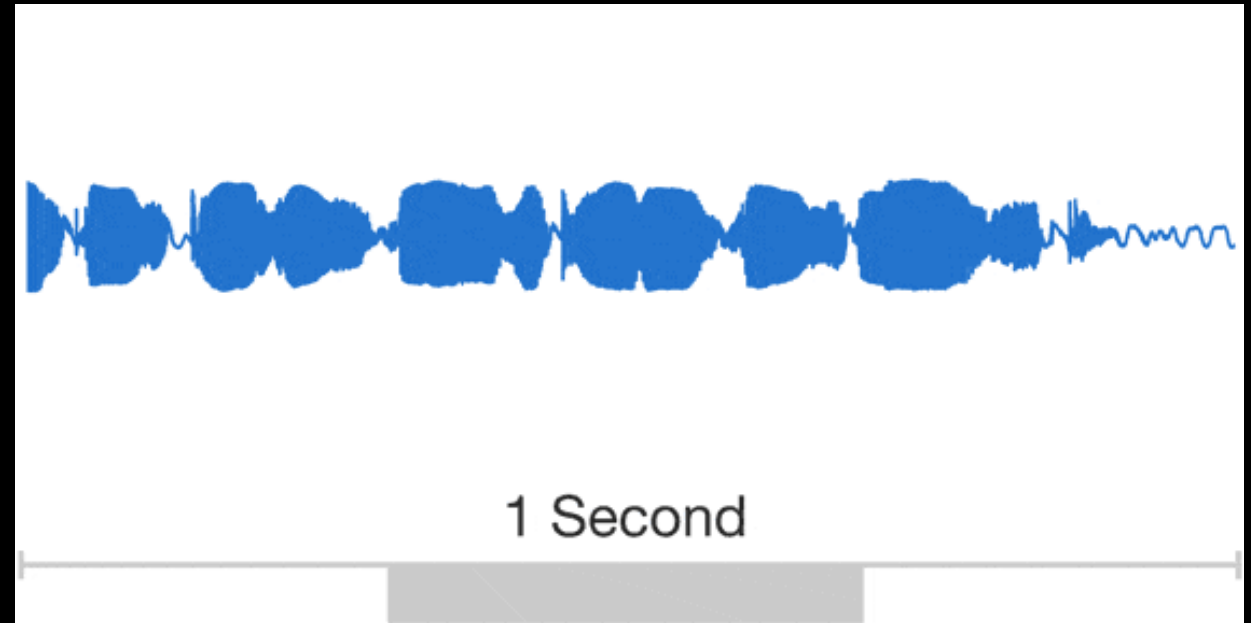# Waveform-based deep learning research at Dolby

