

Text-to-Speech Synthesis Techniques for MIDI-to-Audio Synthesis

Erica Cooper, Xin Wang  , Junichi Yamagishi
National Institute of Informatics, Japan

PLATINUM SPONSOR



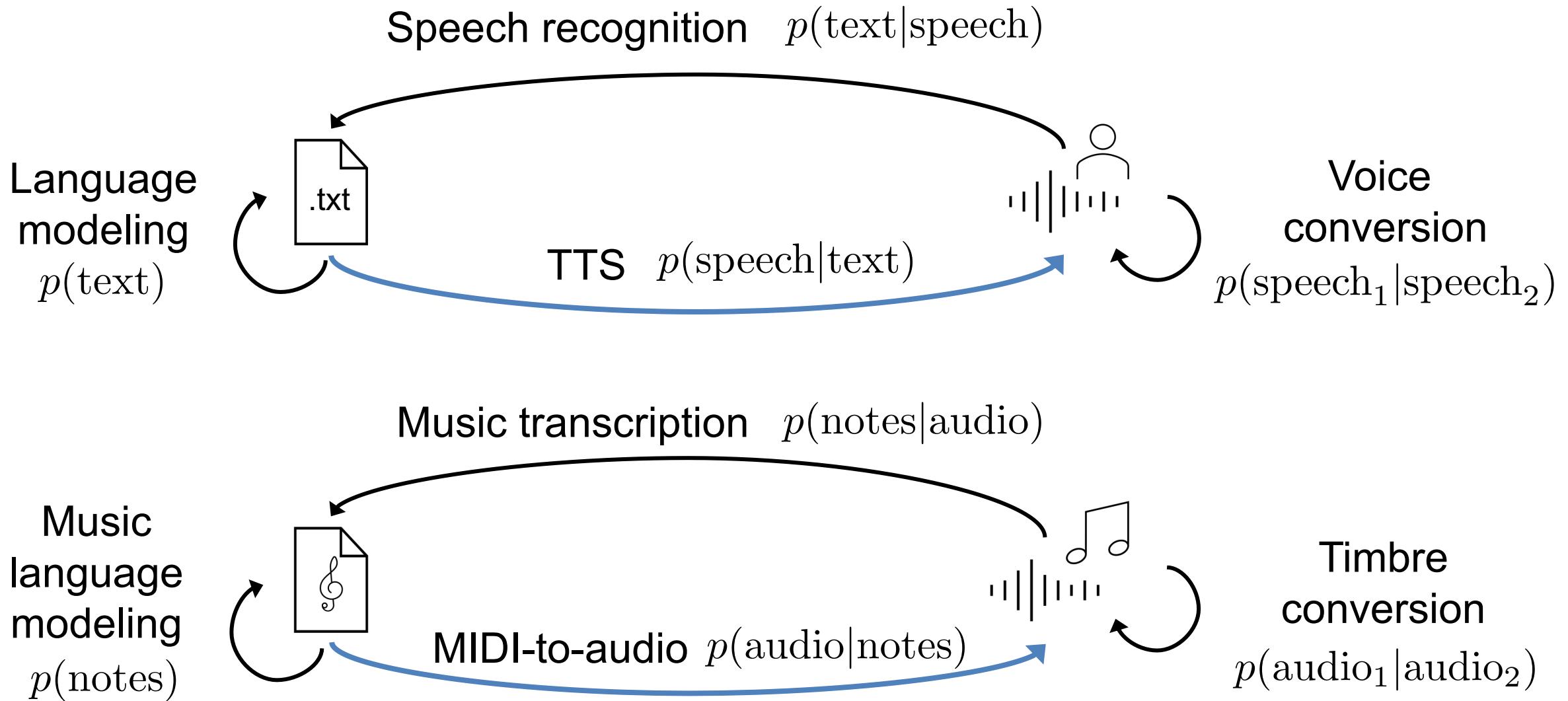
GOLD SPONSORS



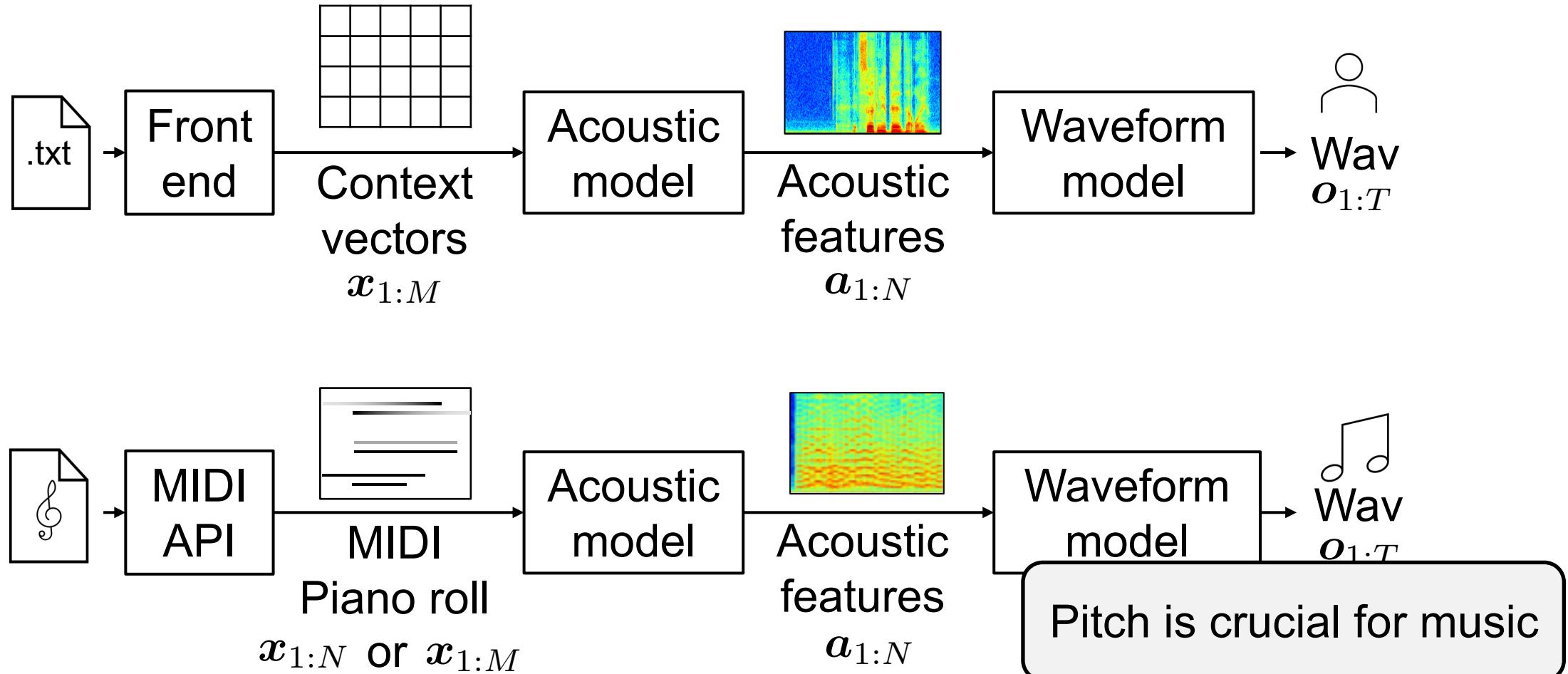
PARTNERS



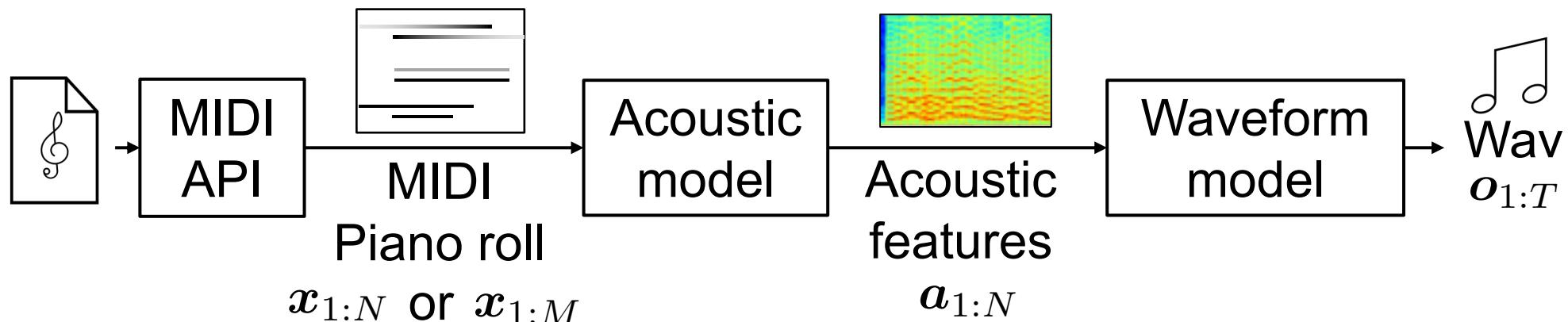
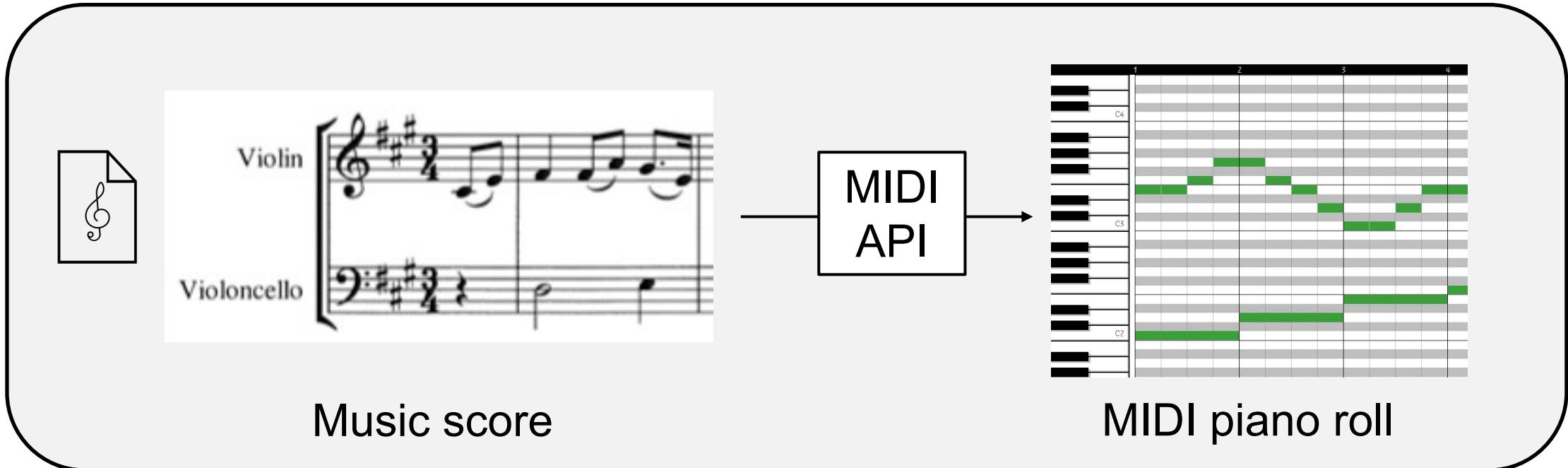
Motivation



