



Audio samples



<https://www.kecl.ntt.co.jp/people/kaneko.takuhiro/projects/waveunetd/>

# *Wave-U-Net Discriminator:*

Fast and Lightweight Discriminator for  
Generative Adversarial Network-Based Speech Synthesis



Takuhiro Kaneko



Hirokazu Kameoka



Kou Tanaka

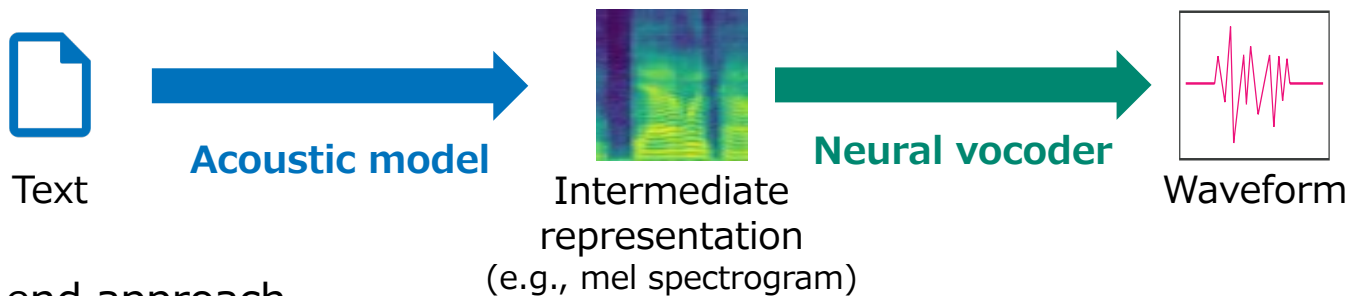


Shogo Seki

ICASSP 2023

## Advancement of speech synthesis

- Two-stage approach



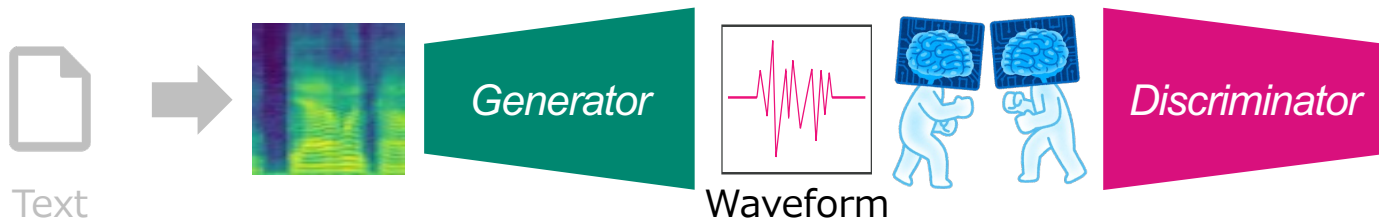
- End-to-end approach



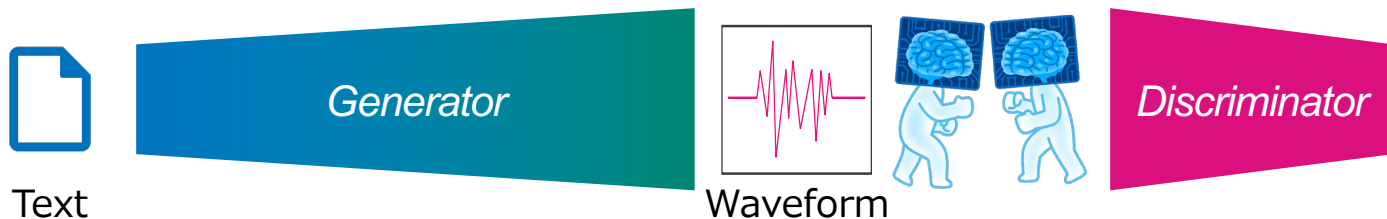
**Common objective: High-quality speech synthesis**

## GAN [Goodfellow+2014]-based speech synthesis

- Two-stage approach (e.g., HiFi-GAN [Kong+2020])