**Jordi Pons** (@jordiponsdotme / www.jordipons.me)
Representing the recent work of the Applied AI team!

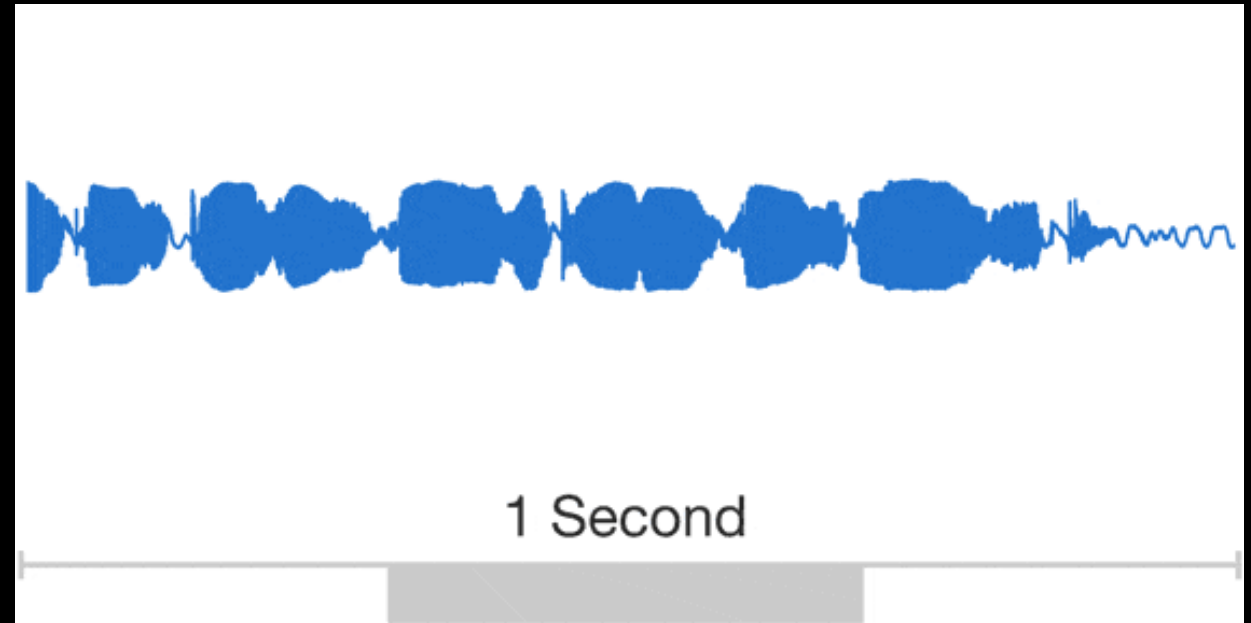# Challenges of Deep Learning in Audio

HIGH DIMENTIONALITY



1 Second

Animation from:
https://deepmind.com/blog/wavenet-generative-model-raw-audio

# Challenges of Deep Learning in Audio

HIGH DIMENTIONALITY

MULTI-LEVEL TEMPORAL
DEPENDANCY



1 Second

Animation from:
https://deepmind.com/blog/wavenet-generative-model-raw-audio
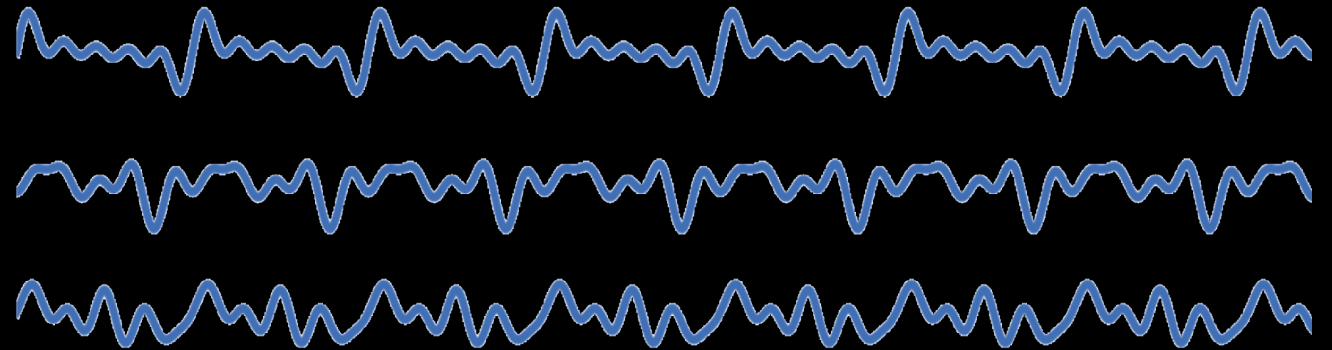
# Challenges of Deep Learning in Audio

