

Audio Synthesis and Bandwidth Extension



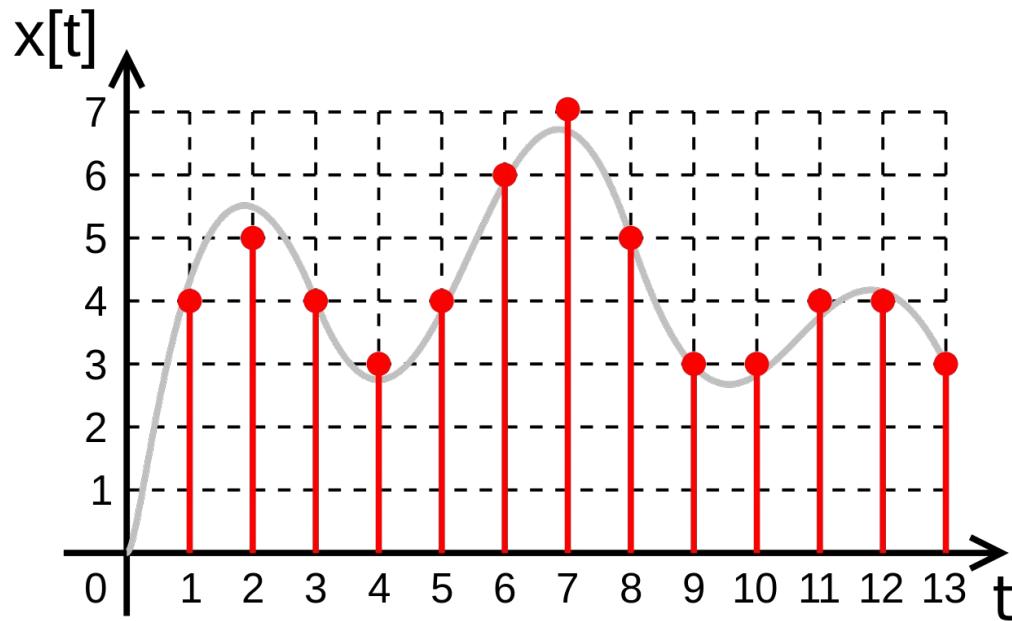
Audio Processing Problems

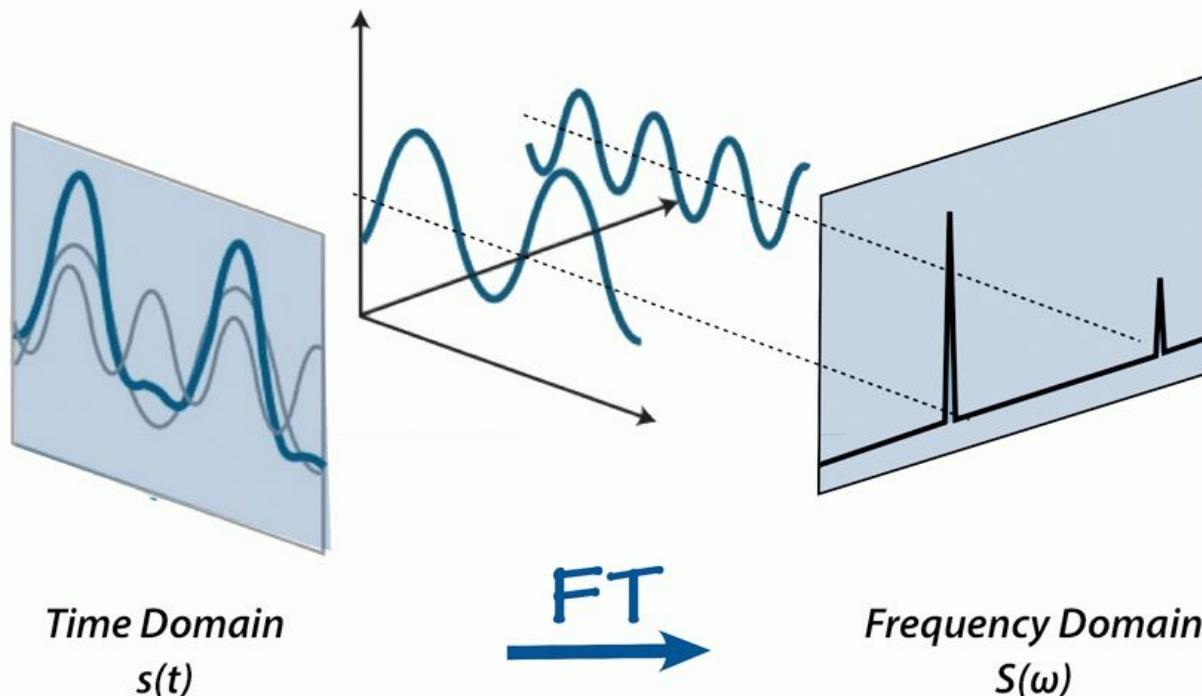
Samsung AI Center - Moscow

- Text To Speech
- Automatic Speech Recognition
- Acoustic Scene Recognition
 - Source Separation
 - Source Localization
- Audio Enhancement
 - Denoising
 - Bandwidth extension
 - Echo cancellation
- ...



Introduction to Signal Processing





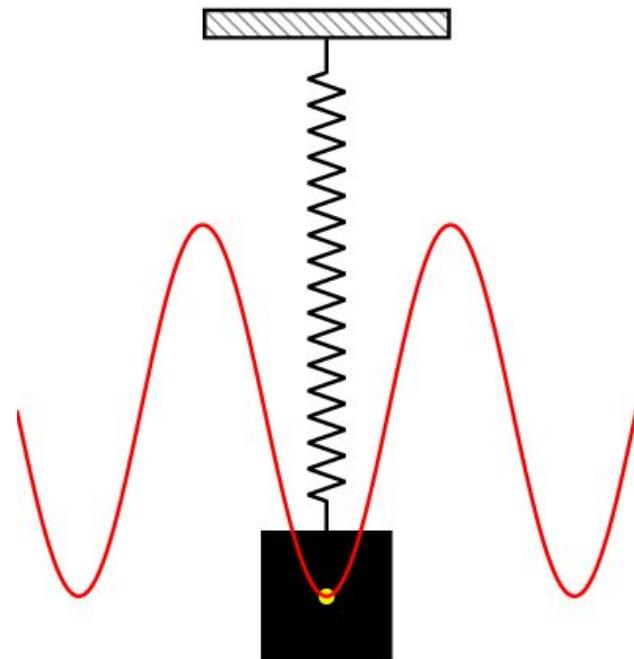
$$\frac{\partial^2 x(t)}{\partial t^2} = -\frac{k}{m}x(t)$$

$$x(t) = A \cos(\phi + \omega t)$$

amplitude phase (cycle) frequency

Fourier Transform

$$x(t) = \sum_k A_k \cos(\phi_k + \omega_k t)$$



$$x(t) = \sum_k A_k \cos(\phi_k + \omega_k t)$$

$$e^{i\omega t} = \cos(\omega t) + i \sin(\omega t)$$

$$x(t) = \sum_k \overbrace{A_k e^{i\phi_k}}^{F(\omega_k)} e^{i\omega_k t}$$

Fourier coefficient

$$x(t) = \sum_k F(\omega_k) e^{i\omega_k t}$$

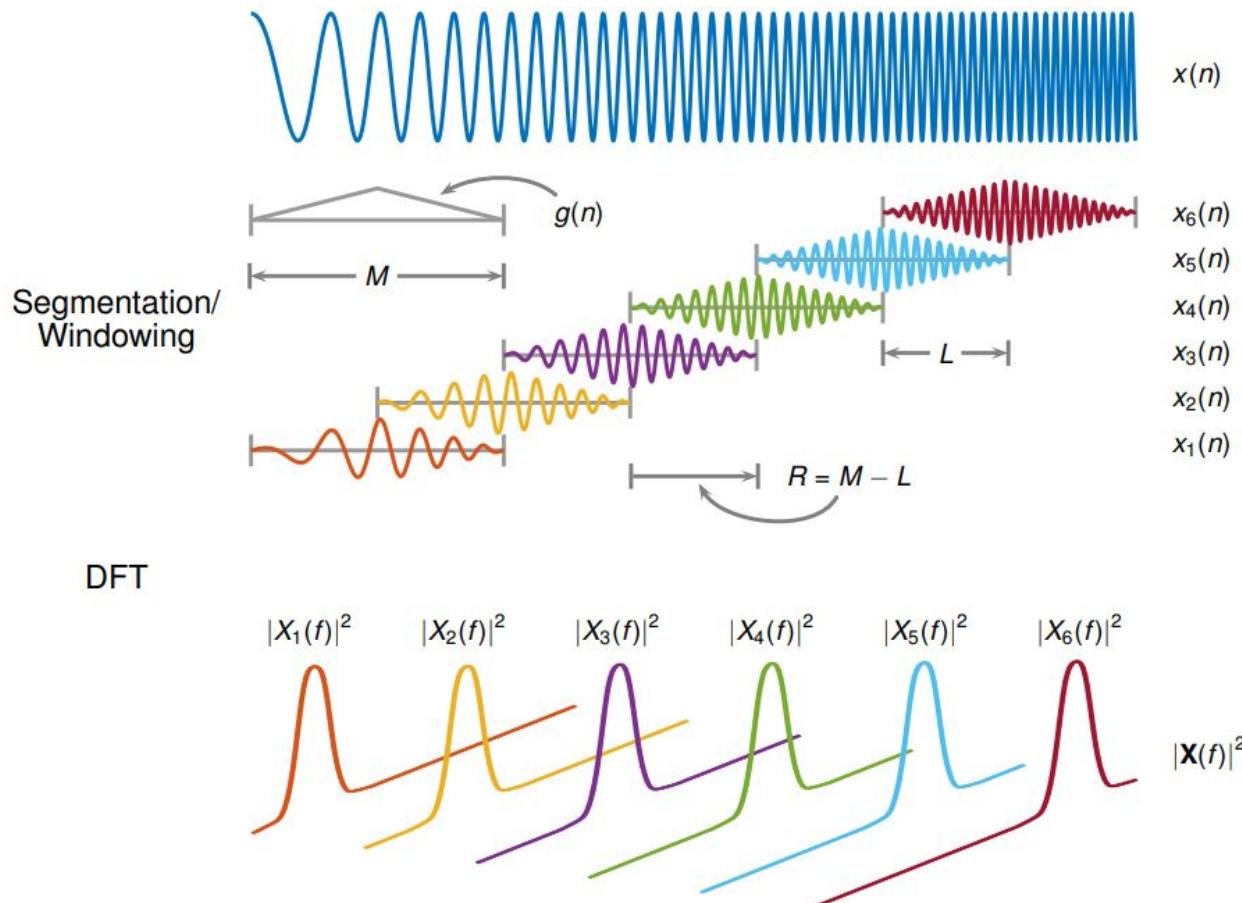
$$F(\omega_k) = \int x(t) e^{-i\omega_k t} dt$$