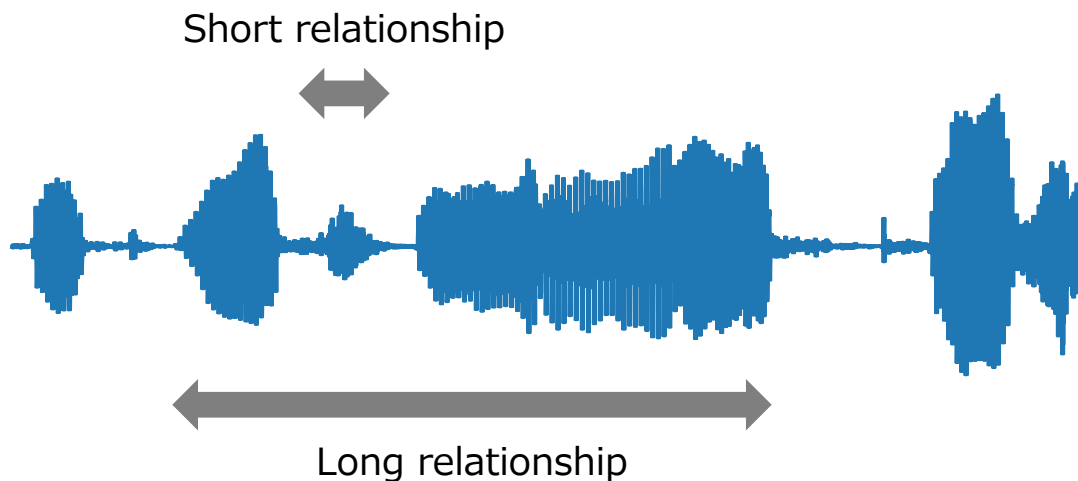
- End-to-end approach (e.g., VITS [Kim+2021])



**Challenge: How to design an adequate discriminator?**

# Challenge of GAN-based approach

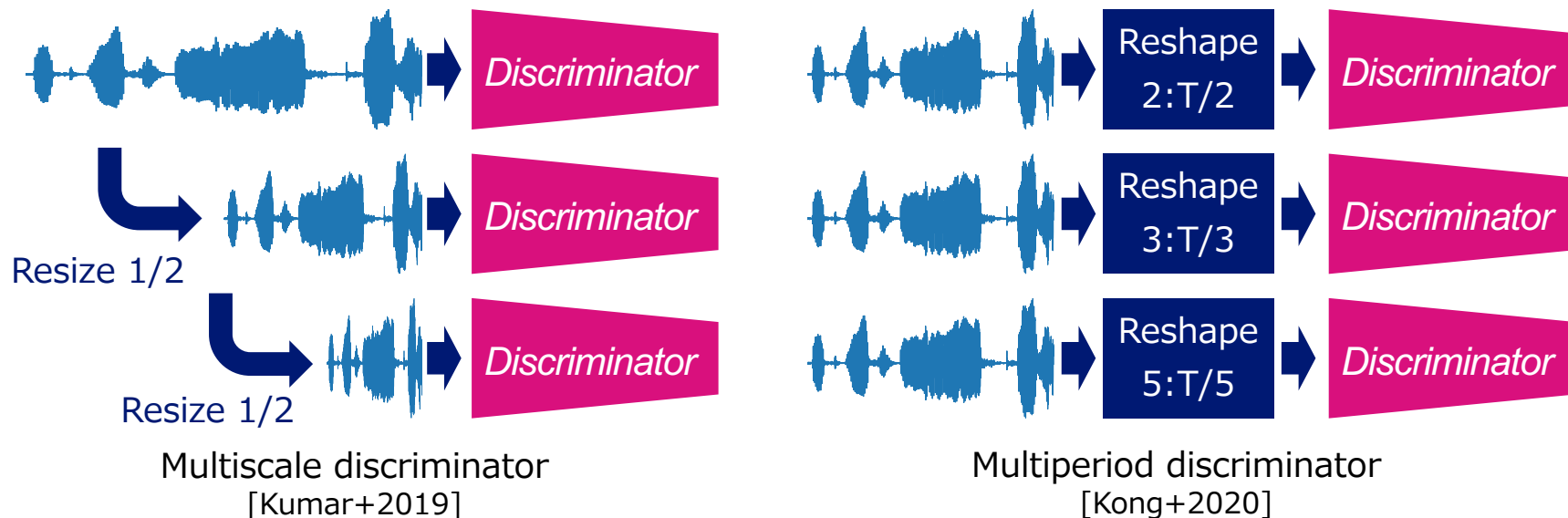
Speech has multilevel (e.g., multiscale) structures



**Discriminator must capture multilevel structures**

# Previous study

An ensemble of discriminators was used



- 😊 Multilevel structures can be captured
- 😞 Model size & computation time increase according to #discriminators