HIGH DIMENTIONALITY
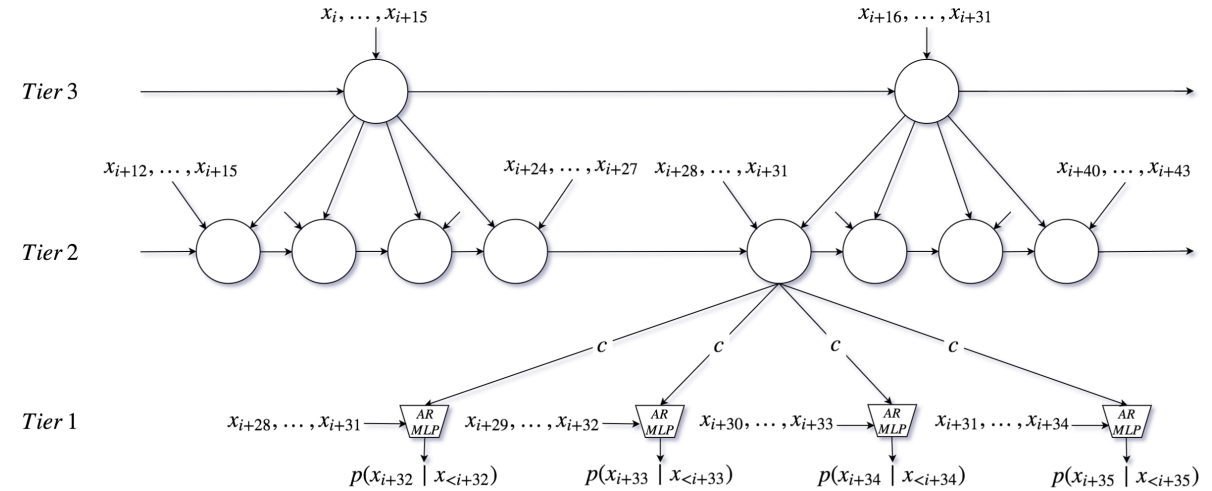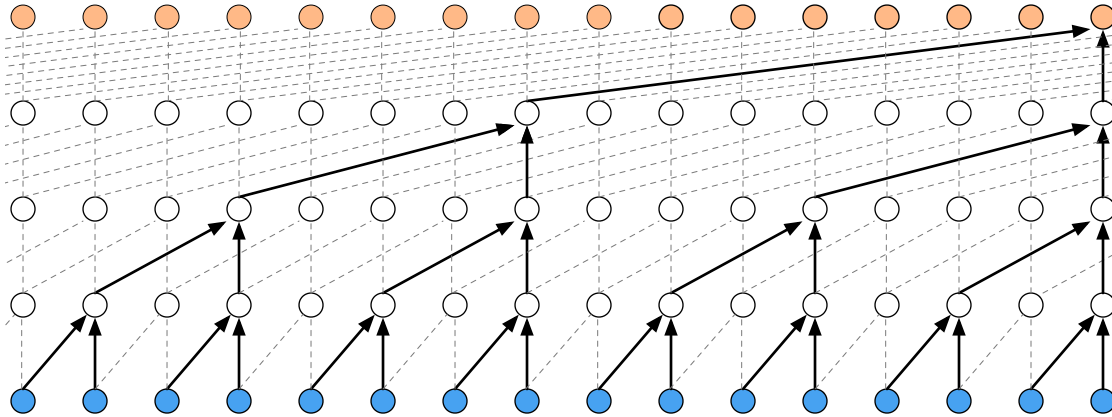
MULTI LEVEL TEMPORAL
DEPENDANCY

PERCEPTION MATTERS



Which of these waves sound different?

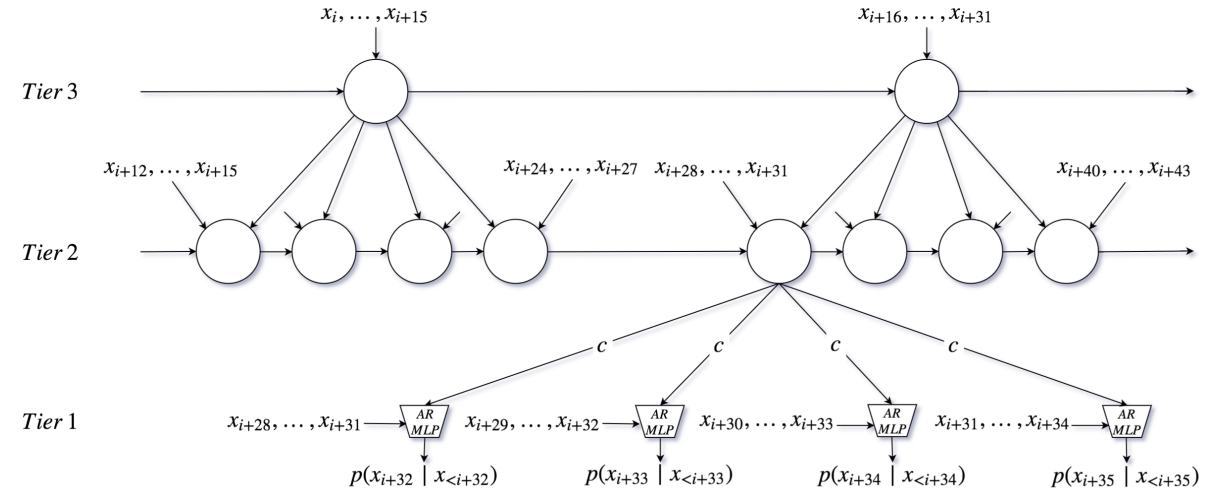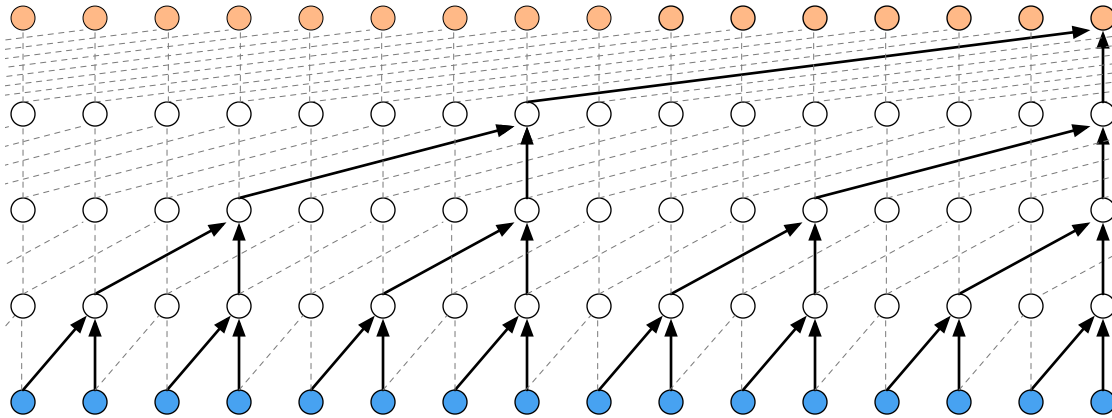Image by Jesse Engel - Problems with WaveNet (DAFx 2019)

