

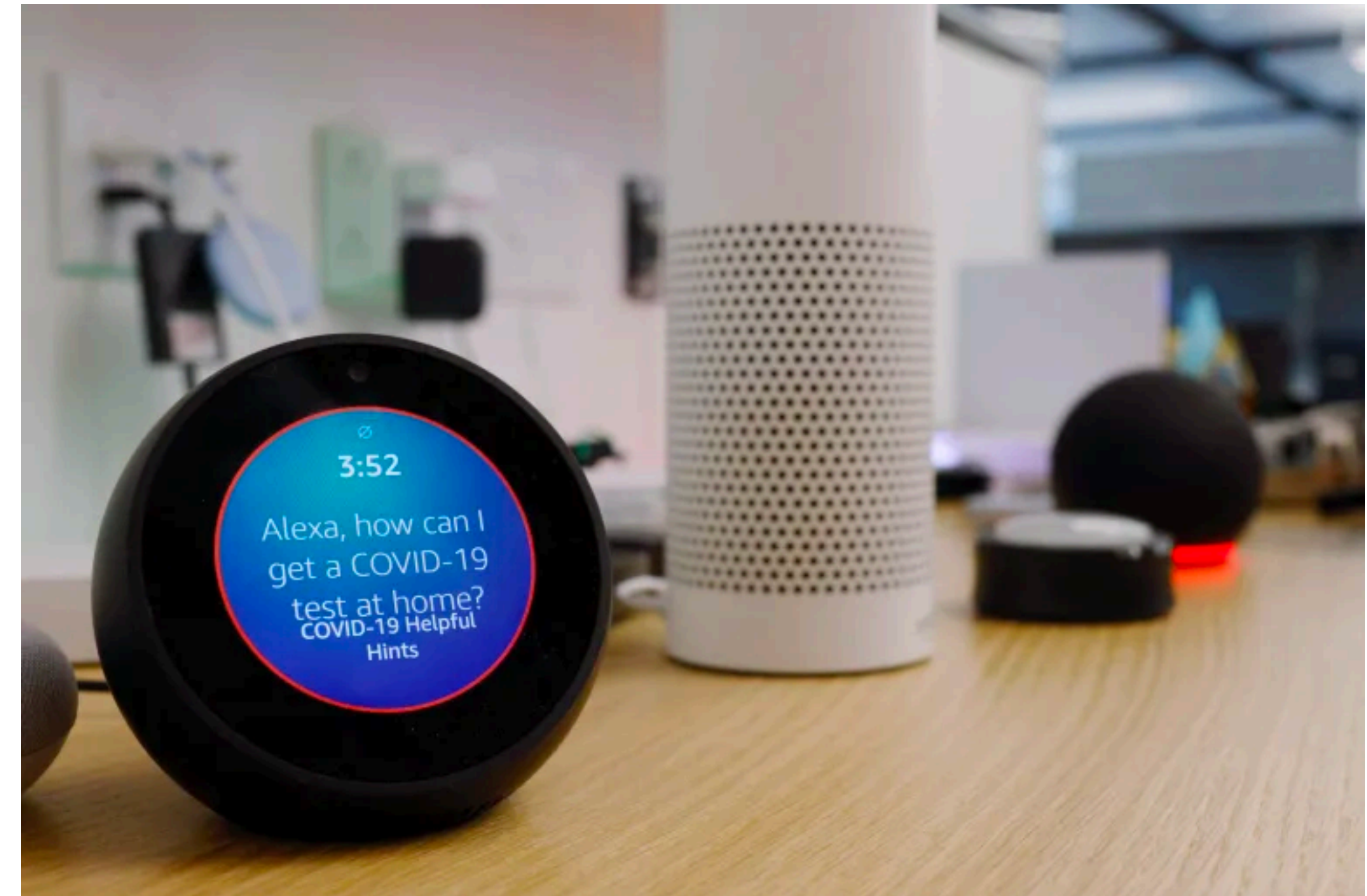
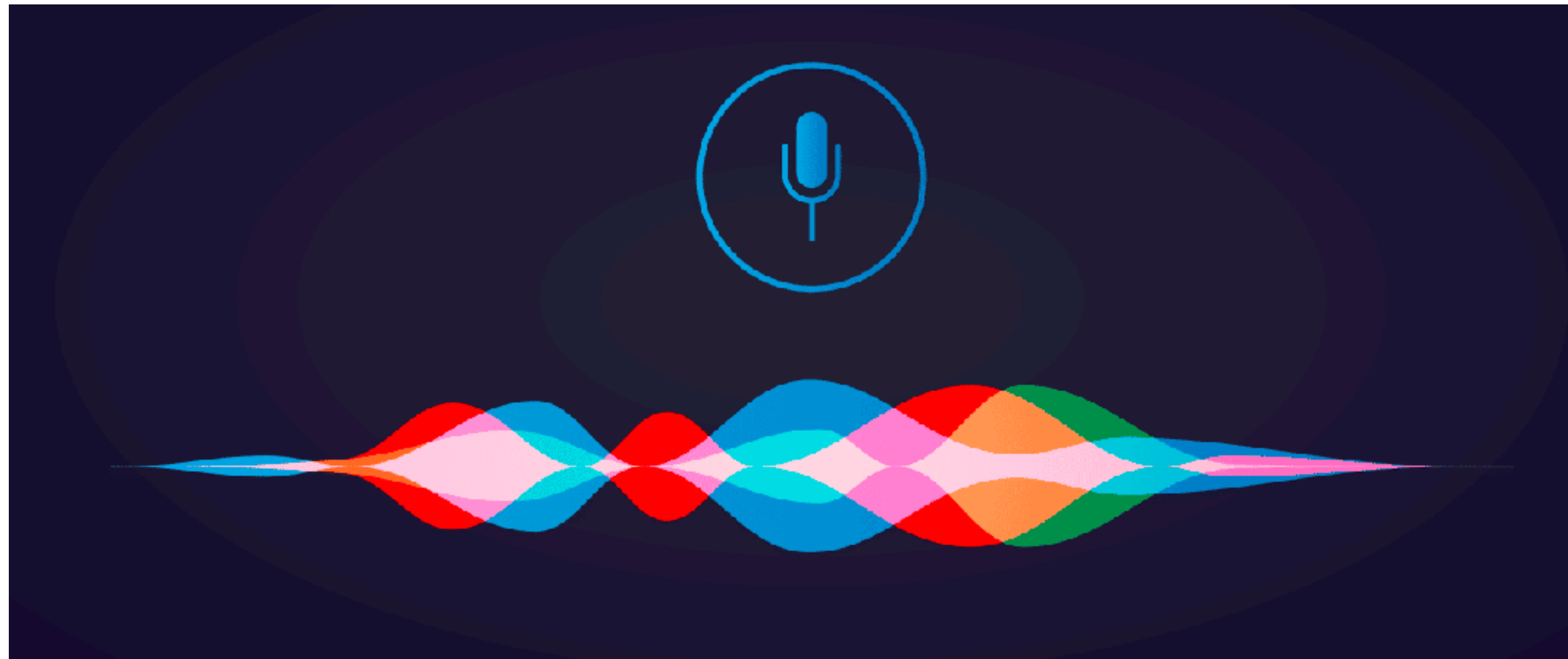
Lecture 18: Text-to-Speech Synthesis

Zhizheng Wu

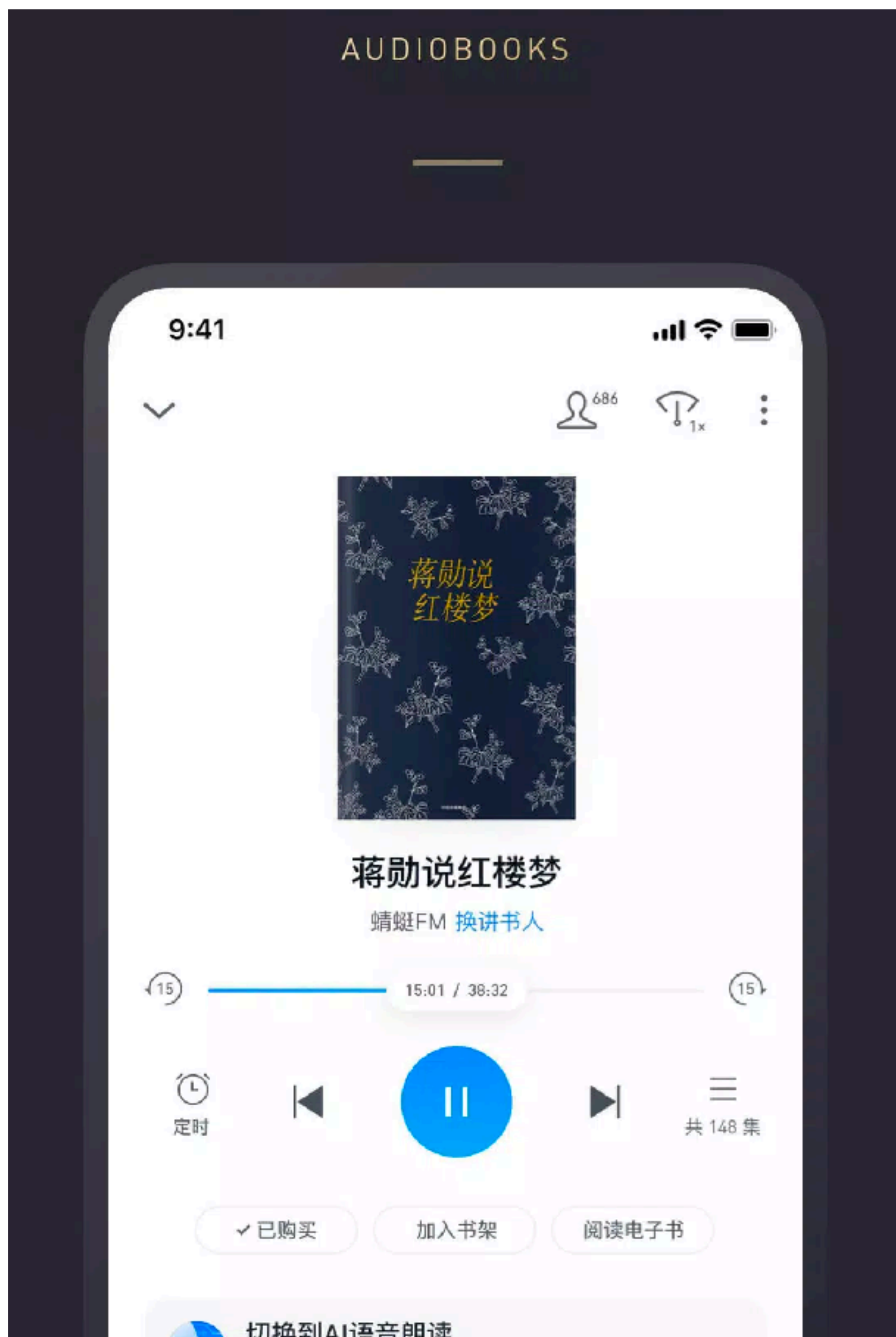
Agenda

- ▶ Applications
- ▶ Overview of text to speech
- ▶ Frontend
- ▶ Acoustic model
- ▶ Waveform generator
- ▶ Tools and readings