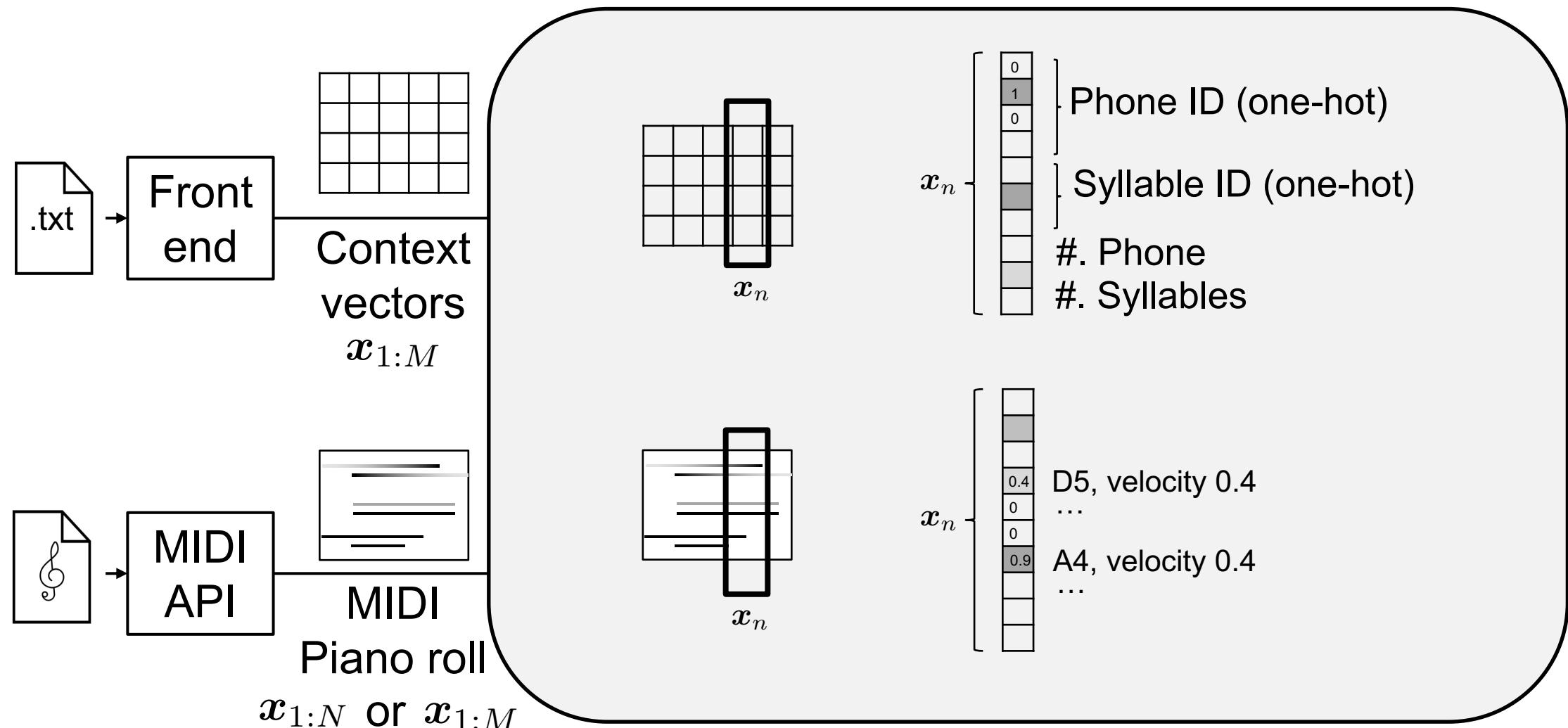
TTS & MIDI-to-Audio Synthesis



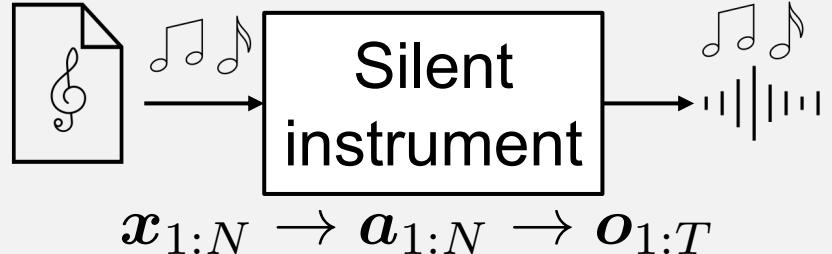
TTS & MIDI-to-Audio Synthesis



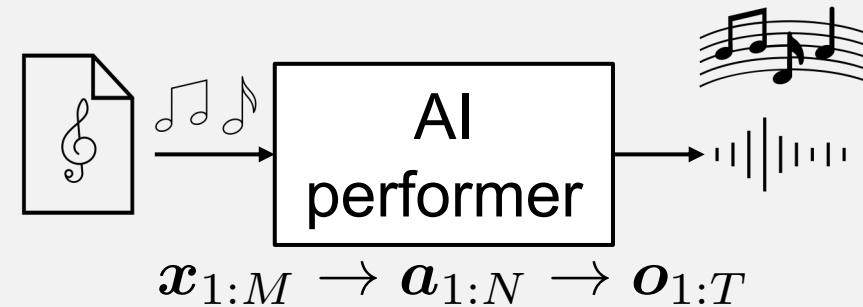
TTS & MIDI-to-Audio Synthesis



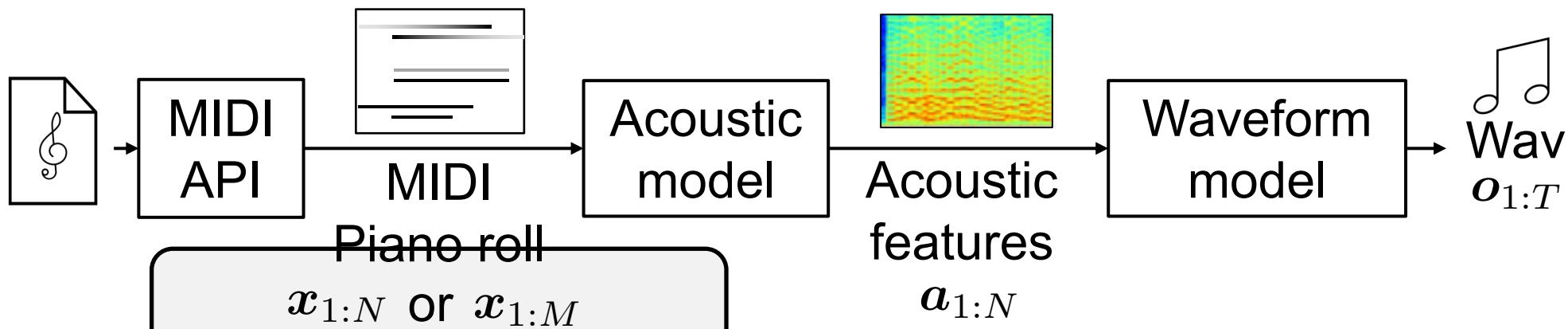
TTS & MIDI-to-Audio Synthesis



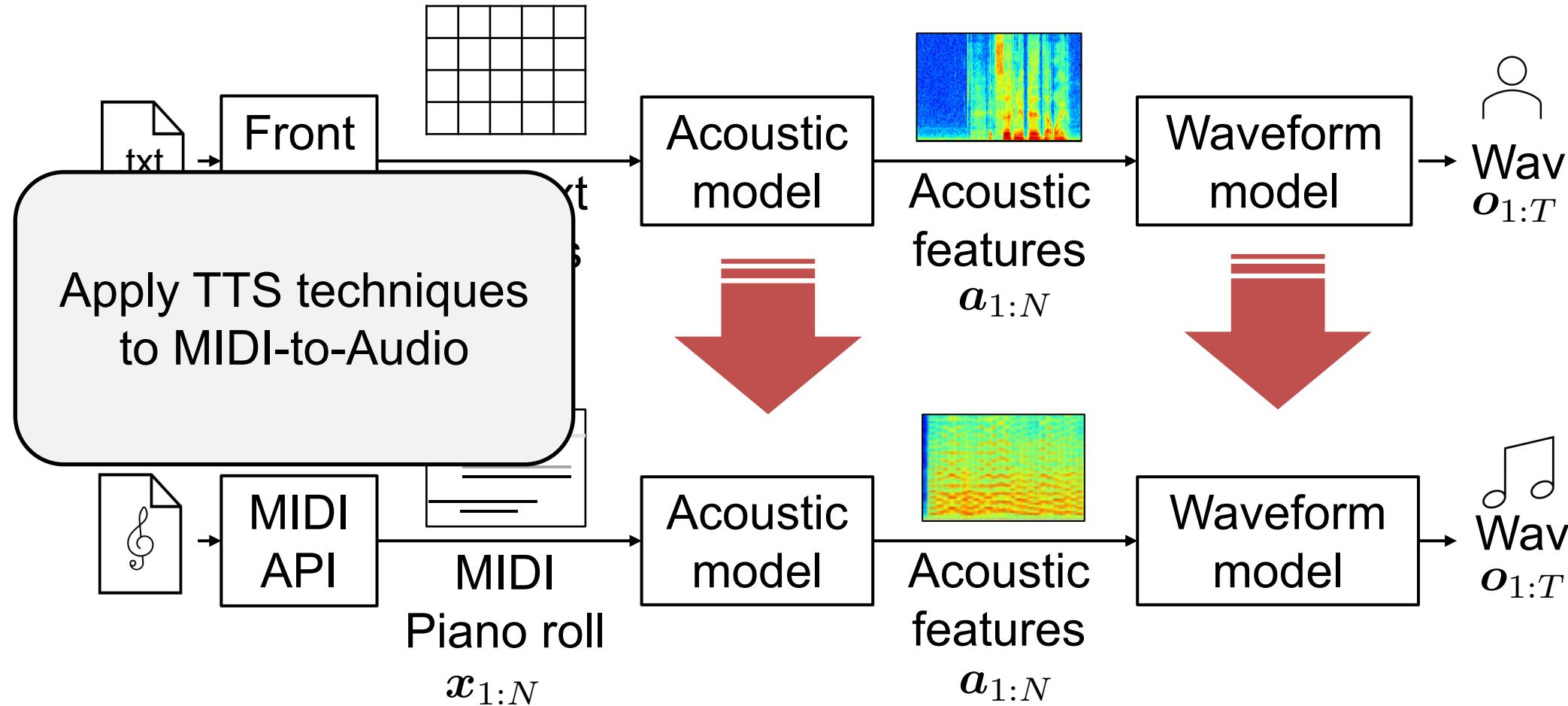
Aligned



Not completely aligned

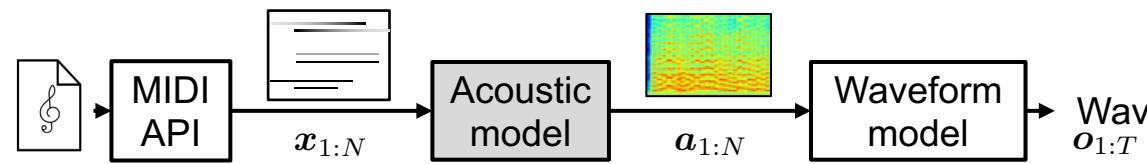


TTS & MIDI-to-Audio Synthesis

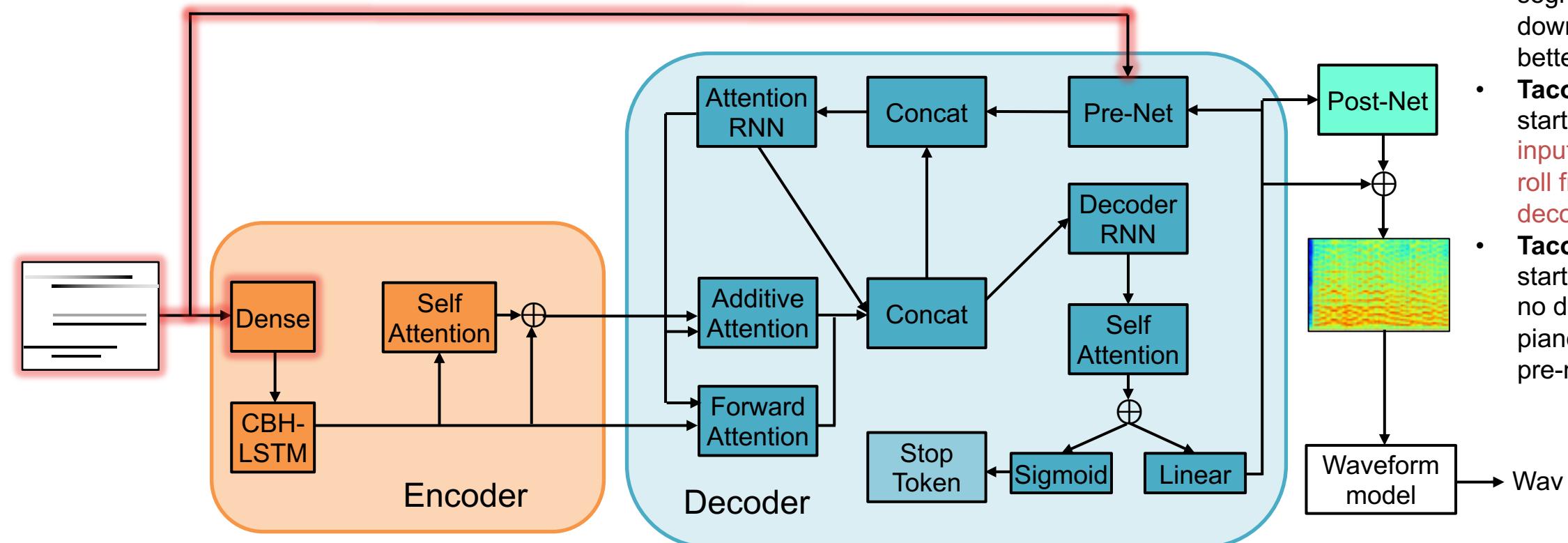