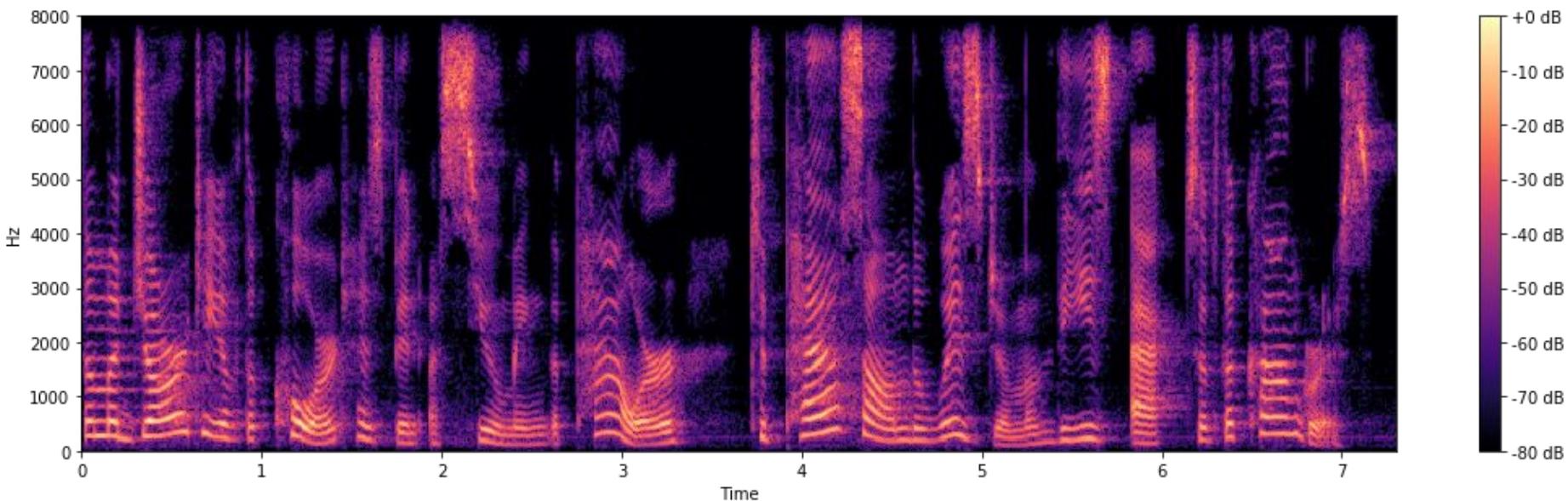
Time Duration		
Finite	Infinite	
Discrete FT (DFT) $X(k) = \sum_{n=0}^{N-1} x(n)e^{-j\omega_k n}$ $k = 0, 1, \dots, N-1$	Discrete Time FT (DTFT) $X(\omega) = \sum_{n=-\infty}^{+\infty} x(n)e^{-j\omega n}$ $\omega \in (-\pi, +\pi)$	discr. time n
Fourier Series (FS) $X(k) = \int_0^P x(t)e^{-j\omega_k t} dt$ $k = -\infty, \dots, +\infty$	Fourier Transform (FT) $X(\omega) = \int_{-\infty}^{+\infty} x(t)e^{-j\omega t} dt$ $\omega \in (-\infty, +\infty)$	cont. time t
discrete freq. k	continuous freq. ω	