Applications



Applications



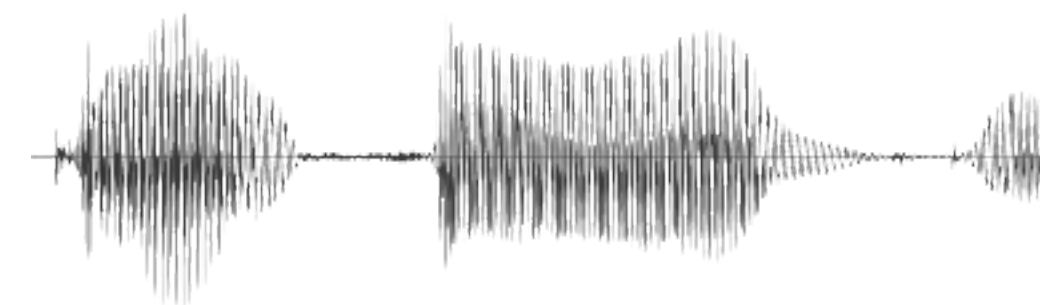
Applications



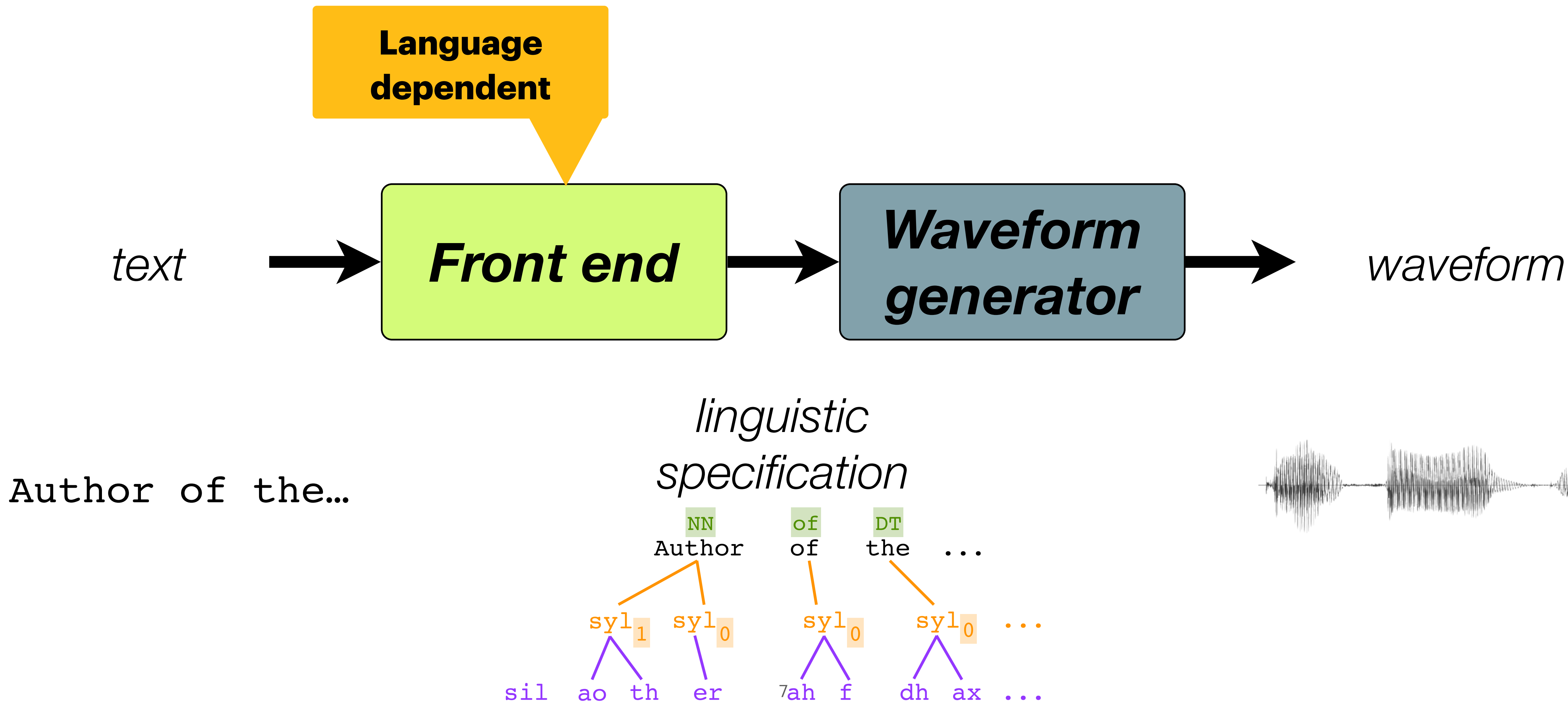
The end-to-end problem we want to solve



Author of the..



The two-stage pipeline



The three-stage pipeline



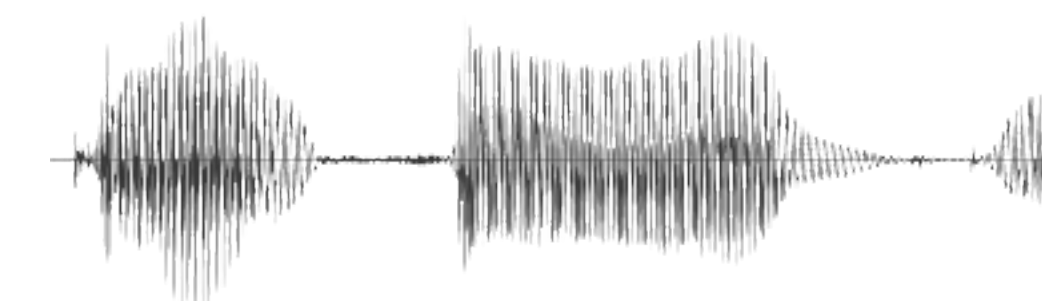
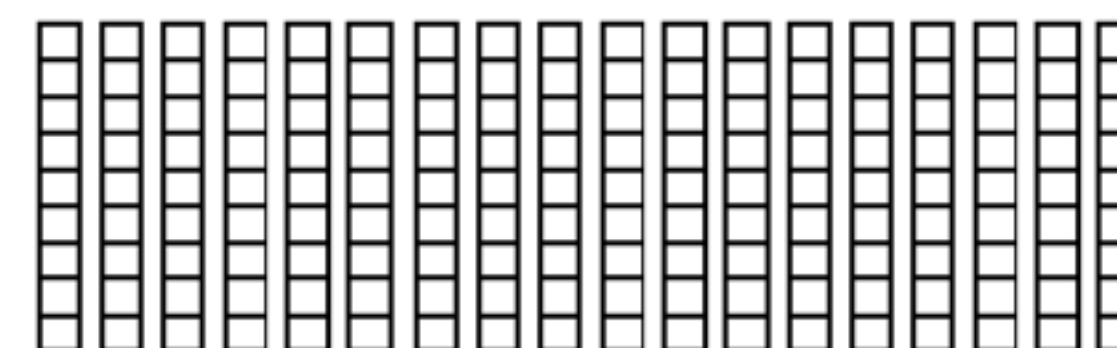
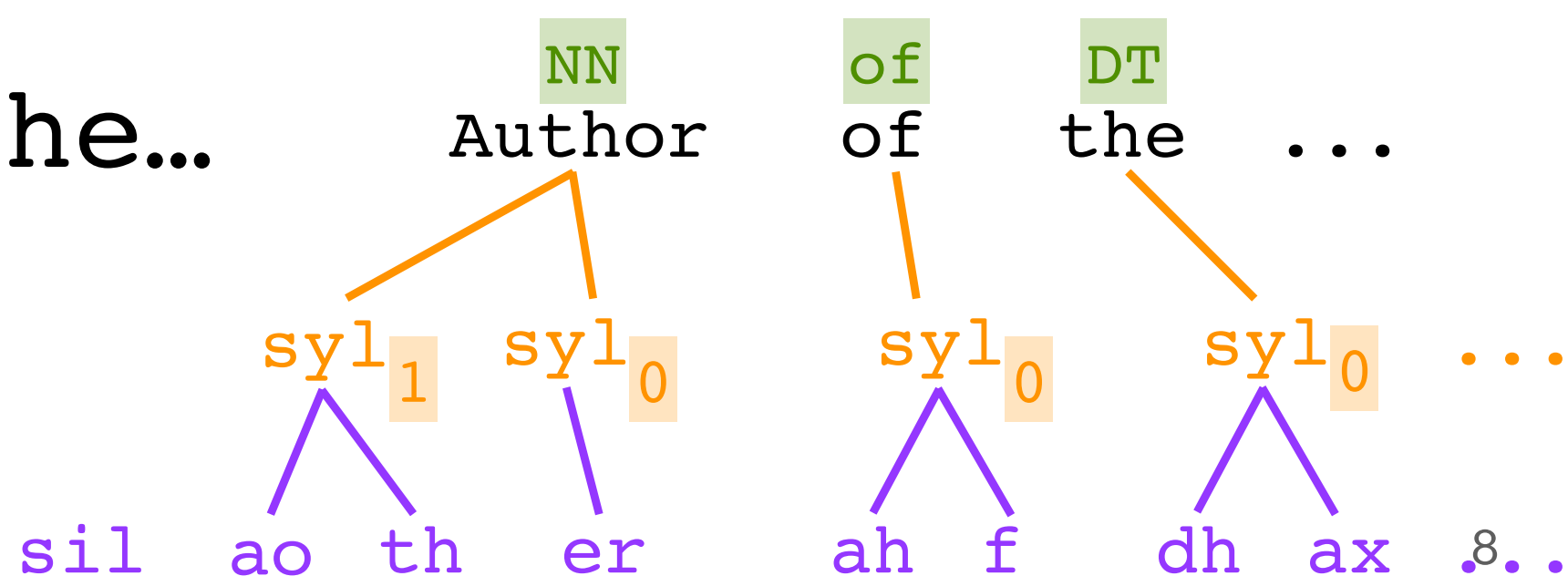
text

*linguistic
specification*

acoustic features

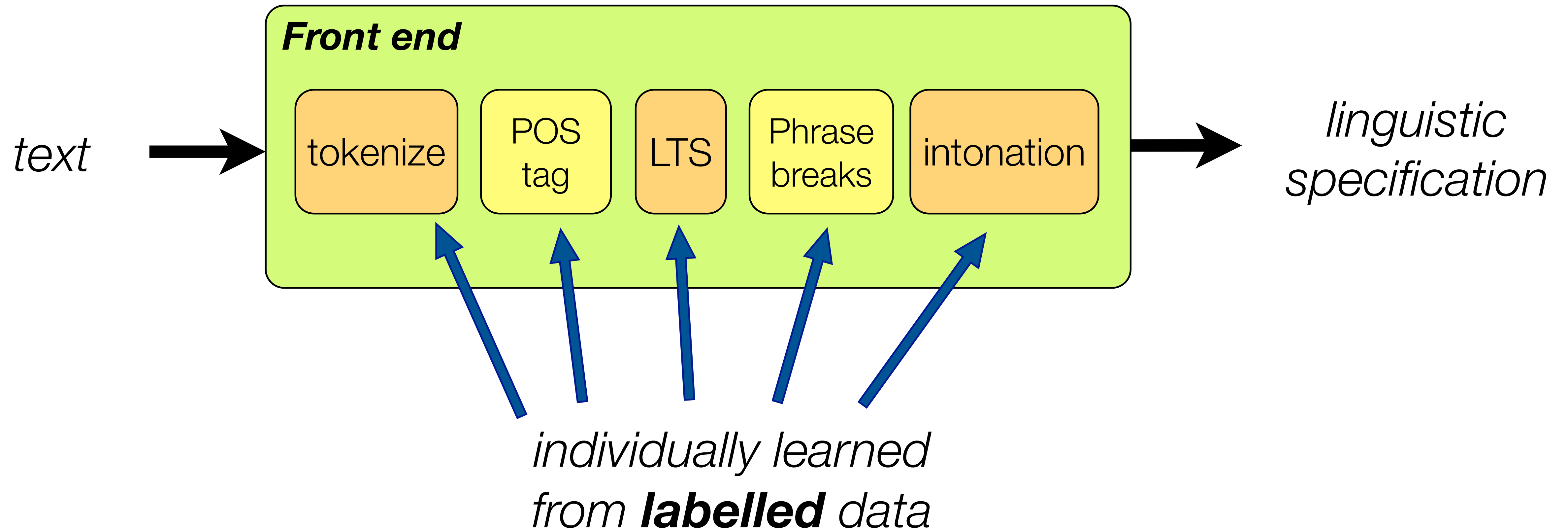
waveform

Author of the...



