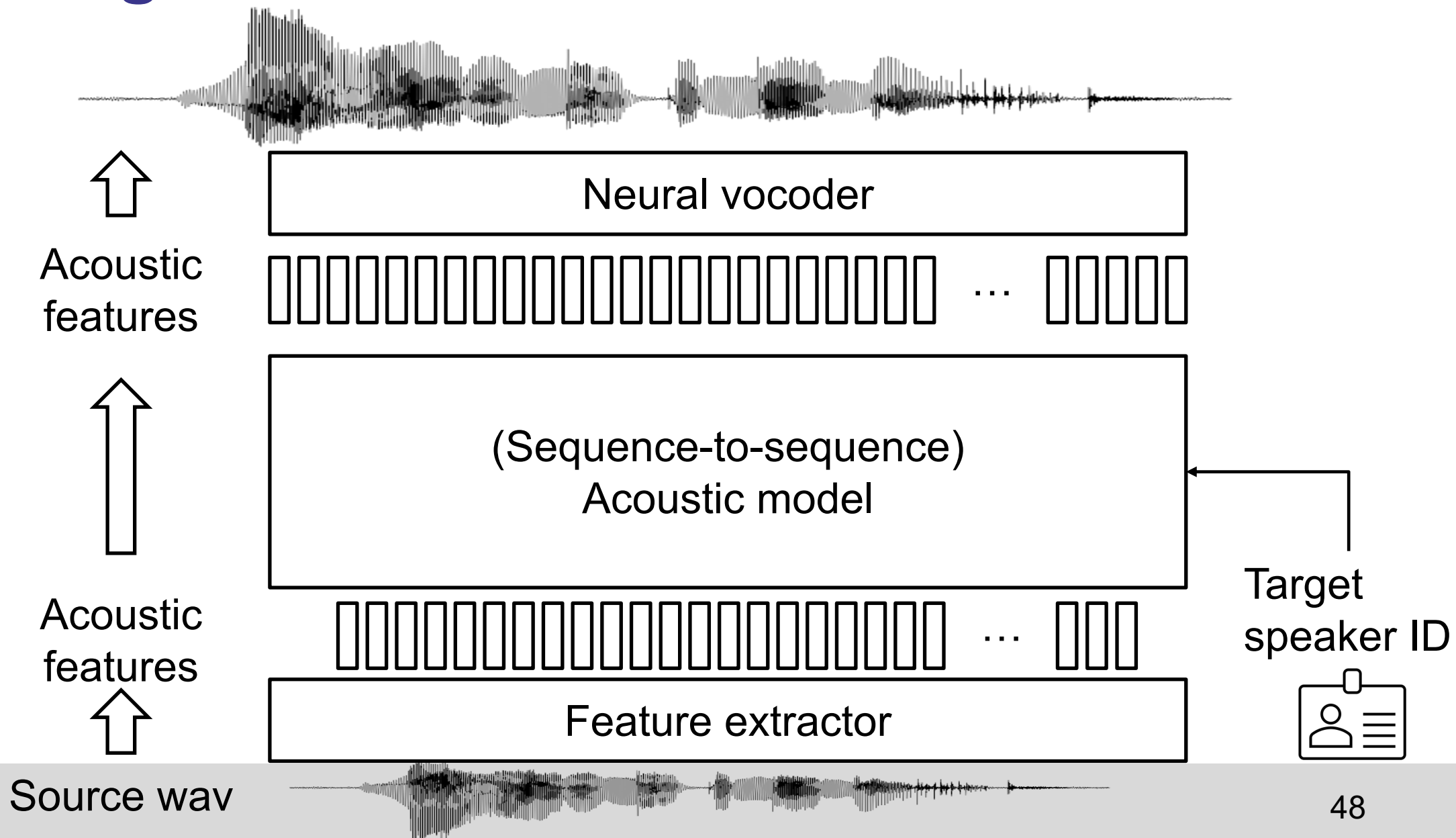


# Building TTS

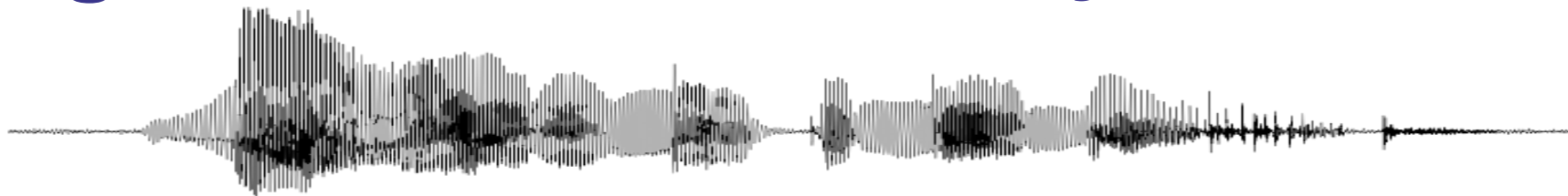
A10 in ASVspoof 2019



# Building VC



# Building TTS/VC is not that easy



Personal opinion

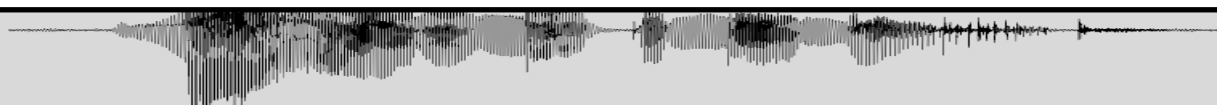
Expert knowledge  
Interpretability

Softwares available  
Computation power

Classifical methods

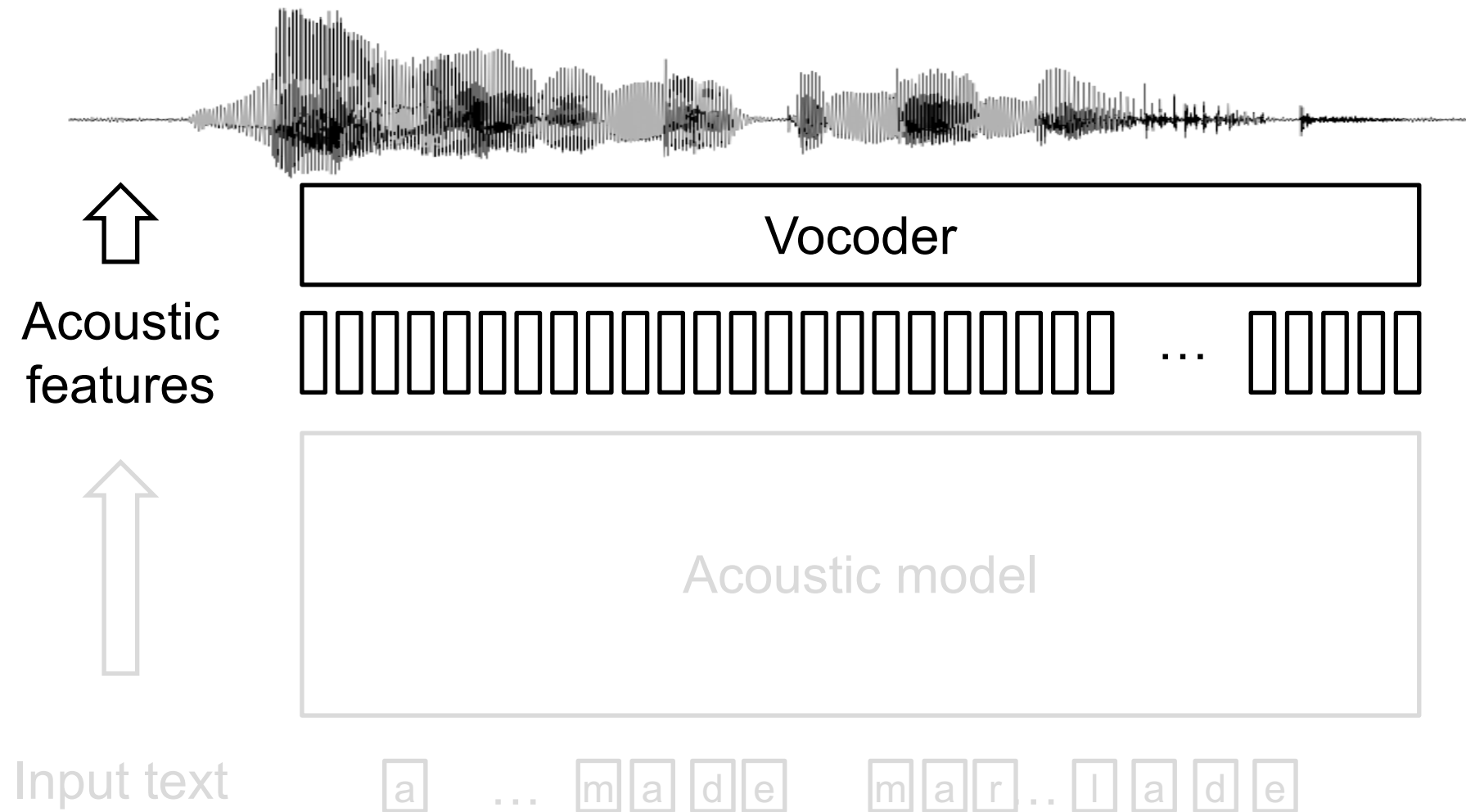


End-to-end methods

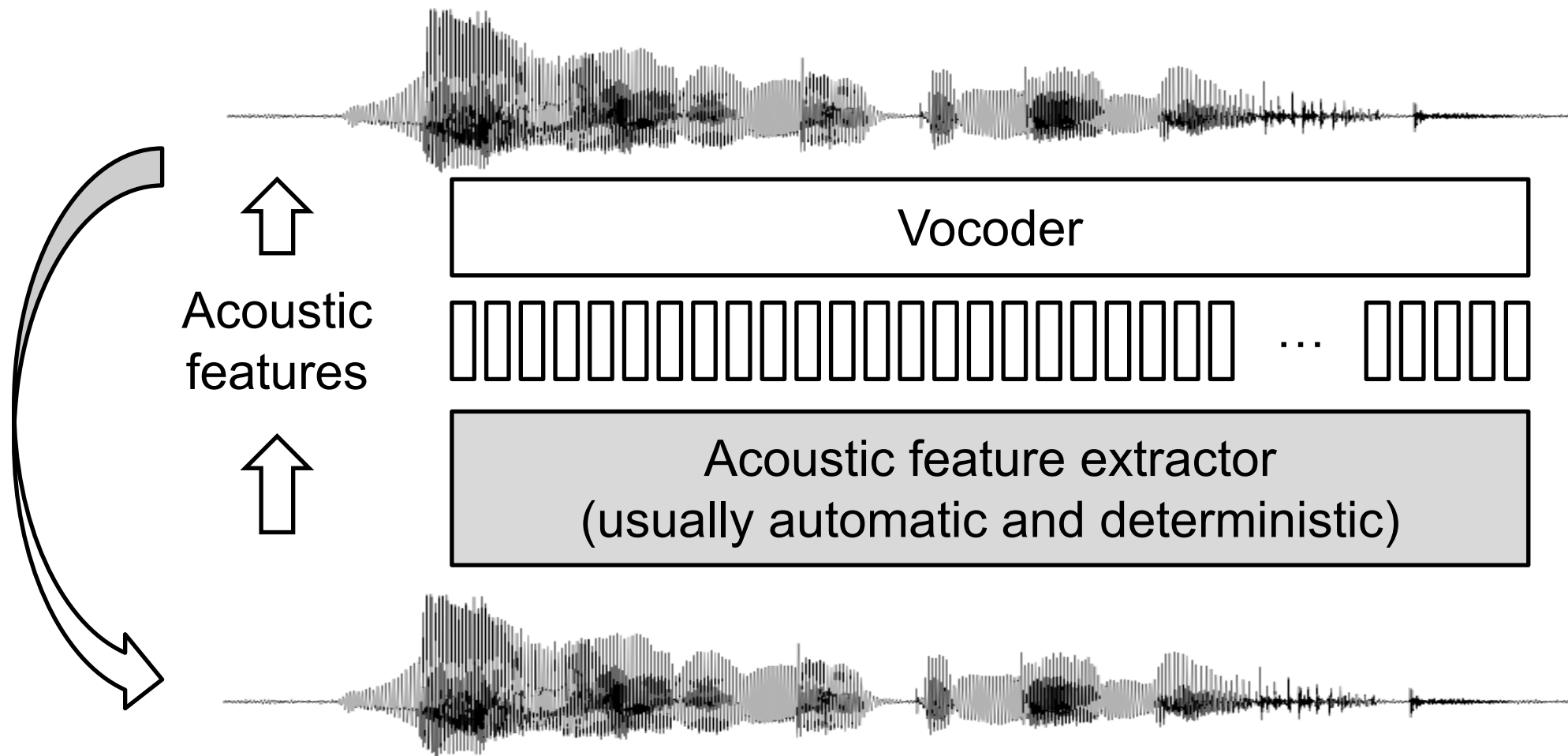


# Alternative method – vocoding

# Vocoding



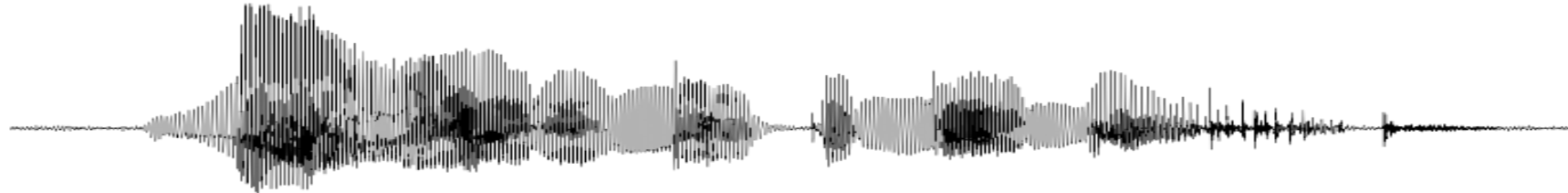
# Vocoding



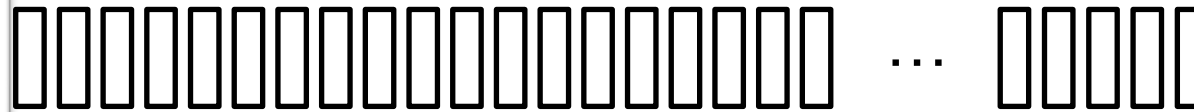
Copy-synthesis, analysis-by-synthesis, copy-resynthesis, ...



