

Waveform generation for text-to-speech synthesis using pitch-synchronous multi-scale generative adversarial networks

Lauri Juvela¹, Bajibabu Bollepalli¹, Junichi Yamagishi^{2,3}, Paavo Alku¹

¹Aalto University, Finland ²National Institute of Informatics, Japan

³The Centre for Speech Technology Research, University of Edinburgh, United Kingdom

1 Introduction

- Neural waveform generators are used in state-of-the-art TTS
- Problem: sample-by-sample synthesis with WaveNet is slow, parallel models are heavy and difficult to train
- Generative models are required to capture the stochastic components in the speech waveform
- GANs are promising, but have suffered from training instability

2 Text-to-speech synthesis system

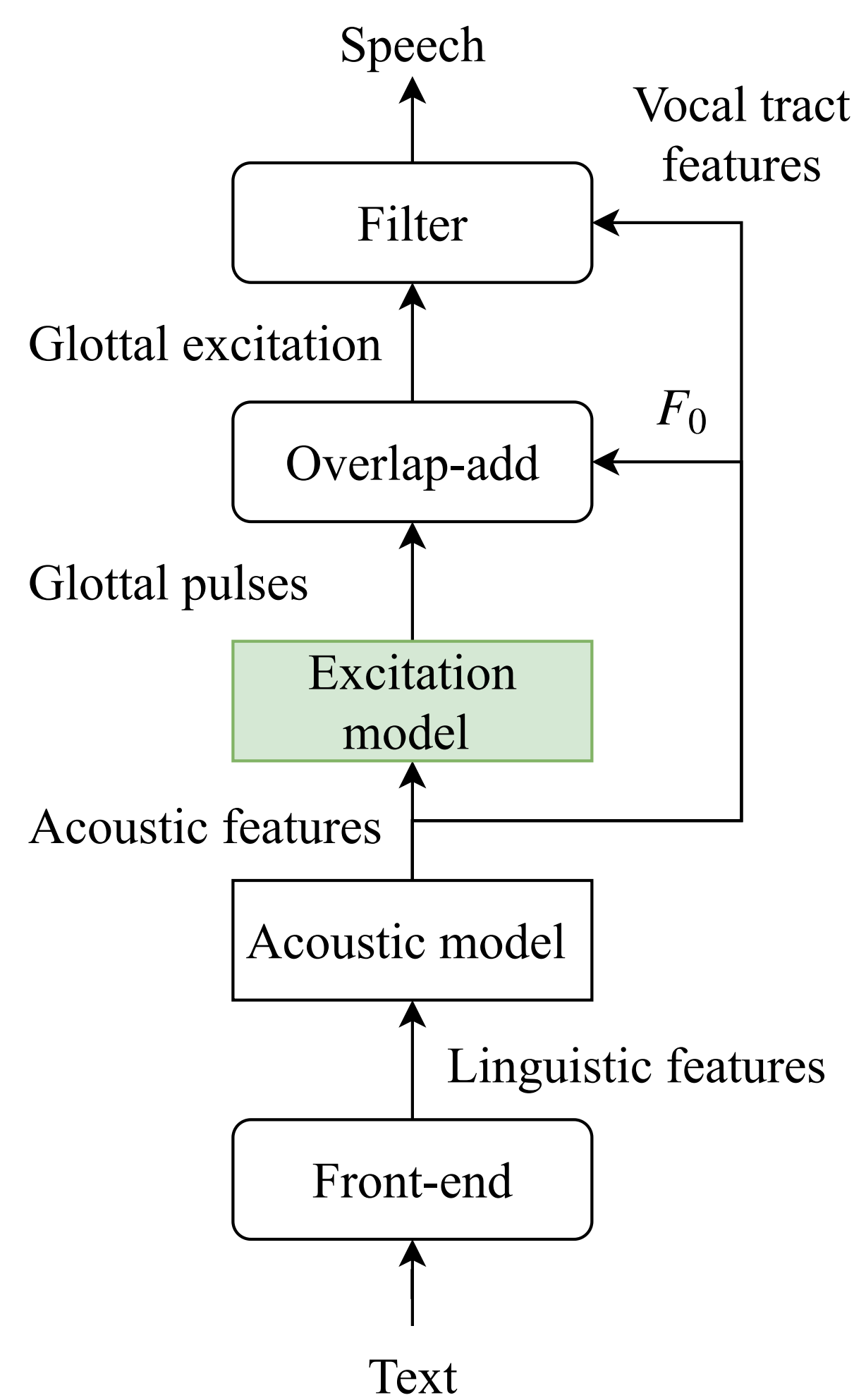
- Some speech processing can make neural vocoding more manageable:

- Use spectral envelope to remove vocal tract resonances (glottal inverse filtering)

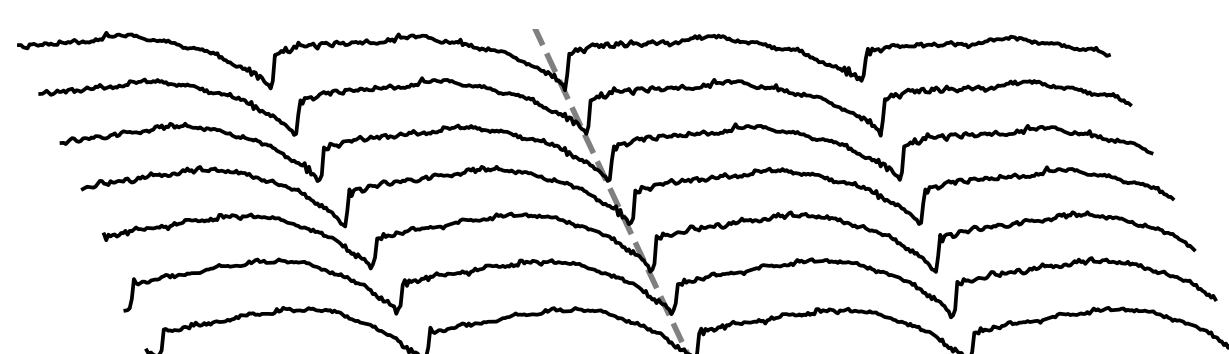
- Use pitch marking to create phase-locked waveform representations

- **Use generative neural nets as excitation model**

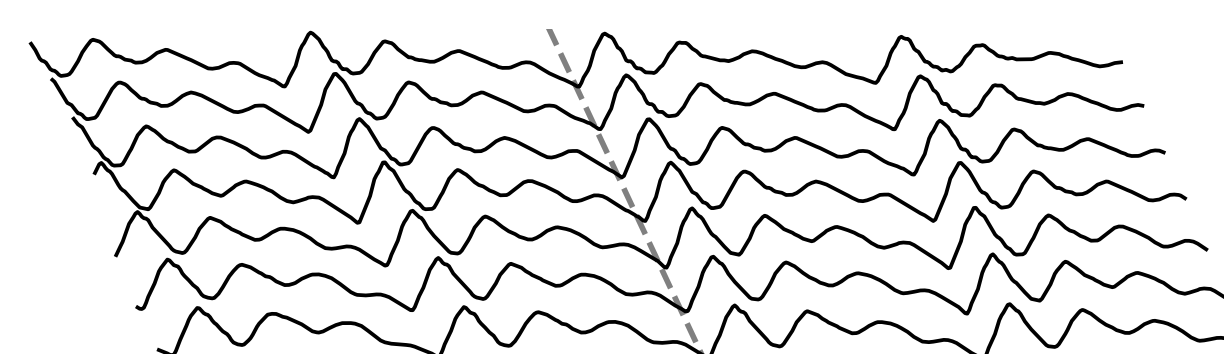
- Front-end and acoustic model are as in conventional neural net based statistical parametric speech synthesis