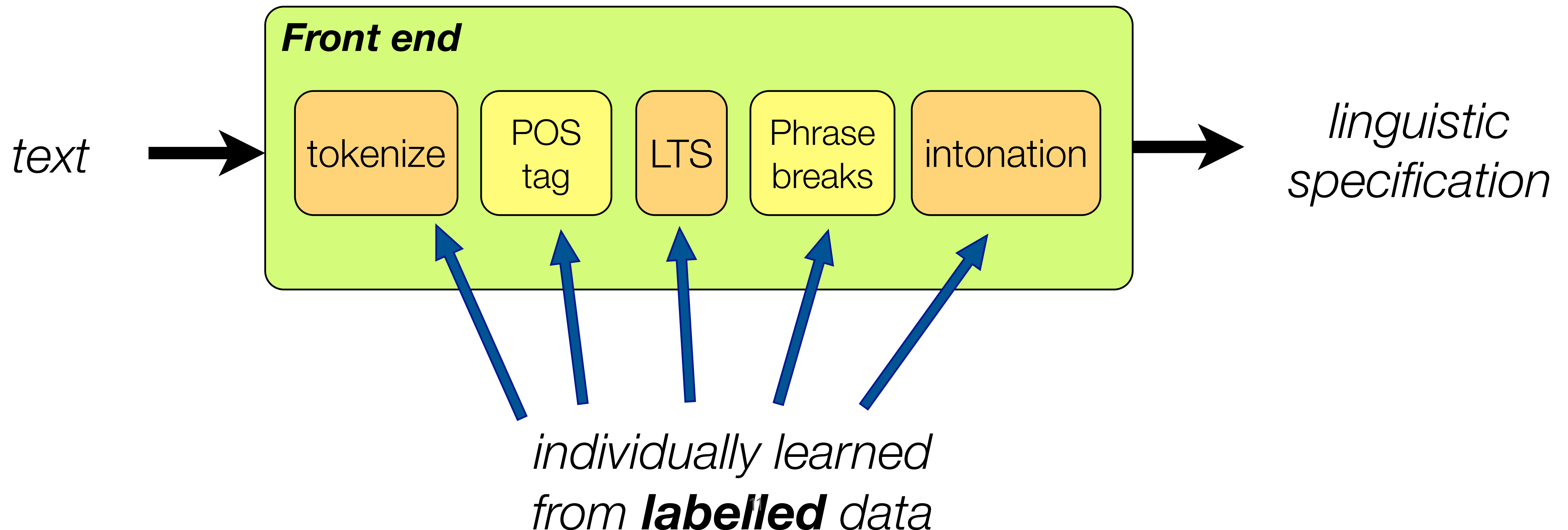
Front end

Front end



Classic front end

- ▶ A chain of **processes**
- ▶ Each process is performed by a **model**
- ▶ These models are independently trained in a **supervised** fashion on annotated data



Neural front end

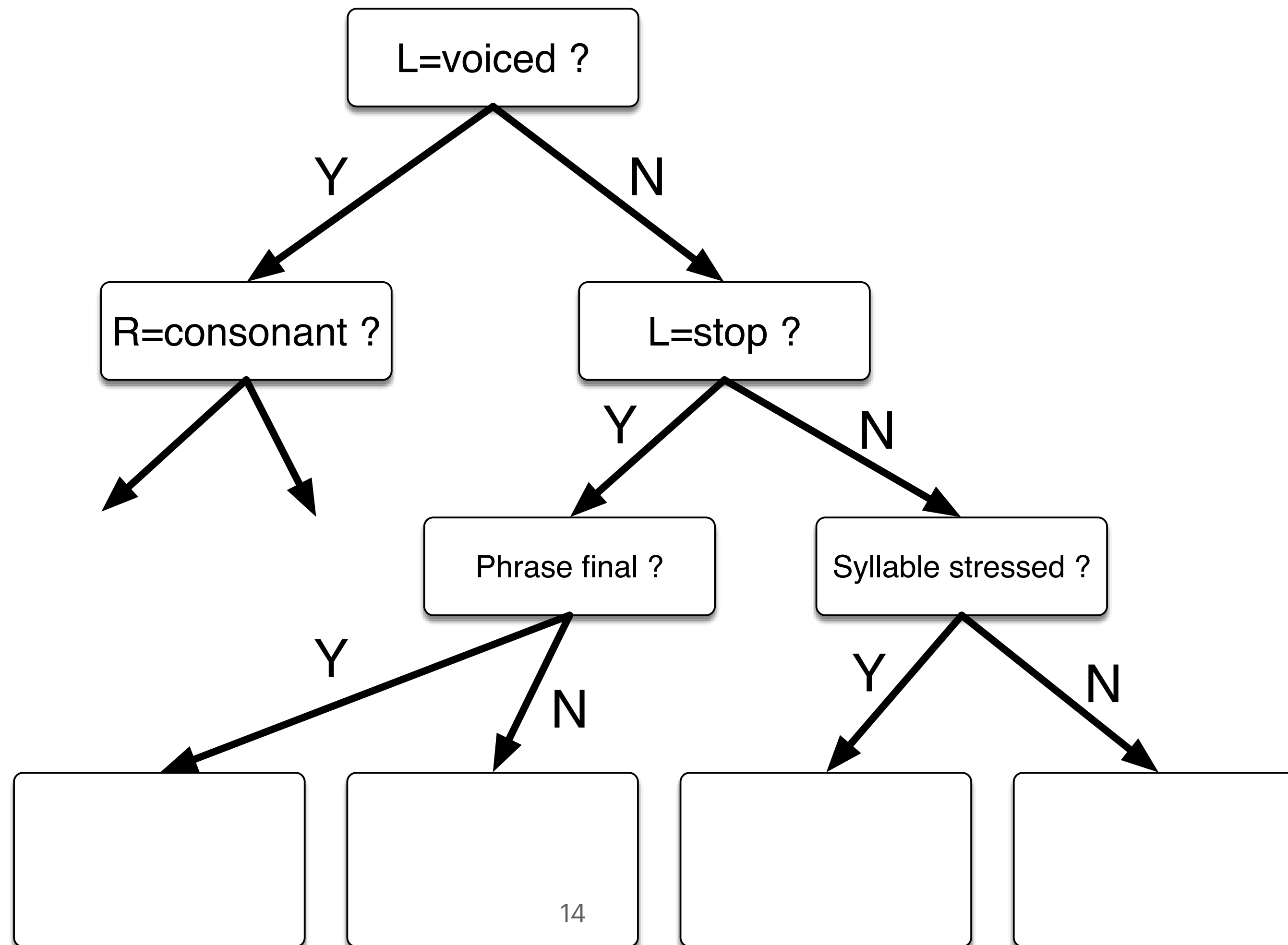
- ▶ Learn by a neural network



Acoustic model

Acoustic model - Decision tree

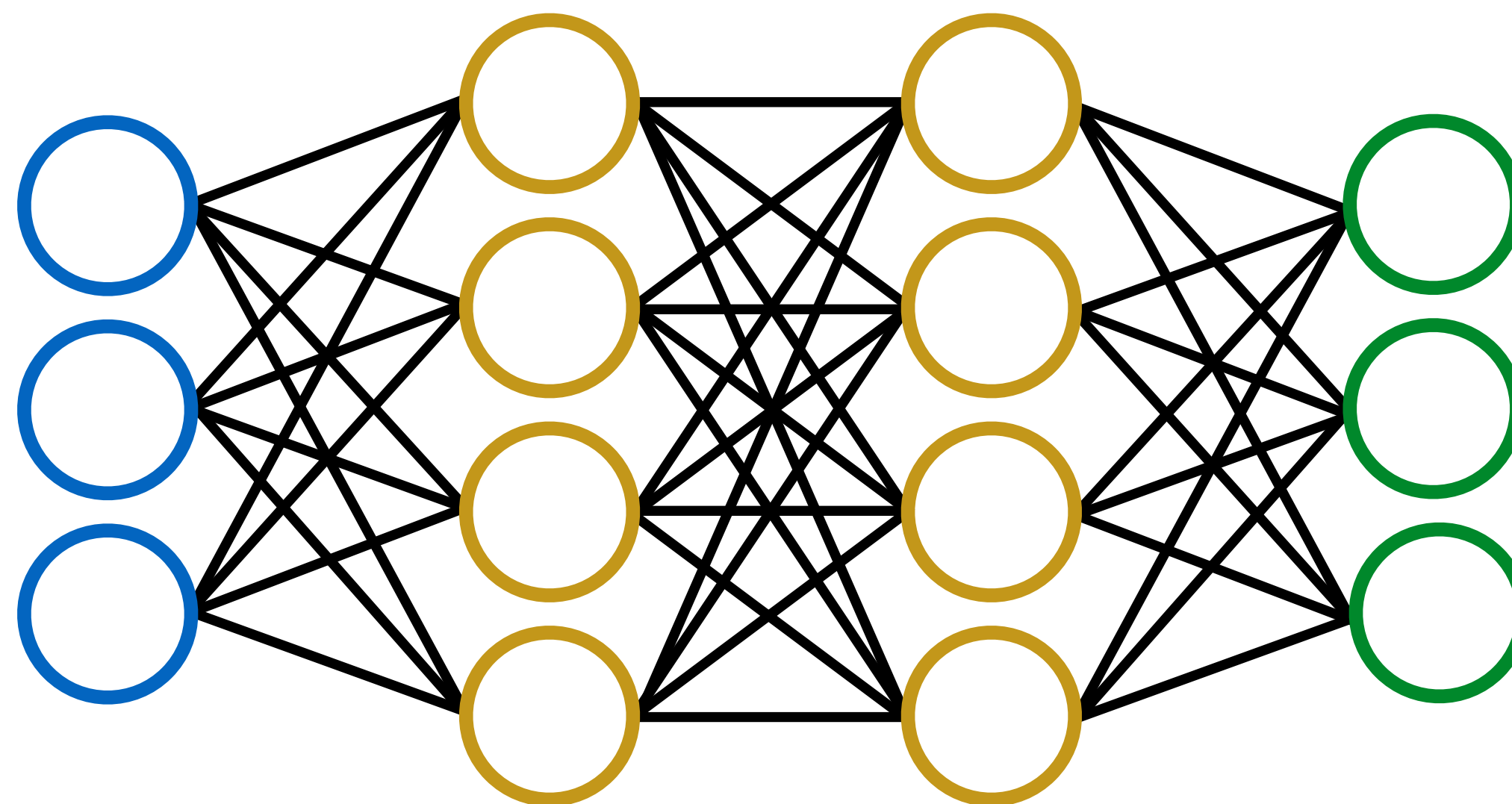
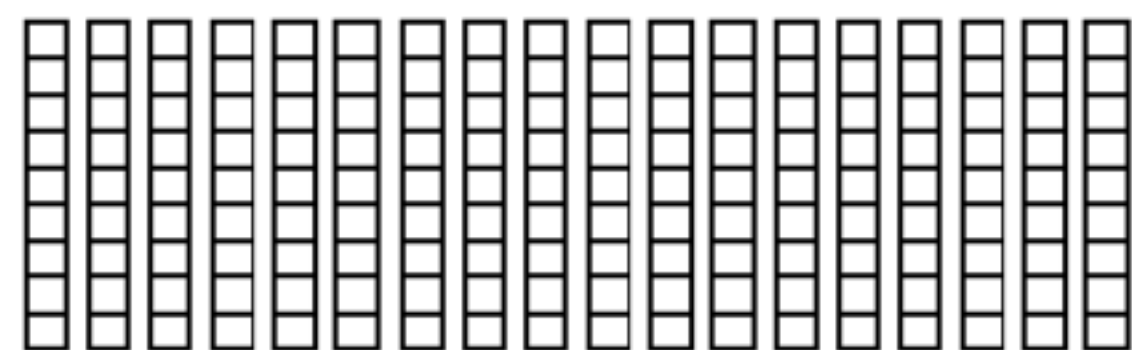
- ▶ Decision tree to group HMM states, which model acoustic feature distribution



