# TEXT-TO-SPEECH SYNTHESIS USING STFT SPECTRA BASED ON LOW-/MULTI-RESOLUTION GENERATIVE ADVERSARIAL NETWORKS

Yuki Saito, **Shinnosuke Takamichi**, and Hiroshi Saruwatari

The University of Tokyo, Japan

## 1. SYNOPSIS

In Text-To-Speech (TTS) synthesis w/ Short-Term Fourier Transform (STFT) spectra:

(1) Proposes Generative Adversarial Networks (GANs)[1]-based training algorithm using low-/multi-resolution spectra.

(2) Demonstrates that the proposed algorithm using GANs w/ low-resolution spectra improves synthetic speech quality.

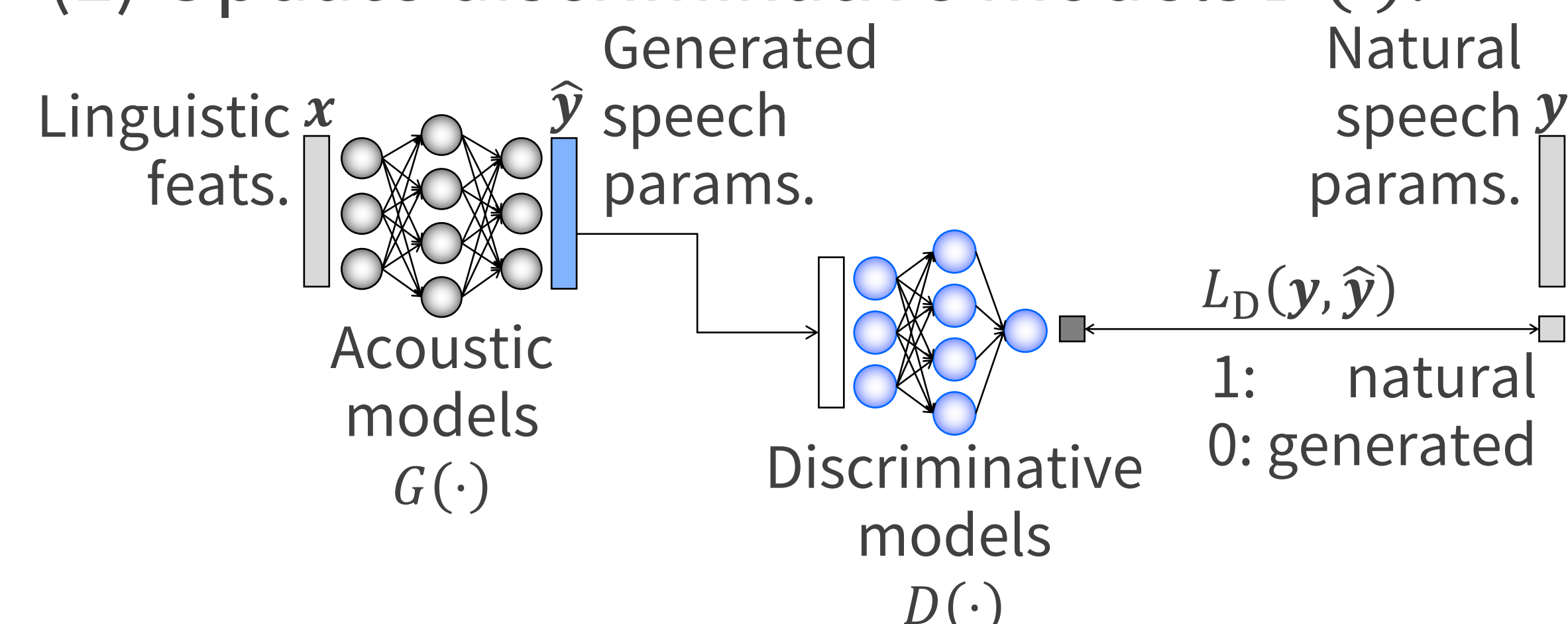## 2. CONVENTIONAL ALGORITHMS

### DNN-based TTS w/ STFT spectra[2]

(1) Generate raw STFT amplitude spectra.