Waveform representation



Glottal excitation

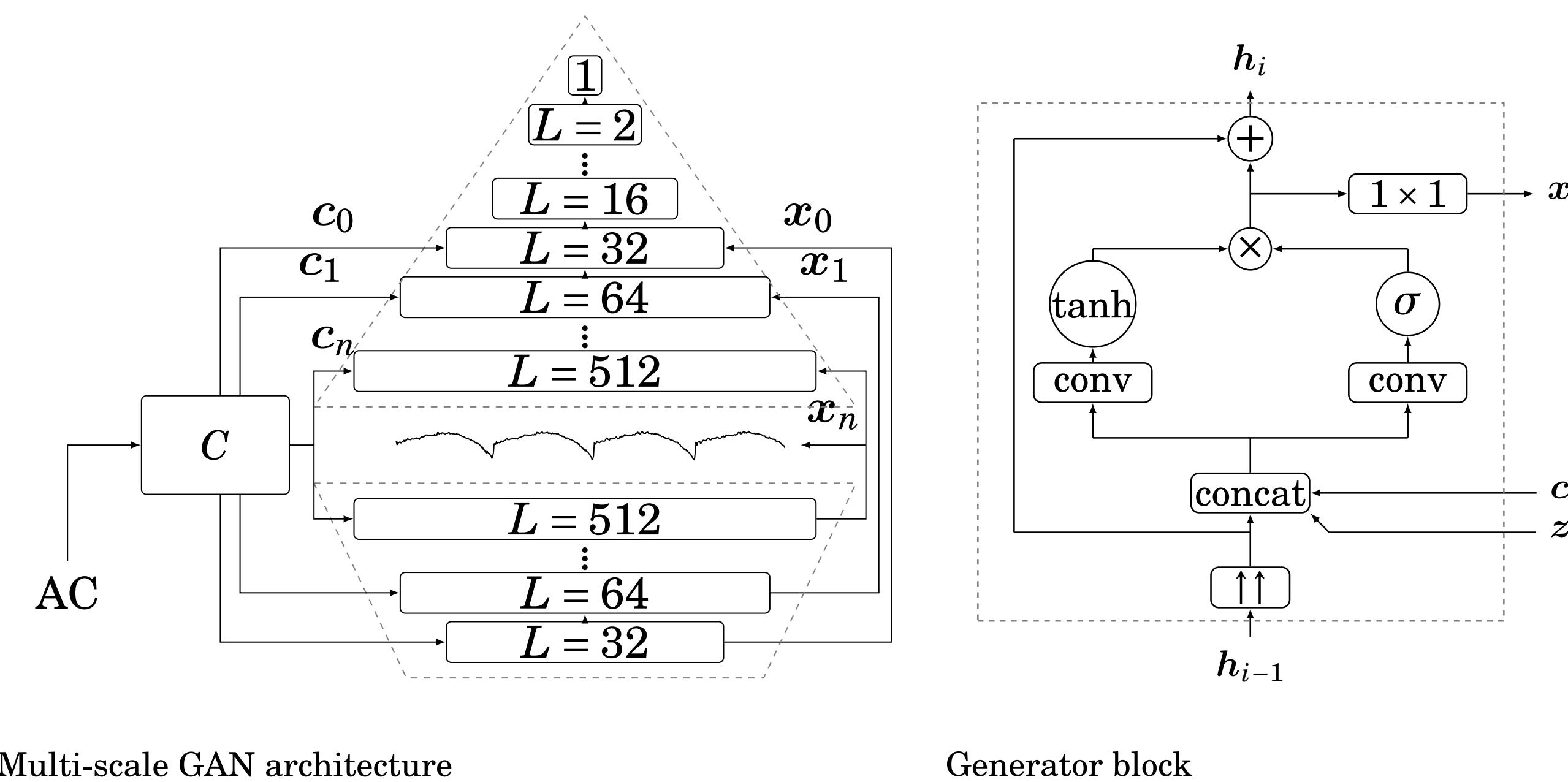


Speech signal

3 GAN waveform model

Multi-scale GAN architecture

Progressive upsampling in Generator



Training criteria

Wasserstein GAN: x is real data and $\hat{x} = G(z, c)$ at all timescales

$$\mathcal{L}_D^W = -\mathbb{E}_{x \sim p_{\mathcal{D}}} [D(x, c)] + \mathbb{E}_{\hat{x} \sim p_G} [D(\hat{x}, c)], \quad (1)$$

with gradient penalty for a smooth Discriminator

$$\mathcal{L}_D^{GP} = \mathbb{E}_{x \sim p_{\mathcal{D}}, \hat{x} \sim p_G} [(\max\{0, \|\nabla_{\hat{x}} D(\tilde{x}, c)\| - 1\})^2], \quad (2)$$

and regularization for large gradients in the data manifold

$$\mathcal{L}_D^{R1} = \mathbb{E}_{x \sim p_{\mathcal{D}}} [\|\nabla_x D(x, c)\|^2]. \quad (3)$$