# Previous study

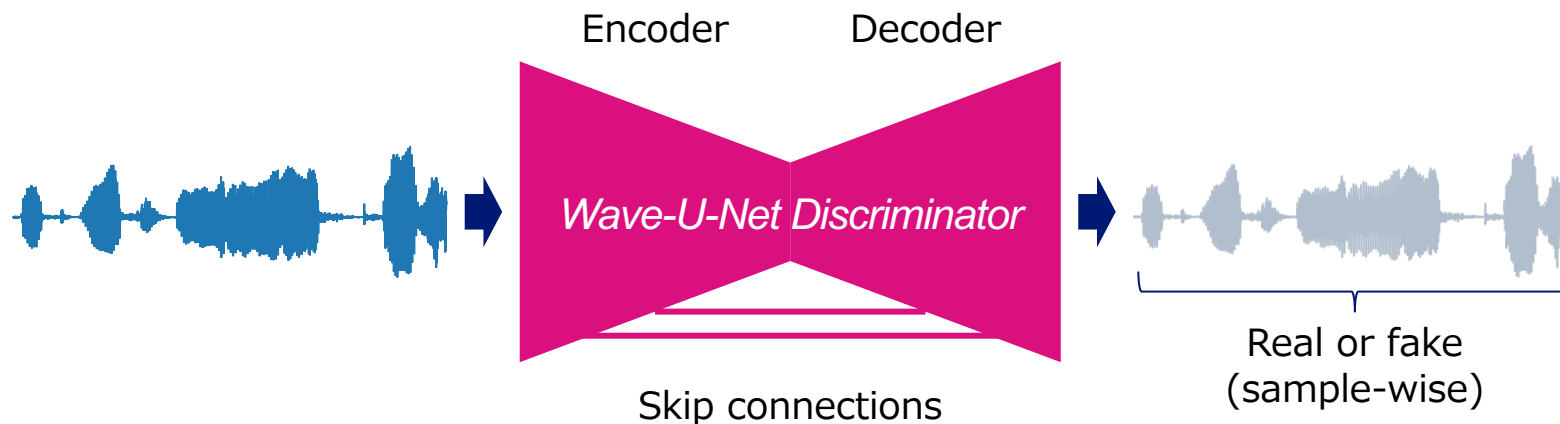
An ensemble of discriminators was used



- 😊 Multilevel structures can be captured
- 😞 Model size & computation time increase according to #discriminators