Methods

□ Acoustic model



1. TTS model: Tacotron (Wang 2017, Yasuda 2019)

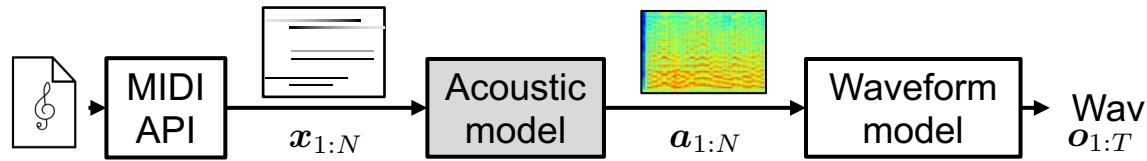


Wang, Y. et al. Tacotron: Towards End-to-End Speech Synthesis. in *Proc. Interspeech* 4006–4010 (2017).

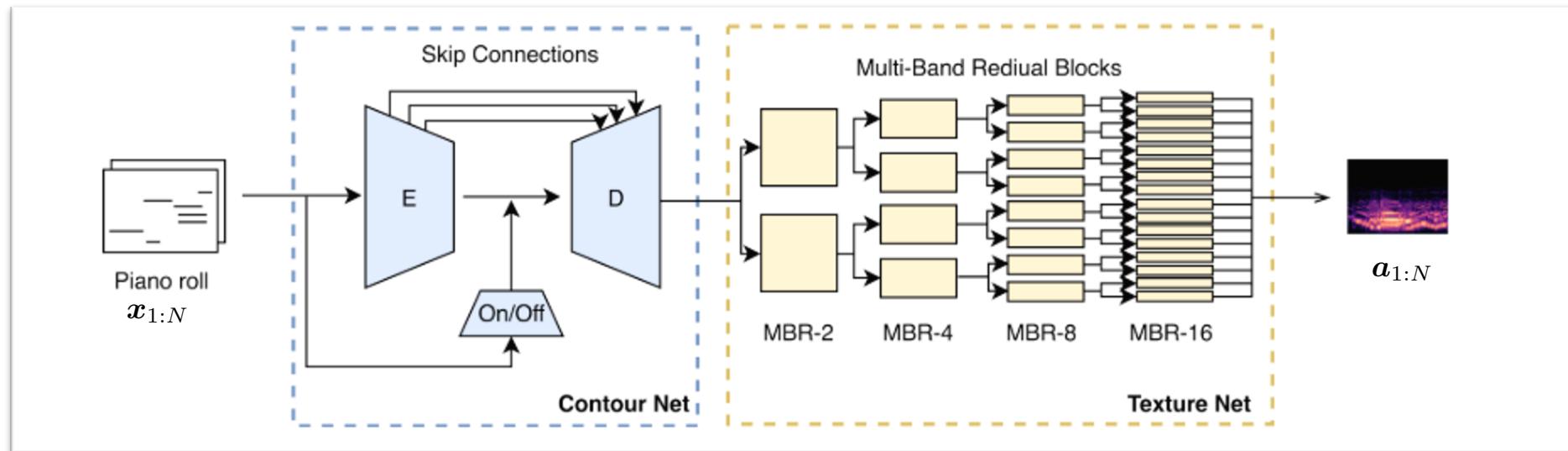
Yasuda, Y., Wang, X., Takaki, S. & Yamagishi, J. Investigation of enhanced Tacotron text-to-speech synthesis systems with self-attention for pitch accent language. in *Proc. ICASSP* 6905–6909 (2019).