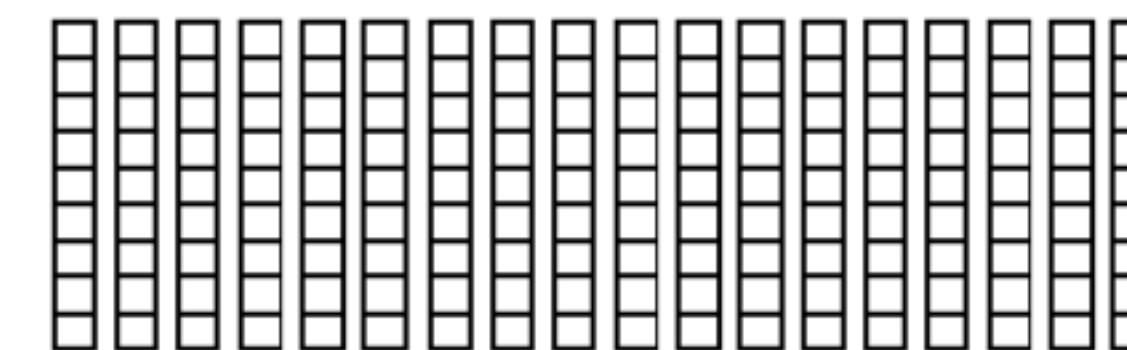
Acoustic model: DNN

- ▶ Feedforward neural network

input
Linguistic features

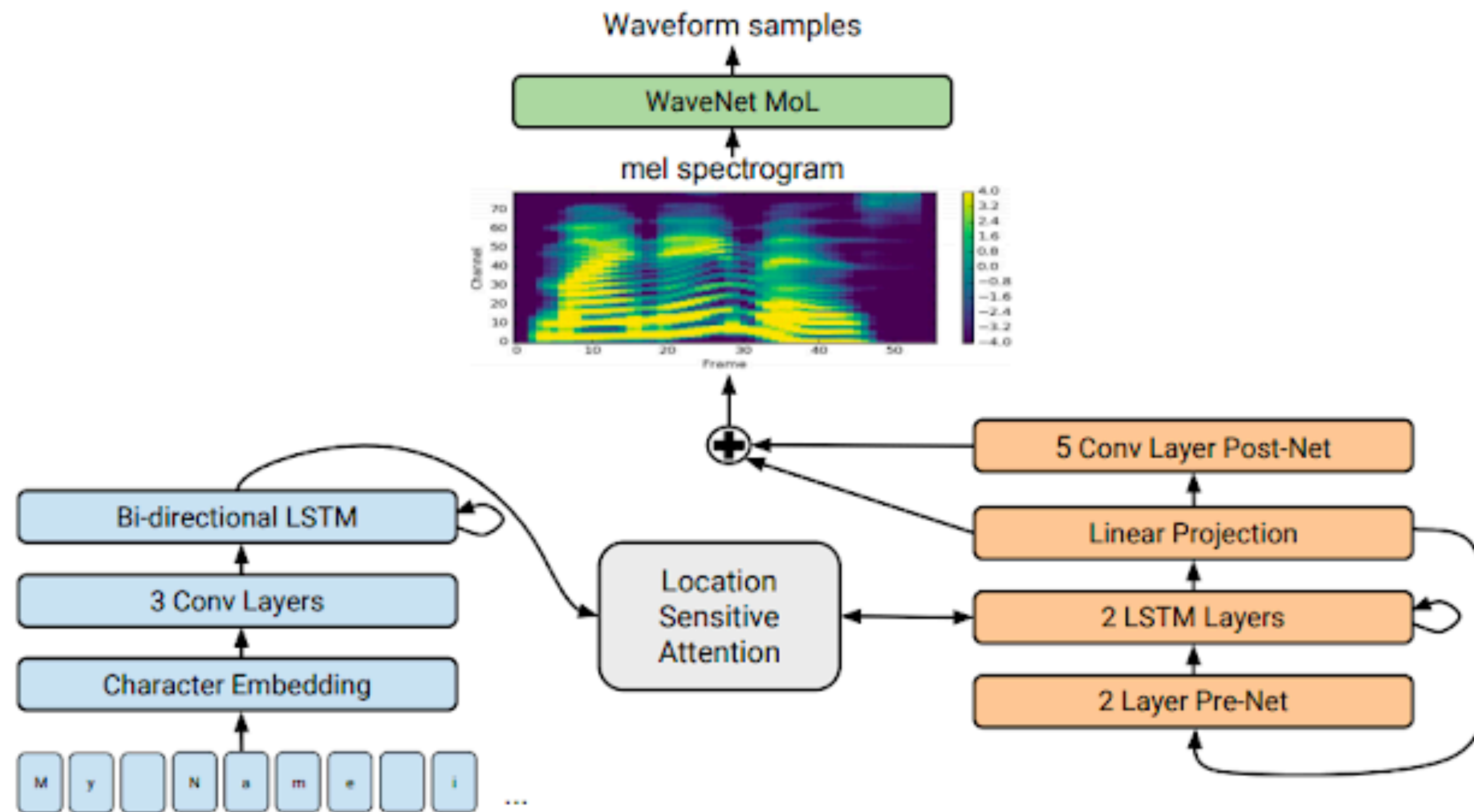


output
acoustic features



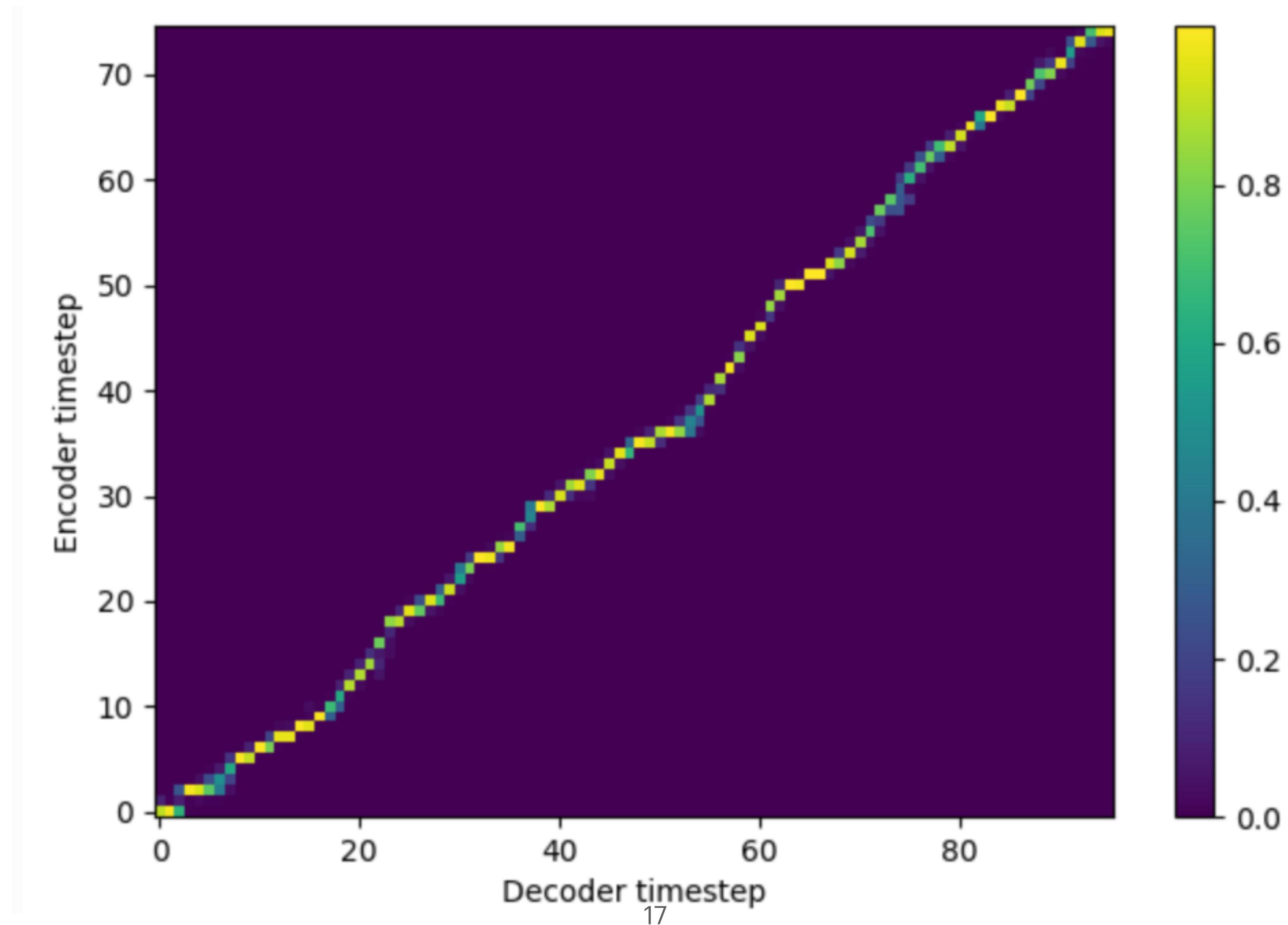
Acoustic model - RNN based

- ▶ Tacotron2: A sequence-to-sequence model based on Recurrent Neural Networks



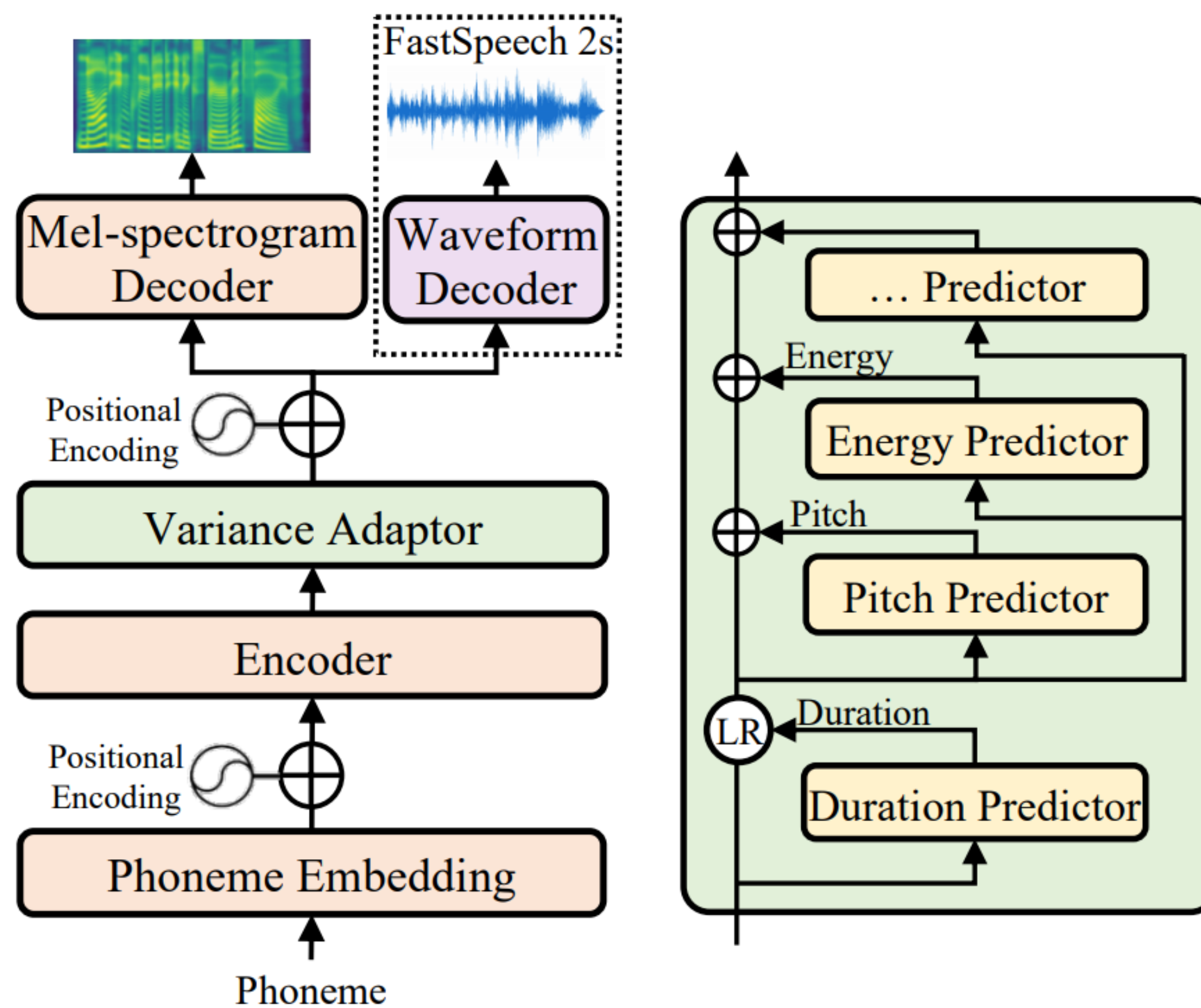
Acoustic model - RNN based

- ▶ Attention



Acoustic model - Transformer based