# Our solution

## Wave-U-Net Discriminator



Multilevel structures can be captured via an encoder and decoder



#discriminators is one → Fast and lightweight

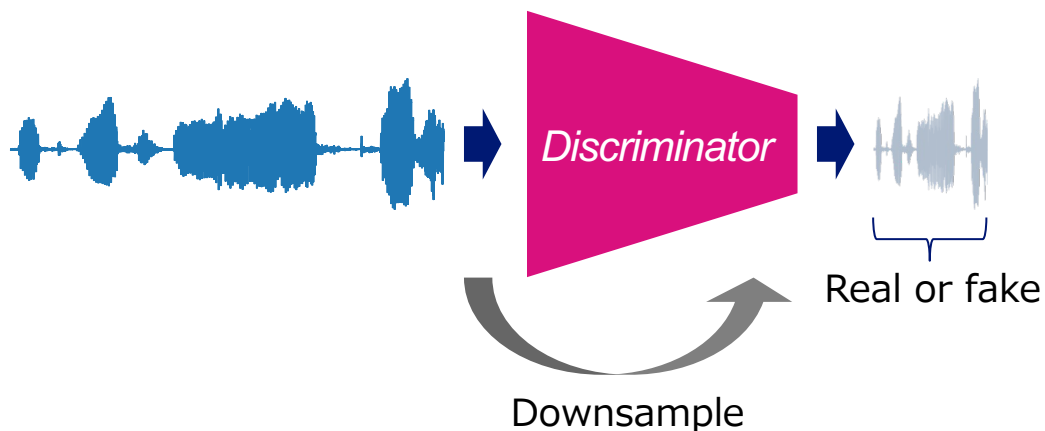


# Method



# Previous discriminator

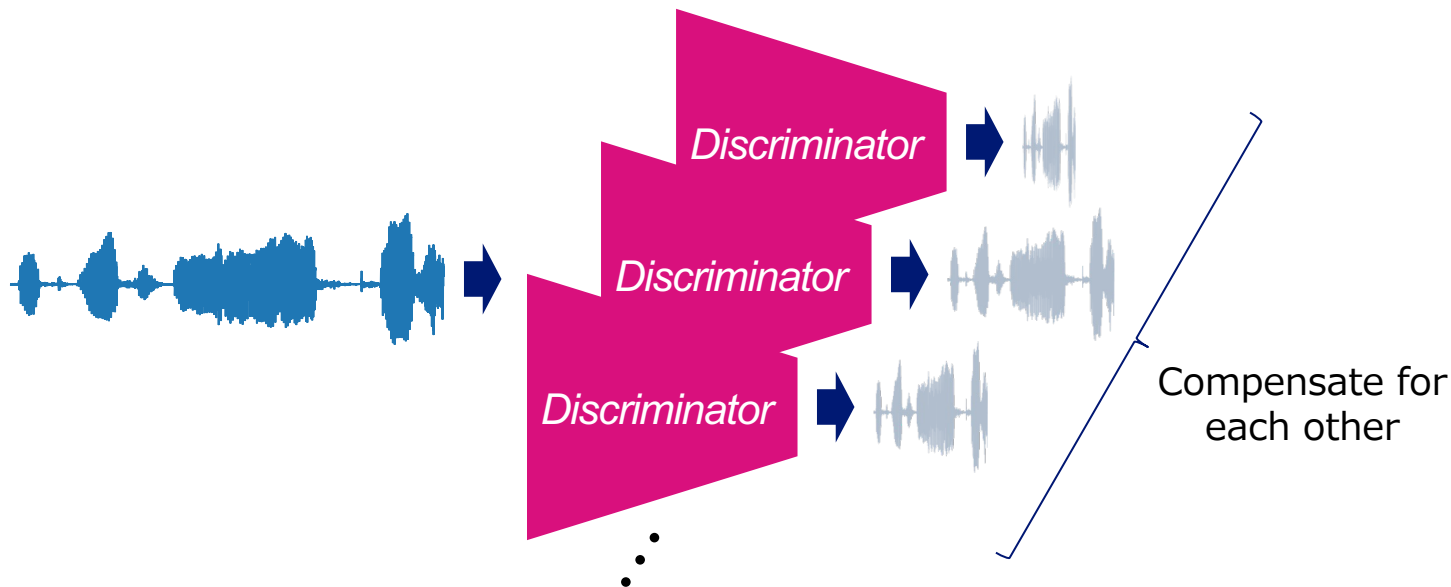
**Encoder architecture** (e.g., MelGAN [Kumar+2019], HiFi-GAN [Kong+2020])