(2) Reconstruct phase spectra by using Griffin and Lim's method[3].

Pros. avoiding a vocoding process

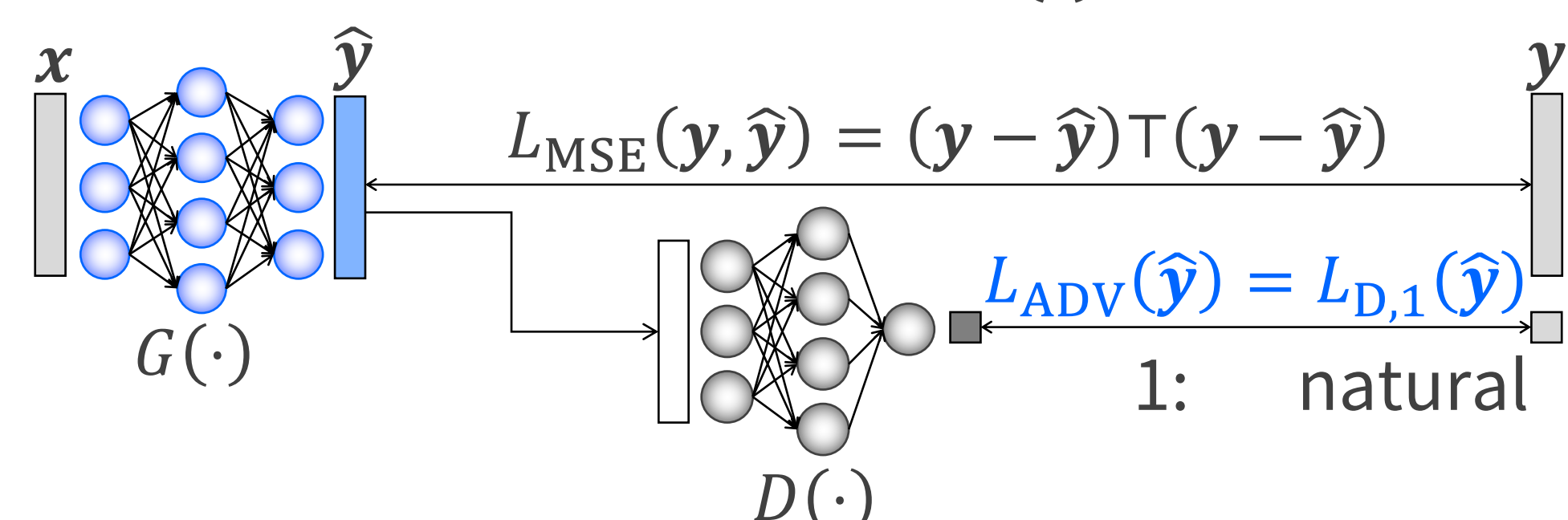Cons. over-smoothing of amplitude spectra

### GAN-based TTS w/ vocoders[4]

(1) Update discriminative models $D(\cdot)$.



Linguistic feats. $x$ → Acoustic models $G(\cdot)$ → Generated speech params. $\widehat{y}$ → Discriminative models $D(\cdot)$

Natural speech params. $y$

$L_D(y, \widehat{y})$

1: natural
0: generated

$$L_D(y, \widehat{y}) = \underbrace{-\sum_t \log D(y_t)}_{\substack{L_{D,1}(y) \\ \text{(loss for natural params.)}}} \underbrace{- \sum_t \log(1 - D(\widehat{y}_t))}_{\substack{L_{D,0}(\widehat{y}) \\ \text{(loss for generated params.)}}}$$

(2) Update acoustic models $G(\cdot)$.