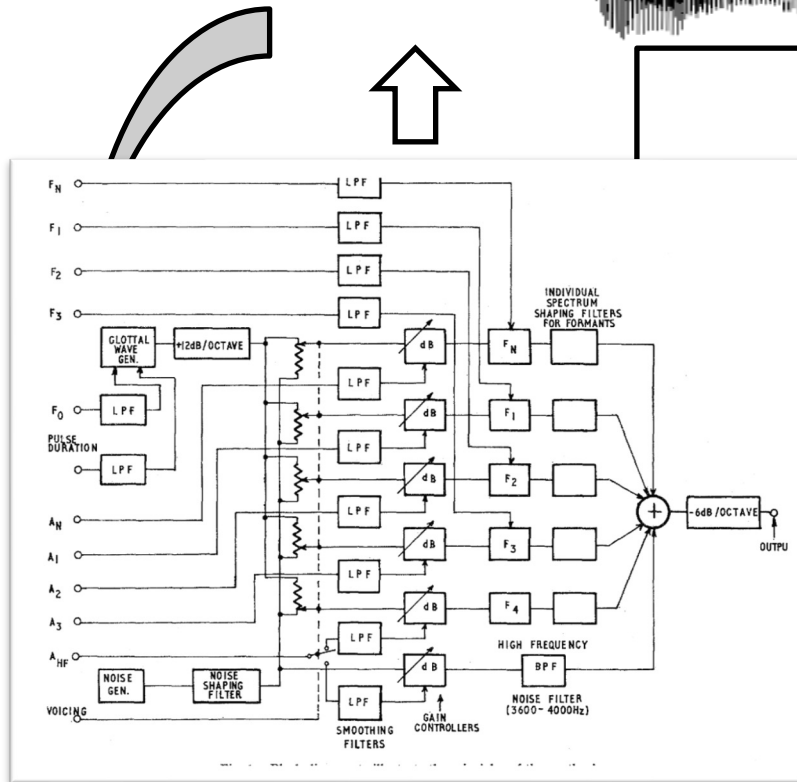
# Vocoding



Formant synthesizer (deterministic)



Acoustic feature extractor  
(usually automatic and deterministic)

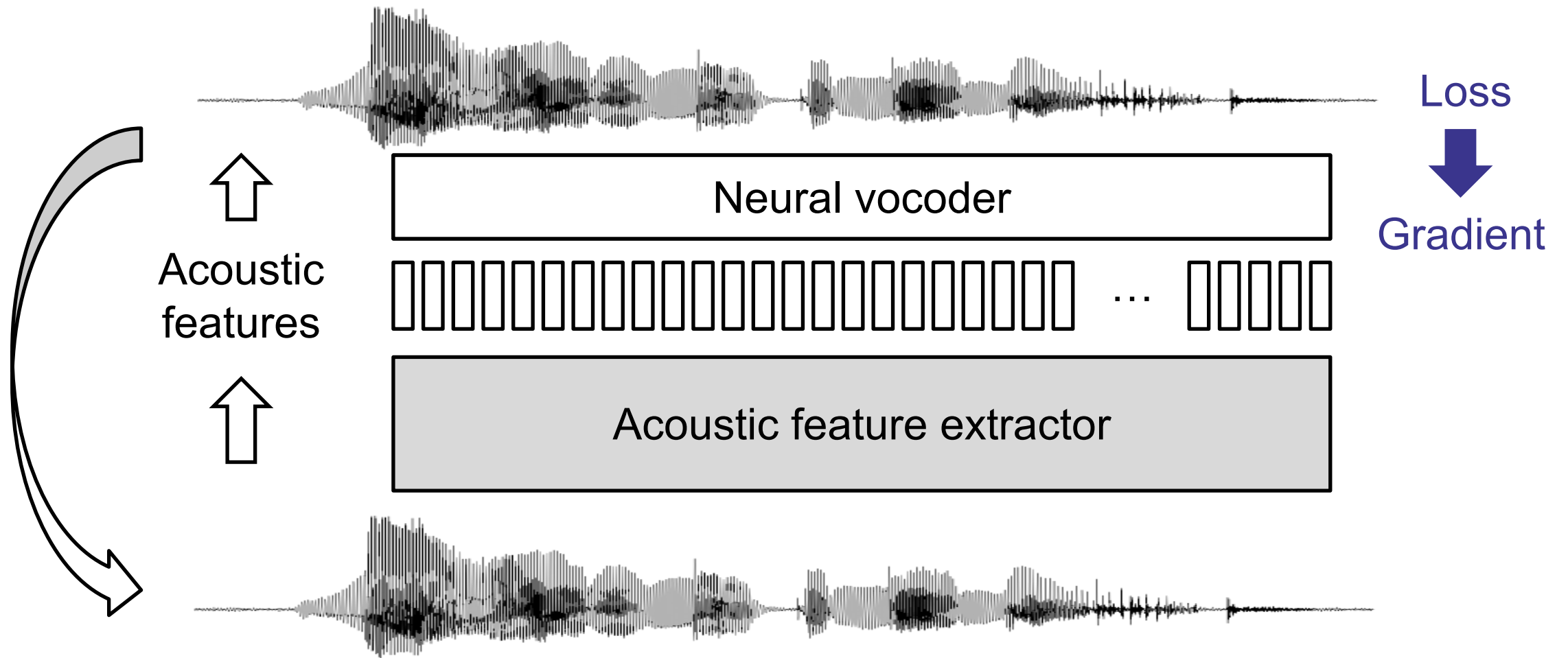


## IV. Synthesis of Copies of Natural Utterances

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, analysis-by-syn

# Vocoding using neural vocoder

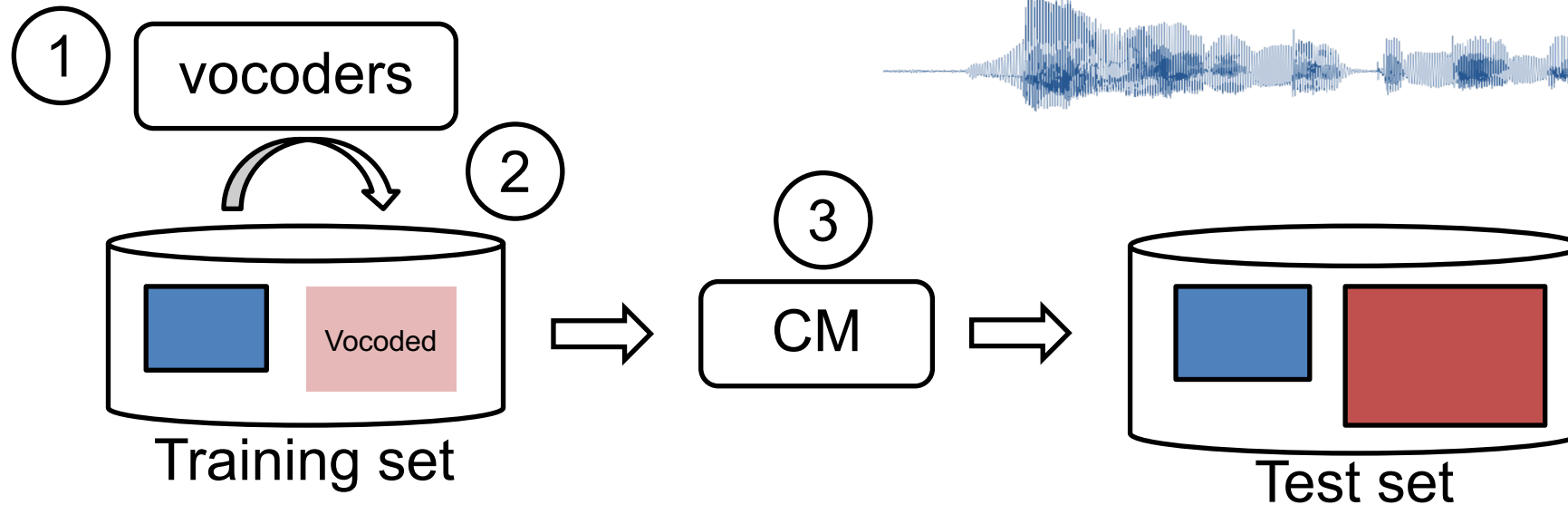
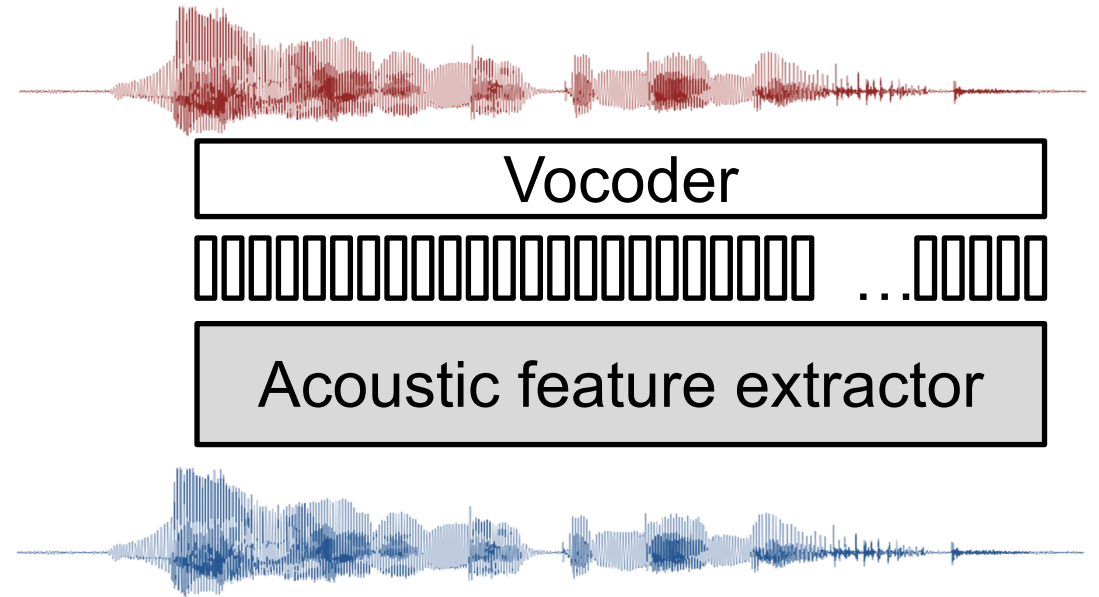


We do copy-synthesis when training the neural vocoders

# Creating vocoded spoofed data

## □ Three steps

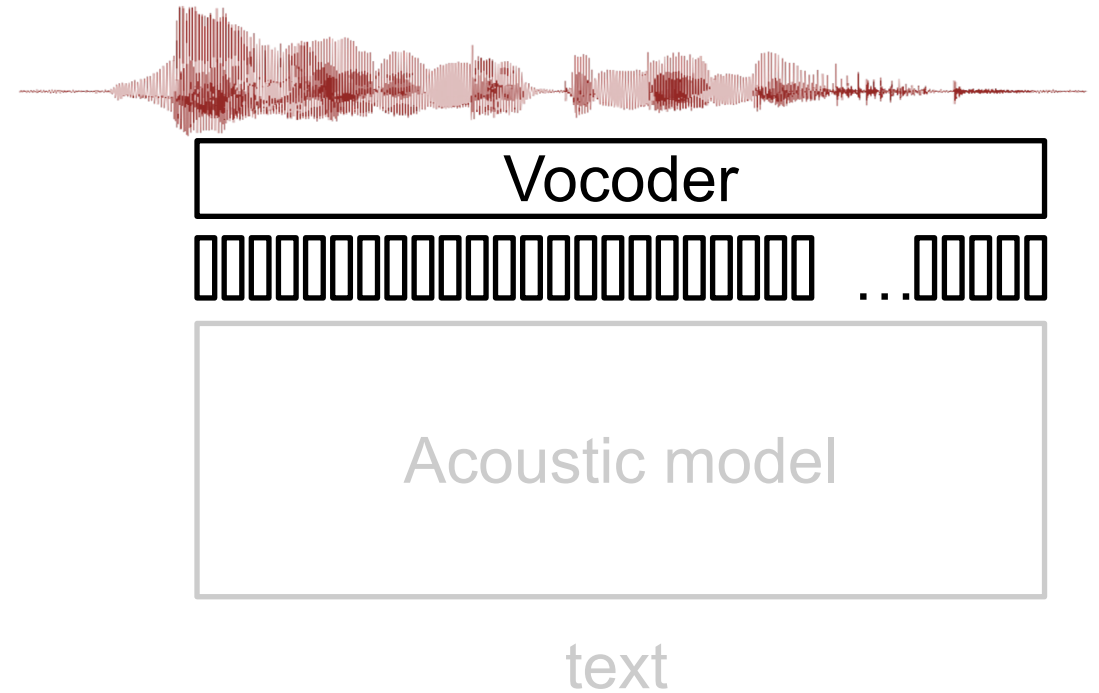
1. Prepare (or training) vocoders
2. Do vocoding on **bona fide** data
3. Train the CM using {**bona fide**, **vocoded spoofed**}