Spectral magnitude matching for the Generator

$$\mathcal{L}_G^{\text{FFT}} = \mathbb{E} [(|\text{FFT}(x_n)| - |\text{FFT}(\hat{x}_n)|)^2]. \quad (4)$$

Sound samples

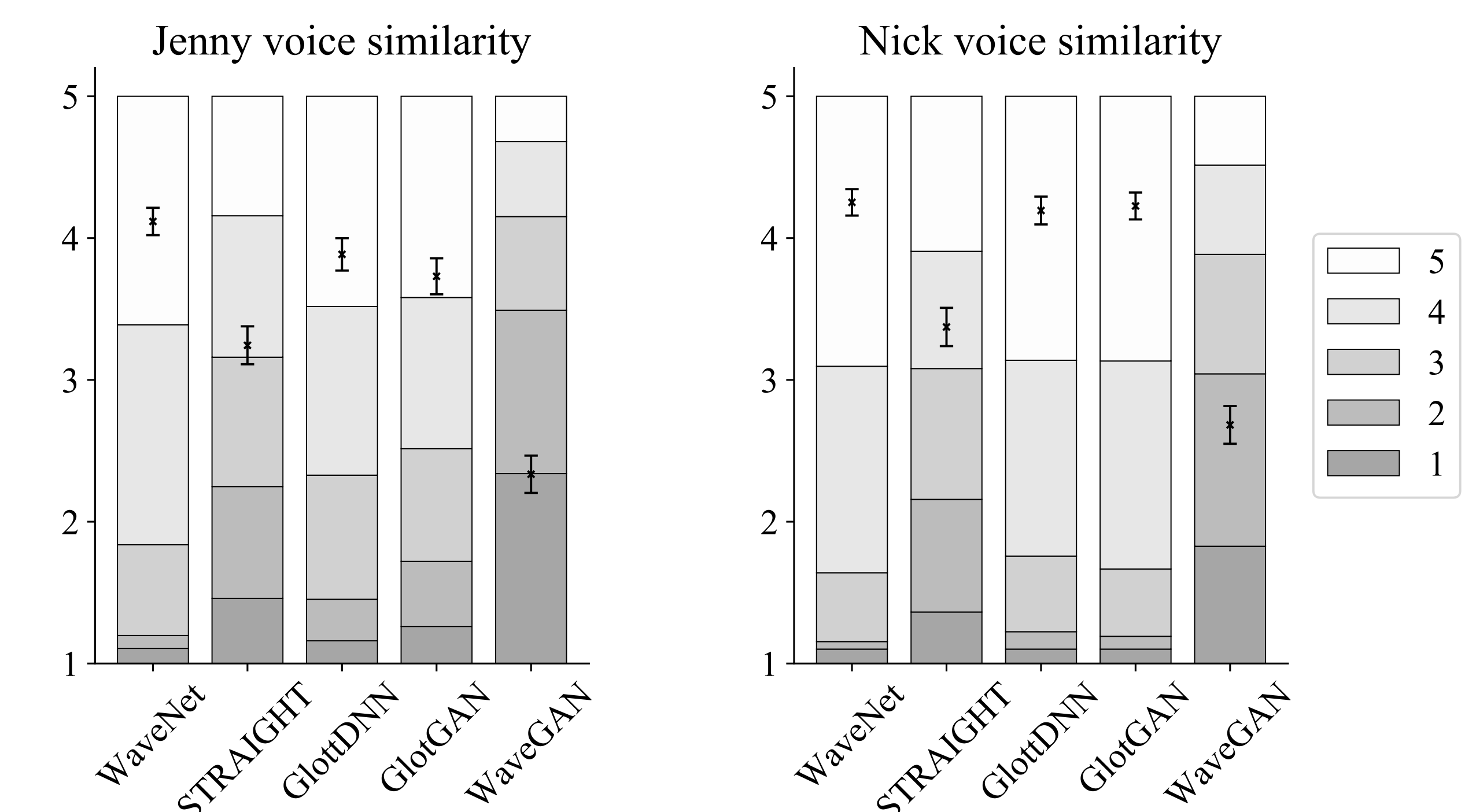
<https://users.aalto.fi/~ljuvela/multiscale-gan/>

Source code (TensorFlow)