Methods

□ Acoustic model



2. Reference model from music field: PerformanceNet (Wang 2019)

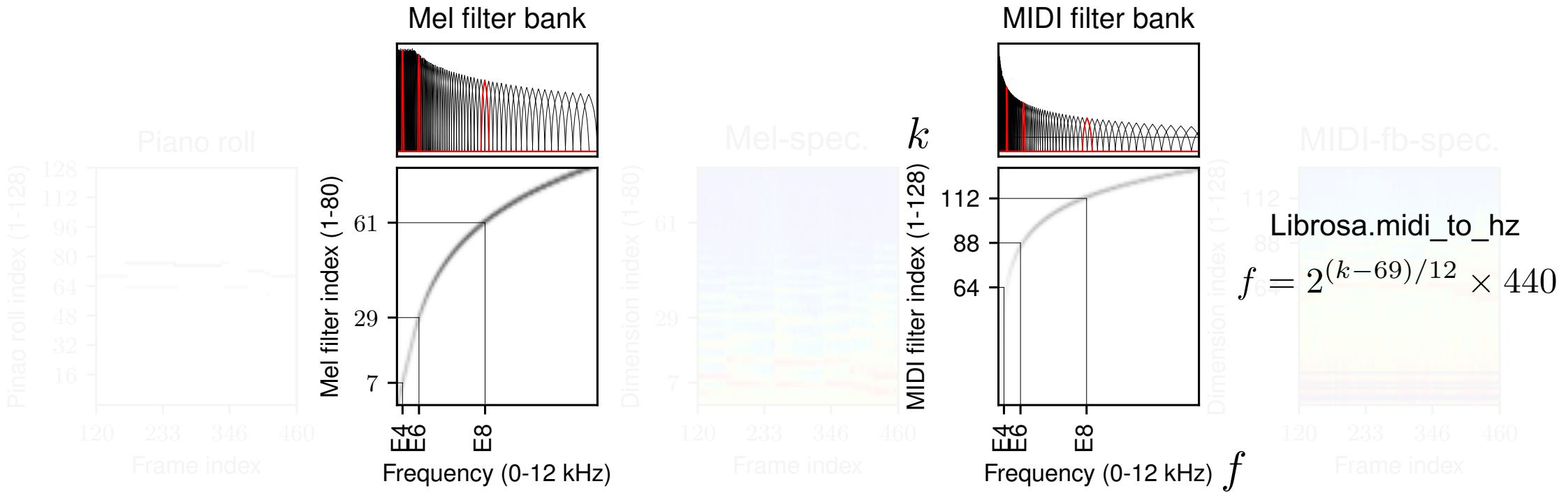
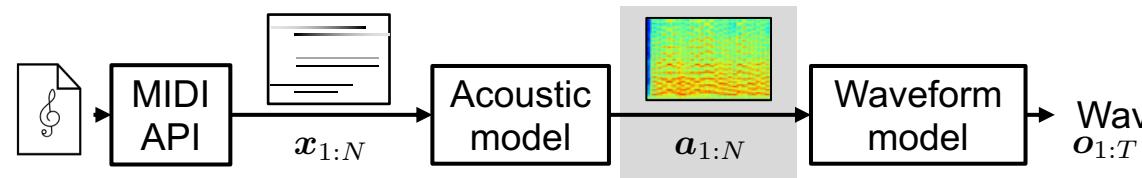


Adopted from figure 3 in (Wang 2019)

Methods

□ Acoustic features

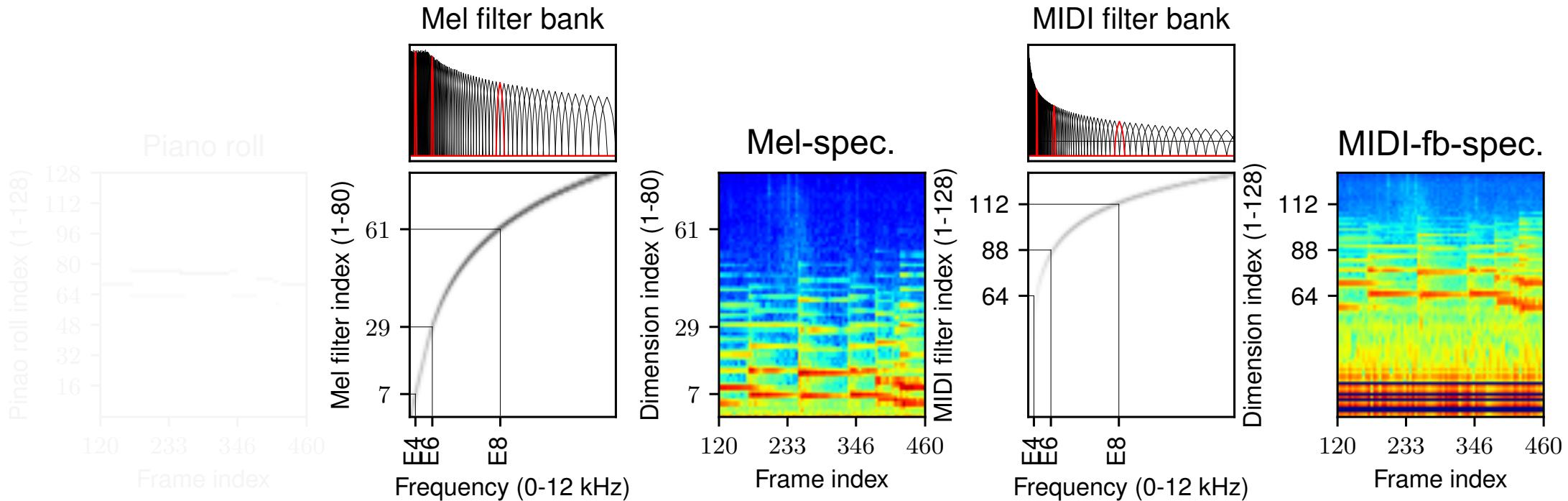
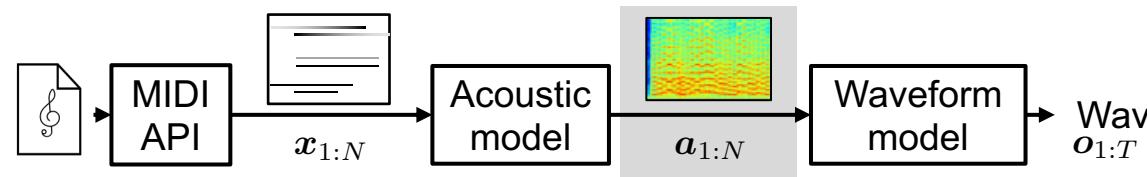
1. Mel-spectrogram
2. MIDI-filterbank-spectrogram



Methods

□ Acoustic features

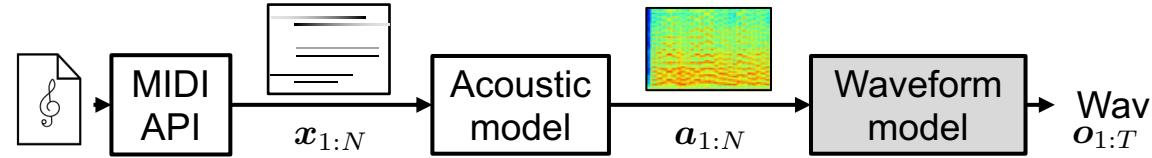
1. Mel-spectrogram
2. MIDI-filterbank-spectrogram



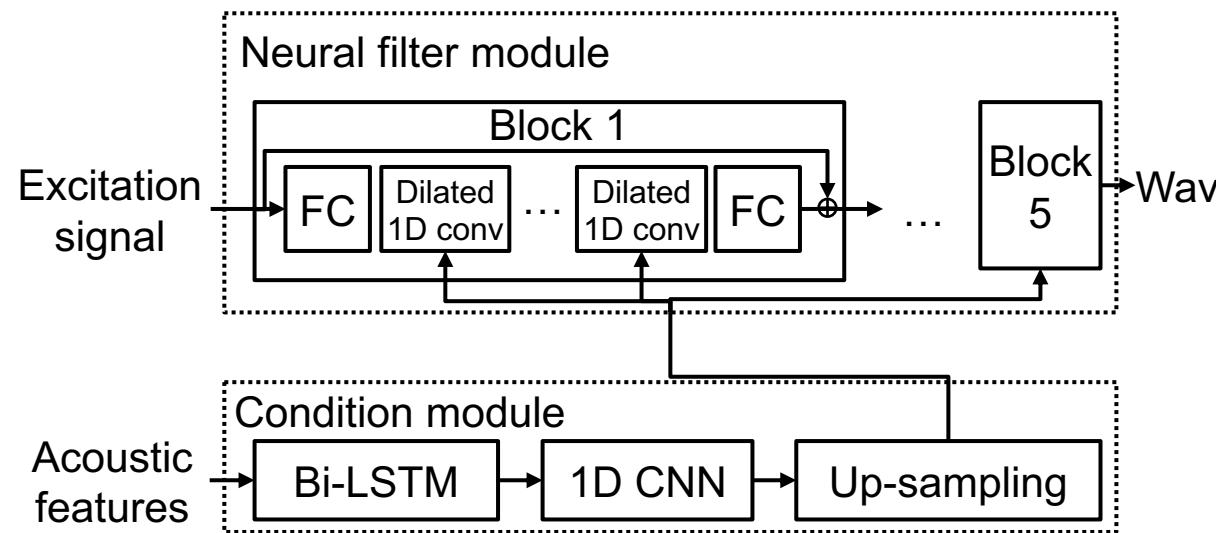
Blank lines in MIDI-fb-spec. due to frequency resolution of FFT (see appendix)