Short-Time Fourier Transform

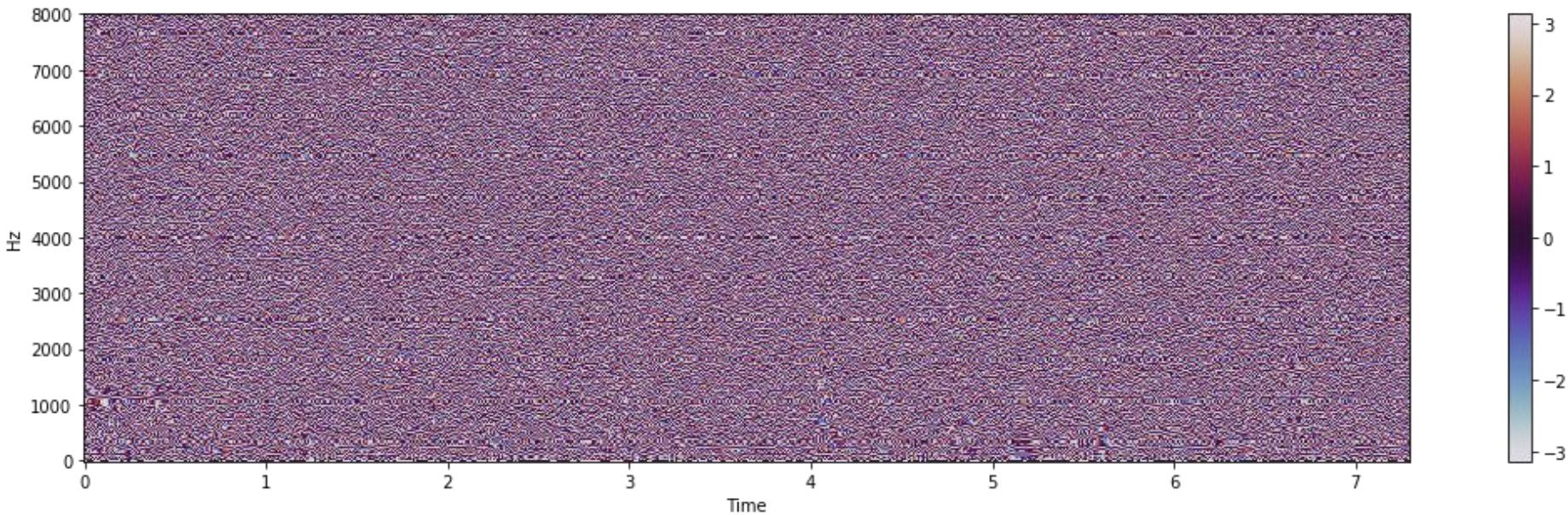
Samsung AI Center - Moscow



Amplitude spectrogram



Phase spectrogram

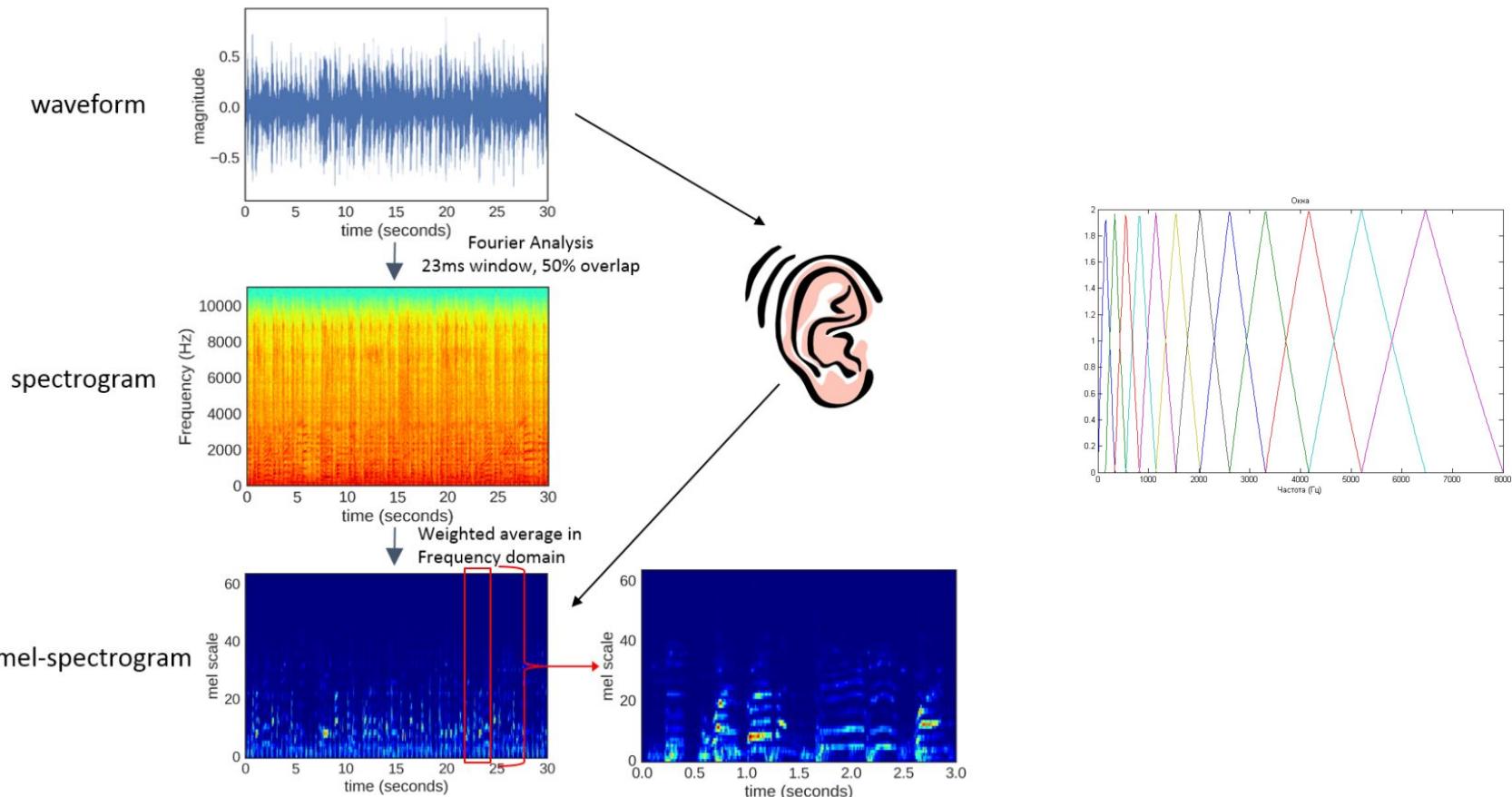


Griffin-Lim

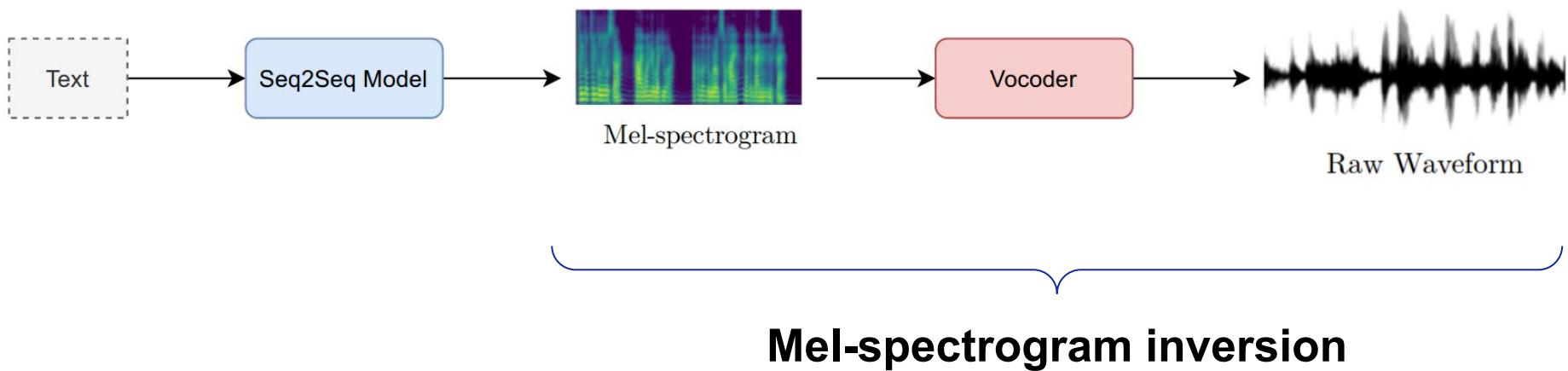
```
def griffin_lim(F, max_iters=32):
    phi = Uniform(0, 2 * pi, size=F.shape)
    for _ in range(max_iters):
        y = ISTFT(F * exp(1j * phi))
        phi = angle(STFT(y))
    return y
```