**Real/fake is determined using the abstracted features**

# Previous discriminator

**Encoder architecture** (e.g., MelGAN [Kumar+2019], HiFi-GAN [Kong+2020])

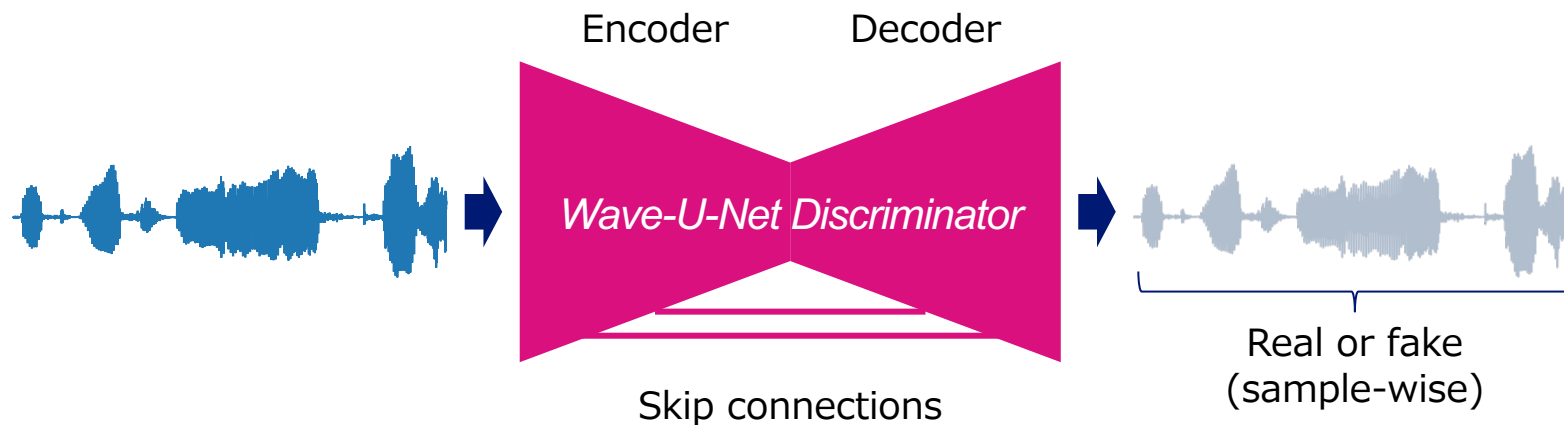


**Real/fake is determined using the abstracted features**

**→ Multiple discriminators are required to capture detailed structures**

# Wave-U-Net discriminator

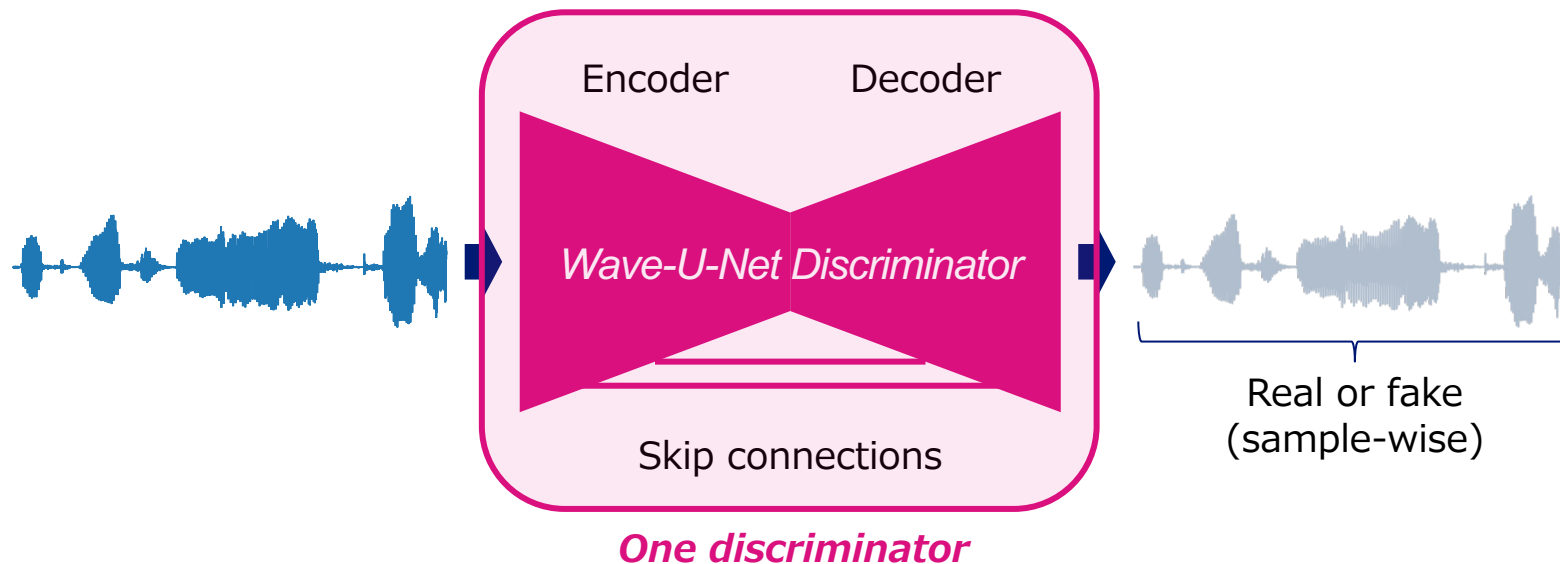
## Encoder-decoder architecture



Real/fake is determined in a sample-wise manner

# Wave-U-Net discriminator

## Encoder-decoder architecture

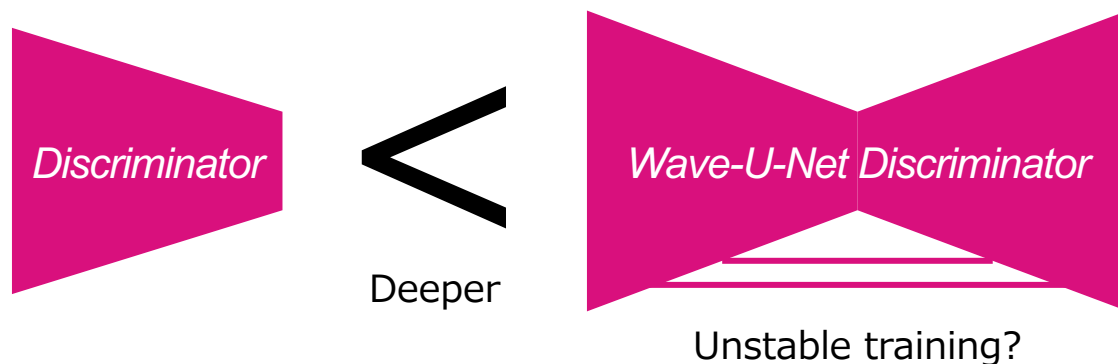


Real/fake is determined in a sample-wise manner

→ One discriminator is sufficient to capture detailed structures

# Challenge in training

## Unstable training of Wave-U-Net discriminator



**Wave-U-Net discriminator is deeper than typical discriminator  
→ Causes unstable training (saturate adversarial losses)**