Methods

□ Waveform model



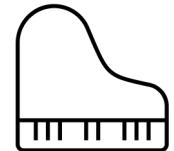
- Based on music neural source-filter (NSF) model (Zhao 2020) but
 - No harmonic-plus-noise structure



Experiments

□ Database: MAESTRO v2.0 (Hawthorne 2018)

- Real **piano** performances in International Piano-e-Competition
- MIDI was recorded simultaneously during performance
 - Aligned audio & piano roll



Data split of MAESTRO

- For experiments:
 - Follow official data split
 - 24kHz, 16 bits PCM

Split	Performances	Duration (hours)	Size (GB)	Notes (millions)
Train	967	161.3	97.7	5.73
Validation	137	19.4	11.8	0.64
Test	178	20.5	12.4	0.76
Total	1282	201.2	121.8	7.13

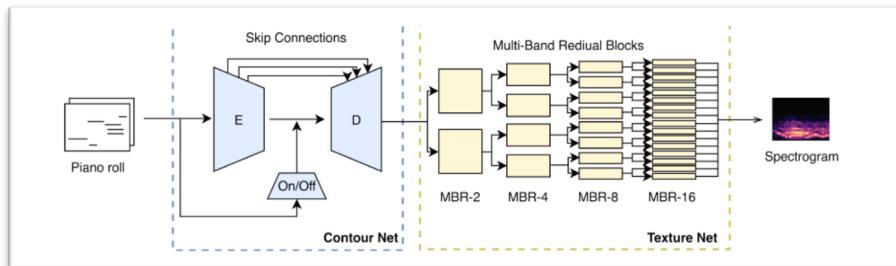
From <https://magenta.tensorflow.org/datasets/maestro#v200>

Experiments

□ Models in comparison



Original PerformanceNet

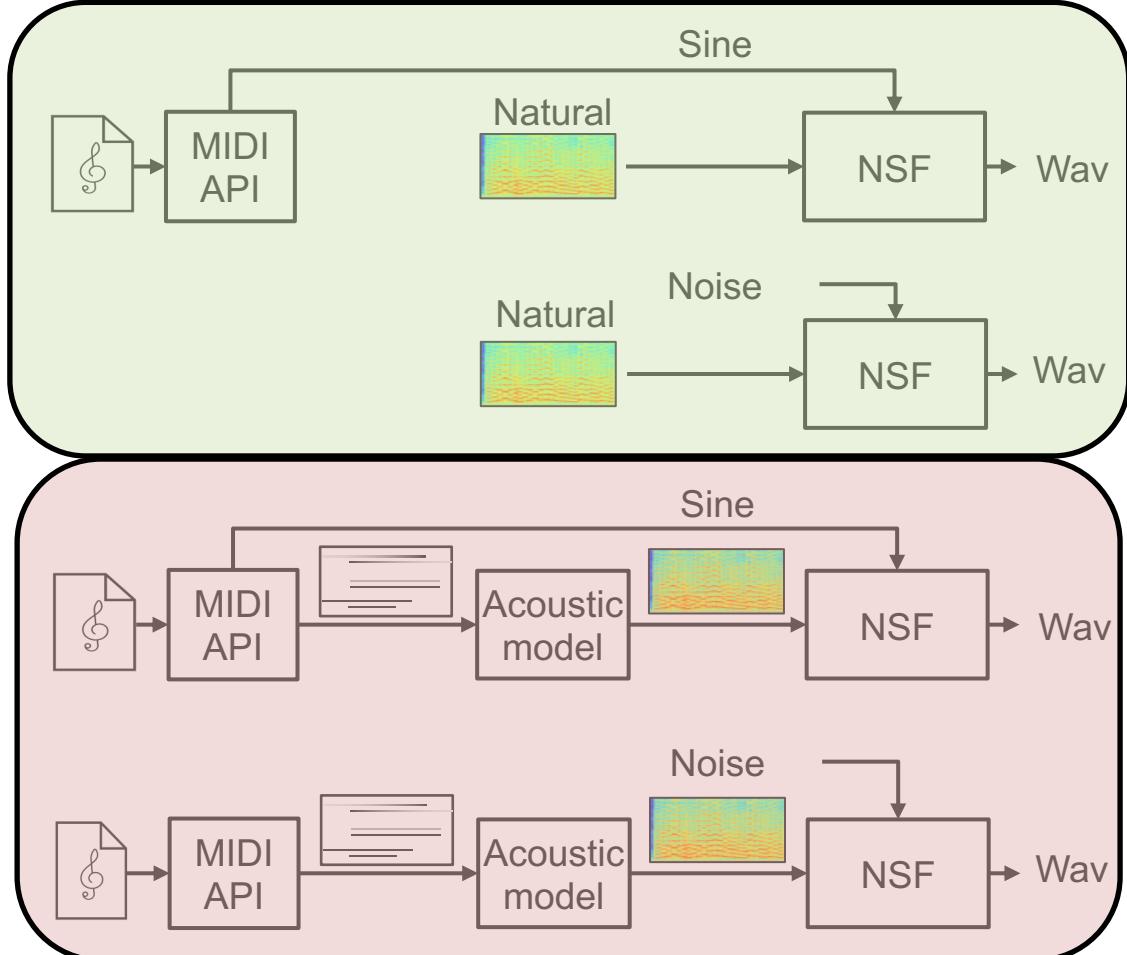


System ID	Acoustic model	Acoustic feature	Excit. signal	Wave. model	Pitch mismatch note	mismatch chord	MOS (mean)
Natural	-	-	-	-	-	-	4.04
Fluidsynth	Sample-based MIDI-to-audio software				5.20	6.77	3.66
Pianoteq	Physical-model MIDI-to-audio software				4.82	5.95	4.25
abs-mfb-sin	-	midi-fb	sine	NSF	-	-	3.87
abs-mfb-noi	-	midi-fb	noise	NSF	-	-	3.77
abs-mel-sin	-	mel-spc.	sine	NSF	-	-	2.72
abs-mel-noi	-	mel-spc.	noise	NSF	-	-	3.81
taco2-mfb-sin	taco2	midi-fb	sine	NSF	4.61	6.34	2.97
taco2-mfb-noi	taco2	midi-fb	noise	NSF	4.66	6.36	3.18
taco3-mfb-sin	taco3	midi-fb	sine	NSF	4.78	6.48	3.19
taco3-mfb-noi	taco3	midi-fb	noise	NSF	4.89	6.53	3.19
taco4-mfb-sin	taco4	midi-fb	sine	NSF	4.86	6.39	2.98
taco4-mfb-noi	taco4	midi-fb	noise	NSF	4.97	6.42	2.95
pfnet-mfb-sin	PFNet	midi-fb	sine	NSF	5.59	7.14	3.10
pfnet-mfb-noi	PFNet	midi-fb	noise	NSF	5.78	7.26	3.05
pfnet-mel-sin	PFNet	mel-spec.	sine	NSF	5.66	7.17	1.82
pfnet-mel-noi	PFNet	mel-spec.	noise	NSF	5.72	7.22	1.83
pfnet-spec-GL	PFNet	spec.	-	GL	5.43	6.98	1.62
midi-sin-nsf	-	-	sine	NSF	-	-	-
midi-noi-nsf	-	-	noise	NSF	-	-	-