# Creating vocoded spoofed data

## □ Assumptions

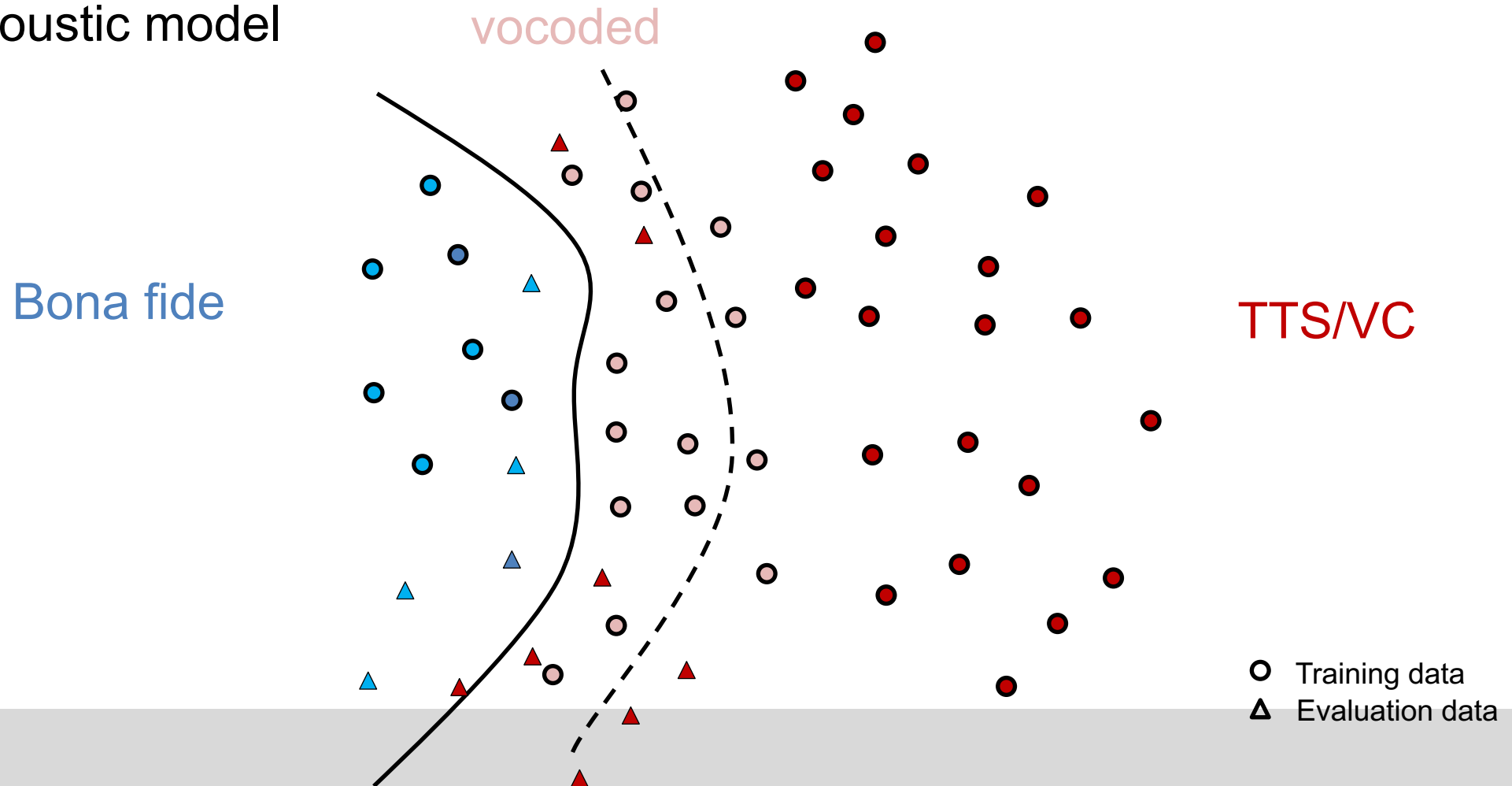
- Vocoding is TTS/VC with a perfect acoustic model
  - x artefacts by the acoustic model
  - ✓ artefacts by the vocoder



# Creating vocoded spoofed data

## □ Assumptions

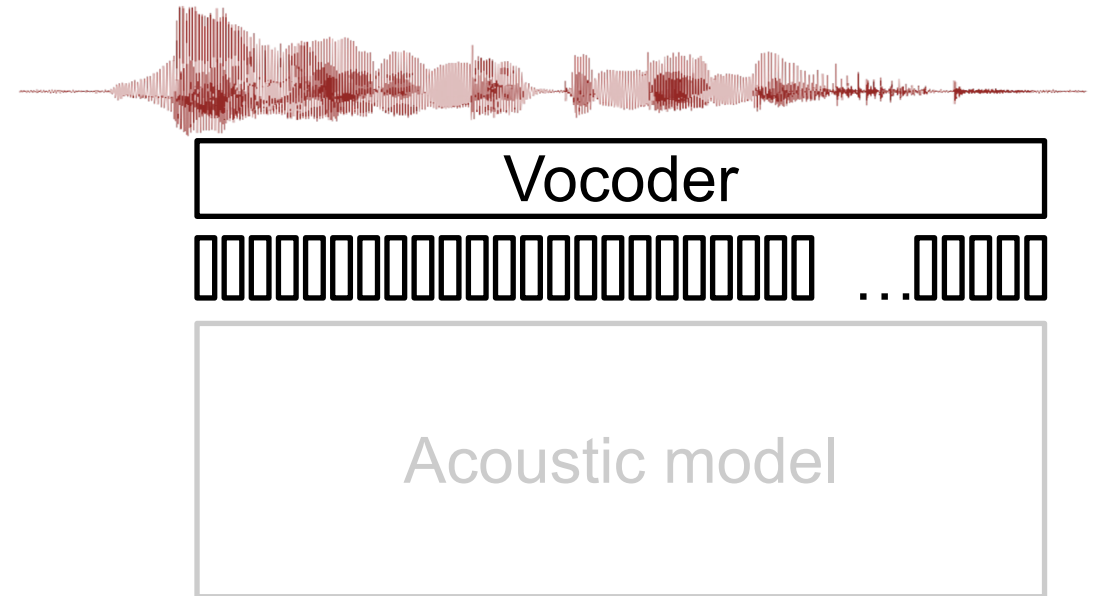
- Vocoding is TTS/VC with a perfect acoustic model



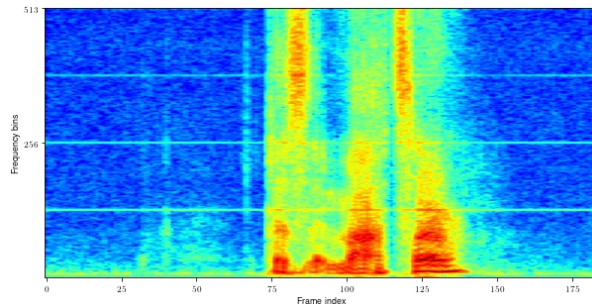
# Creating vocoded spoofed data

## □ Assumptions

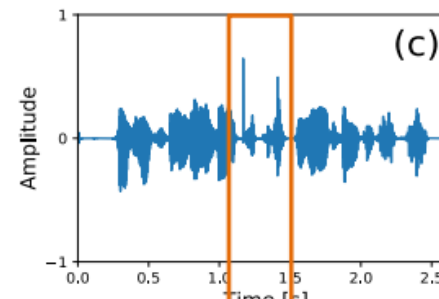
- Vocoding is TTS/VC with a perfect acoustic model
- Actual TTS/VC spoofed data contain artefacts by the vocoder



WaveGlow “bar” (Prenger 2019)



WaveNet “click” (Wu 2018) text

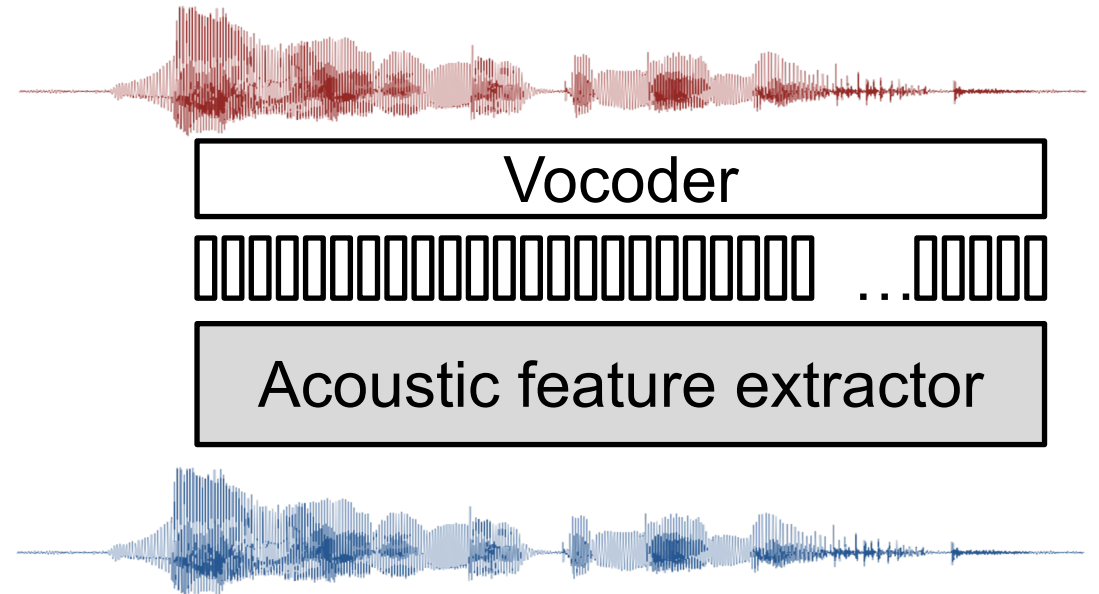


...

# Creating vocoded spoofed data

## □ Potential benefits

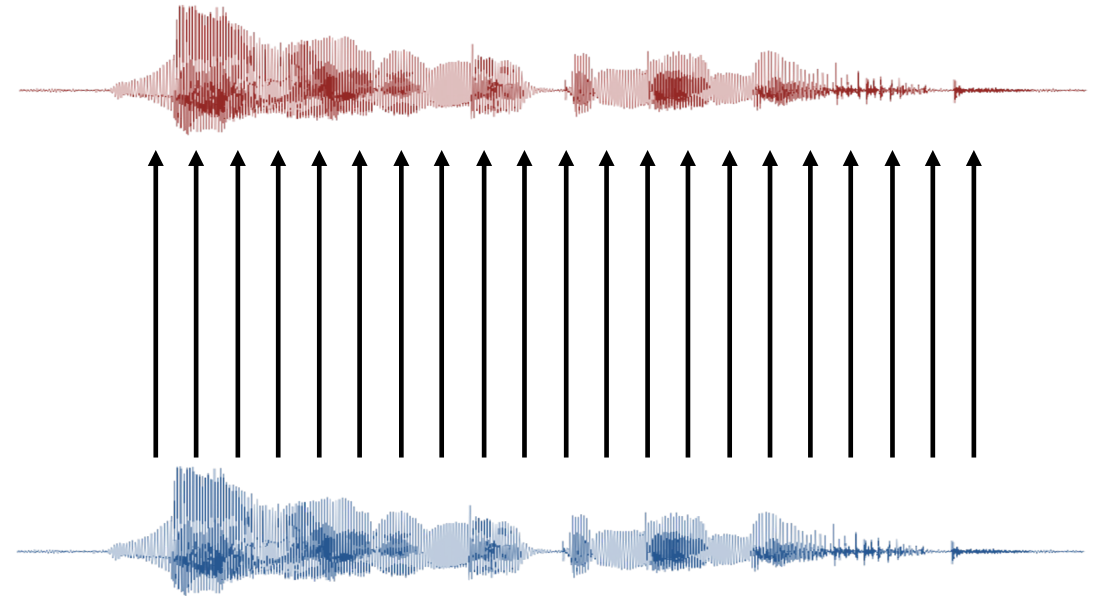
- Preparing vocoders is easier
  - ~~linguistic knowledge~~
  - ~~transcription / annotation~~
  - ~~speaker embedding~~
  - ...