# Techniques for stable training 1

## Careful normalization

- Global normalization

$$b = a / \sqrt{\frac{1}{N} \sum_{i=1}^N (a^i)^2 + \epsilon}$$

Normalized features  $b$  = Original features  $a$

$N$ : Number of features

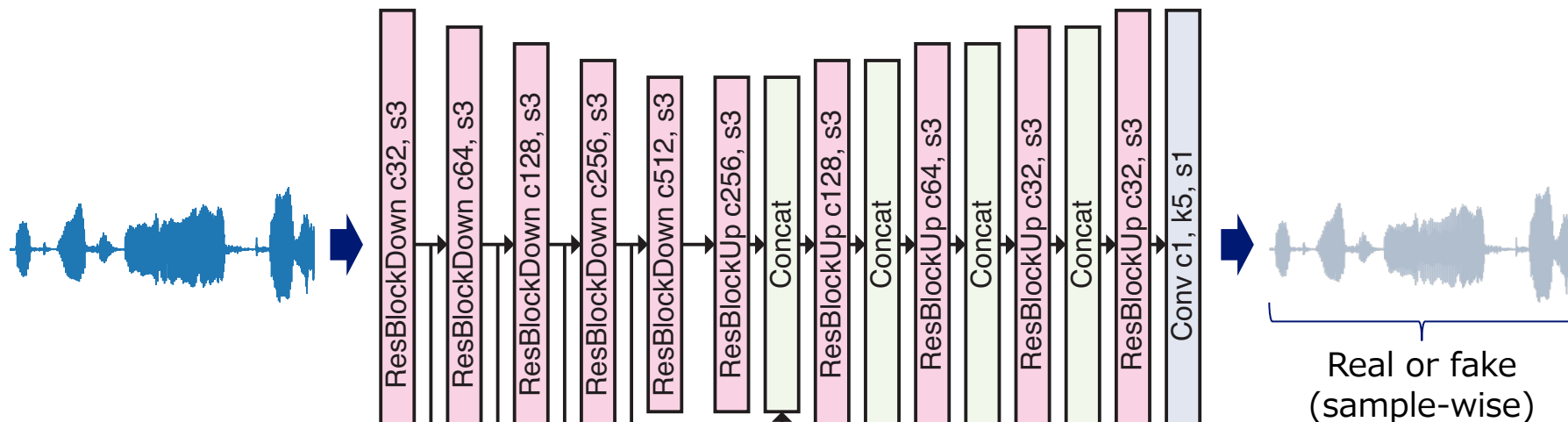
$a^i$ :  $i$ -th feature in  $a$

$\epsilon = 10^{-8}$

Prevents Wave-U-Net Discriminator from restricting itself to specific features

# Techniques for stable training 2

## Introduction of residual connections [He+2016]



Prevents the gradient vanishing problem



# Experiments



## 1. Evaluation on neural vocoders

*Dataset dependency was investigated*

- **Datasets:**
  - › **LJSpeech** [Ito&Johnson2017]: Single English female speaker
  - › **VCTK** [Yamagishi+2016]: Multiple English speakers
  - › **JSUT** [Sonobe+2017]: Single Japanese female speaker
- **Baseline: HiFi-GAN** [Kong+2020]

## 2. Evaluation on *end-to-end TTS*

*Task dependency was investigated*

- **Datasets: LJSpeech** [Ito&Johnson2017]
- **Baseline: VITS** [Kim+2021]

Performance was examined when the original ensemble of discriminators was replaced with a Wave-U-Net discriminator

## Speech quality

- **Subjective metric: MOS** ↑
  - › Mean opinion score on naturalness
- **Objective metric: cFW2VD** ↓ [Kaneko+2022]
  - › Distance between real and synthesized speech in wav2vec 2.0 [Baevski+2020]

## Training speed

- **Time (s/batch)** ↓
  - › Time required for a discriminator to process real and synthesized speech in a batch