Open-source code

Experiments

□ Models in comparison



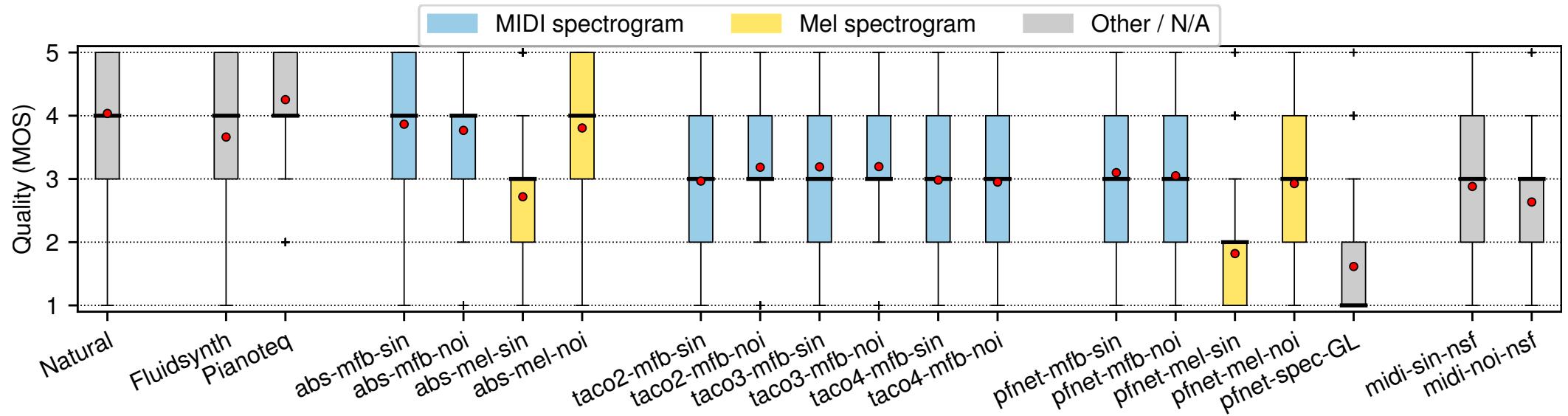
System ID	Acoustic model	Acoustic feature	Excit. signal	Wave. model	Pitch mismatch note	Pitch mismatch chord	MOS (mean)
Natural	-	-	-	-	-	-	4.04
Fluidsynth	Sample-based MIDI-to-audio software						
Pianoteq	Physical-model MIDI-to-audio software						
abs-mfb-sin	-	midi-fb	sine	NSF	-	-	3.87
abs-mfb-noi	-	midi-fb	noise	NSF	-	-	2.77
abs-mel-sin	-	mel-spc.	sine	NSF	-	-	2.72
abs-mel-noi	-	mel-spc.	noise	NSF	-	-	3.81
taco2-mfb-sin	taco2	midi-fb	sine	NSF	4.61	6.34	2.97
taco2-mfb-noi	taco2	midi-fb	noise	NSF	4.66	6.36	3.18
taco3-mfb-sin	taco3	midi-fb	sine	NSF	4.78	6.48	3.19
taco3-mfb-noi	taco3	midi-fb	noise	NSF	4.89	6.53	3.19
taco4-mfb-sin	taco4	midi-fb	sine	NSF	4.86	6.39	2.98
taco4-mfb-noi	taco4	midi-fb	noise	NSF	4.41	6.21	3.02
pfnet-mfb-sin	PFNet	midi-fb	sine	NSF	5.59	7.14	3.10
pfnet-mfb-noi	PFNet	midi-fb	noise	NSF	5.78	7.26	3.05
pfnet-mel-sin	PFNet	mel-spec.	sine	NSF	5.66	7.17	1.82
pfnet-mel-noi	PFNet	mel-spec.	noise	NSF	5.74	7.25	2.93
pfnet-spec-GL	PFNet	spec.	-	GL	5.43	6.98	1.62
midi-sin-nsf	-	-	sine	NSF	-	-	-
midi-noi-nsf	-	-	noise	NSF	-	-	-

Copy-synthesis

MIDI-to-audio

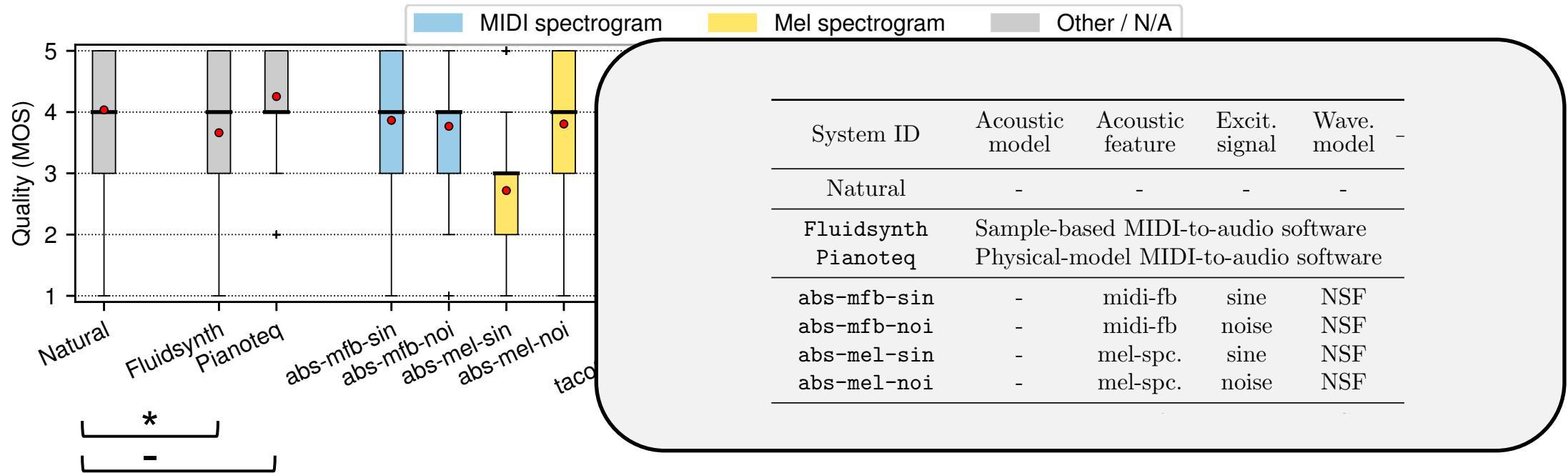
Experiments

- ☐ Listening test: crowd-sourcing, 224 amateur participants



Experiments

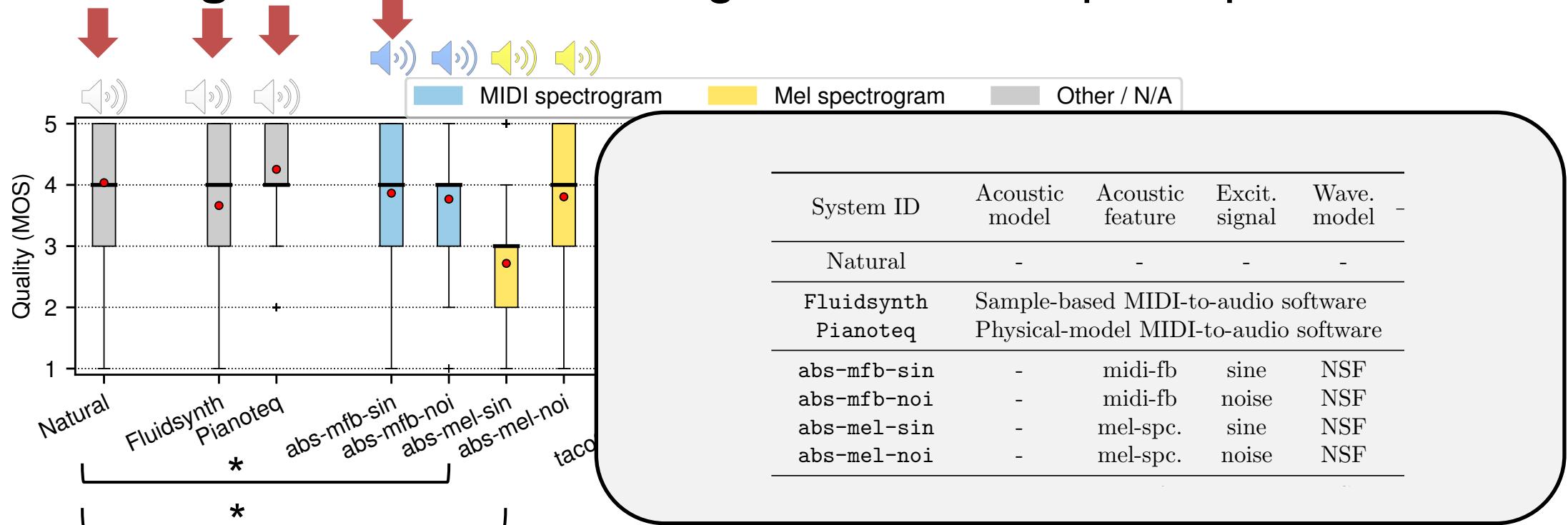
☐ Listening test: crowd-sourcing, 224 amateur participants



Mann-Whitney-U, Holm-Boferroni correction: “**” statistical significance at alpha=0.05, “-” otherwise