# Creating vocoded spoofed data

## □ Potential benefits

- Preparing vocoders is easier
  - ~~linguistic knowledge~~
  - ~~transcription / annotation~~
  - ~~speaker embedding~~
  - ...
- **Bona fide** and **vocoded** waveforms are aligned in time





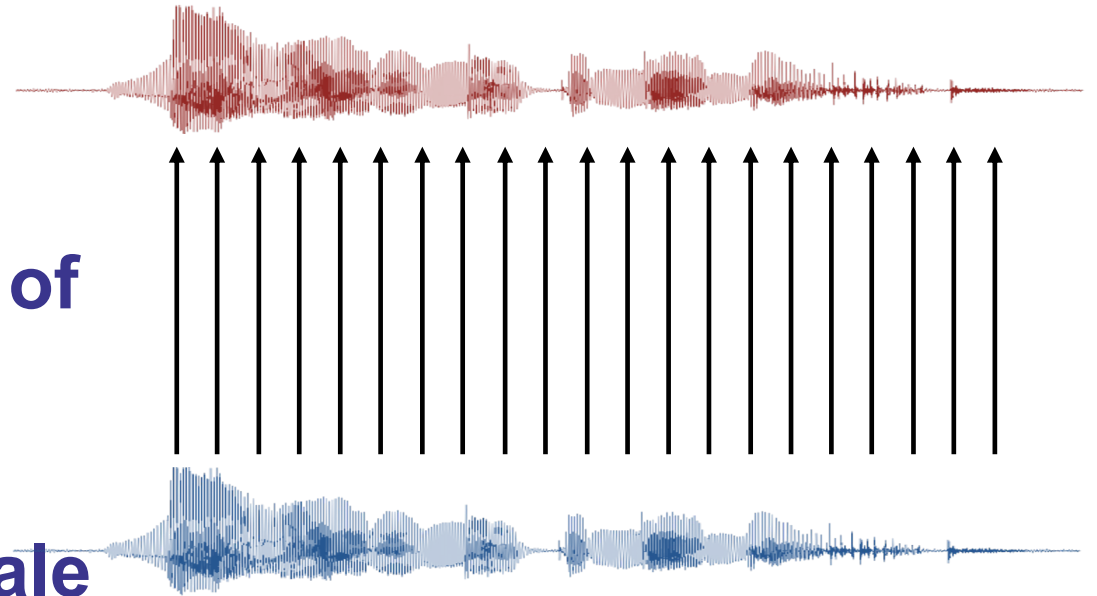
# Questions

## □ Which vocoders? How to train?

- pre-trained?
- fine tuning?

## □ How to exploit the aligned pair of {bona fide, **vocoded spoofed**}?

## □ Improvement of using large scale vocoded data?



# Vocoded spoofed data is a not new idea

## ■ 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.

## ■ 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.

# Question 1

- Which vocoders? How to train?

# Which vocoders?

## Options

Not necessary

- ~~▪ Digital-signal-processing (DSP)~~
- ~~▪ Autoregressive DNN~~
- Non-autoregressive DNN+DSP

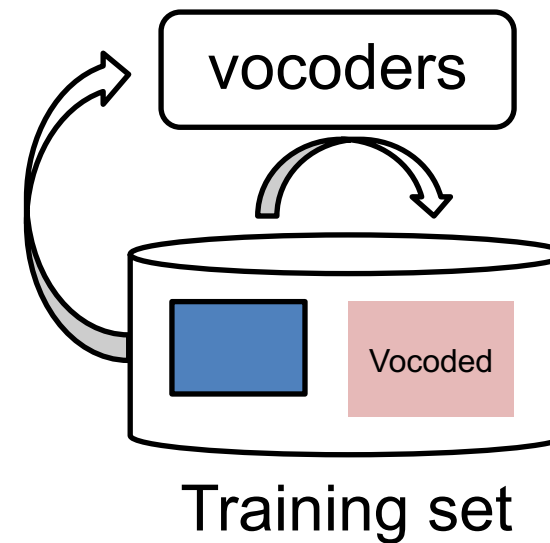
## Constraints

- Fast generation speed, much faster than real-time speed

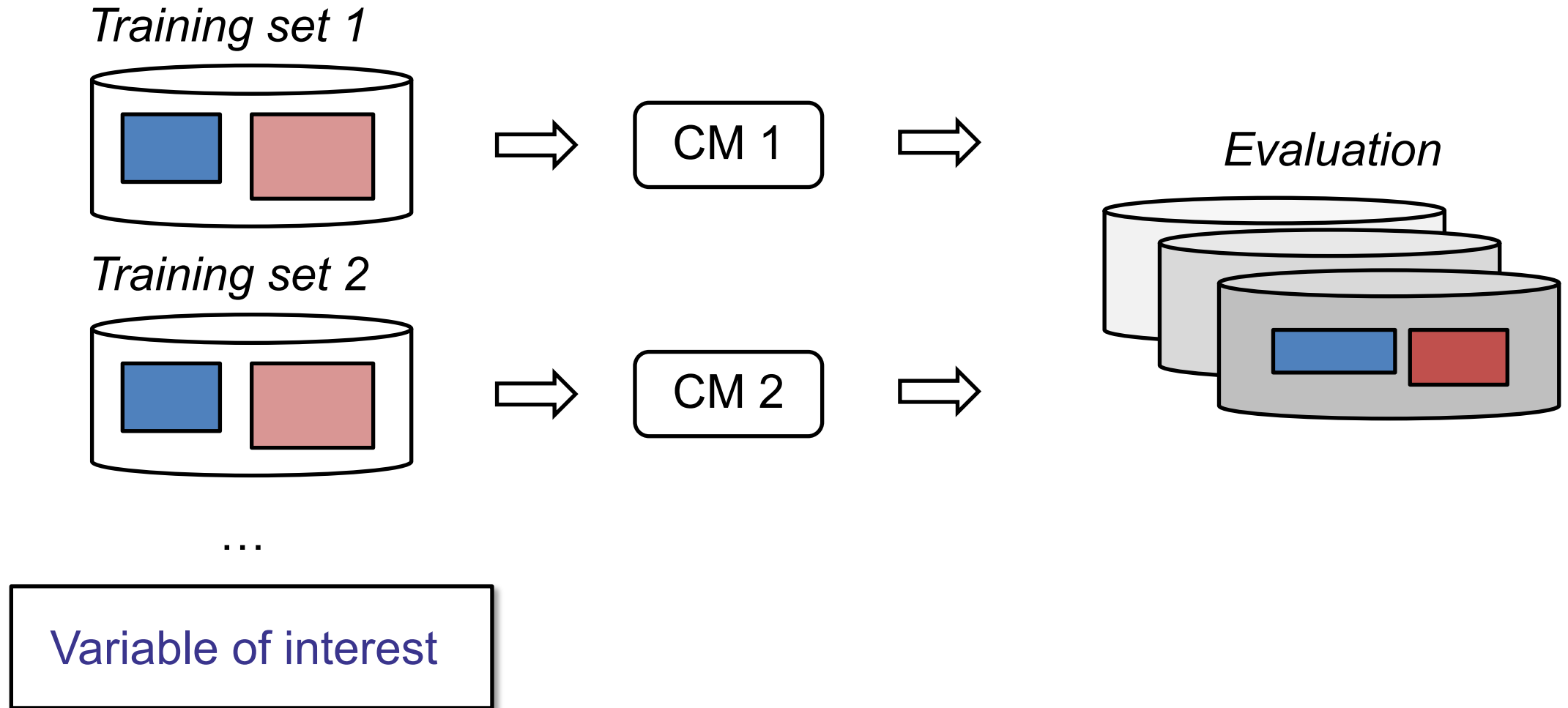
# How to train?

## Options

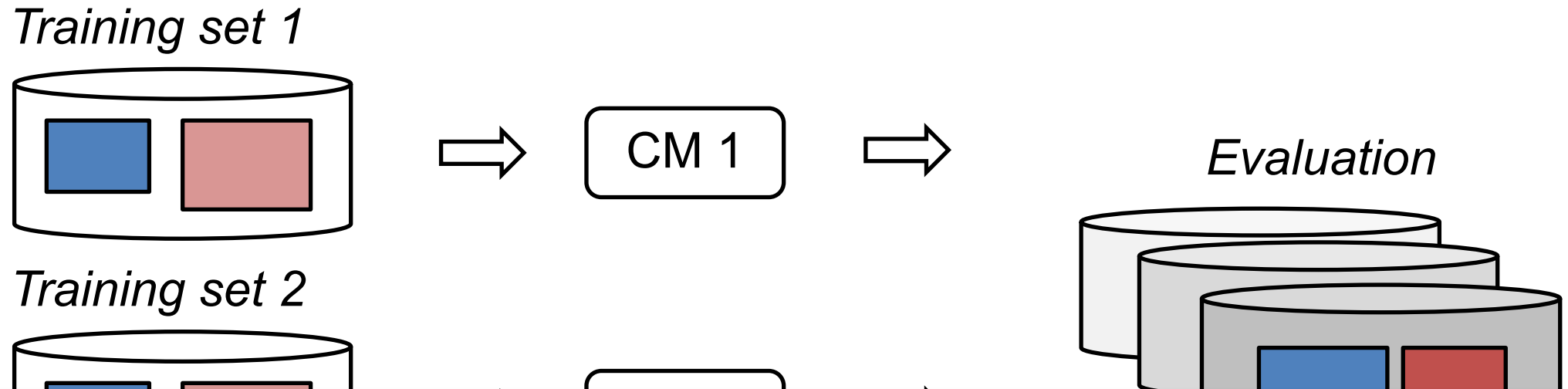
- Pre-trained, off-the-shelf
- Pre-trained, fine-tuning
- Trained from scratch



# Experiment



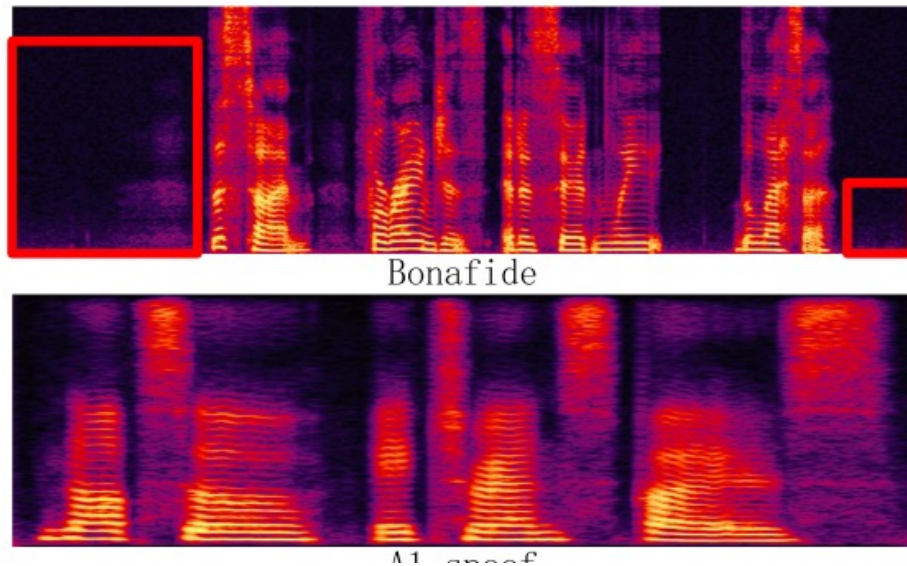
# Experiment



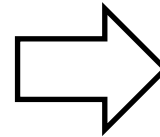
- CM: Wav2vec2.0 front end (Baevski 2020) + pooling + linear output (Wang 2022)
- Evaluation on
  - ASVspoof 2019 LA test set, 2021 LA & DF eval sets
  - ASVspoof 2019 LA test set w/o non-speech, 2021 LA & DF hidden track
  - WaveFake (Frank 2021) , In-the-Wild (Müller 2022)