$x$ → $G(\cdot)$ → $\widehat{y}$

$L_{MSE}(y, \widehat{y}) = (y - \widehat{y})^{\top}(y - \widehat{y})$

$y$

$D(\cdot)$

$L_{ADV}(\widehat{y}) = L_{D,1}(\widehat{y})$

1: natural

$$L_G(y, \widehat{y}) = L_{MSE}(y, \widehat{y}) + \omega_D \frac{\mathbb{E}_{\widehat{y}}[L_{MSE}]}{\mathbb{E}_{\widehat{y}}[L_{ADV}]} L_{ADV}(\widehat{y})$$

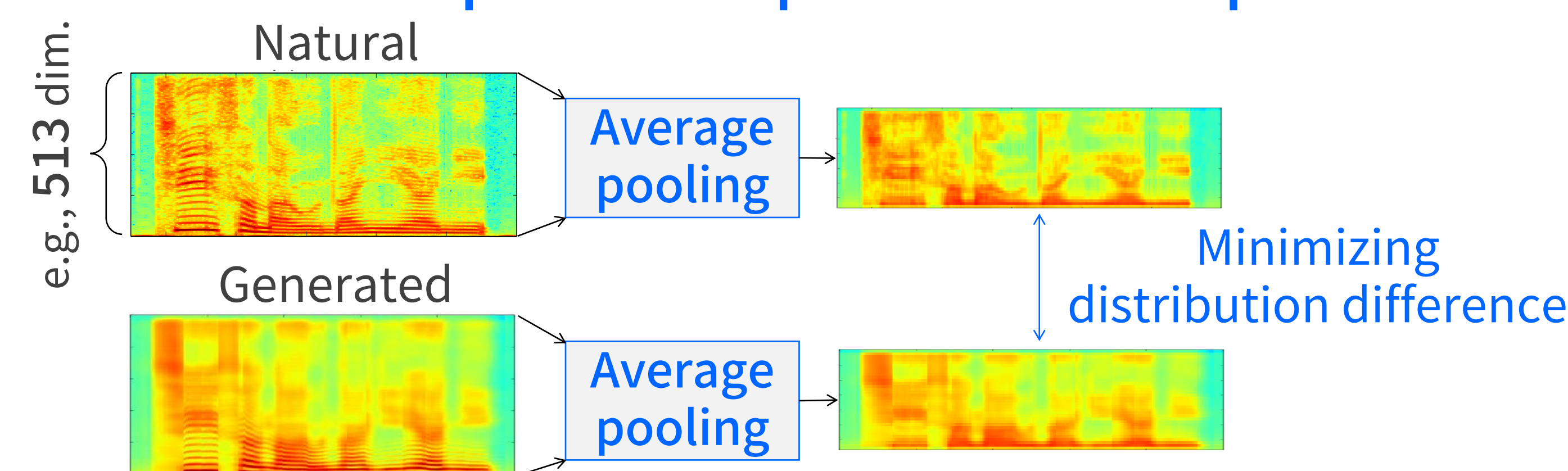Minimizing approx. JS-divergence
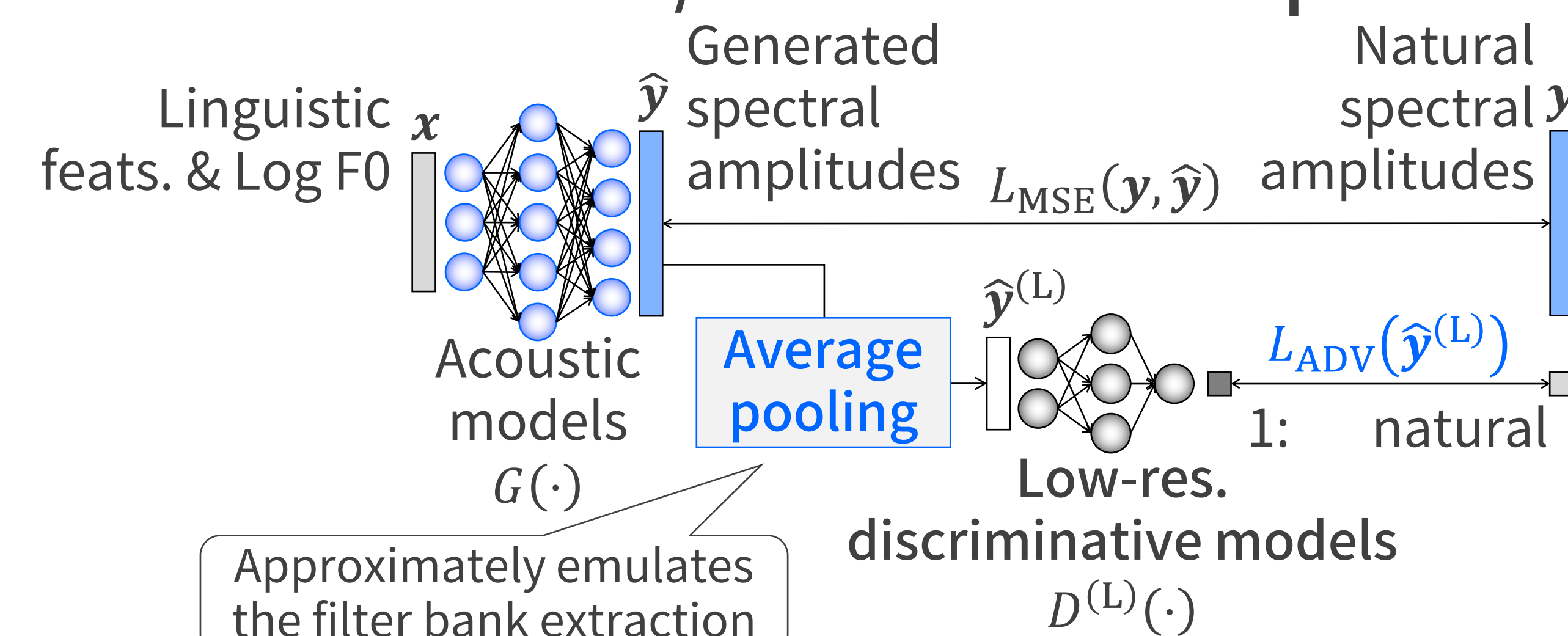
## 3. PROPOSED ALGORITHMS