# Generating Audio Waveforms



WaveNet: A generative model for raw audio (Google DeepMind)

SampleRNN: Multirate RNN based generative model (MILA)
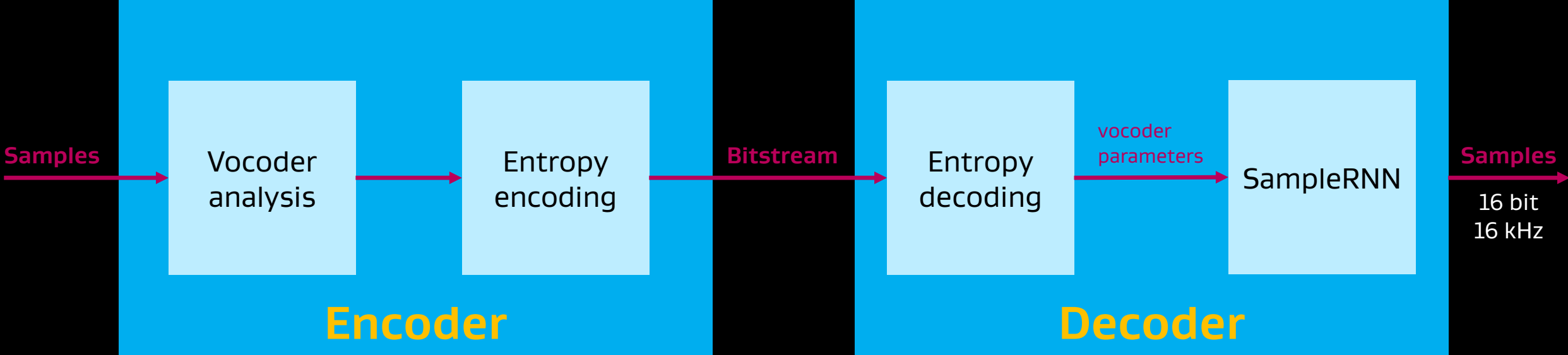
# Generating Audio Waveforms



**WaveNet: A generative model for raw audio (Google DeepMind)**



**SampleRNN: Multirate RNN based generative model (MILA)**

## Low Bitrate Speech Coding

- Wavenet Based Low Rate Speech Coding (Google) *W. Bastiaan Kleijn, Felicia S. C. Lim, Alejandro Luebs, Jan Skoglund, Florian Stimberg, Quan Wang, Thomas C. Walters*

- High-quality speech coding with SampleRNN (Dolby) *Janusz Klejsa, Per Hedelin, Cong Zhou, Roy Fejgin, Lars Villemoes*