Mel-spectrogram

Samsung AI Center - Moscow

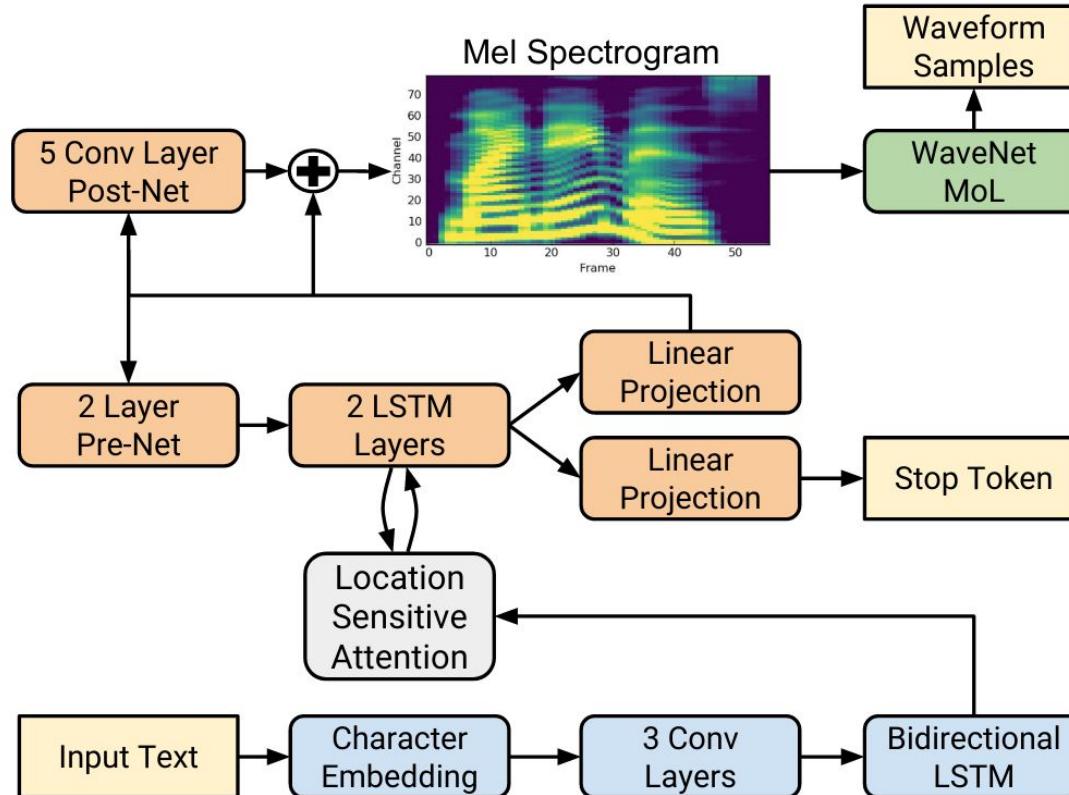


Audio Processing Problems



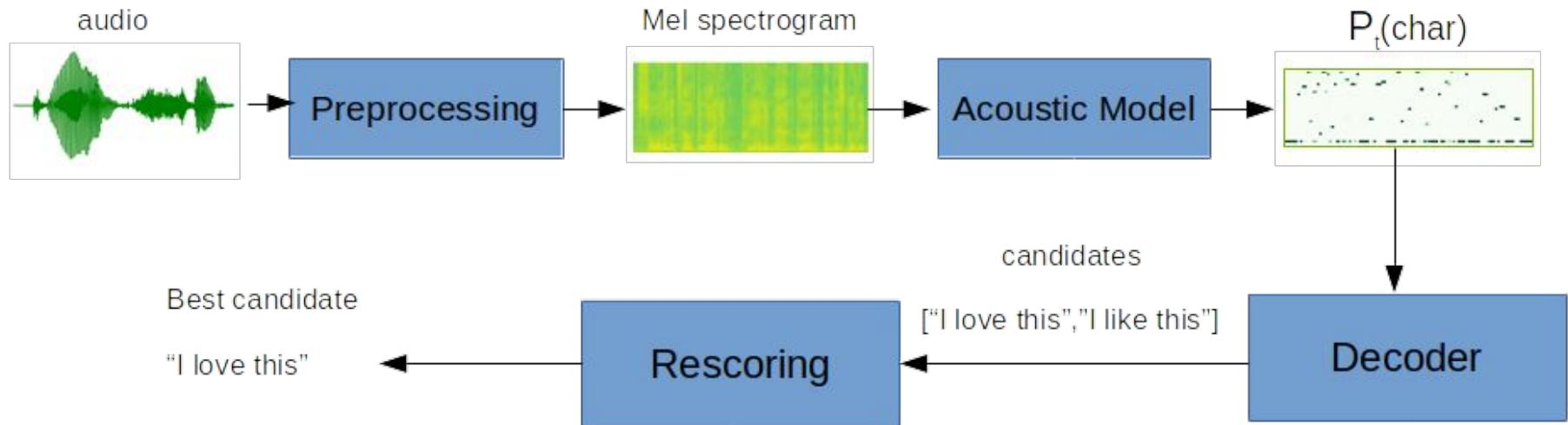
TTS Example - Tacotron 2

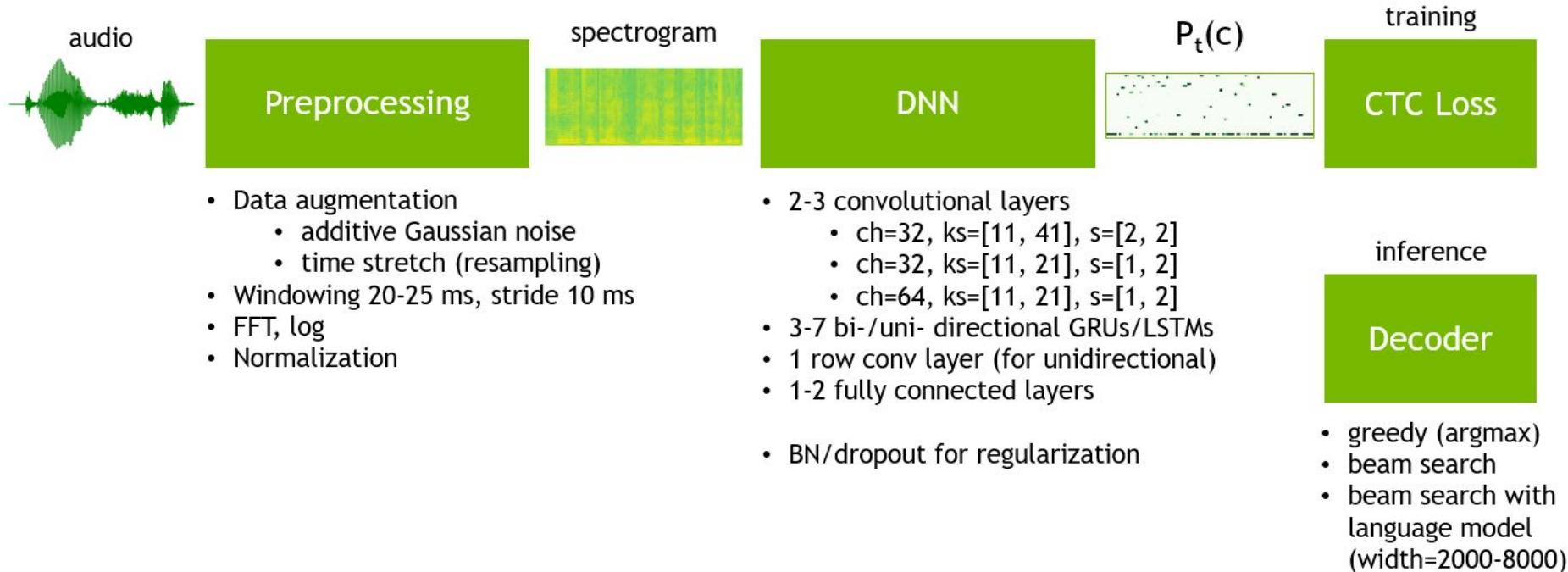
Samsung AI Center - Moscow



Automatic Speech Recognition (ASR)