Experiments

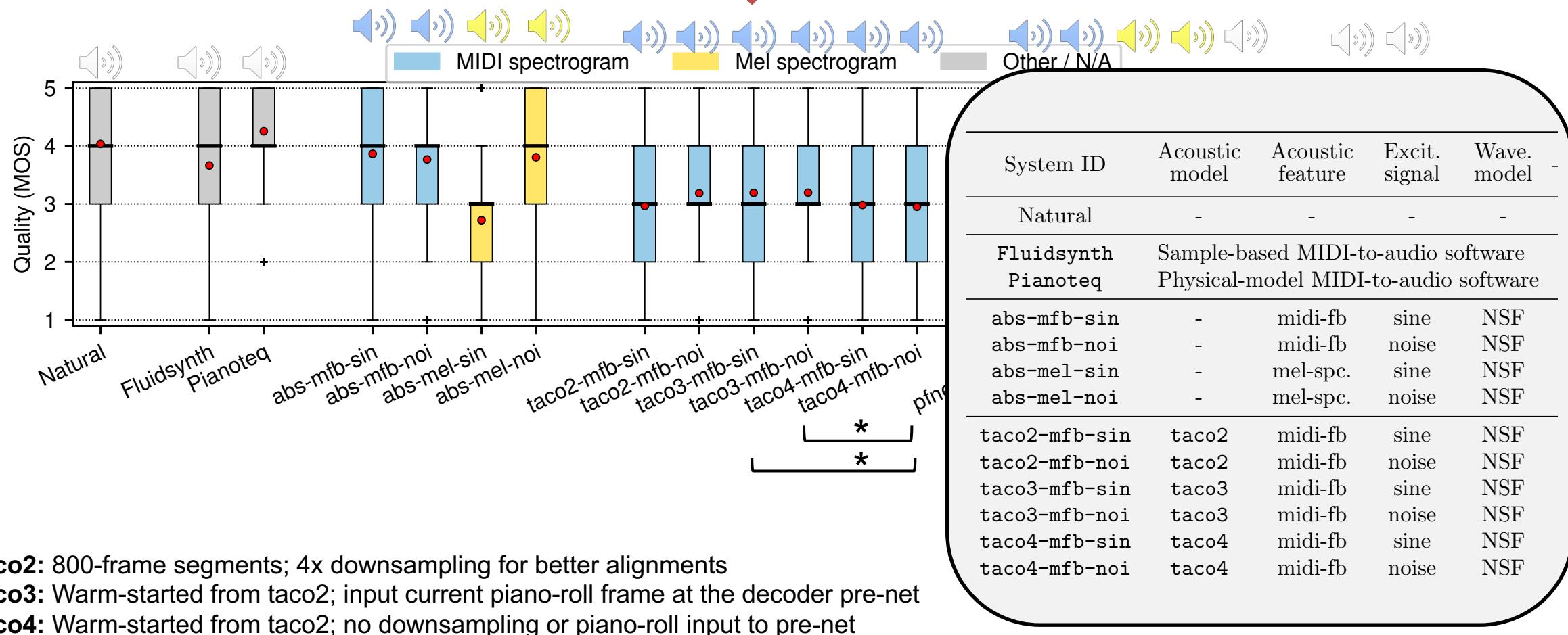
☐ Listening test: crowd-sourcing, 224 amateur participants



- Neural waveform model works well in copy-synthesis condition

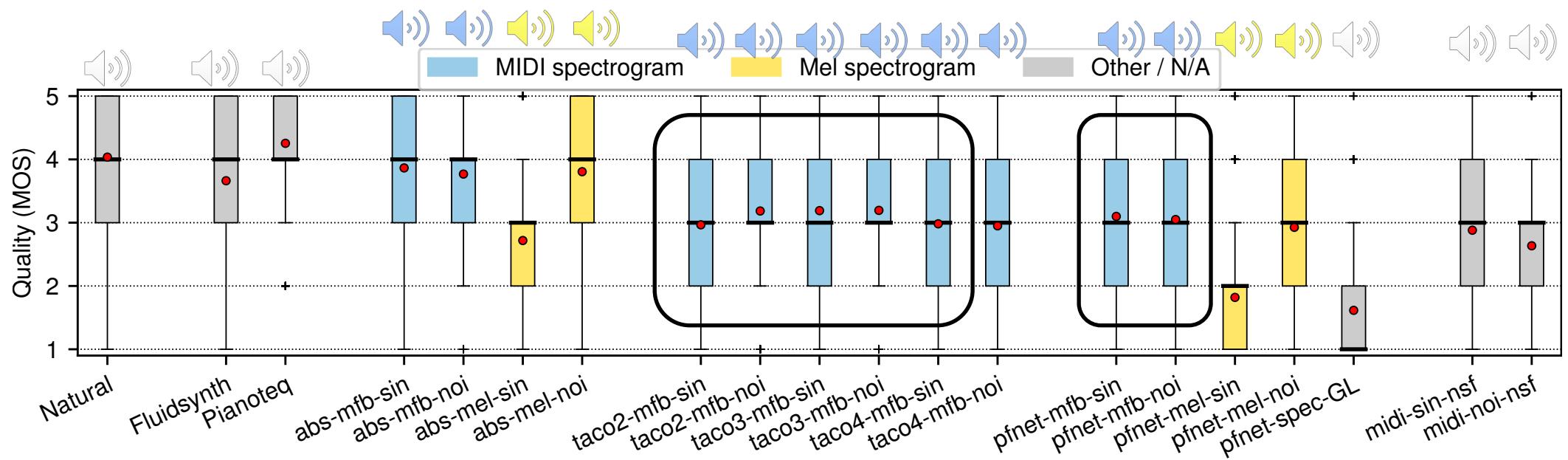
Experiments

☐ Listening test: crowd-sourcing, 224 amateur participants



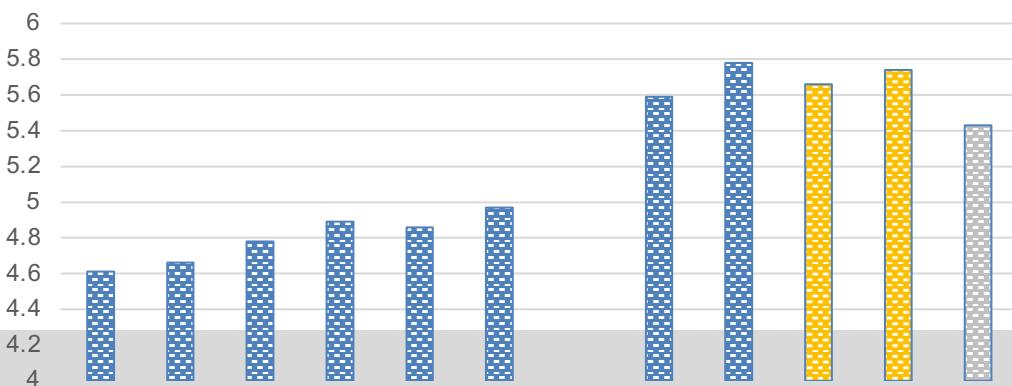
Experiments

☐ Listening test: crowd-sourcing, 224 amateur participants



☐ Pitch distortion

- Natural audio as target
- Evaluated on single notes
- The lower the better



Messages

□ TTS & MIDI-to-audio

- Techniques can be shared: acoustic model, waveform model
- Performance bottleneck on acoustic model (and waveform model)

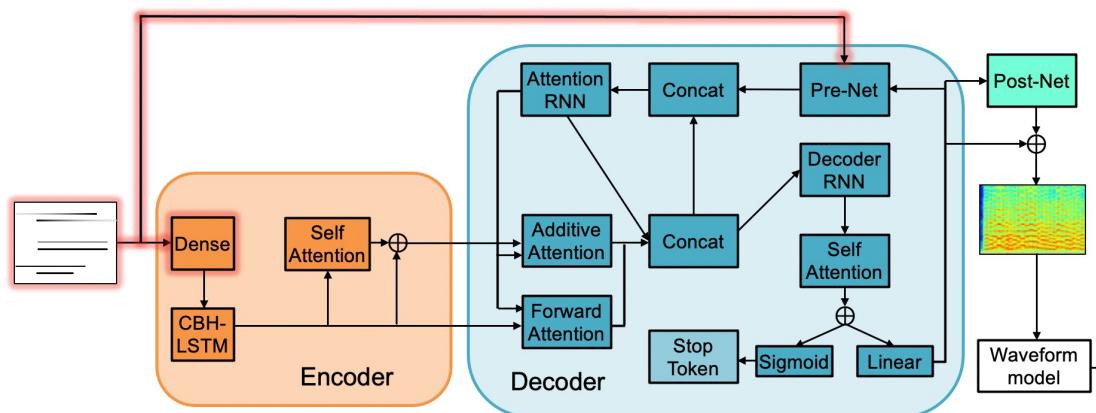
□ On waveform modeling

- Physical-model performs well but lacks reverberation effect
- Sample-based model replies on the sample database
- Non-AR waveform model is OK in copy-synthesis
 - Reverberation is captured
 - Noise excitation is OK

Messages

□ On acoustic model

- Obtaining good alignments for longer input sequences is challenging
- Inputting the piano-roll frame to the decoder prenet helps improve alignments
 - Acceptable for perfectly-aligned performance-MIDI
 - Have to consider other strategies for non-aligned score-MIDI



Thank you!