# Coding Scheme

# What's in the conditioning?
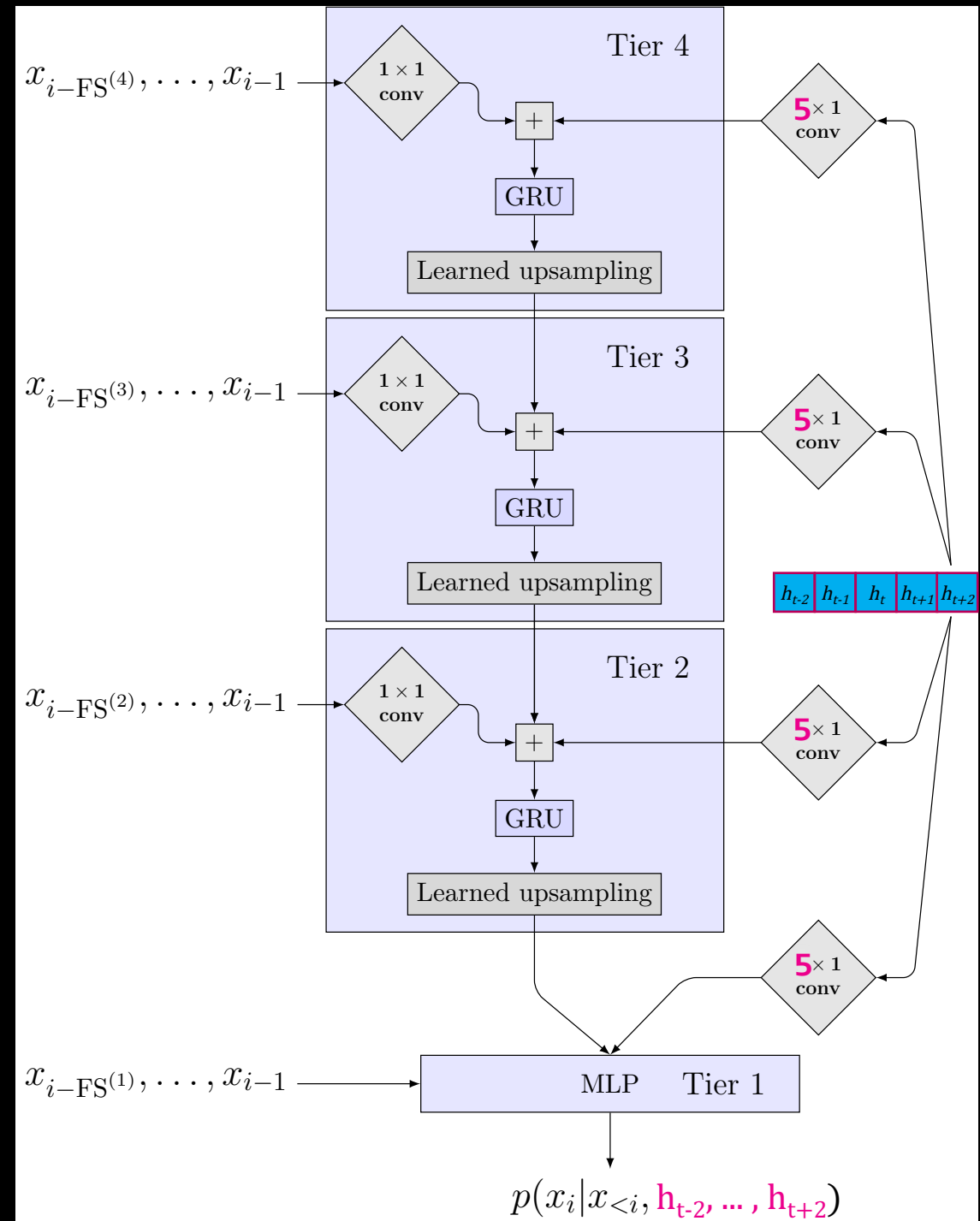
Quantized vocoder parameters:

- Pitch

- LPC filter coefficients

- RMS level of residual

- Voicing level per band (6 bands)