- ▶ FastSpeech2: parallel generation and not depending on the location attention

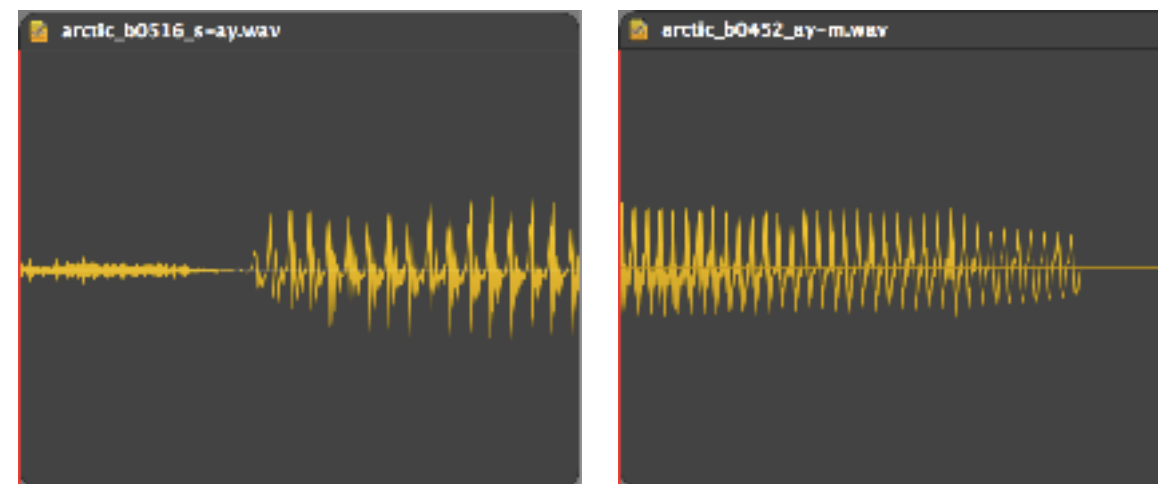


(a) FastSpeech 2

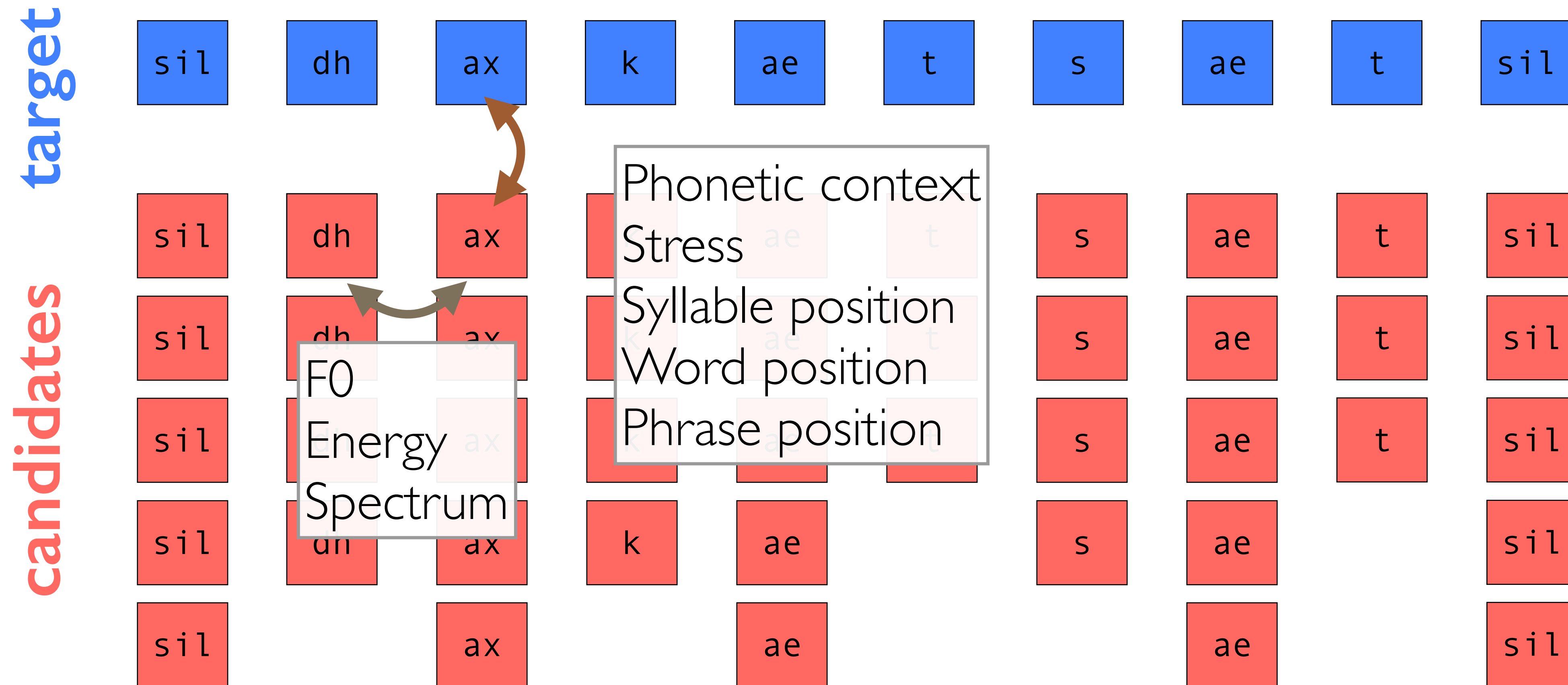
(b) Variance adaptor

Waveform generator

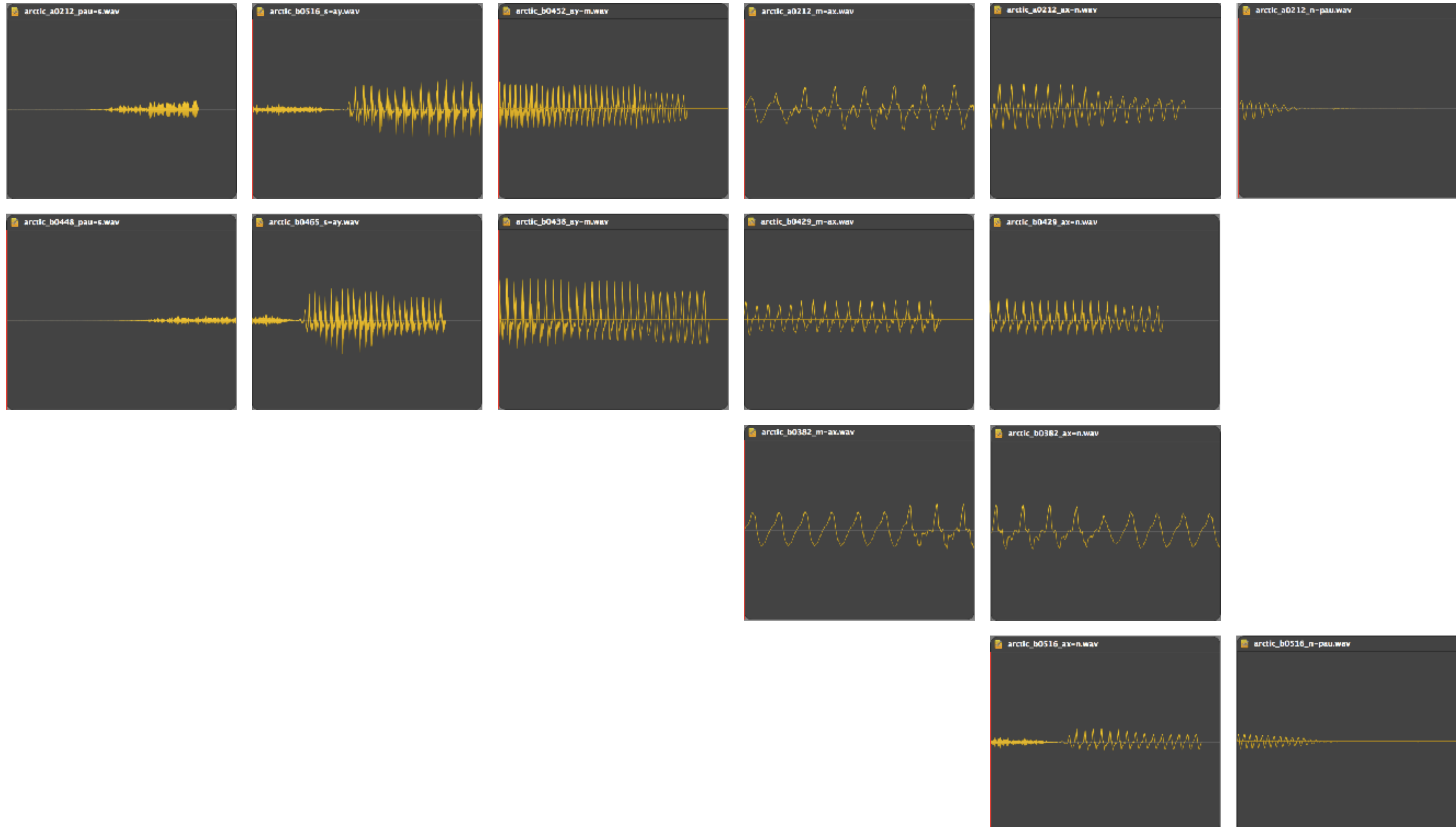
Waveform generator: Waveform concatenation



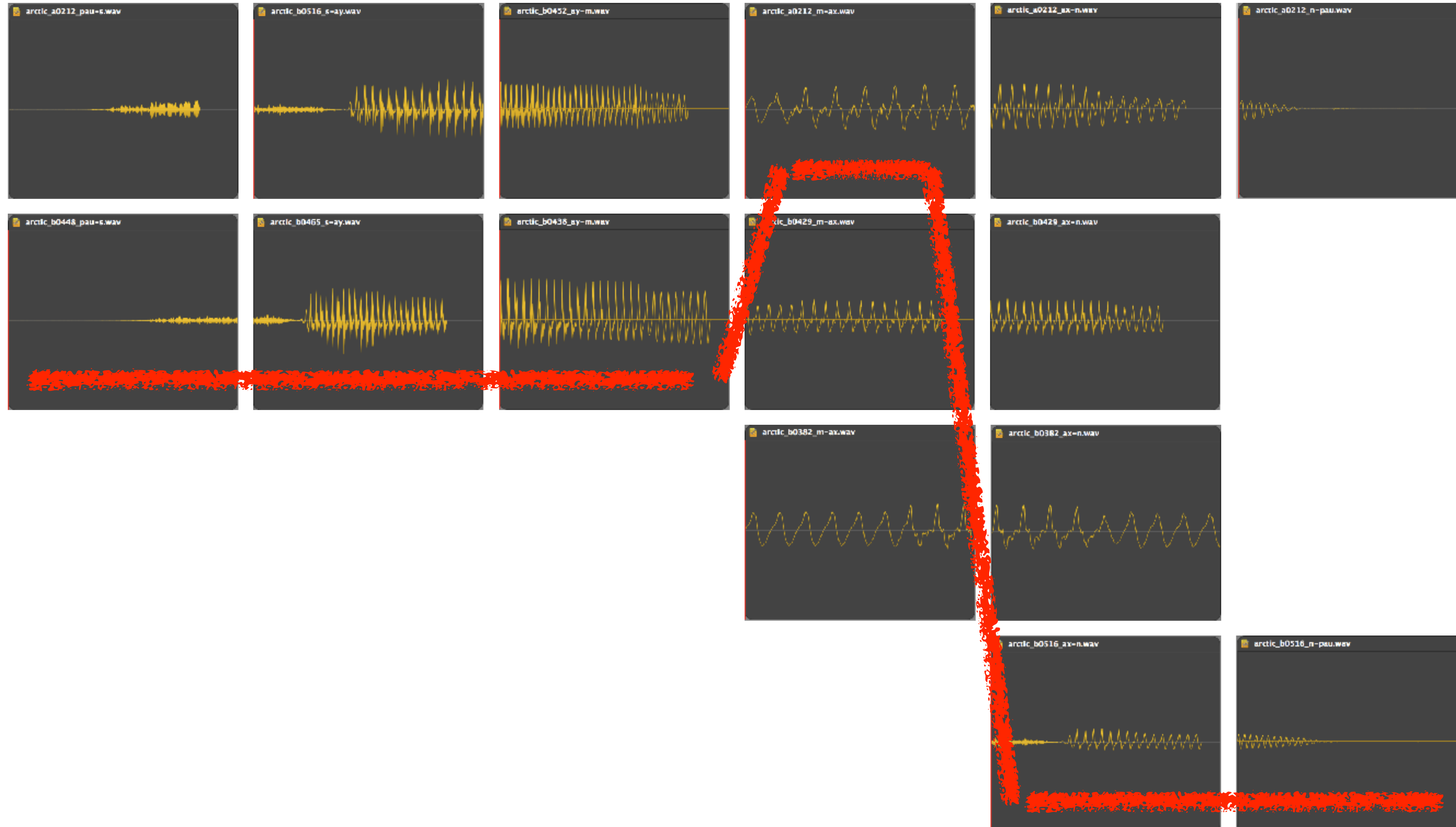
Classical unit selection (drawn here with phone units) - target and join costs