## Model size

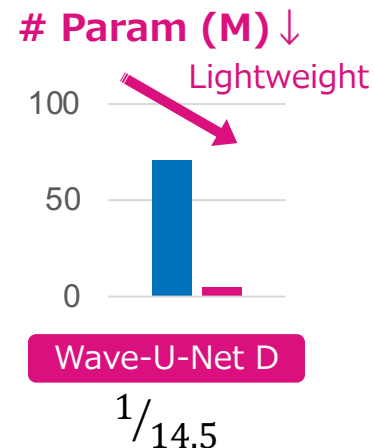
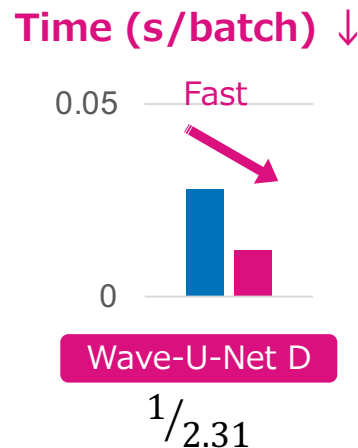
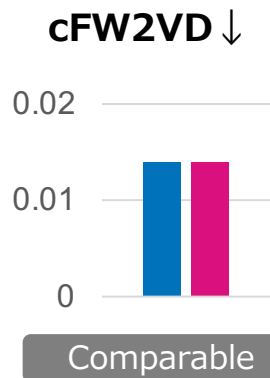
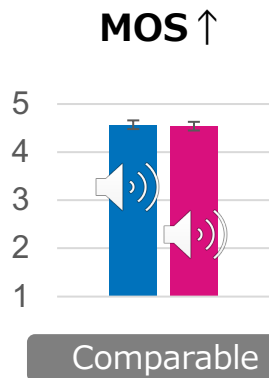
- **# Param (M)** ↓
  - › Number of parameters of a discriminator



# Evaluation on neural vocoders

# Results 1/3

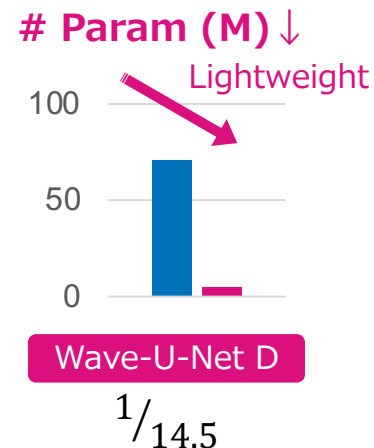
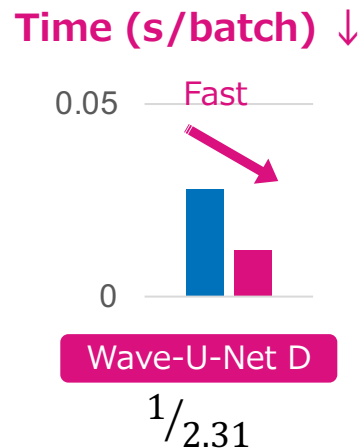
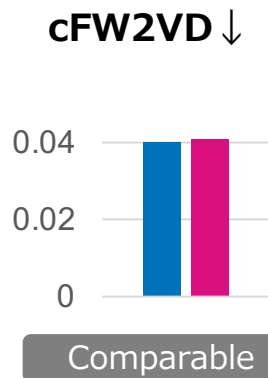
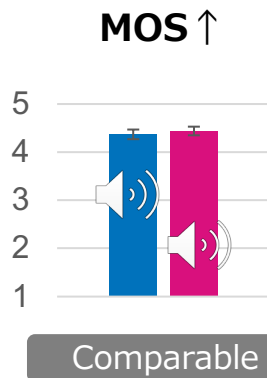
## Evaluation on neural vocoder in LJSpeech