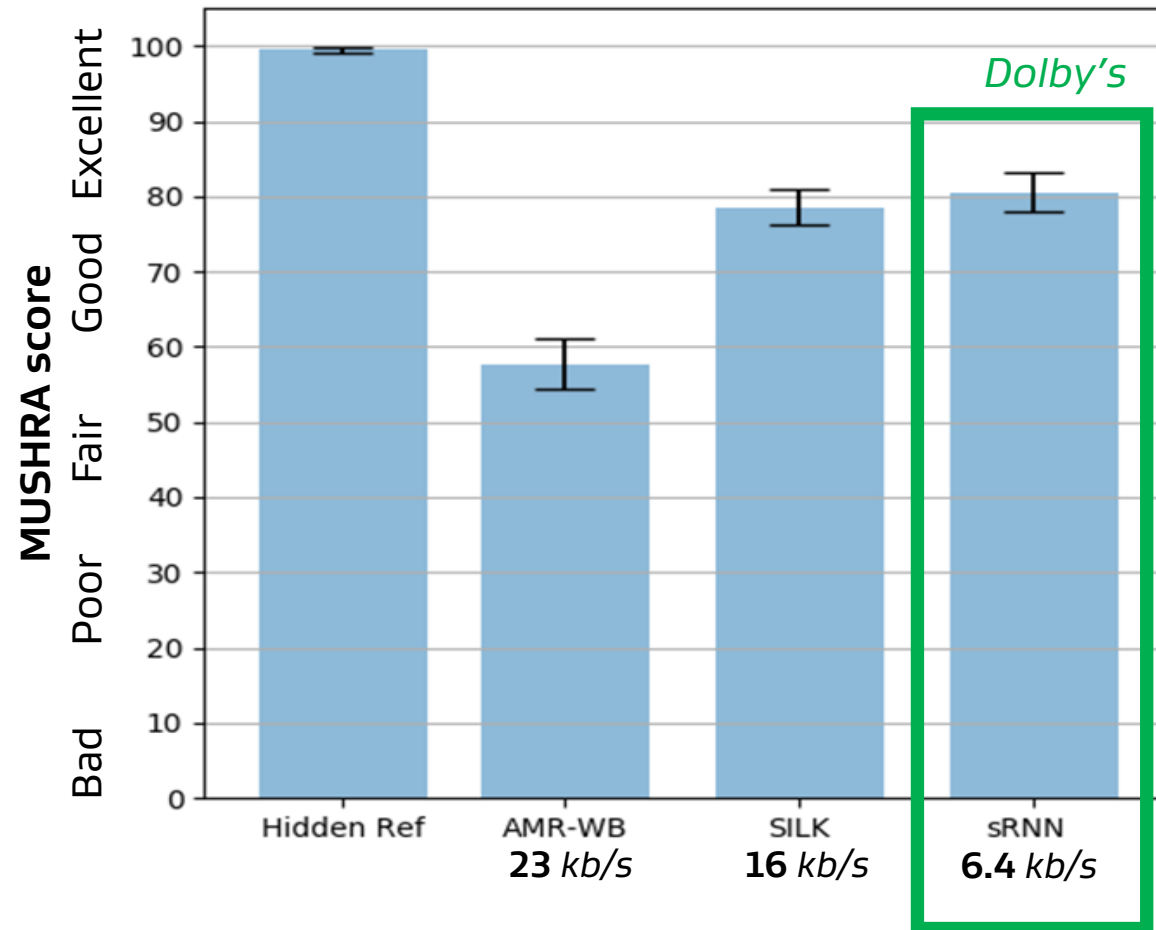
The vocoder is based on:

Per Hedelin, "A sinusoidal LPC vocoder," in 2000 IEEE Workshop on Speech Coding. Proceedings. Meeting the Challenges of the New Millennium (Cat. No.00EX421), Sept 2000, pp. 2–4