# Experiment

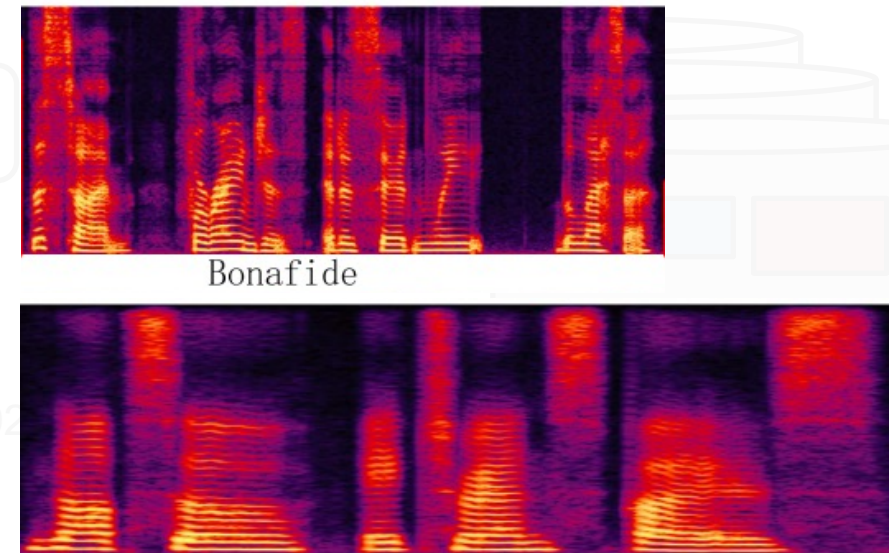
Original test trials



CM



Non-speech trimmed test trials



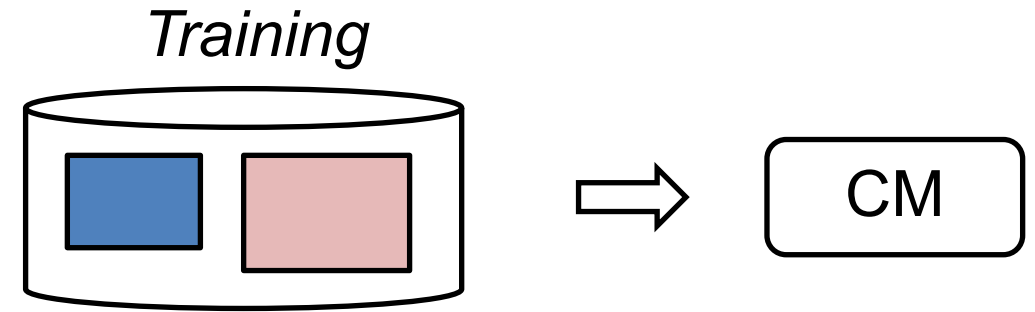
– ASVspoof 2019 LA test set, 2021 LA&DF eval track

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

I personally recommend using both versions of ASVspoof test sets

# Experiment

## □ CM training sets in comparison



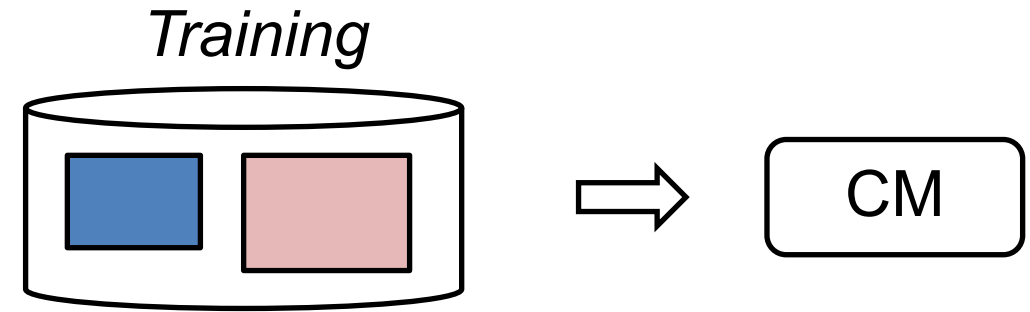
ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data
LA19trn	20	2,580	22,800	-	-	-
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -
Voc.v1	20 same as LA19trn	2,580	10,320	HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS /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

- LA19trn: ASVspoof 2019 LA training set (bona fide + **TTS/VC**)
- WFtrn: WaveFake English subset (bona fide + **vocoded**)
- Voc.v\*: ASVspoof 2019 LA training set bona fide data + **vocoded**



# Experiment

## □ CM training sets in comparison



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

### How to train the vocoder

- LA19trn: ASVspoof 2019 LA training set (bona fide + **TTS/VC**)
- WFtrn: WaveFake English subset (bona fide + **vocoded**)
- Voc.v\*: ASVspoof 2019 LA training set bona fide data + **vocoded**

# Experiment results

😊 Low EER

☹️ High EER

□ EER (% , mean of three runs)

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
ASVspoof 2019 LA	→ LA19eval	2.98	44.48	5.78	5.32	8.74	4.36
ASVspoof 2021 LA	→ LA21eval	7.53	41.57	26.30	17.98	19.29	24.39
ASVspoof 2021 DF	→ 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	9.52
	LA21hid	28.80	27.60	28.30	19.49	17.62	21.43
	DF21hid	23.62	26.18	22.01	13.92	13.50	16.99
	WaveFake	15.76	-	39.27	34.05	17.10	10.89
	InWild	26.65	19.98	41.06	36.46	22.26	19.45
Single EER threshold	→ Pooled	14.24	-	36.57	39.95	19.39	16.35

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

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
Test sets	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
	LA19etrim	15.56	31.62	23.29	16.16	14.99	9.52
	LA21hid	28.80	27.60	28.30	19.49	17.62	21.43
	DF21hid	23.62	26.18	22.01	13.92	13.50	16.99
	WaveFake	15.76	-	39.27	34.05	17.10	10.89
	InWild	26.65	19.98	41.06	36.46	22.26	19.45
	Pooled	14.24	-	36.57	39.95	19.39	16.35

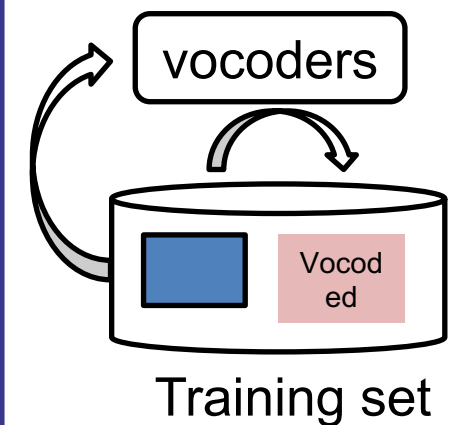
ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data
LA19trn	20	2,580	22,800	-	-	-
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -
Voc.v2	20			HiFiGAN, NSF-HiFiGAN,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

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
Test sets	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
	LA19etrim	15.56	31.62	23.29	10.41	17.68	21.11
	LA21hid	28.80	27.60	28.30	13.92	13.50	16.99
	DF21hid	23.62	26.18	22.01	11.05	17.10	10.89
	WaveFake	15.76	-	39.27	31.05	17.10	10.89
	InWild	26.65	19.98	41.06	36.46	22.26	19.45
	Pooled	14.24	-	36.57	39.95	19.39	16.35

**Vocoders pre-trained by ESPNet (Hayashi 2020)**

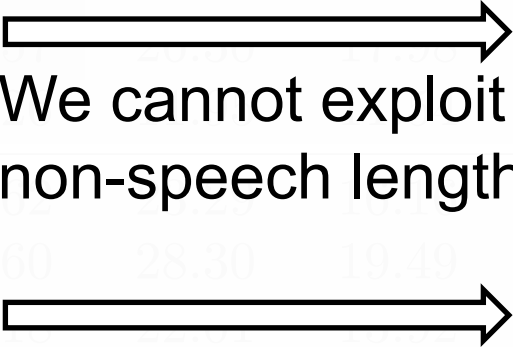
ID	#. Spr.	#. Bona.	#. Spoof.	Vocoder type	Implementation	Vocoder train/fine-tune data
LA19trn	20	2,580	22,800	-	-	-
WFtrn	1	3,930	15,720	HiFiGAN, MB-MelGAN, PWG, WaveGlow	ESPNet toolkit	LJSpeech / -
Voc.v1				HiFiGAN, MB-MelGAN, PWG, StyleMelGAN	ESPNet toolkit	LibriTTS / -
Voc.v2	20			HiFiGAN, NSF-HiFiGAN,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

		pretrained			trained from scratch		finetuned
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	
Test sets	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
	LA19etrim	15.56	31.62	23.29	16.16	14.99	9.52
	LA21hid	28.80	27.60	28.30	19.49	17.62	21.43
	DF21hid	23.62	26.18	22.01	13.92	13.50	16.99
	WaveFake	15.76	-	39.27	34.05	17.10	10.89
	InWild	26.65	19.98	41.06	36.46	22.26	19.45
	Pooled	14.24	-	36.57	39.95	19.39	16.35



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

		Training set					
		LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
Test sets	LA19eval	2.98	14.18	5.78	5.32	8.71	4.36
	LA21eval	7.53	14.18	5.78	5.32	19.29	24.39
	DF21eval	6.67	24.39	14.18	5.78	9.71	13.31
	LA19etrim	15.56	31.00	14.18	5.78	14.99	9.52
	LA21hid	28.80	27.60	28.30	19.49	17.62	21.43
	DF21hid	23.62	26.65	14.18	5.78	13.50	16.99
	WaveFake	15.76	-	-	-	17.20	10.89
	InWild	26.65	19.98	41.06	30.46	22.26	19.45
	Pooled	14.24	-	36.57	39.95	19.39	16.35


  
**We cannot exploit non-speech length**  
**Reasonably good**

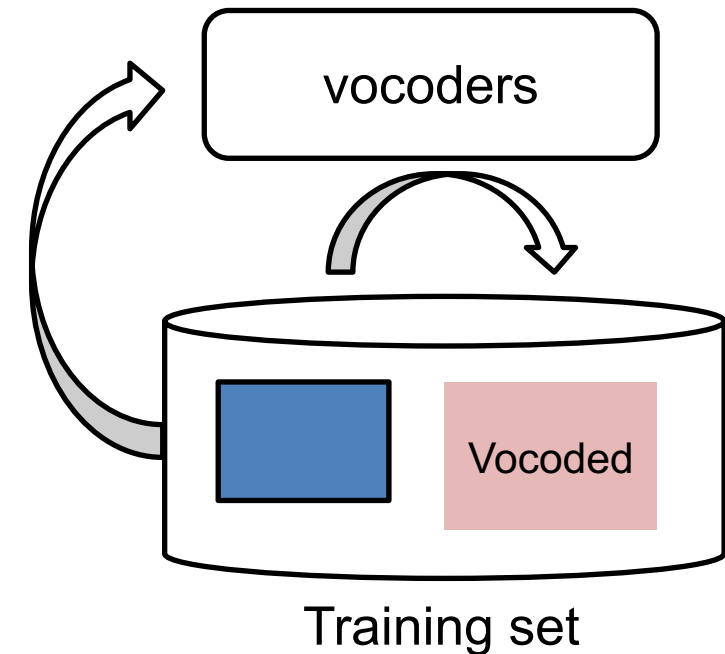
# Vocoding data are not useless

## □ Which vocoders?

- practical choice – non-autoregressive neural vocoders
- more analysis later

## □ How to train?

- expose vocoder to the data to be vocoded



# Vocoding data are not useless

## ❑ CAUTION !

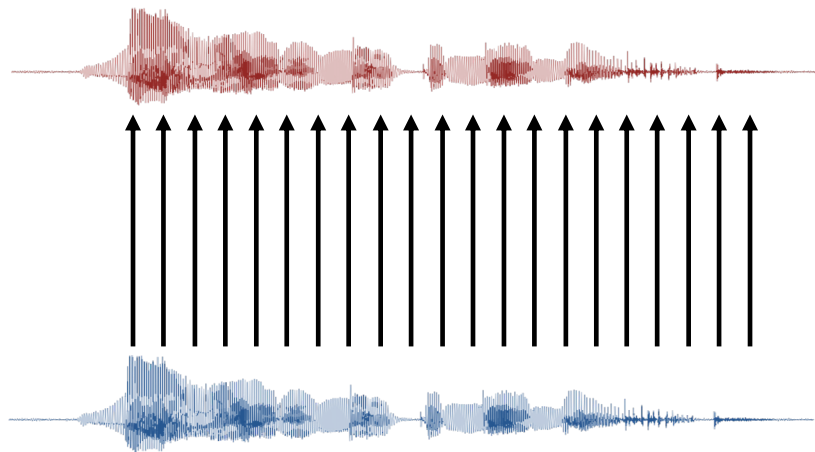
- Vocoder overlap with test sets?
  - partially in ASVspoof DF 2021
  - unknown in In-the-wild

- Other CMs performed poorly  
<https://arxiv.org/abs/2210.10570>

	LA19 trn	WF trn	Voc. v1	Voc. v2	Voc. v3	Voc. v4
LA19trn	0.10	41.69	14.25	40.72	22.83	28.08
LA15eval	32.42	23.82	49.89	28.34	23.90	37.60
LA19eval	3.32	49.80	23.52	40.79	37.71	29.68
LA21eval	23.38	57.18	36.43	62.69	49.17	49.84
DF21eval	29.45	47.40	40.27	52.76	47.22	51.76
LA19etrim	21.18	40.10	33.23	41.97	40.94	40.56
LA21hid	39.27	49.53	44.59	49.07	44.18	49.52
DF21hid	35.86	45.12	43.05	48.52	45.18	49.81
WaveFake	45.85	-	22.17	21.15	14.38	26.16
InWild	72.19	91.28	84.33	45.93	61.38	28.60
Pooled	37.68	-	43.90	51.16	46.26	51.35



- How to exploit the aligned pair of {bona fide, vocoded spoofed}?