Tacotron code: <https://github.com/nii-yamagishilab/self-attention-tacotron>

NSF code: <https://github.com/nii-yamagishilab/project-NN-Pytorch-scripts>

Samples: <https://nii-yamagishilab.github.io/samples-xin/main-midi2audio.html>

PLATINUM SPONSOR



GOLD SPONSORS

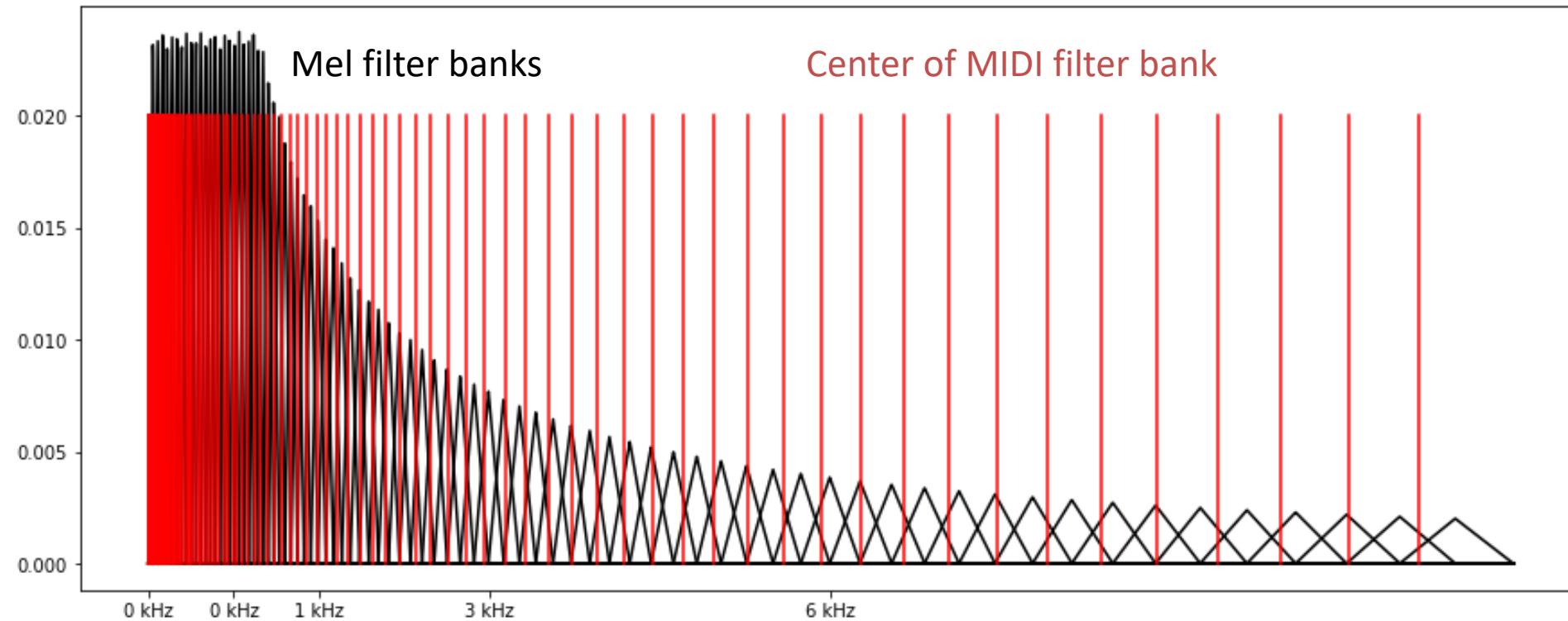


PARTNERS



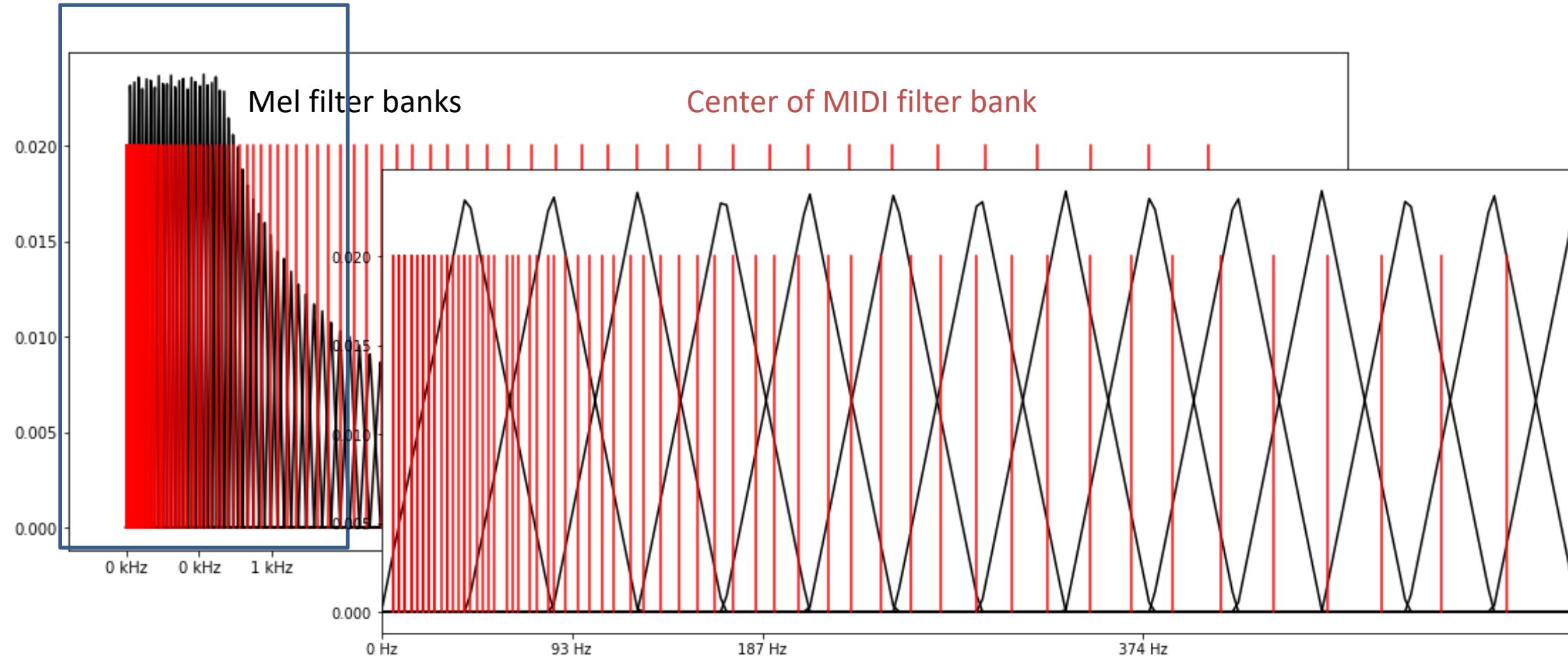
Appendix

□ On MIDI filterbank



Appendix

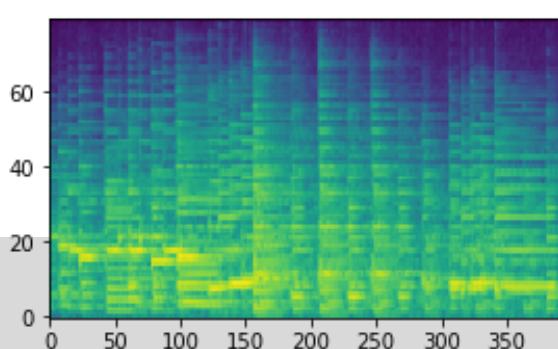
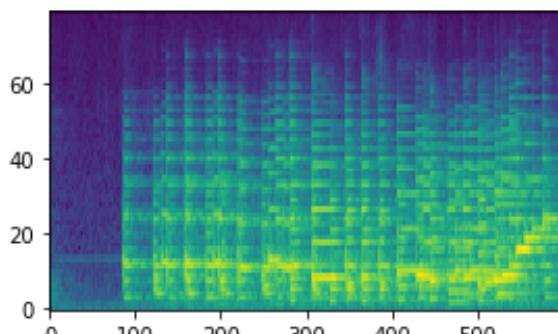
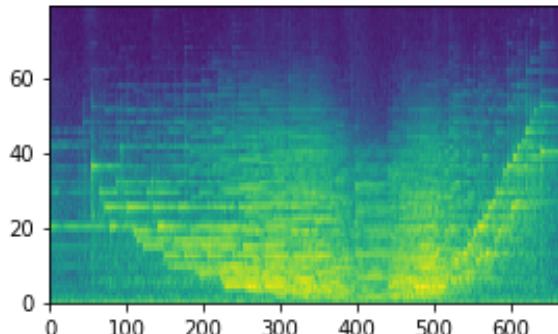
□ On MIDI filterbank



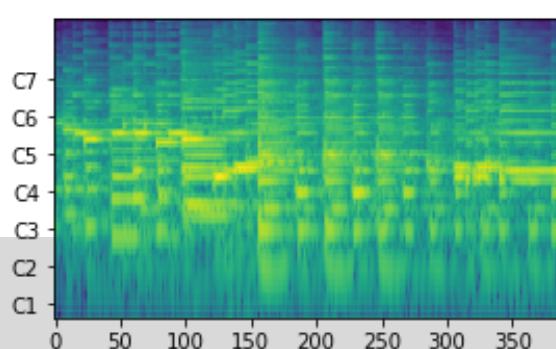
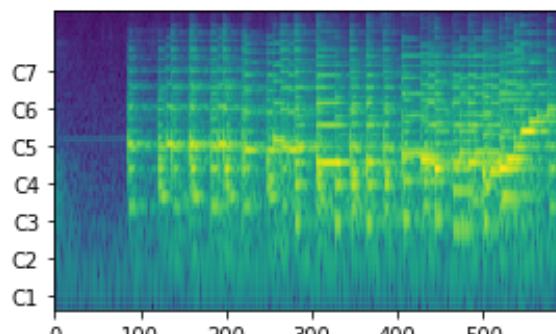
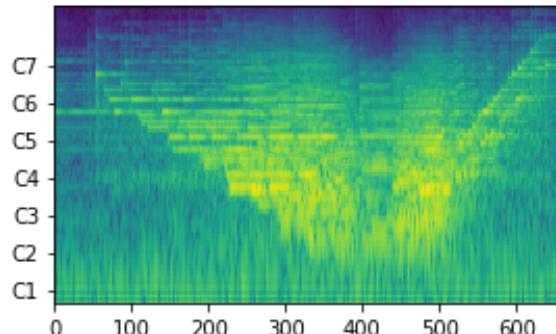
Appendix

□ On MIDI filterbank

Mel-spectrogram



MIDI-centered filter-bank



CQT