# Conditional SampleRNN

- By itself, SampleRNN can only 'babble'
  → we need conditioning

- 4-tier configuration
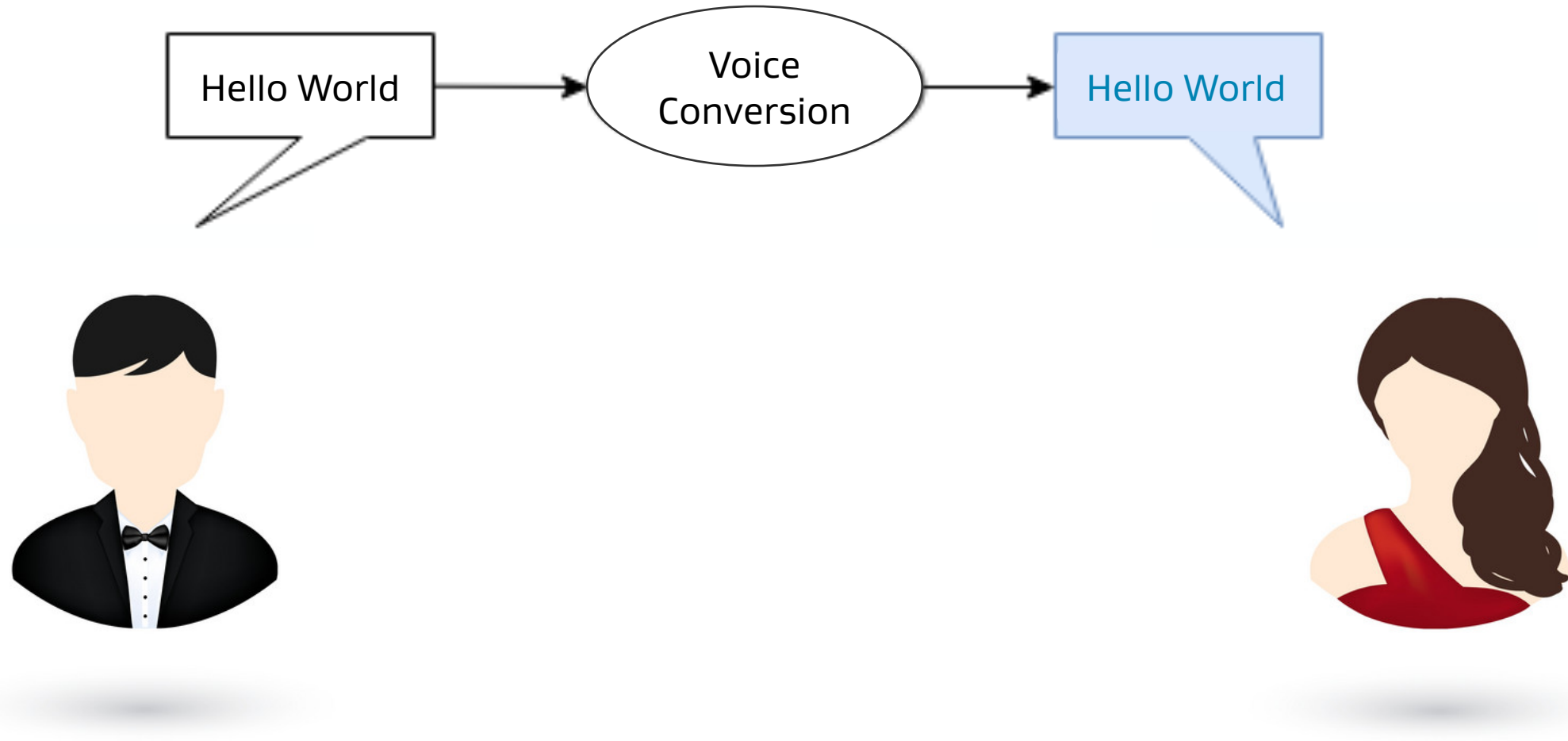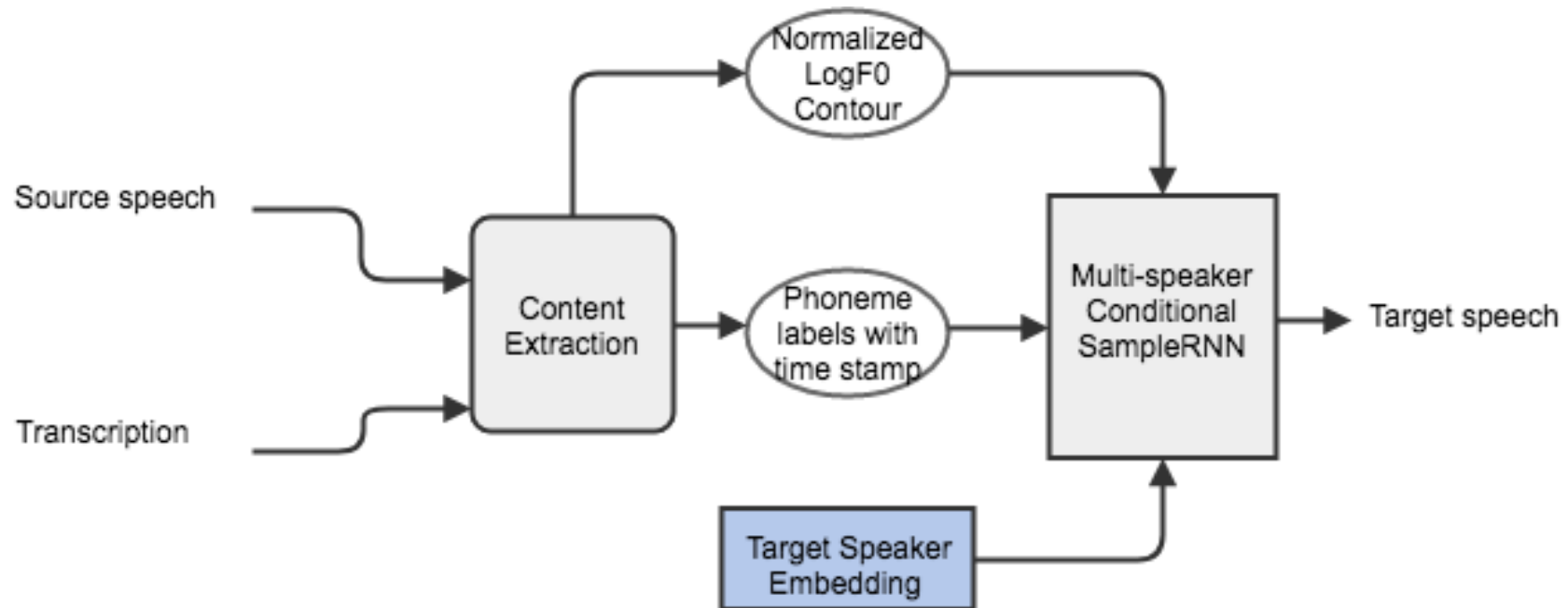
- Conditioning with lookahead