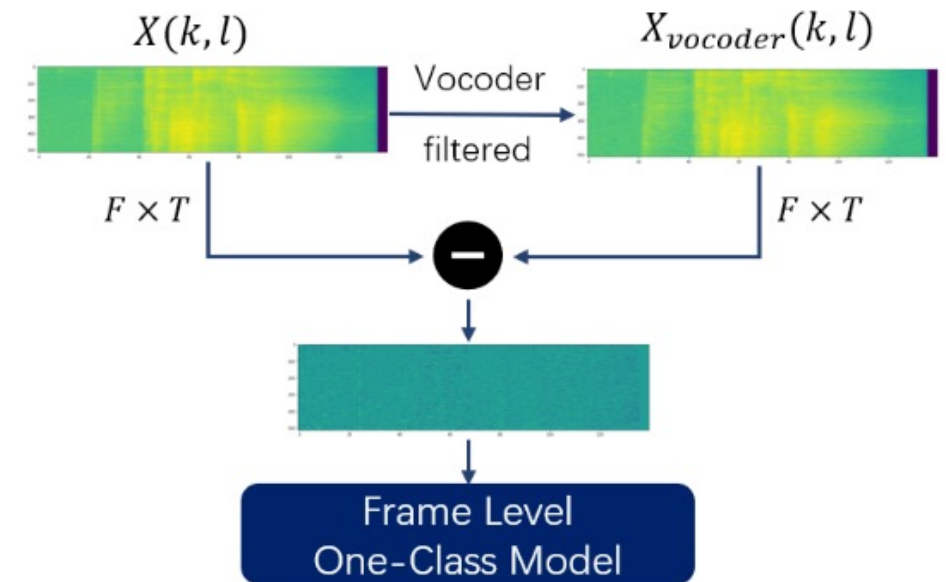
# Contrastive learning

## □ Use contrastive features (Wang 2021)

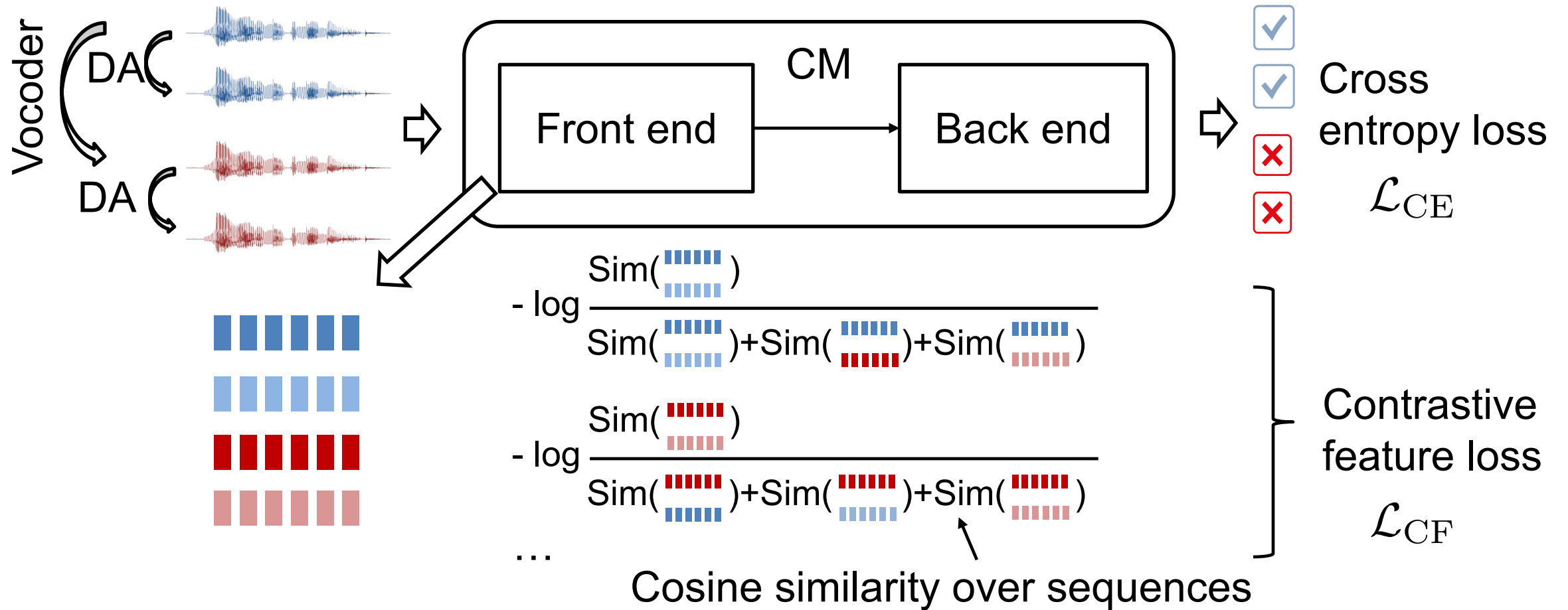
- Vocoder is needed during testing

## □ Use contrastive learning

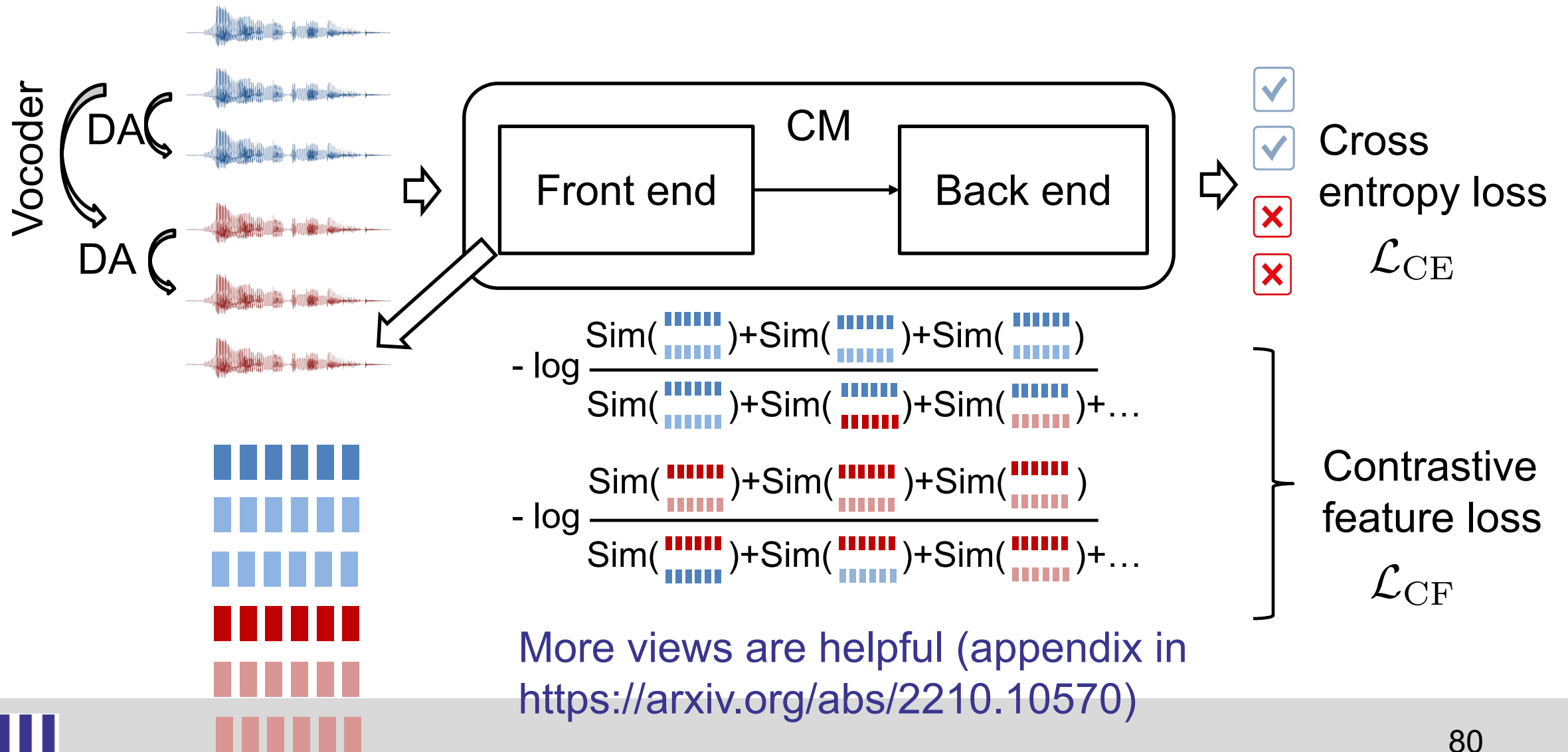
- Supervised contrastive loss (Khosla 2020)
- Vocoder is NOT needed during testing
- What is needed:
  - an additional CM training loss
  - data augmentation (DA) to create multi-view, e.g., RawBoost (Tak 2022)



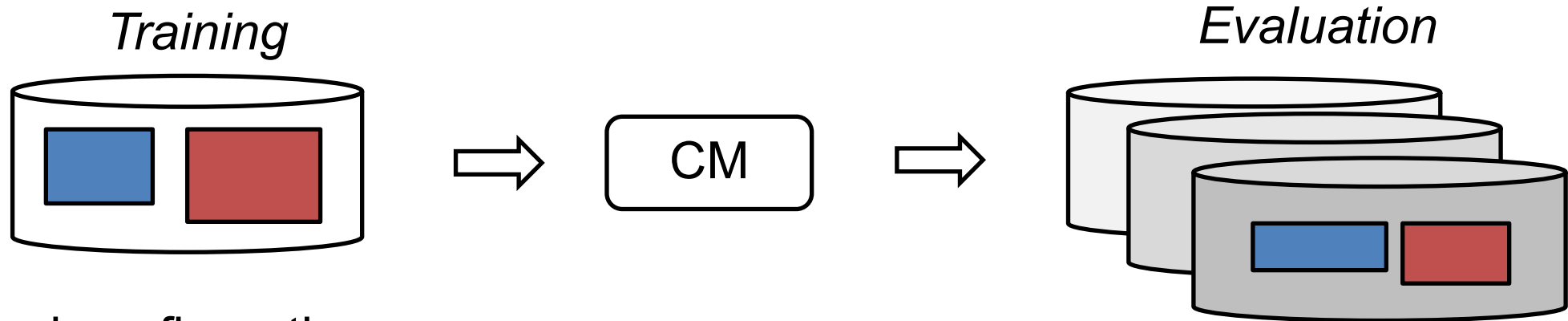
# Contrastive learning



# Contrastive learning



# Experiment



- Fixed configuration
  - training data **Voc.v4** (ASVspoof 2019 LA trn bona fide + **vocoded**)
  - Wav2vec-based CM
  - multiple test sets
- Variable of interest: how CM is trained

# Experiment 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	Voc. v4	
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	

← +Contrastive loss

# Experiment 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	Voc. v4
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	3.79
	LA21hid	28.80	21.43	21.18	19.37	26.98	14.57
	DF21hid	23.62	16.99	13.64	14.29	16.85	7.78
	WaveFake	15.76	10.89	26.37	6.87	24.62	2.50
	InWild	26.65	19.45	16.17	12.08	17.07	7.55
	Pooled	14.24	16.35	13.12	13.13	13.68	13.15

best

② vs ④

RawBoost is useful

④ vs ⑦

Contrastive feature loss is useful

# Experiment results

from Experiment I      control groups      best

Training criterion	$\mathcal{L}_{CE}$				$\mathcal{L}_{CE} + \mathcal{L}_{CF}$			
	×		RawBoost		RawBoost			
Data augmentation	LA19 trn	Voc. v4	LA19 trn	Voc. v4	LA19 trn	Voc. v4	Voc. v4	
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	

vocoded data  
+ contrastive learning is  
helpful



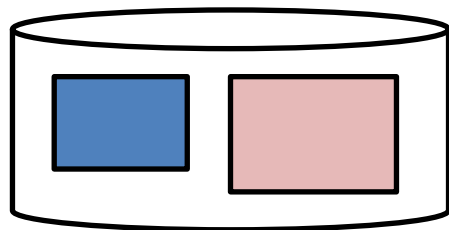
# Analysis

# Can we detect TTS/VC w/ unseen vocoders?

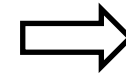
## □ Options

- ~~DSP vocoders~~
- ~~Neural autoregressive (AR) vocoders~~
- Non-autoregressive DNN+DSP