**Wave-U-Net discriminator reduces computation time & model size while retaining speech quality**

# Results 2/3

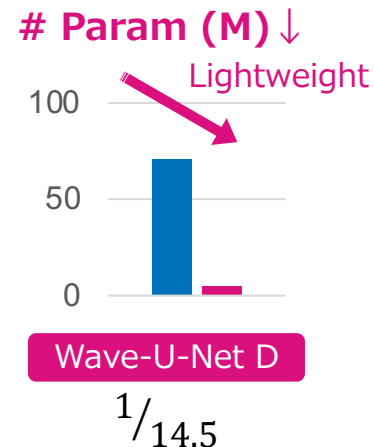
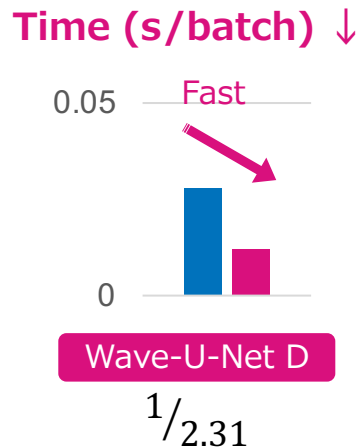
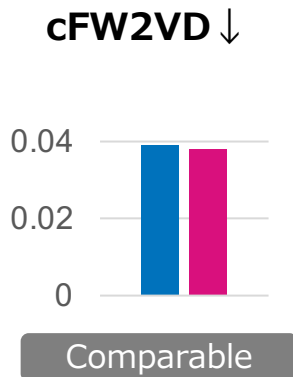
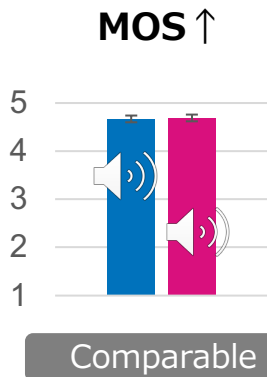
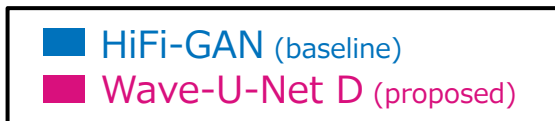
## Evaluation on neural vocoder in VCTK



**Wave-U-Net discriminator reduces computation time & model size while retaining speech quality**

# Results 3/3

## Evaluation on neural vocoder in JSUT

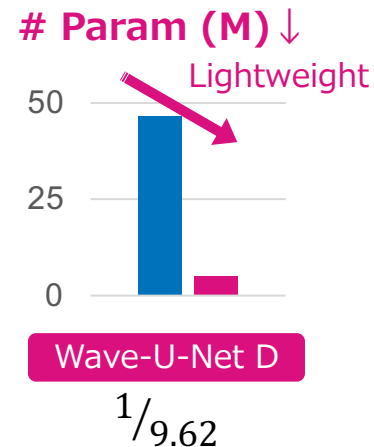
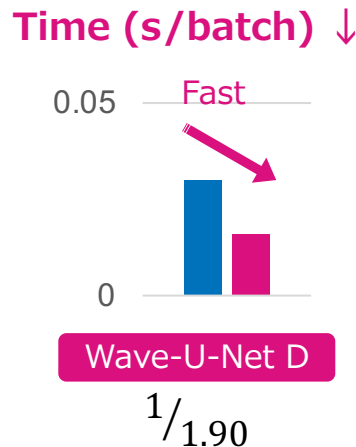
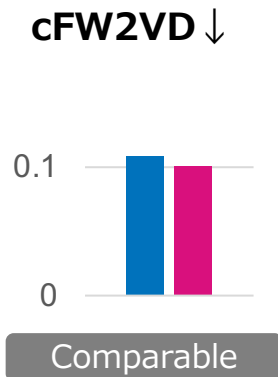
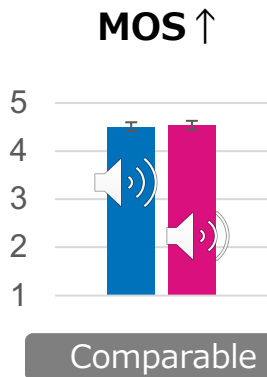
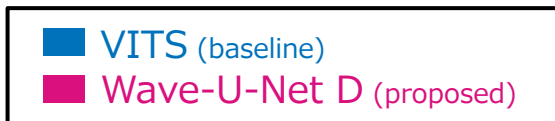


**Wave-U-Net discriminator reduces computation time & model size while retaining speech quality**



# Evaluation on end-to-end TTS

## Evaluation on end-to-end TTS



**Wave-U-Net discriminator reduces computation time & model size while retaining speech quality**





# Conclusion

# Conclusion

## Objective

- Make a discriminator **faster & more lightweight**

## Proposal

- *Wave-U-Net Discriminator*

## Experiments

- Make a discriminator **faster & more lightweight** while retaining speech quality

## Future work

- Application to **other tasks**
  - › Singing speech synthesis, emotional speech synthesis, ...

### Wave-U-Net Discriminator



### Audio samples



<https://www.kecl.ntt.co.jp/people/kaneko.takuhiro/projects/waveunetd/>