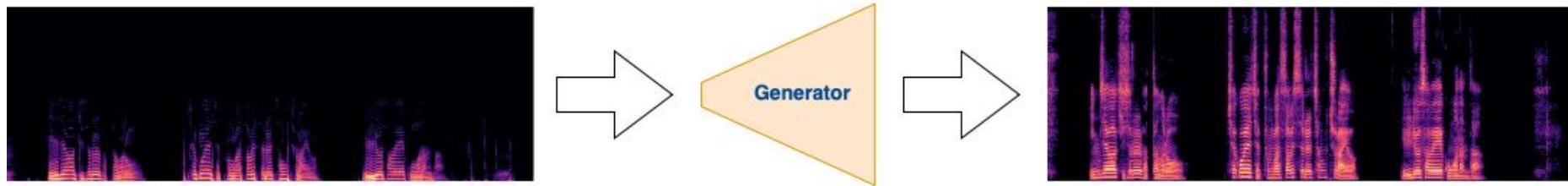
Samsung AI Center - Moscow





Bandwidth Extension

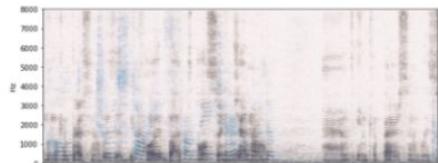
Samsung AI Center - Moscow



Multi-modal Speech Enhancement

Samsung AI Center - Moscow

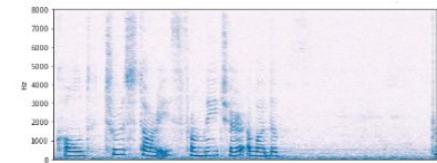
Noisy speech



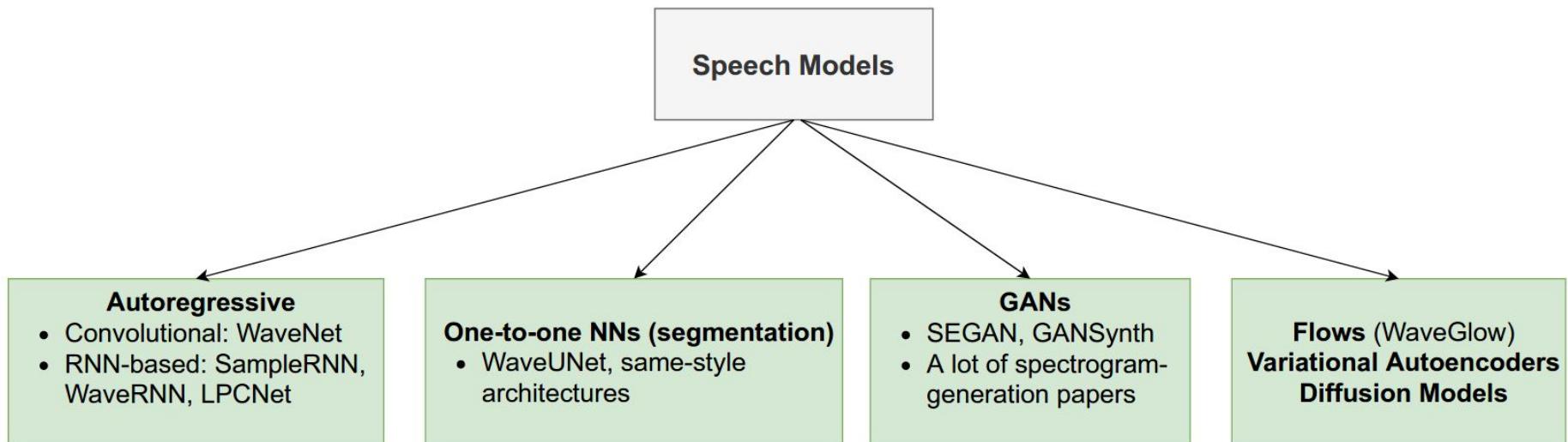
Accelerometer

Generator

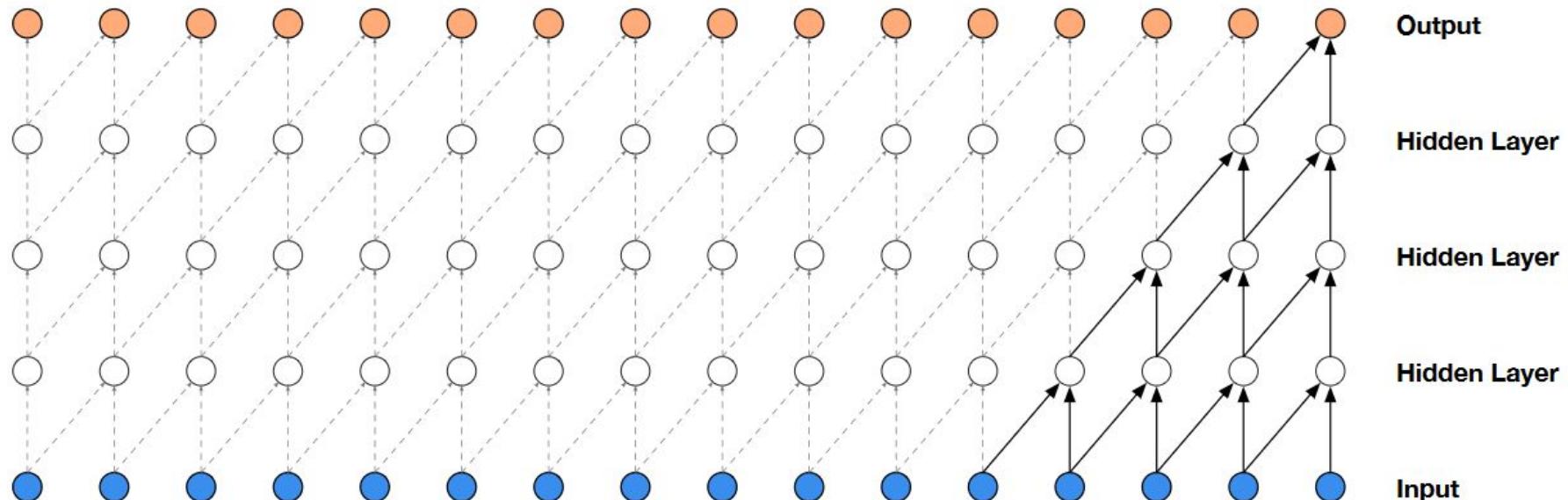
Clean speech



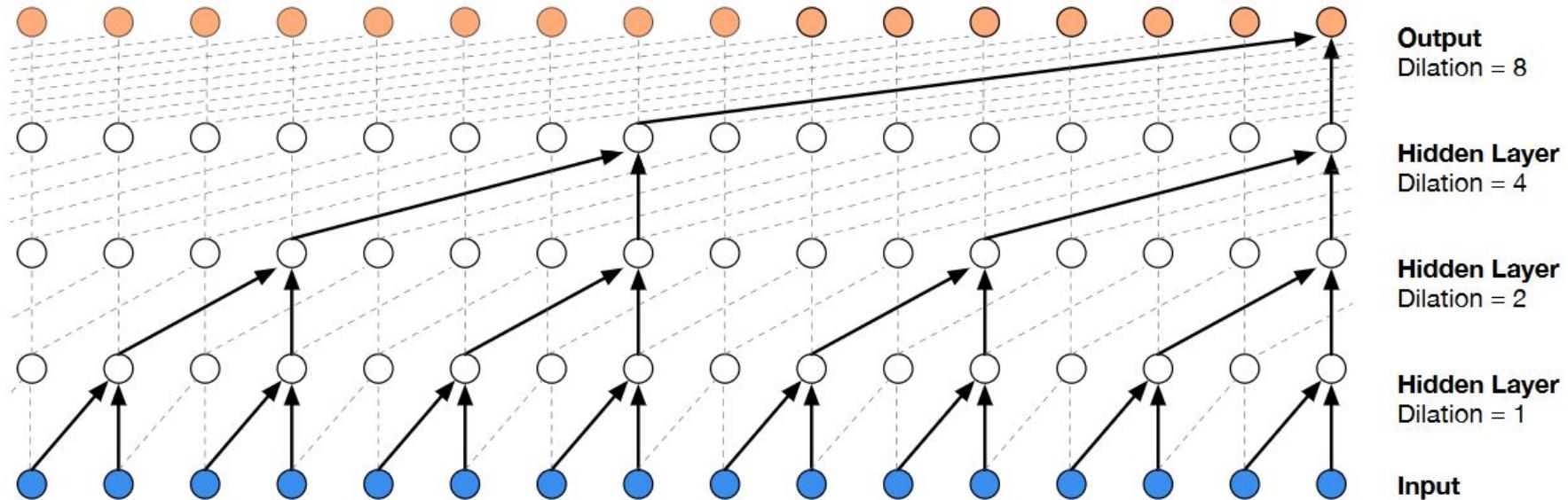
Architectures



Causal Convolutions have small receptive field

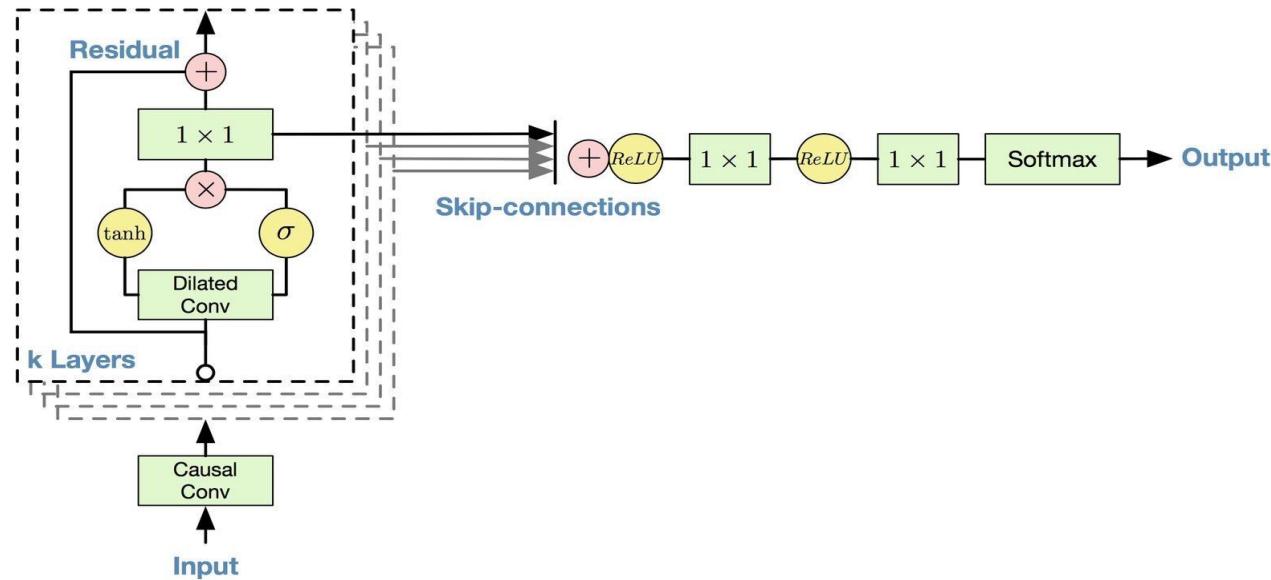


WaveNet uses causal *dilated* convolutions



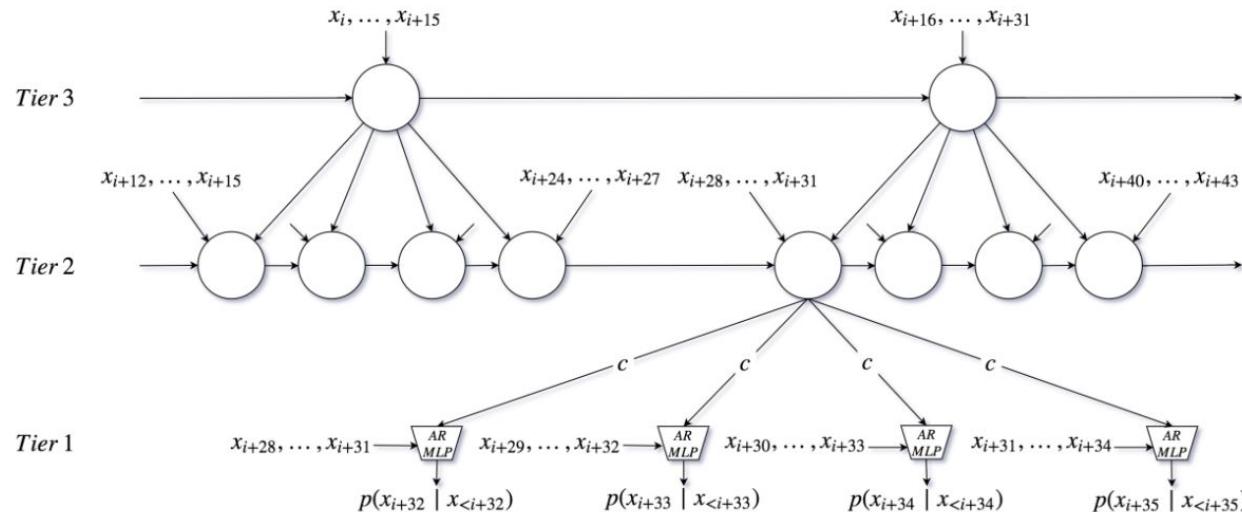
Dilations grows only up to some limit: 1, 2, 4, ..., 512, 1, 2, 4, ..., 512, 1, 2, 4, ..., 512