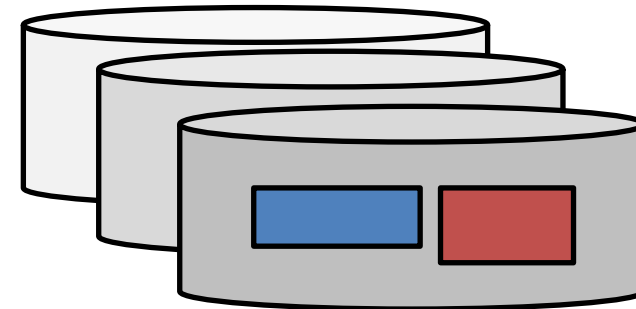
Vocoded data from non-AR vocoders



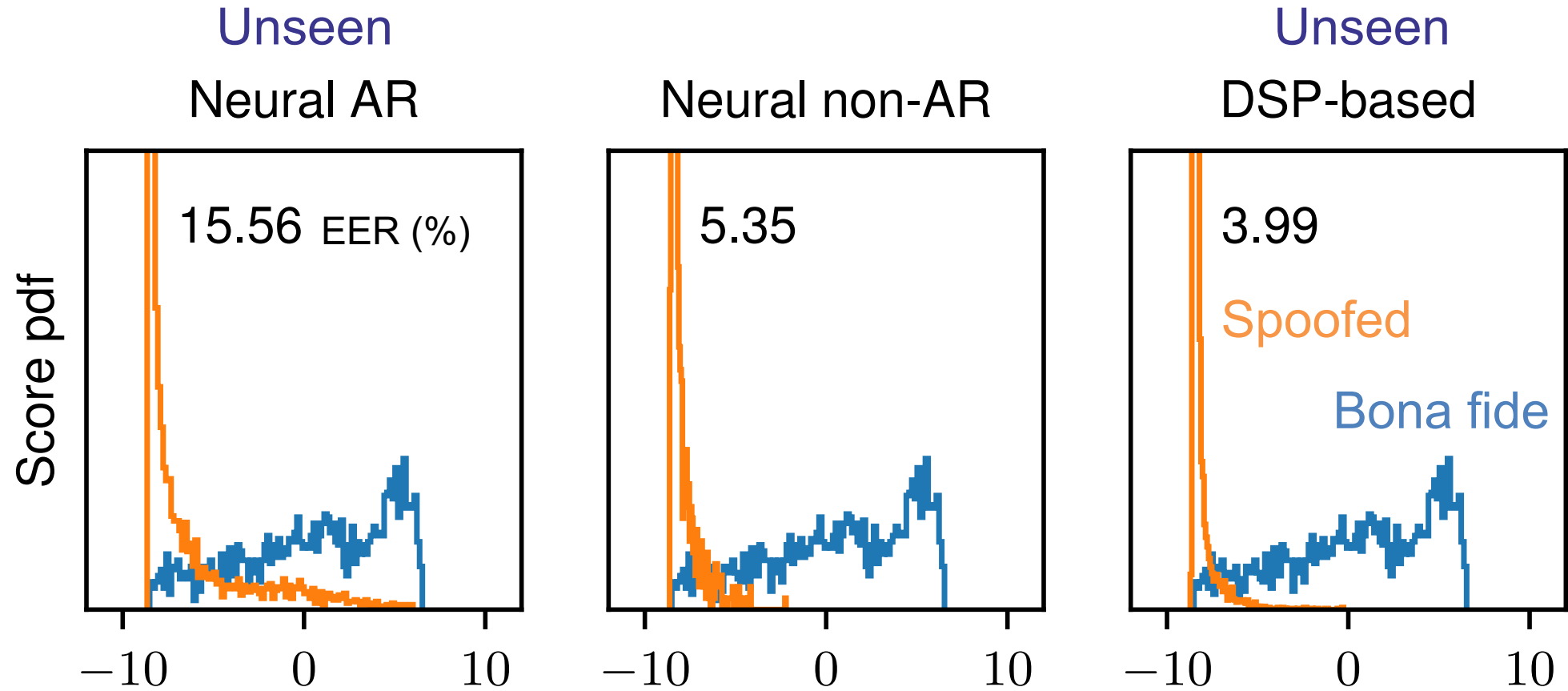
CM



Spoofed data from TTS/VC using Neural AR / DSP vocoders

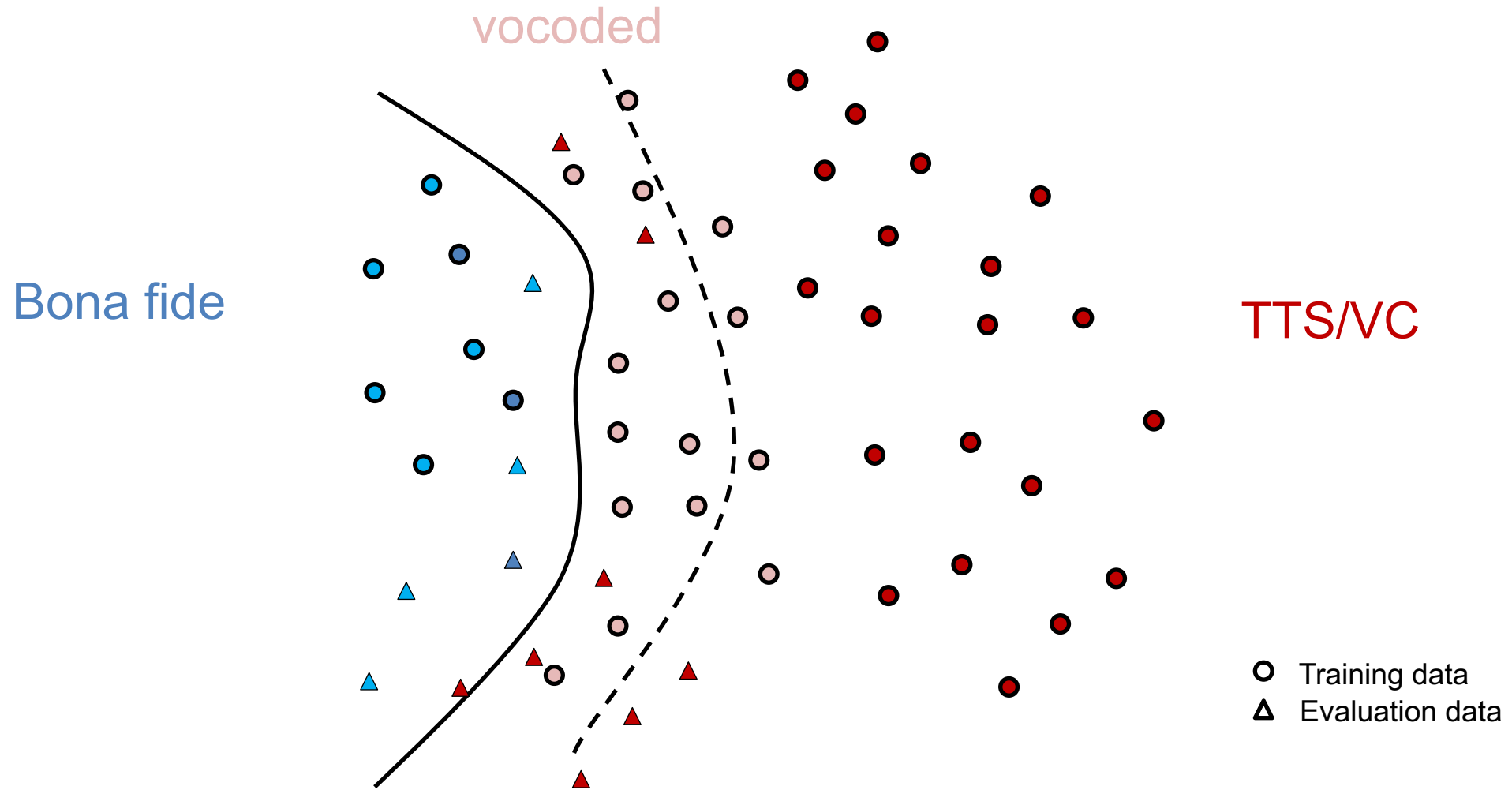


# Can we detect TTS/VC w/ unseen vocoders?



Results on ASVspoof 2021 DF test set

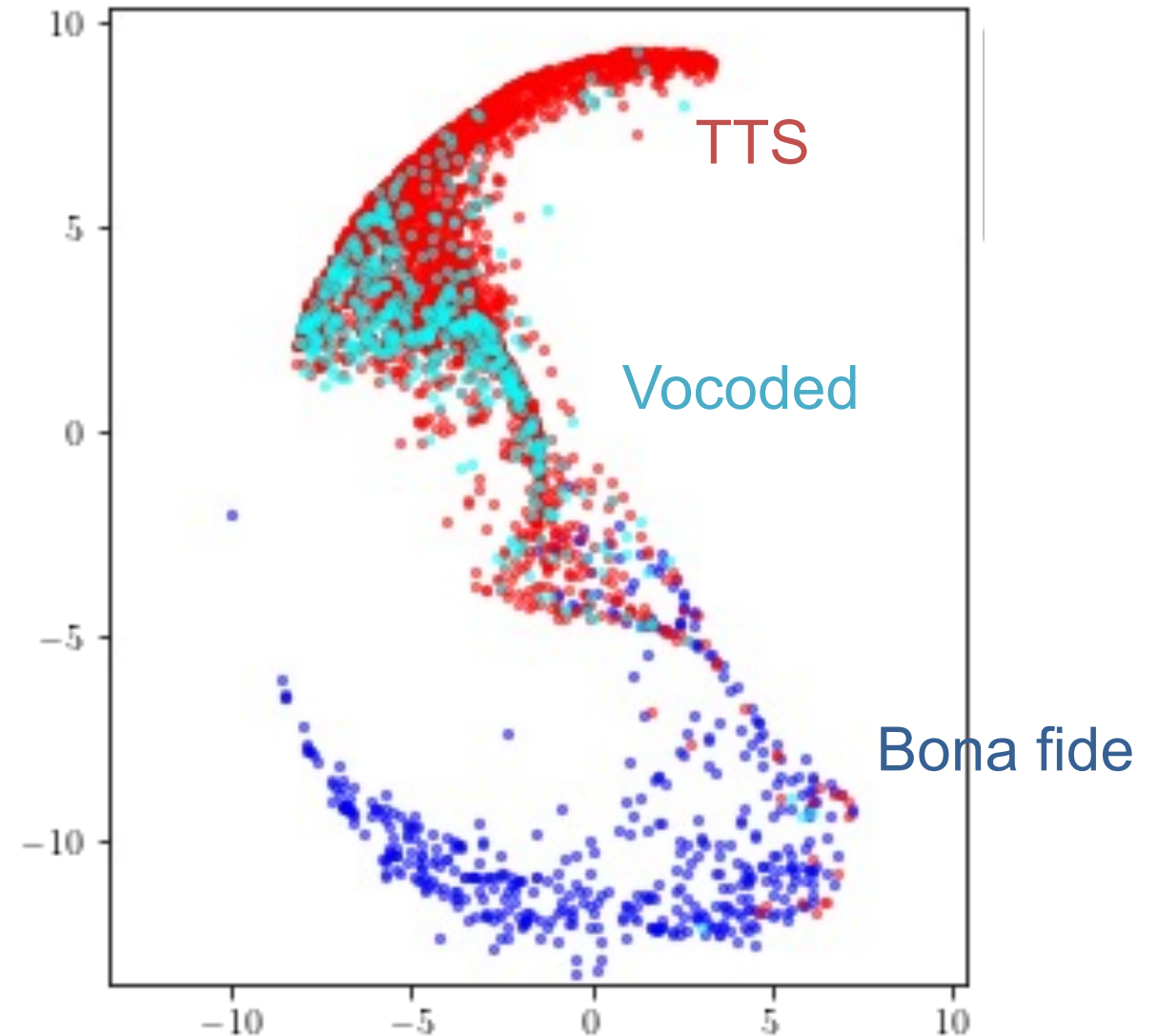
# Is vocoded data closer to bona fide data?



# Is vocoded data closer to bona fide data?

Yes

- bonafide: LJspeech data
- vocoded: full-band MelGAN
- TTS: full-band MelGAN + FastSpeech2 / Tacotron2



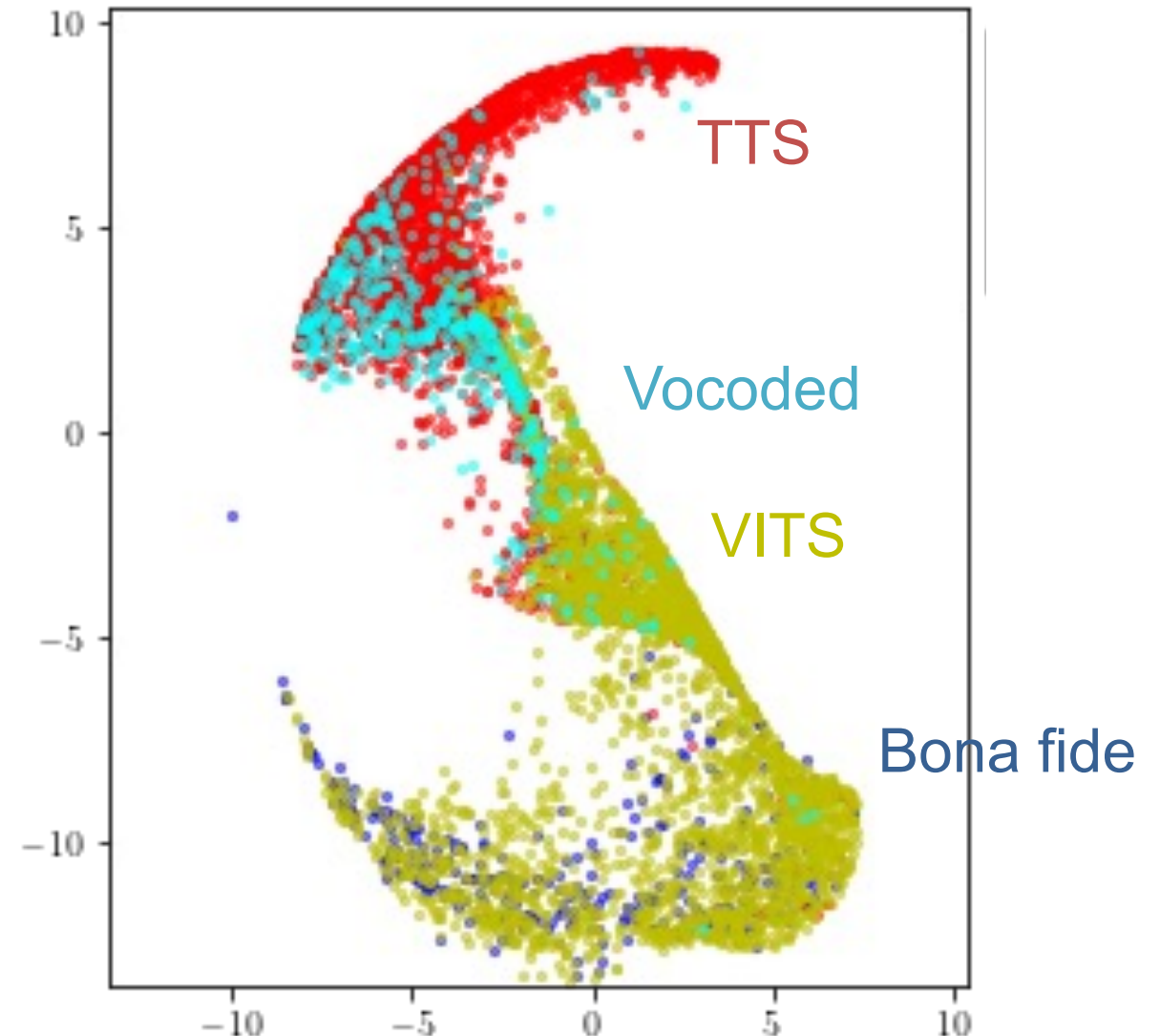
# Is vocoded data closer to bona fide data?