# Listening Tests

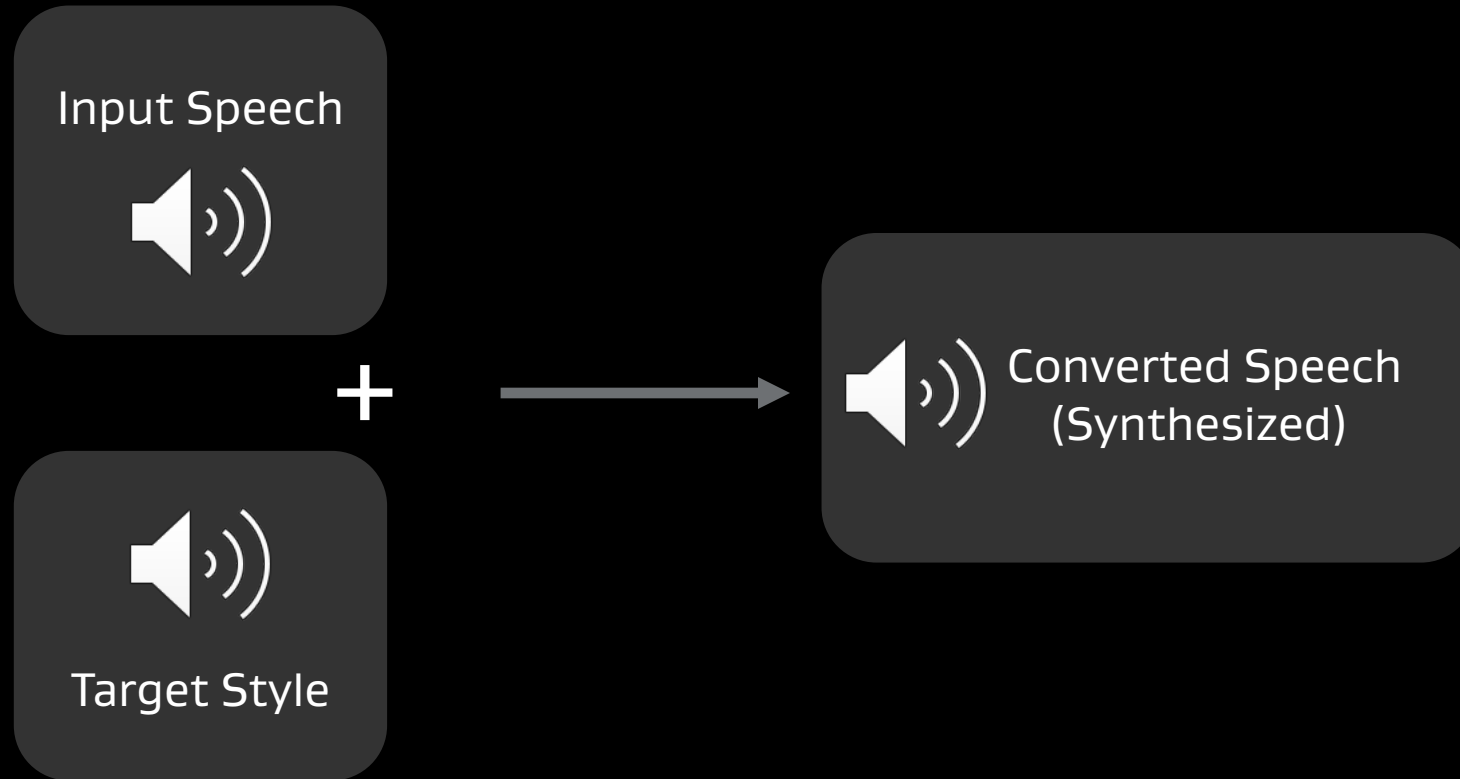High quality speech at 2.5x lower bitrate than SOTA codecs

# Voice Conversion

# Voice Conversion : Demo

# End-to-end Learning Audio Research

## Voice Conversion with Conditional SampleRNN @ Interspeech 2018
*Cong Zhou, Michael Horgan, Vivek Kumar, Cristina Vasco, Dan Darcy*

## High-quality speech coding with SampleRNN @ ICASSP 2019
*Janusz Klejsa, Per Hedelin, Cong Zhou, Roy Fejgin, Lars Villemoes*

**Jordi Pons** (jpons@dolby.com /@jordiponsdotme / www.jordipons.me)