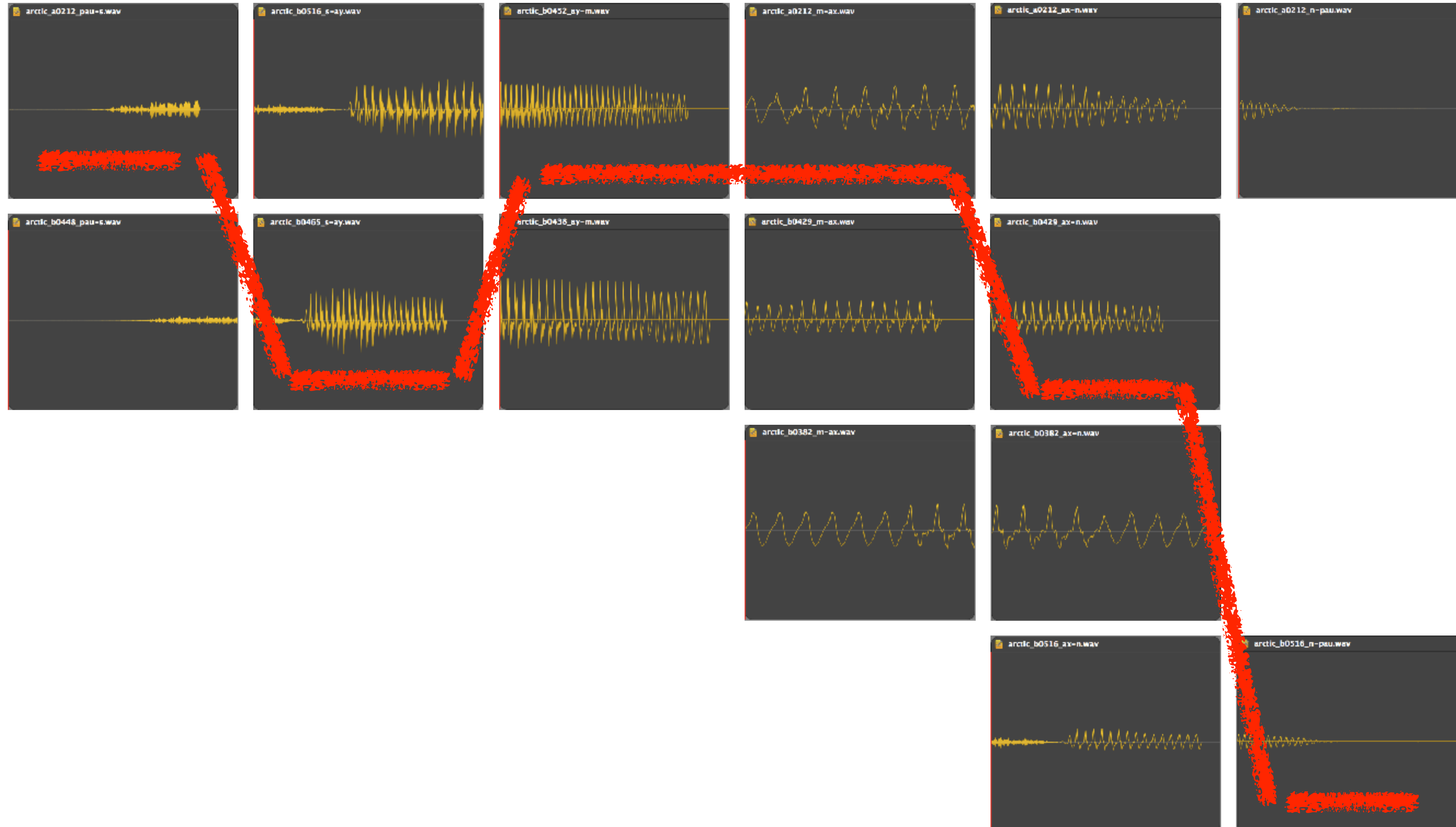
“Simon”



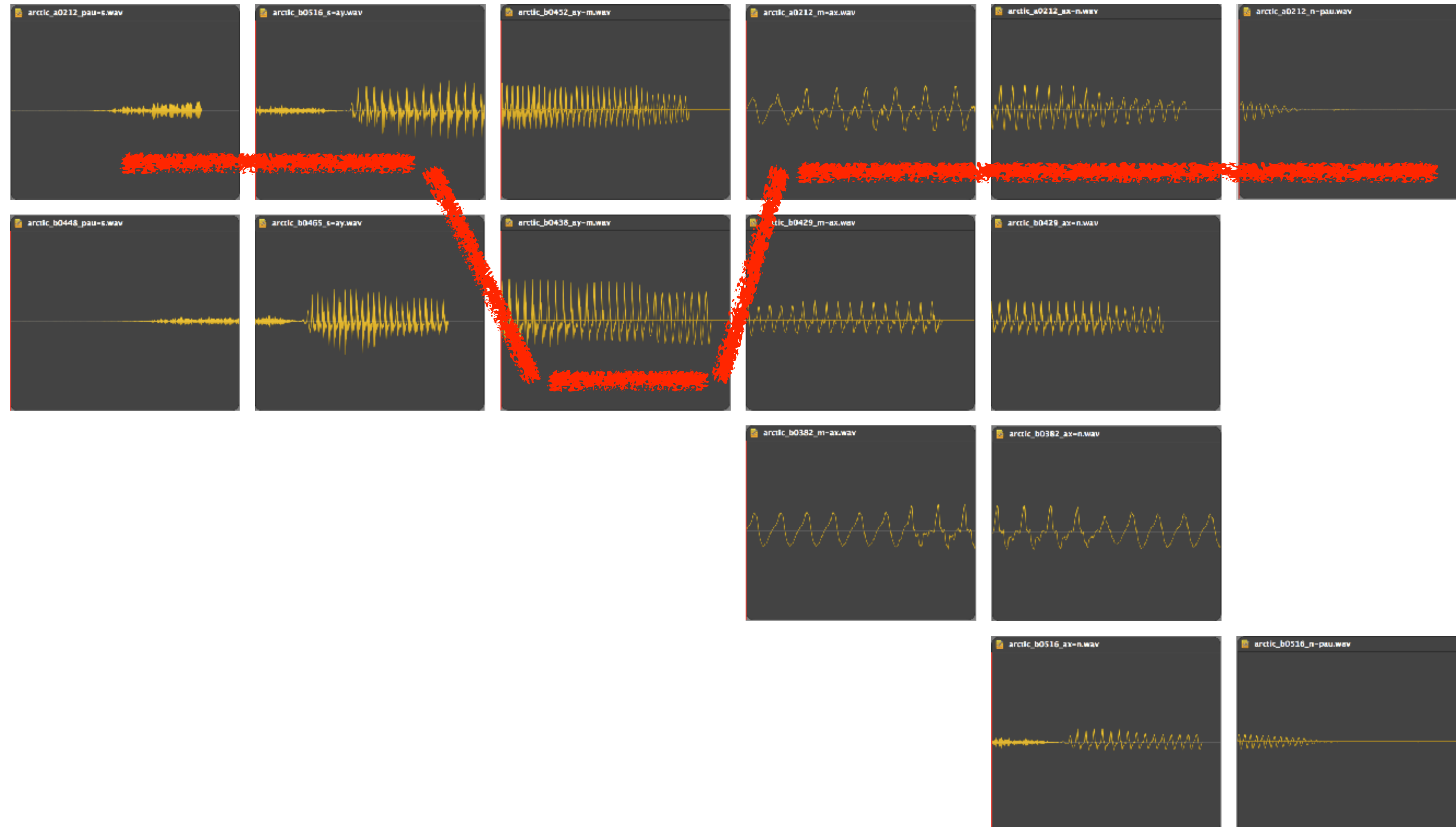
“Simon”



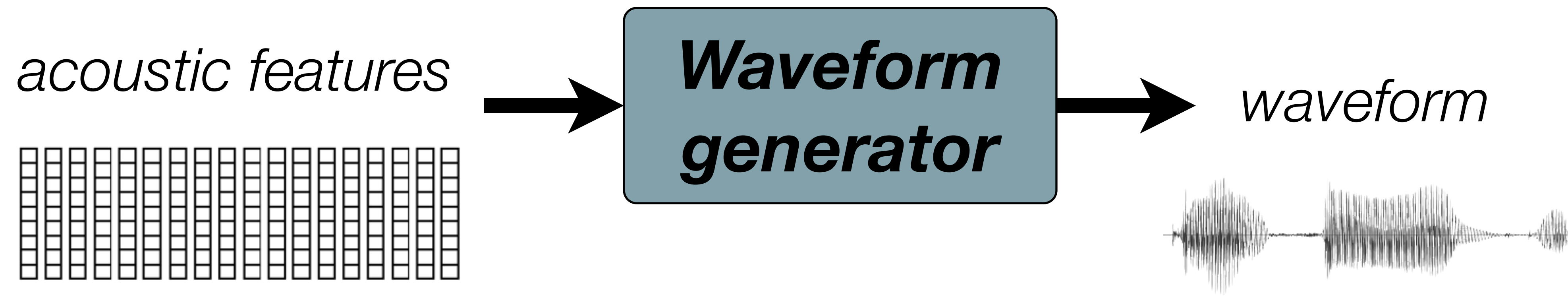
“Simon”



“Simon”