## Yes

- bonafide: LJspeech data
- vocoded: full-band MelGAN
- TTS: full-band MelGAN + FastSpeech2 / Tacotron2

## but not always

- end-to-end TTS: VITS



# Limitations

## Generalization?

- ✓ DSP-based vocoders

- ? Neural AR vocoders

## Generalization to end-to-end TTS?

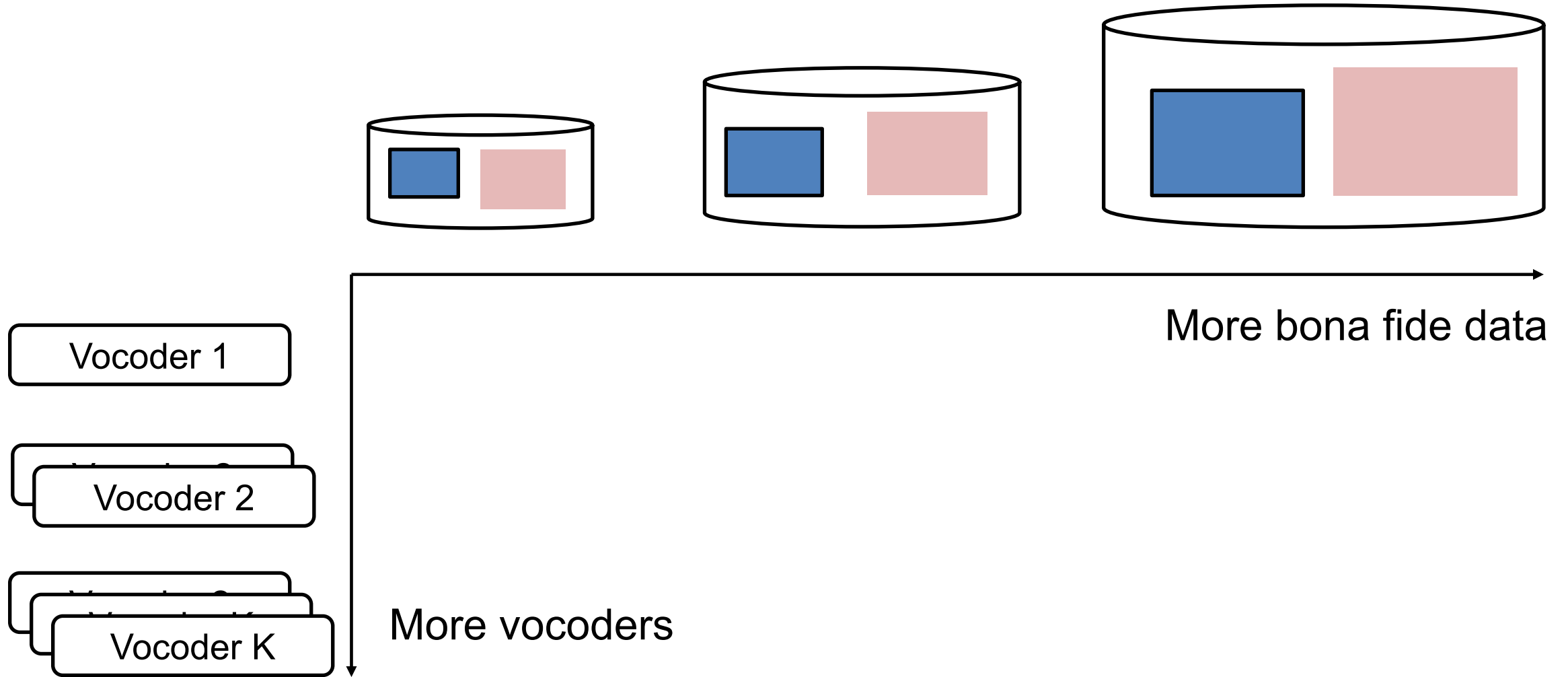
- X VITS

# Question 3

- Benefit of large-scale vocoded data?

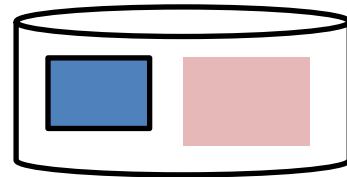


# Get much more data



# Get much more data

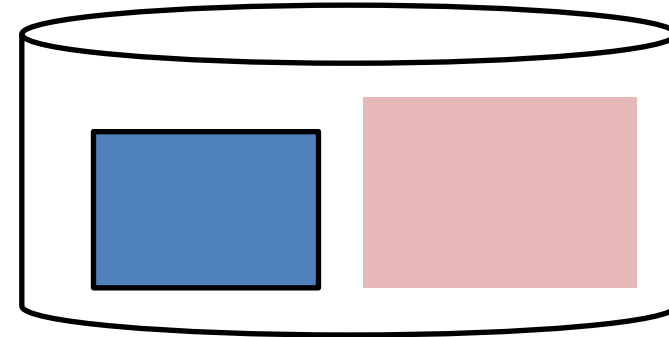
## ☐ Vocoding VoxCeleb2



ASVspoof 2019 LA

Bona fide: 2.42 hours

Vocoded: 2.42 x 4



VoxCeleb2 dev

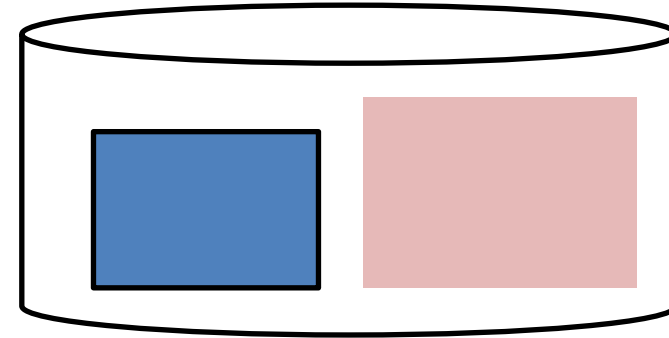
2360 hours

2360 x 4

# Get much more data using VoxCeleb2

## ☐ Usage 1: train CM as usual

- Practical issue – data size is too large
  - random sampling data to train CM



VoxCeleb2 dev

# Get much more data using VoxCeleb2

the best in Experiment II

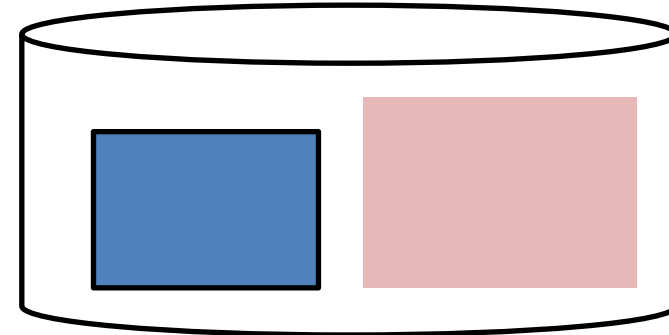


Training data size	LA19 trn.	Voc.v4	Vocoded VoxCeleb2 dev					
	-	×1	×0.3	×0.6	×1.2	×2.4	×6.0	×12.0
LA19eval	0.21	2.21	8.20	7.04	5.40	6.53	5.40	5.63
LA21eval	3.30	17.90	20.19	16.73	14.33	18.10	17.44	17.84
DF21eval	4.12	5.04	7.49	5.41	5.39	6.00	5.65	5.64
LA19etrim	9.00	3.79	6.17	5.53	5.16	5.25	5.22	5.14
LA21hid	26.98	14.57	13.98	12.34	11.47	11.64	11.37	11.43
DF21hid	16.85	7.78	11.02	9.71	9.90	10.05	9.99	10.04
WaveFake	24.62	2.50	14.94	10.39	8.38	5.52	4.88	4.94
InWild	17.07	7.55	16.12	15.63	14.19	13.43	13.32	13.77
Pooled	13.68	11.27	13.52	11.79	9.98	9.01	8.36	8.37

# Get much more data using VoxCeleb2

## ☐ Usage 1: train CM as usual

- Practical issue – data size is too large
  - random sampling data to train CM



VoxCeleb2 dev

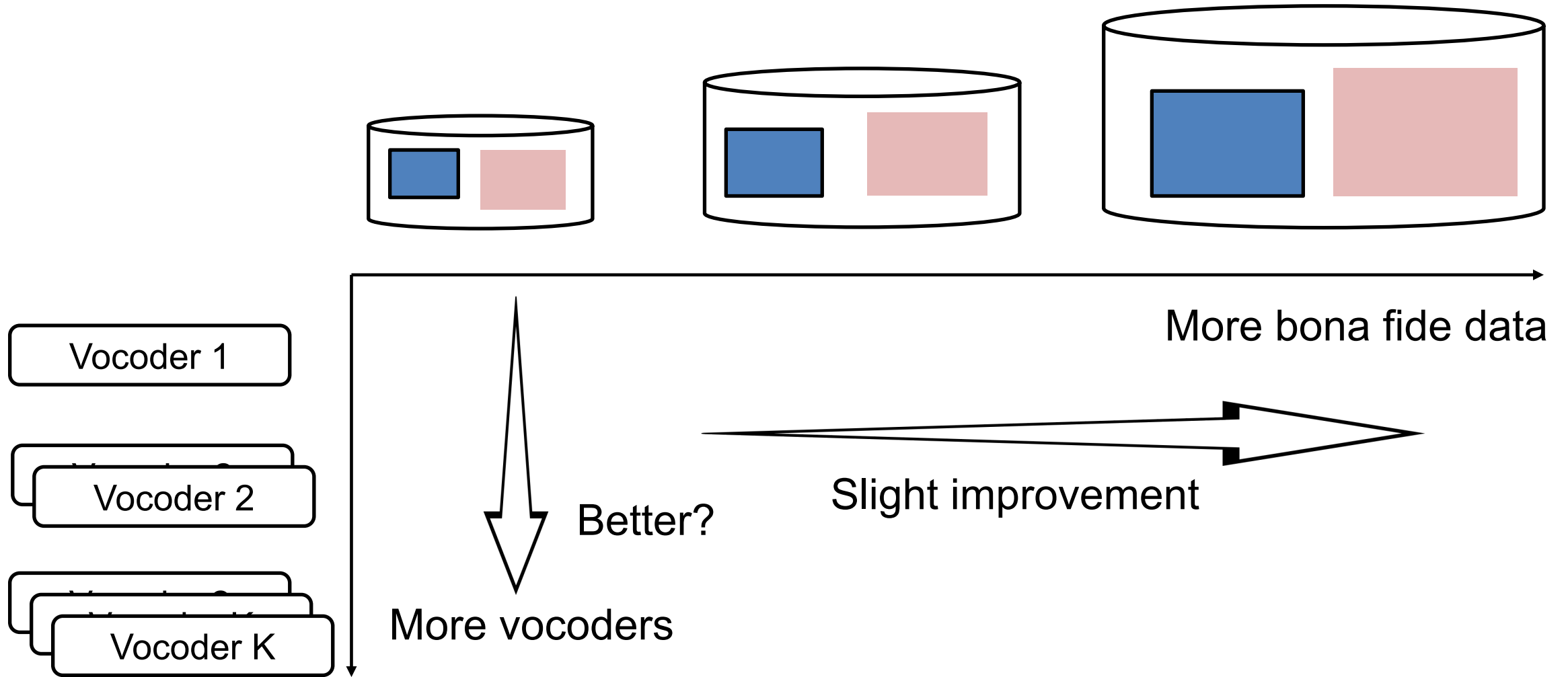
## ☐ Usage 2: train CM feature extractor

- wav2vec2.0 ..
- ...

Limited improvement

<https://arxiv.org/abs/2309.06014>

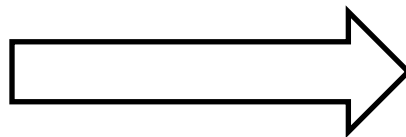
# Get much more data



# Summary



Vocoding



Contrastive learning