WaveNet: one block structure

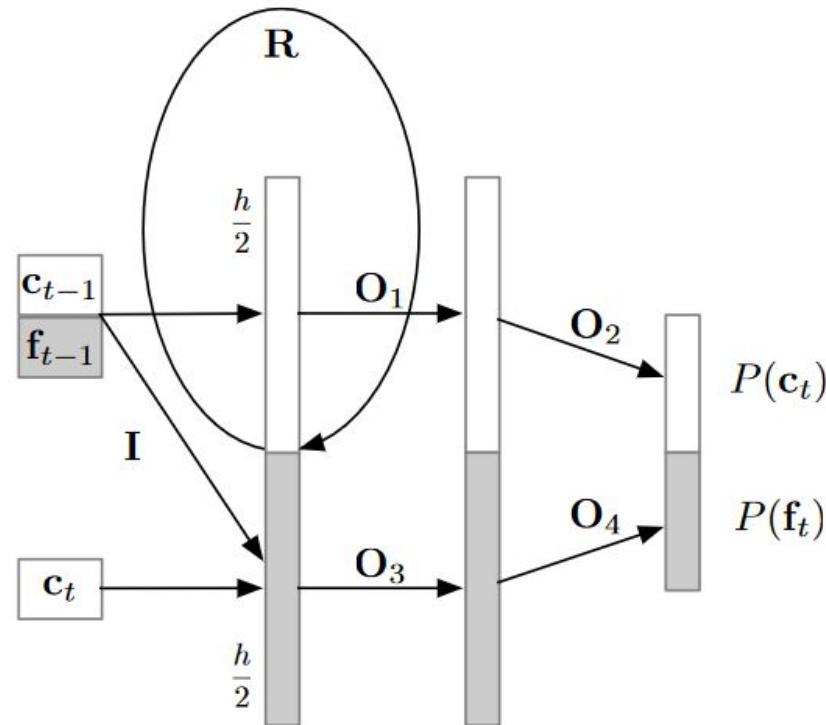


Multiscale RNNs

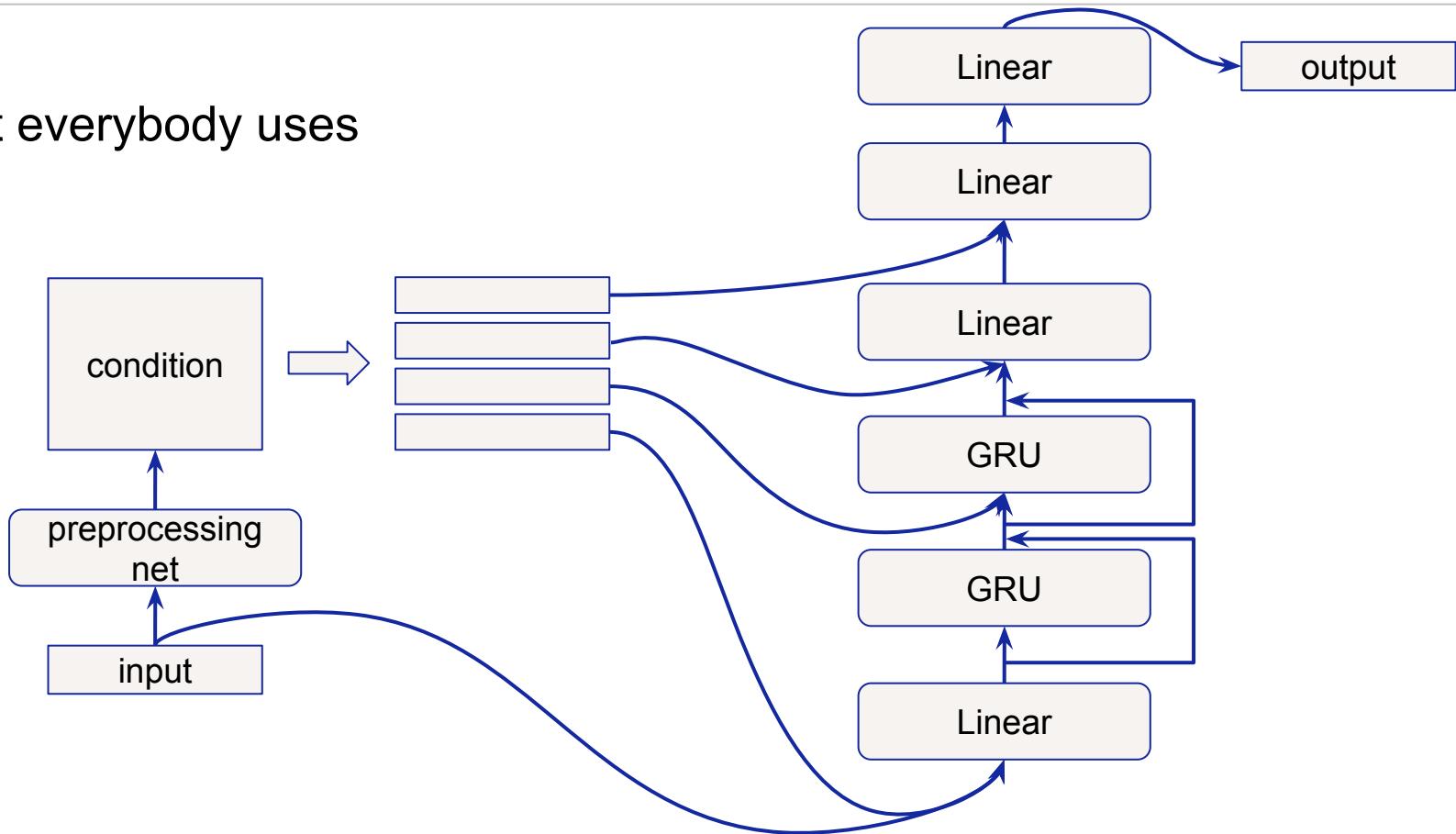
neither popular nor very effective



Original net: audio 8bit -> 16bit



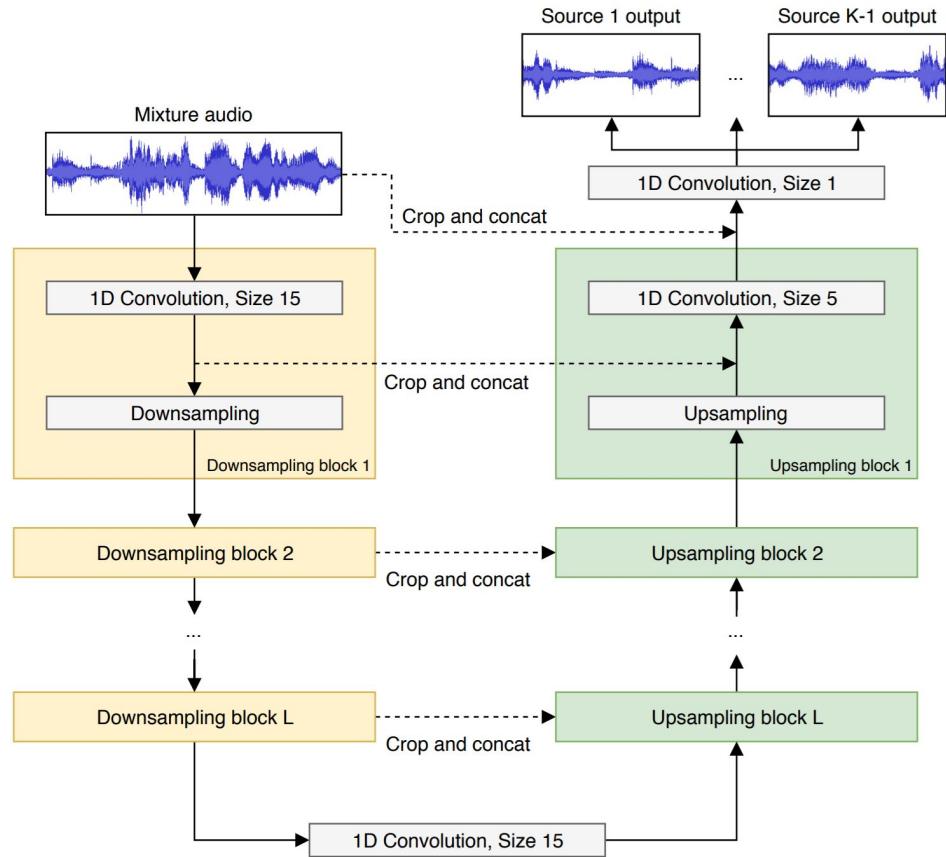
What everybody uses



Pros:

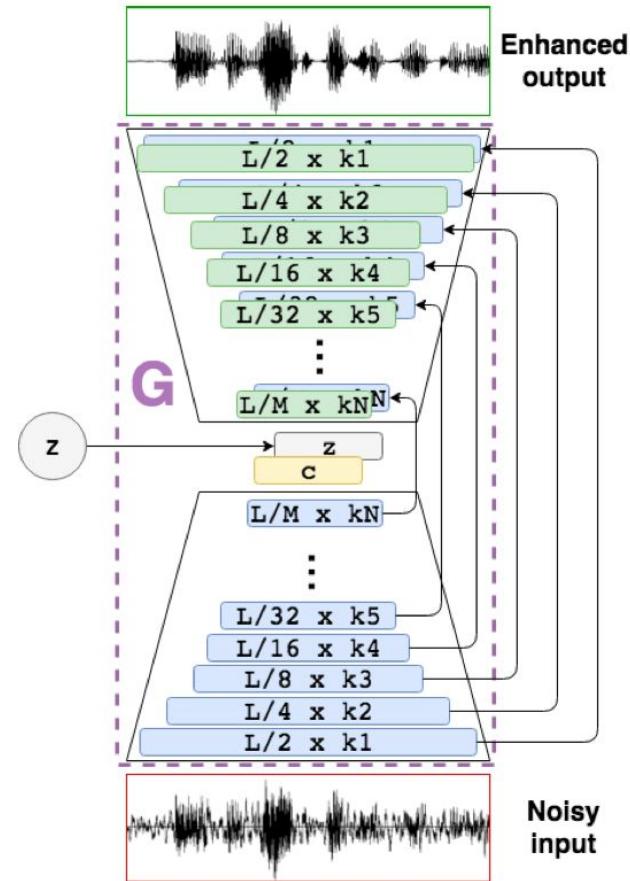
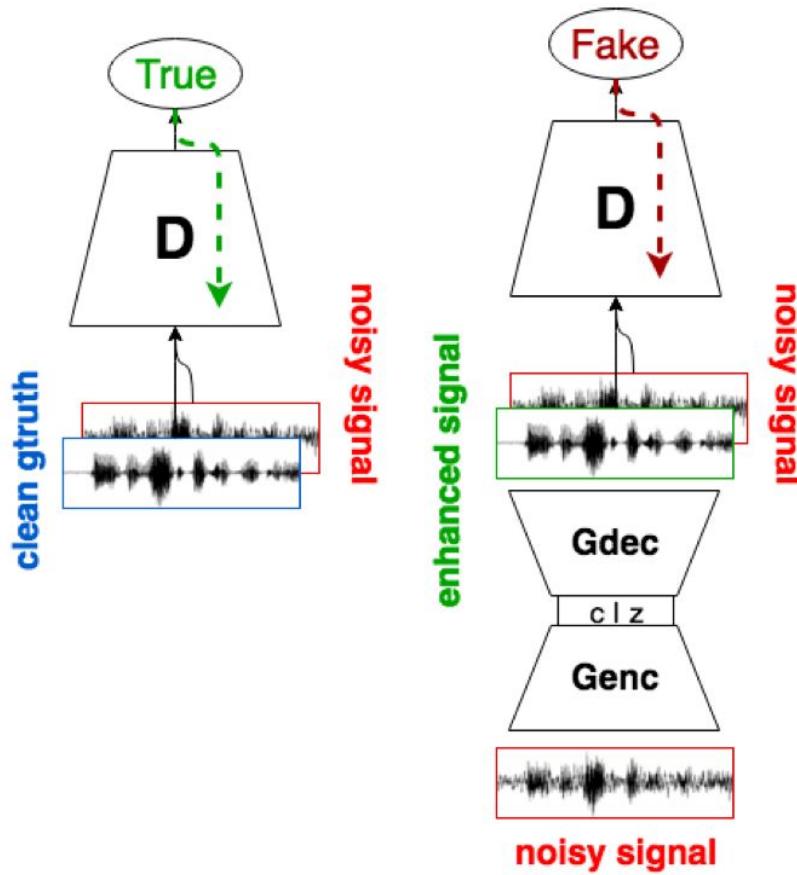
- + lightweight
- + operates on several scales
- + fully convolutional

Cons:



GANs (*)