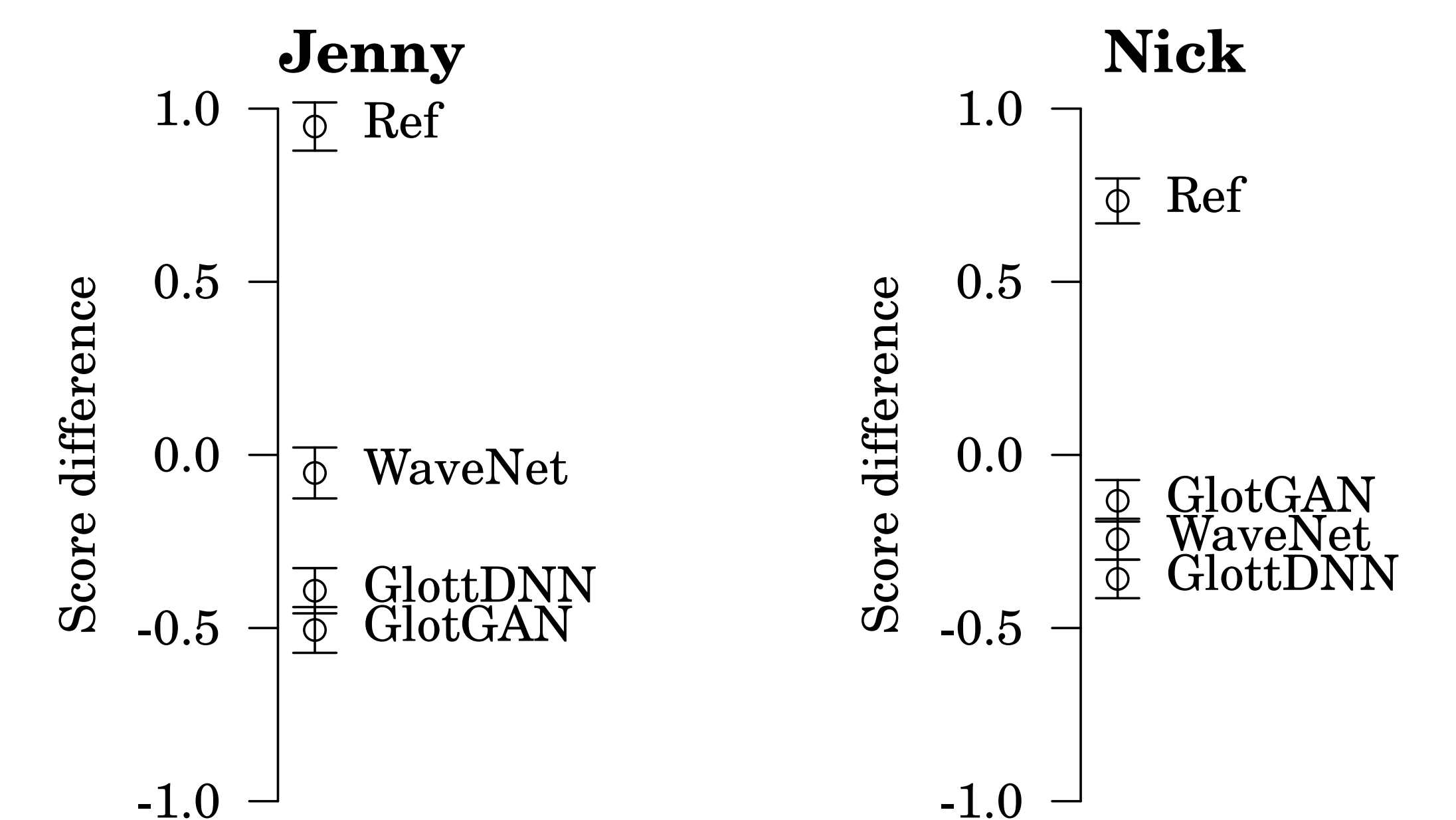
<https://github.com/ljuvela/multiscale-GAN>

4 Listening tests

- TTS systems trained on Jenny (4.7h) and Nick (1.8h) datasets
- Difference mean opinion score (DMOS) test for speaker similarity



- Category comparison rating (CCR) test to evaluate quality



- **The proposed method can match a WaveNet vocoder in TTS quality**