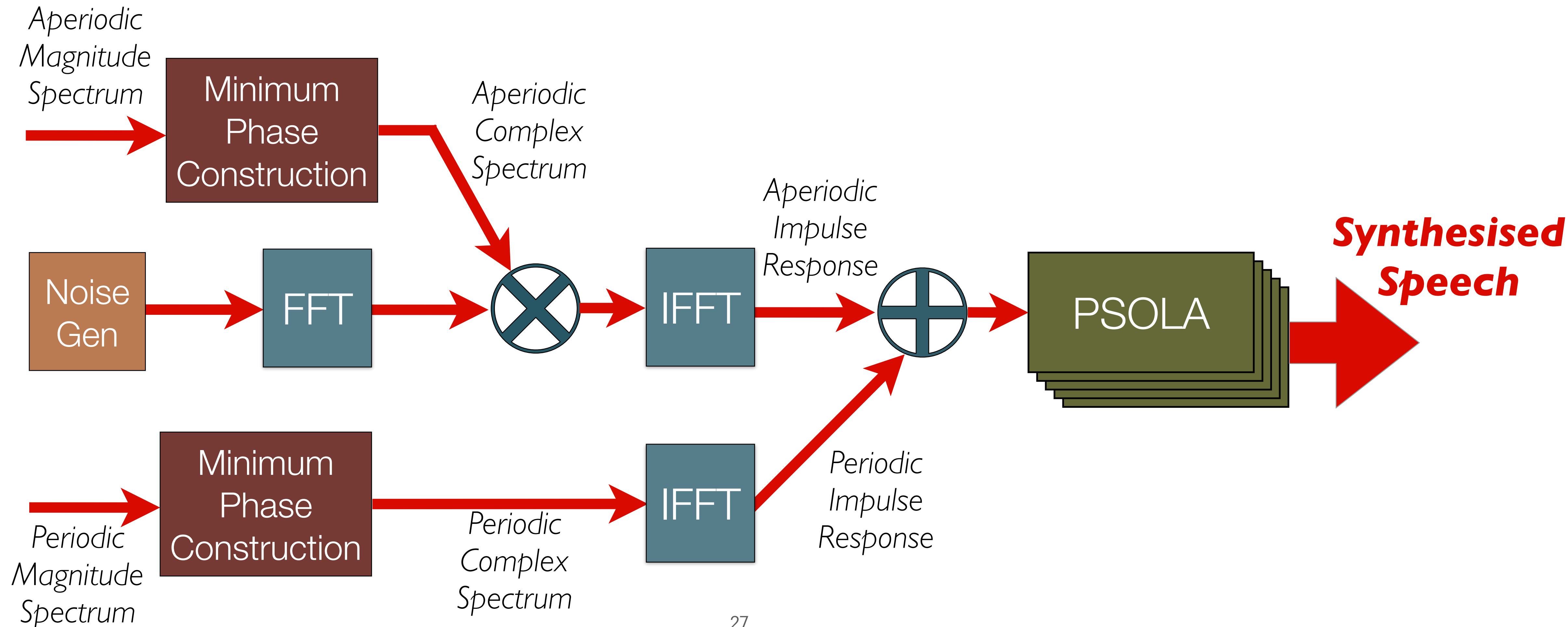


Waveform generator: Vocoder



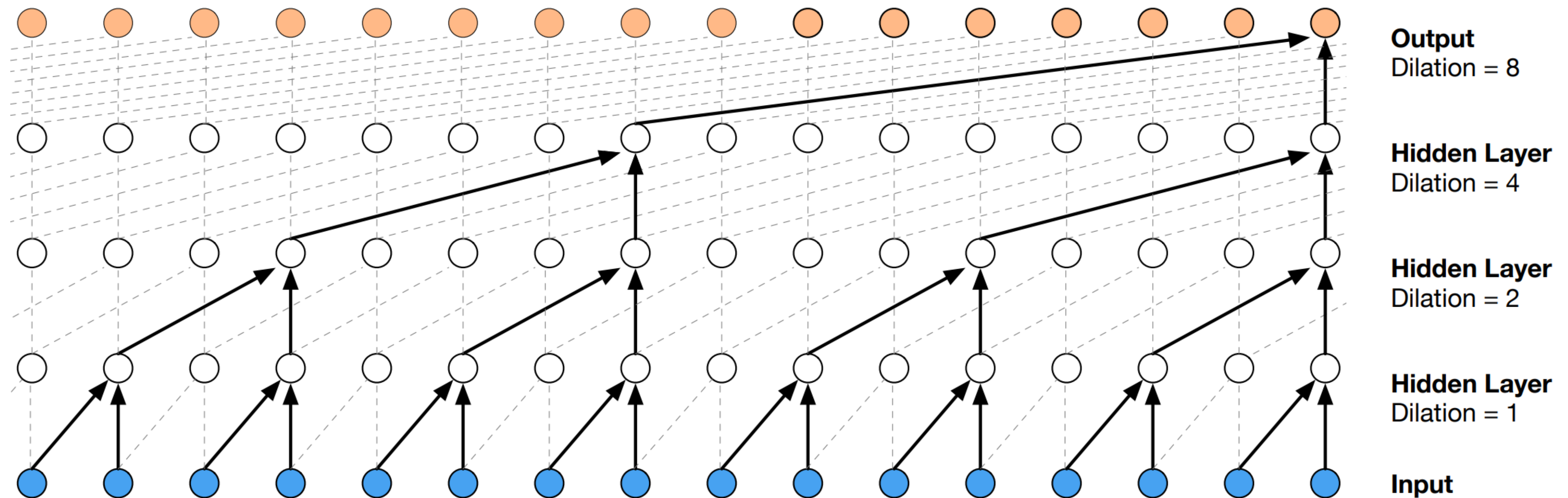
Vocoder - Signal processing based

► WORLD vocoder



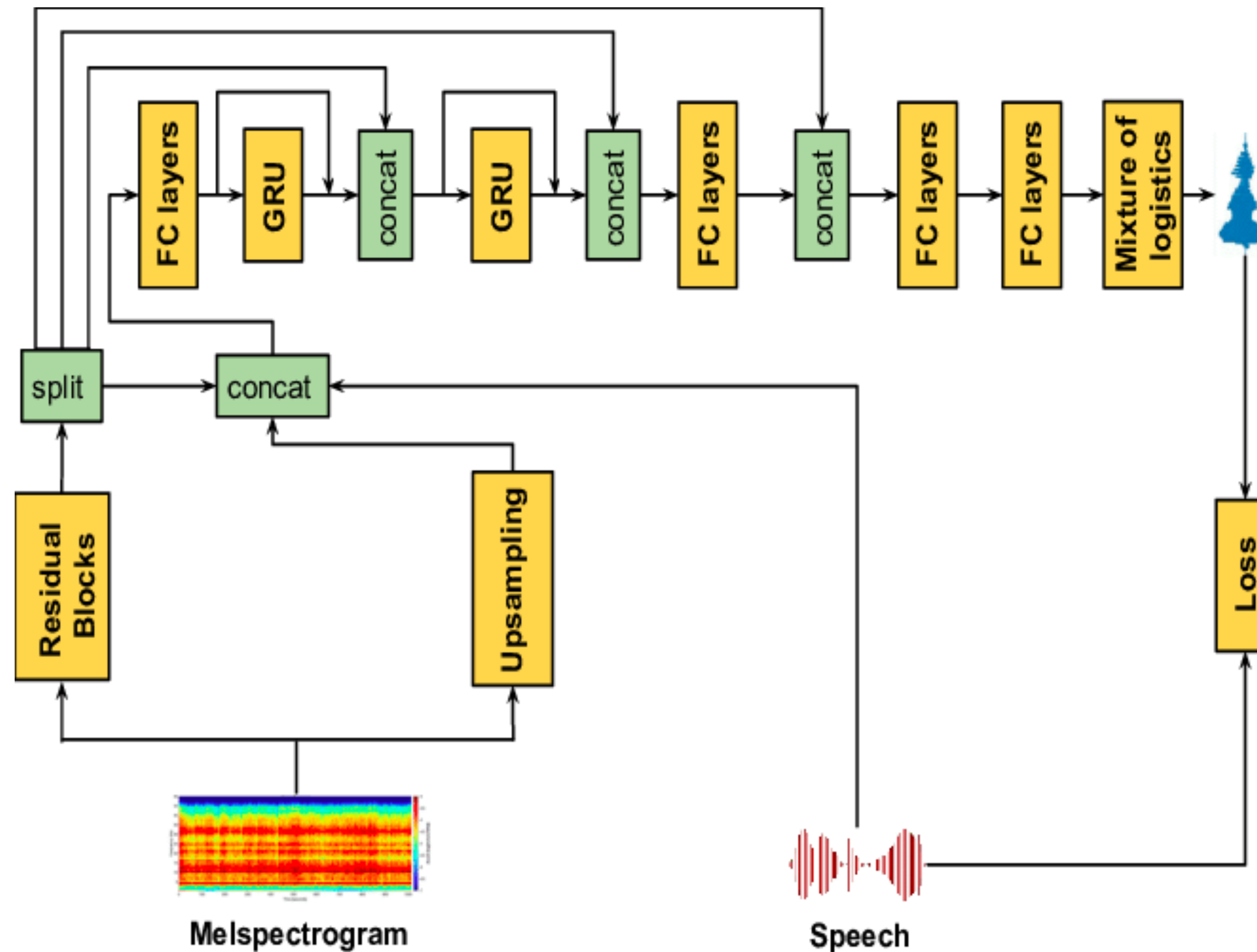
Vocoder: Autoregressive

- ▶ WaveNet: autoregressive model with dilated causal convolution



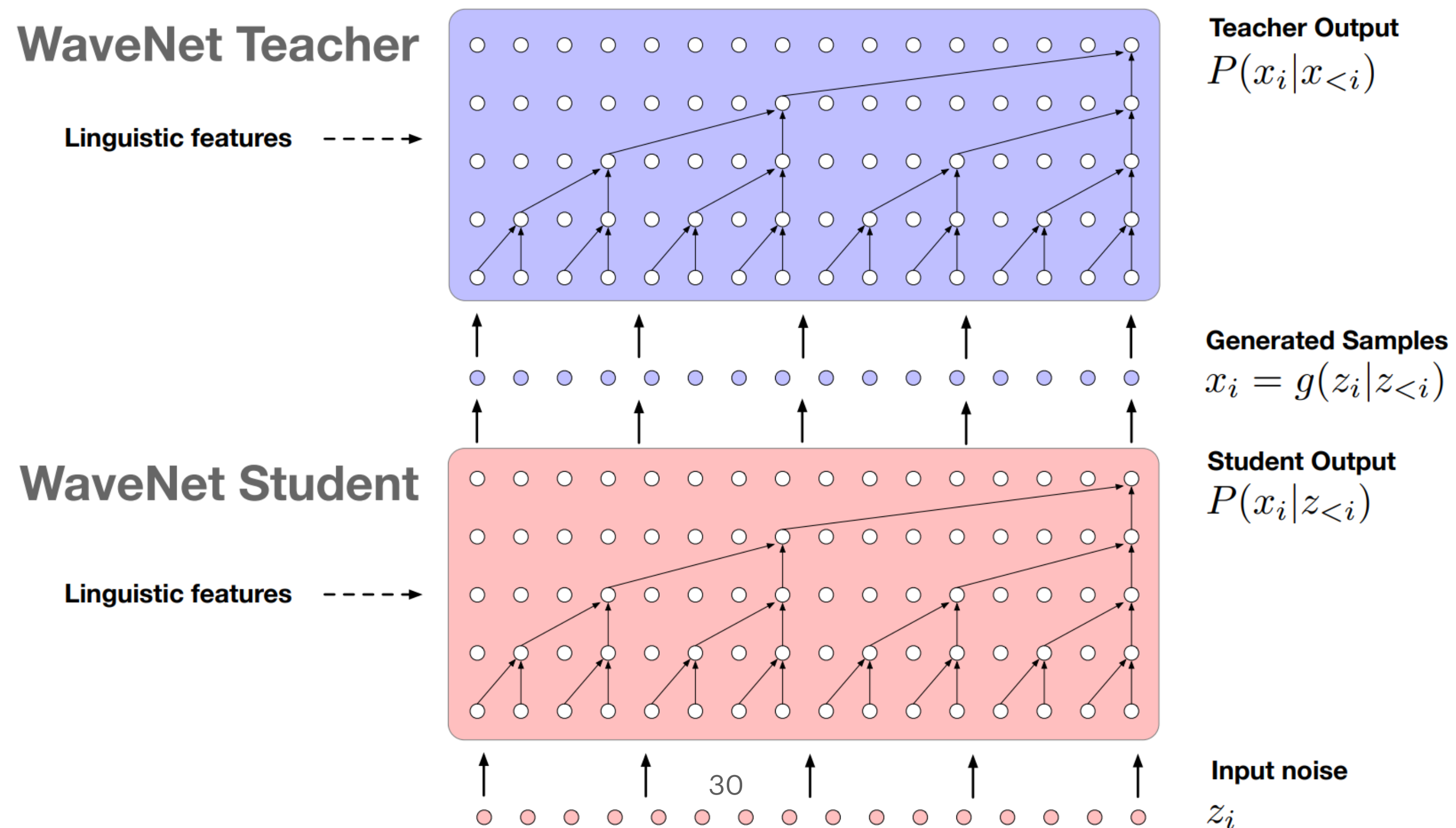
Vocoder: Autoregressive

- ▶ WaveRNN: autoregressive model with RNN



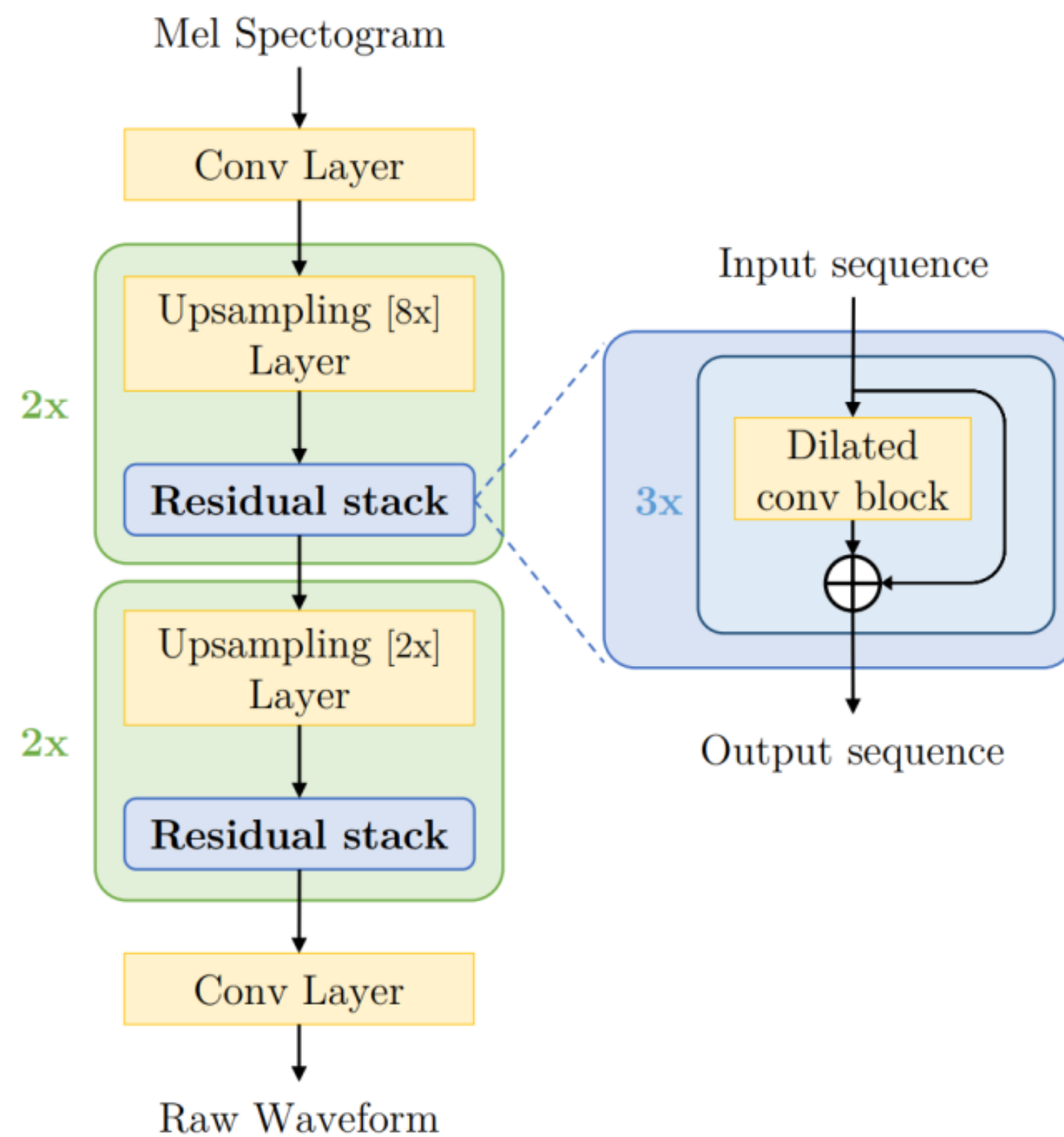
Vocoder: Flow based

- ▶ AF (autoregressive flow) and IAF (inverse autoregressive flow)
 - Parallel inference of IAF student
 - Parallel training of AF teacher