### Motivation

Amplitude spectra: high-dimensional features including spectral & excitation characteristics.

GAN-based TTS[3]: effective for spectral envelopes (i.e., mel-cepstral coefficients).

Assumption:

low-res. spectra ≒ spectral envelopes



Natural → Average pooling

Generated → Average pooling

Minimizing distribution difference

e.g., 513 dim.

### GAN-based TTS w/ low-res. STFT spectra



Linguistic feats. & Log F0 $x$ → Acoustic models $G(\cdot)$ → Generated spectral amplitudes $\widehat{y}$

$L_{MSE}(y, \widehat{y})$

Natural spectral $y$ amplitudes

Average pooling → $\widehat{y}^{(L)}$ → Low-res. discriminative models $D^{(L)}(\cdot)$

$L_{ADV}(\widehat{y}^{(L)})$

1: natural

Approximately emulates the filter bank extraction

$$L_G^{(Low)}(y, \widehat{y}) = L_{MSE}(y, \widehat{y}) + \omega_D^{(L)} \frac{\mathbb{E}_{\widehat{y}}[L_{MSE}]}{\mathbb{E}_{\widehat{y}^{(L)}}[L_{ADV}]} L_{ADV}(\widehat{y}^{(L)})$$

### GAN-based TTS w/ multi-res. STFT spectra

$$L_G^{(Multi)}(y, \widehat{y}) = L_{MSE}(y, \widehat{y}) + \omega_D^{(L)} \frac{\mathbb{E}_{\widehat{y}}[L_{MSE}]}{\mathbb{E}_{\widehat{y}^{(L)}}[L_{ADV}]} L_{ADV}(\widehat{y}^{(L)}) + \omega_D \frac{\mathbb{E}_{\widehat{y}}[L_{MSE}]}{\mathbb{E}_{\widehat{y}}[L_{ADV}]} L_{ADV}(\widehat{y})$$

## 4. EXPERIMENTAL EVALUATION

### Experimental conditions

| | |
|---|---|
| Dataset | 4,007 utterances taken from Japanese female speaker (subset of JSUT[5] corpus, 3,808/199 for training/evaluation) |
| STFT analysis | Frame length: 400, shift length: 80, FFT length: 1024, analysis window: Hamming |
| Average pooling | Zero padding: 6, pooling width $w$: {14, 30, 70}, stride: $w/2$ |
| $\omega_D$ and $\omega_D^{(L)}$ | 1.0 |
| Dims. of $x/y/y^{(L)}$ | 444 (linguistic feats, durations, Log F0, UV)/513/{74, 34, 14} (dim. of $y^{(L)}$ was changed in accordance with $w$) |
| DNN architectures | Feed-Forward (see our paper) |