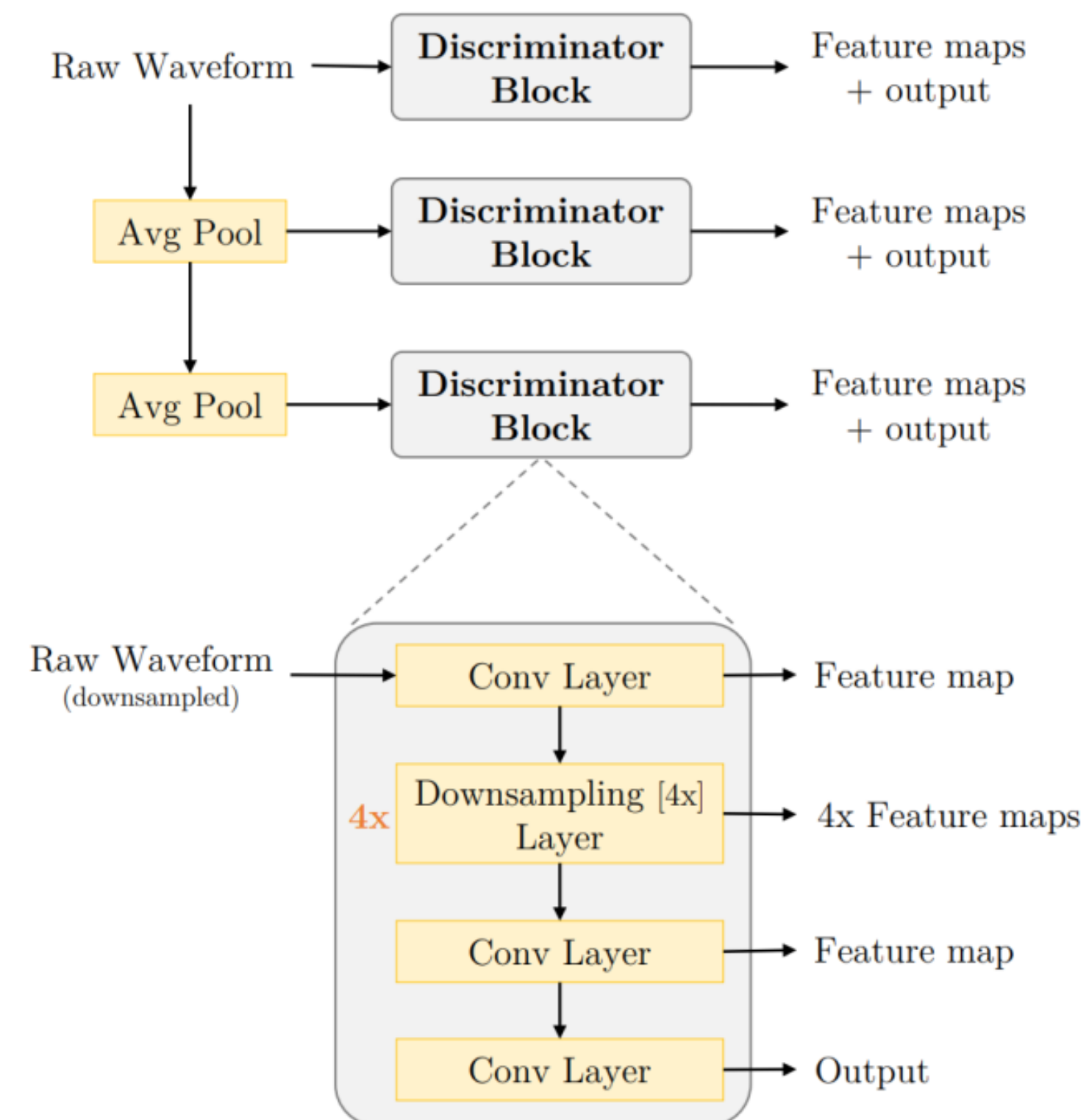


Vocoder: GAN based

- ▶ MelGAN: Generator + Discriminator



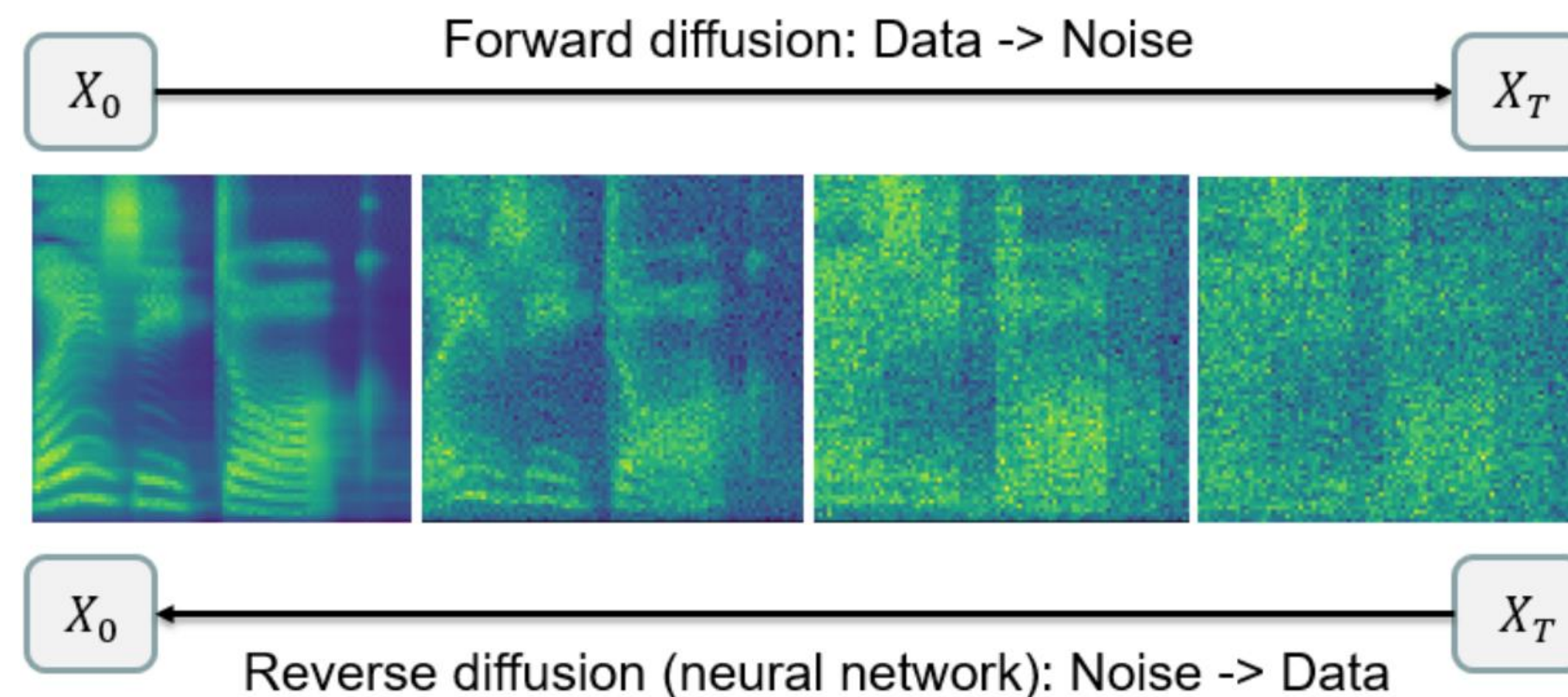
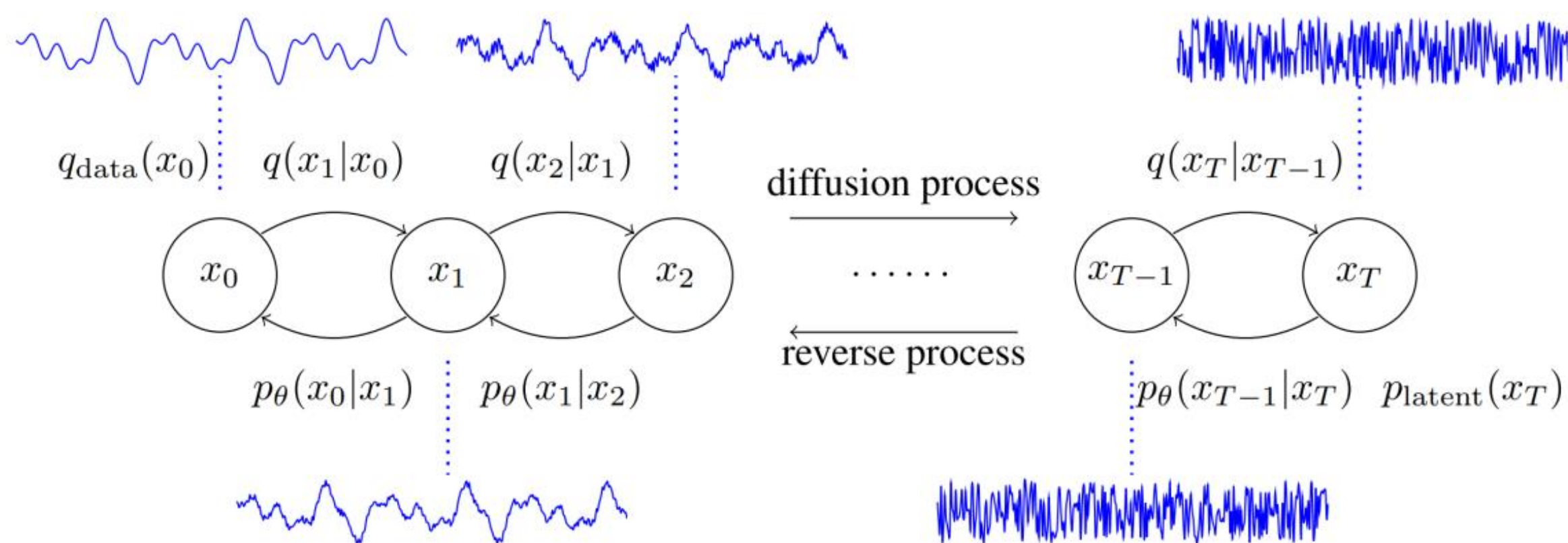
(a) Generator



(b) Discriminator

Vocoder: Diffusion based

- ▶ Diffusion probabilistic model
 - Forward process: diffusion
 - Reverse process: denoising



Tools

- ▶ TTS open-source
 - <https://github.com/coqui-ai/TTS>
 - <https://github.com/espnet/espnet>

- ▶ Acoustic models
 - Tacotron2: <https://github.com/NVIDIA/tacotron2>
 - FastSpeech2: <https://github.com/ming024/FastSpeech2>

- ▶ Vocoder
 - <https://github.com/coqui-ai/TTS/tree/dev/TTS/vocoder>
 - <https://github.com/NVIDIA/BigVGAN>

Readings

- ▶ Interspeech 2022 TTS tutorial
 - https://github.com/tts-tutorial/interspeech2022/blob/main/INTERSPEECH_Tutorial_TTS.pdf
- ▶ Text-to-Speech Synthesis
 - <https://www.cambridge.org/core/books/texttospeech-synthesis/D2C567CEF939C7D15B2F1232992C7836>