### Subjective evaluation (preference AB tests)

MSE: minimizing $L_{MSE}(y, \widehat{y})$[2]

ADV: minimizing $L_G(y, \widehat{y})$[4]

ADV-Low: minimizing $L_G^{(Low)}(y, \widehat{y})$ (proposed)

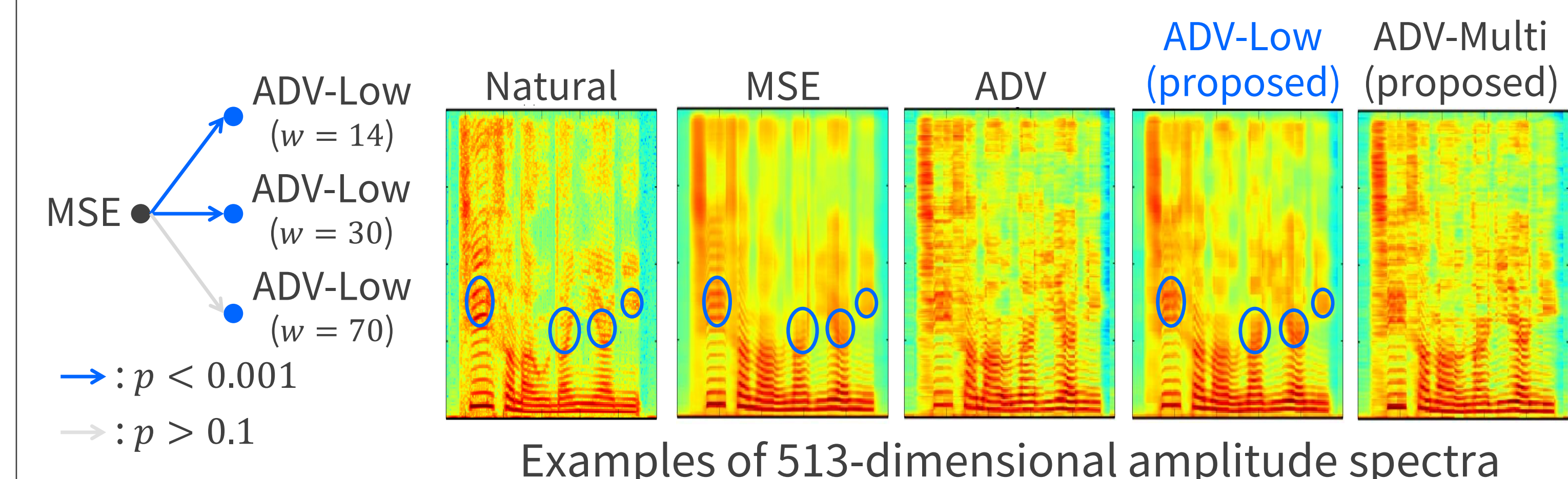ADV-Multi: minimizing $L_G^{(Multi)}(y, \widehat{y})$ (proposed)

Results:

(1) ADV & ADV-Multi didn't improve speech quality.
⇒ Because of difficulty in minimizing distribution differences in original resolution

(2) ADV-Low improved speech quality.
⇒ Effect of the spectral envelopes compensation



MSE → ADV-Low ($w = 14$), ADV-Low ($w = 30$), ADV-Low ($w = 70$)

→ : $p < 0.001$
→ : $p > 0.1$

Natural | MSE | ADV | ADV-Low (proposed) | ADV-Multi (proposed)

Examples of 513-dimensional amplitude spectra

[References]
[1] Goodfellow et al., *Proc. NIPS,* 2014. [2] Takaki et al., *Proc. INTERSPEECH*, 2017. [3] Griffin et al., *IEEE Trans. on ASLP*, 1984. [4] Saito et al., *IEEE/ACM Trans. on ASLP*, 2018. [5] Sonobe et al., *arXiv:1711.00354,* 2017.