Samsung AI Center - Moscow

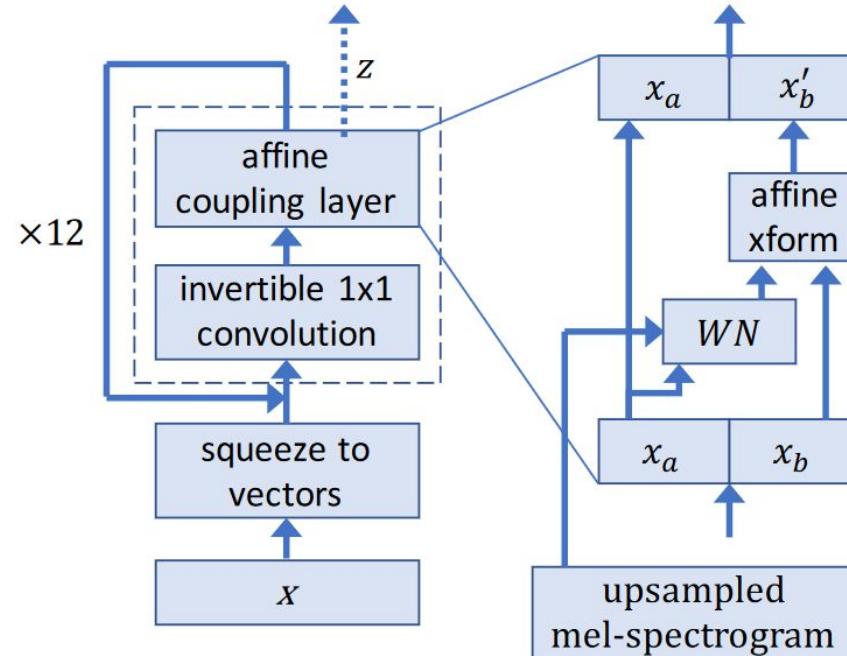


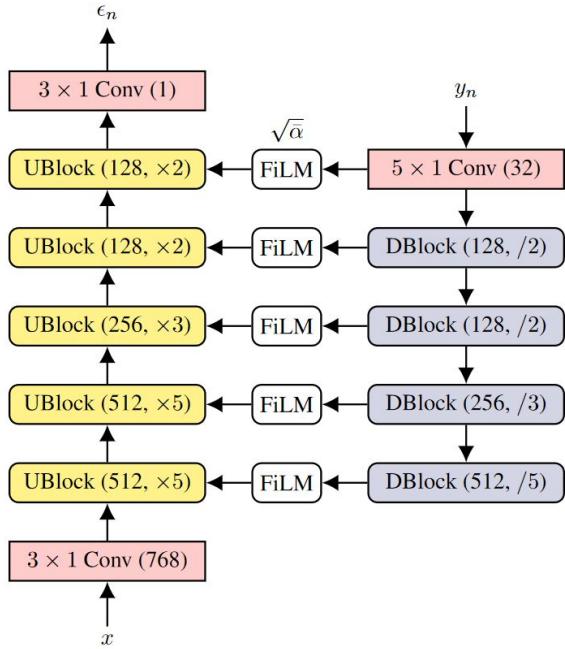
Pros:

- + highly parallel
- + fast on GPUs

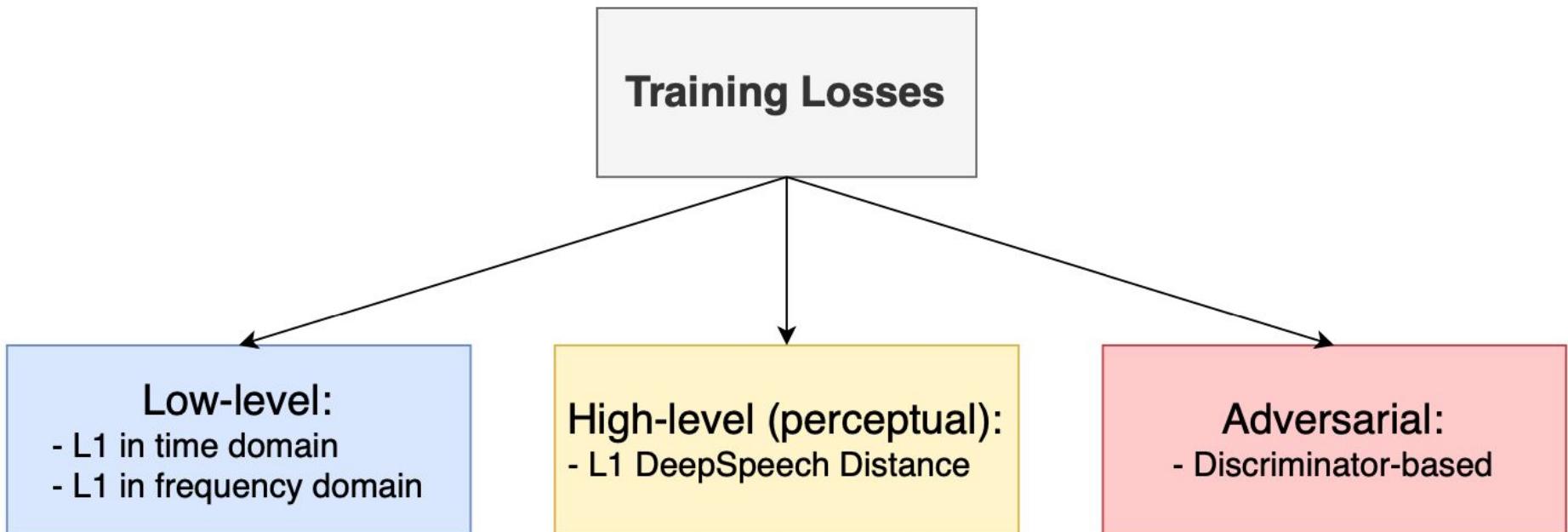
Cons:

- slow on CPUs



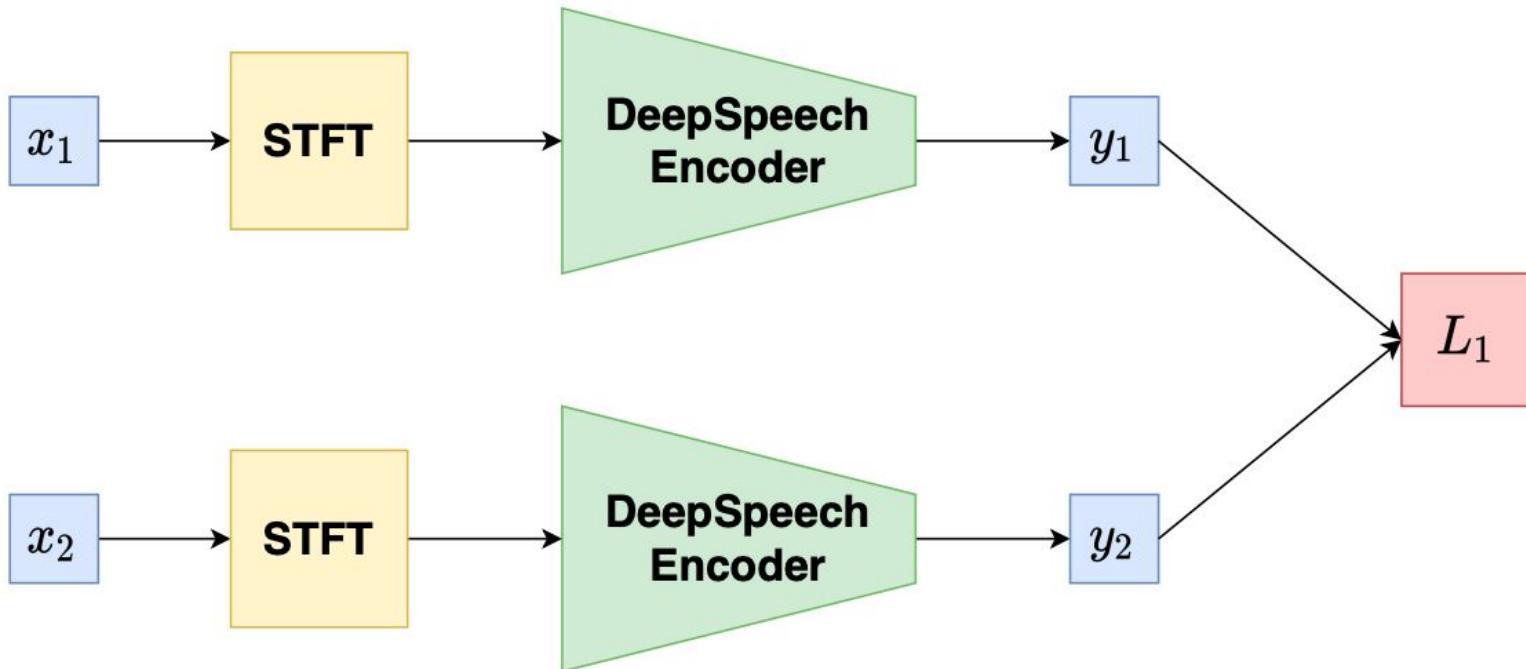


Loss Functions

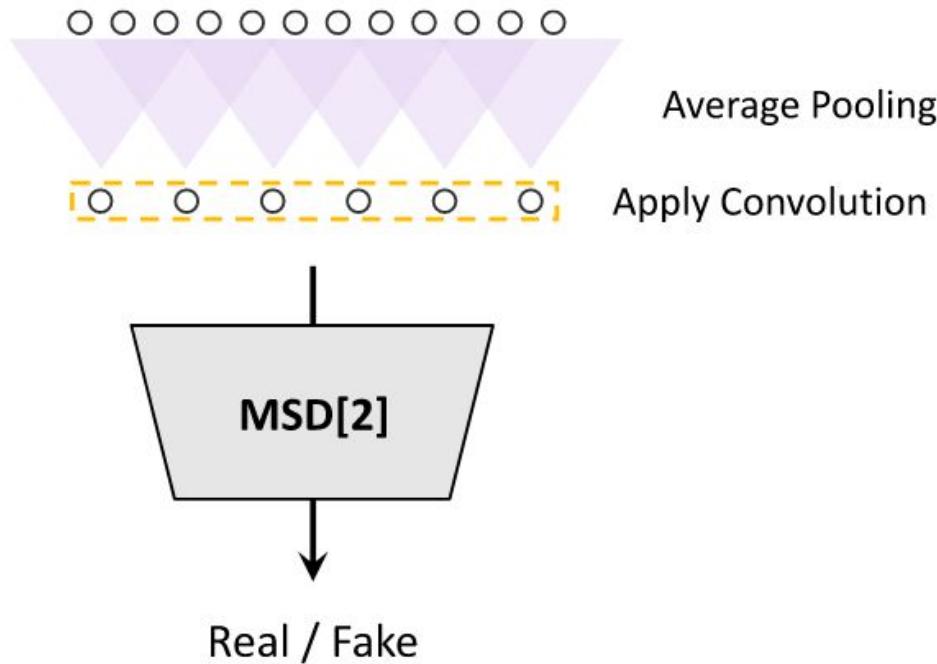


$$L_1(y, \hat{y}) = \|y - \hat{y}\|_1$$

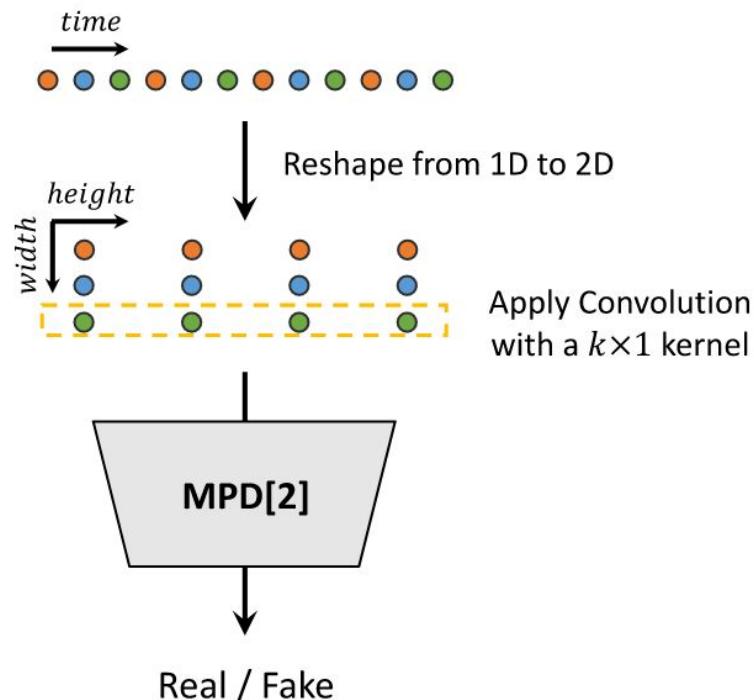
$$L_1^{stft}(y, \hat{y}) = \|\log(1 + |STFT(y)|) - \log(1 + |STFT(\hat{y})|)\|_1$$



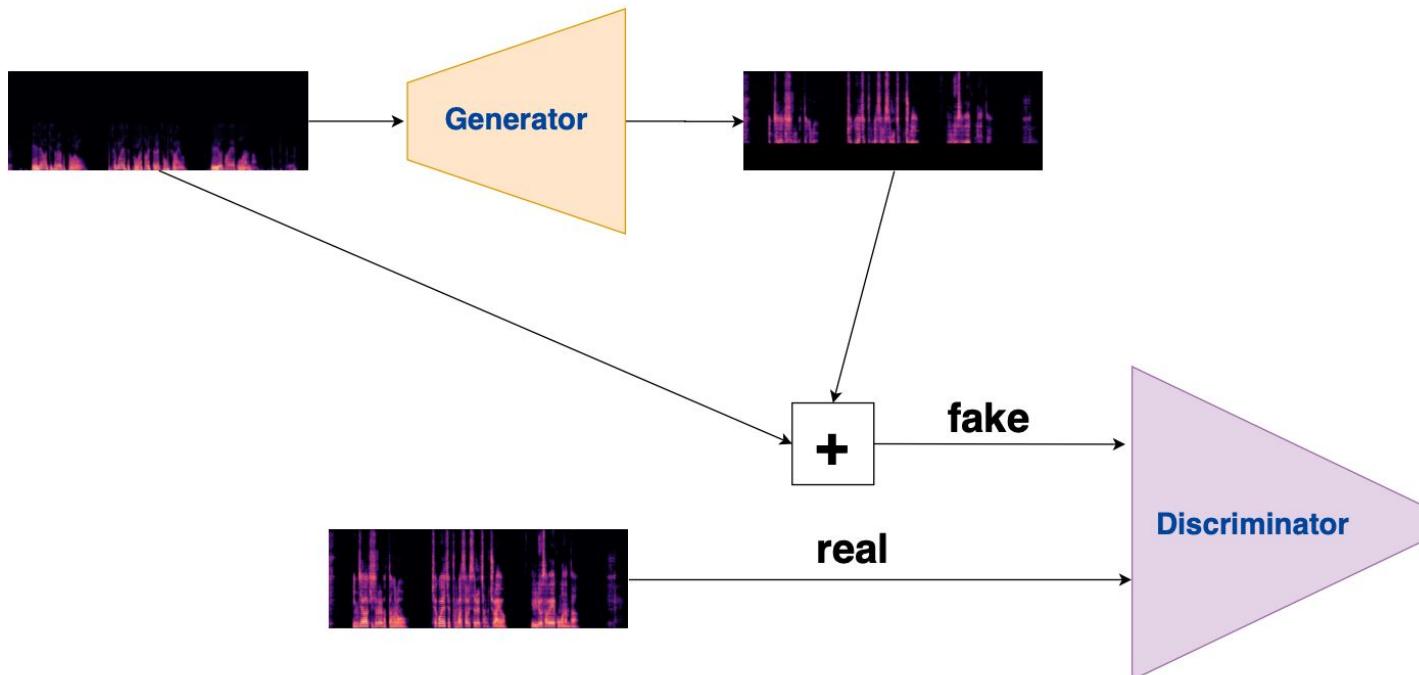
Multi-Scale Discriminators (MSD)



Multi-Period Discriminators (MPD)



Spectrogram-based Discriminators



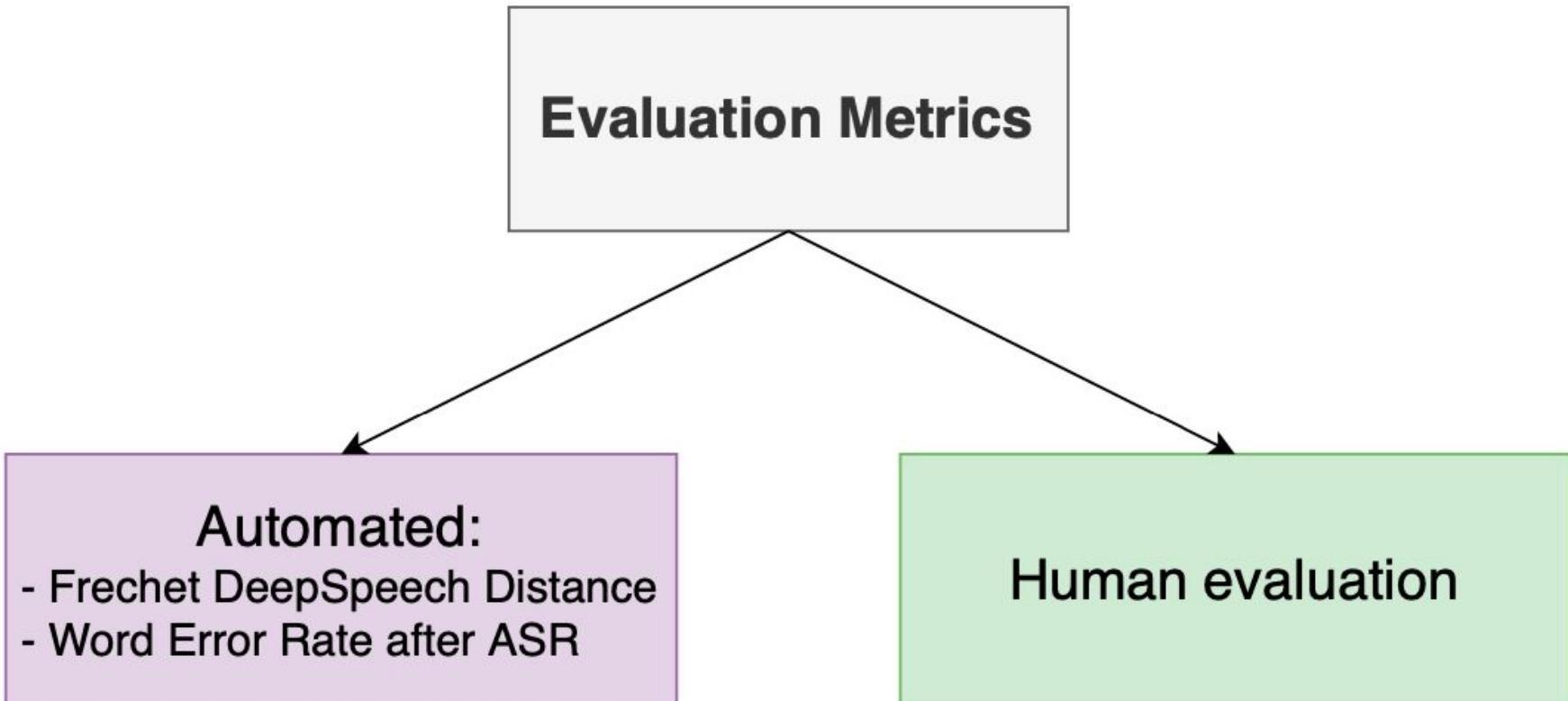
Hinge Loss:

$$\min_{D_k} \mathbb{E}_x \left[\min(0, 1 - D_k(x)) \right] + \mathbb{E}_{s,z} \left[\min(0, 1 + D_k(G(s, z))) \right], \forall k = 1, 2, 3$$

$$\min_G \mathbb{E}_{s,z} \left[\sum_{k=1,2,3} -D_k(G(s, z)) \right]$$

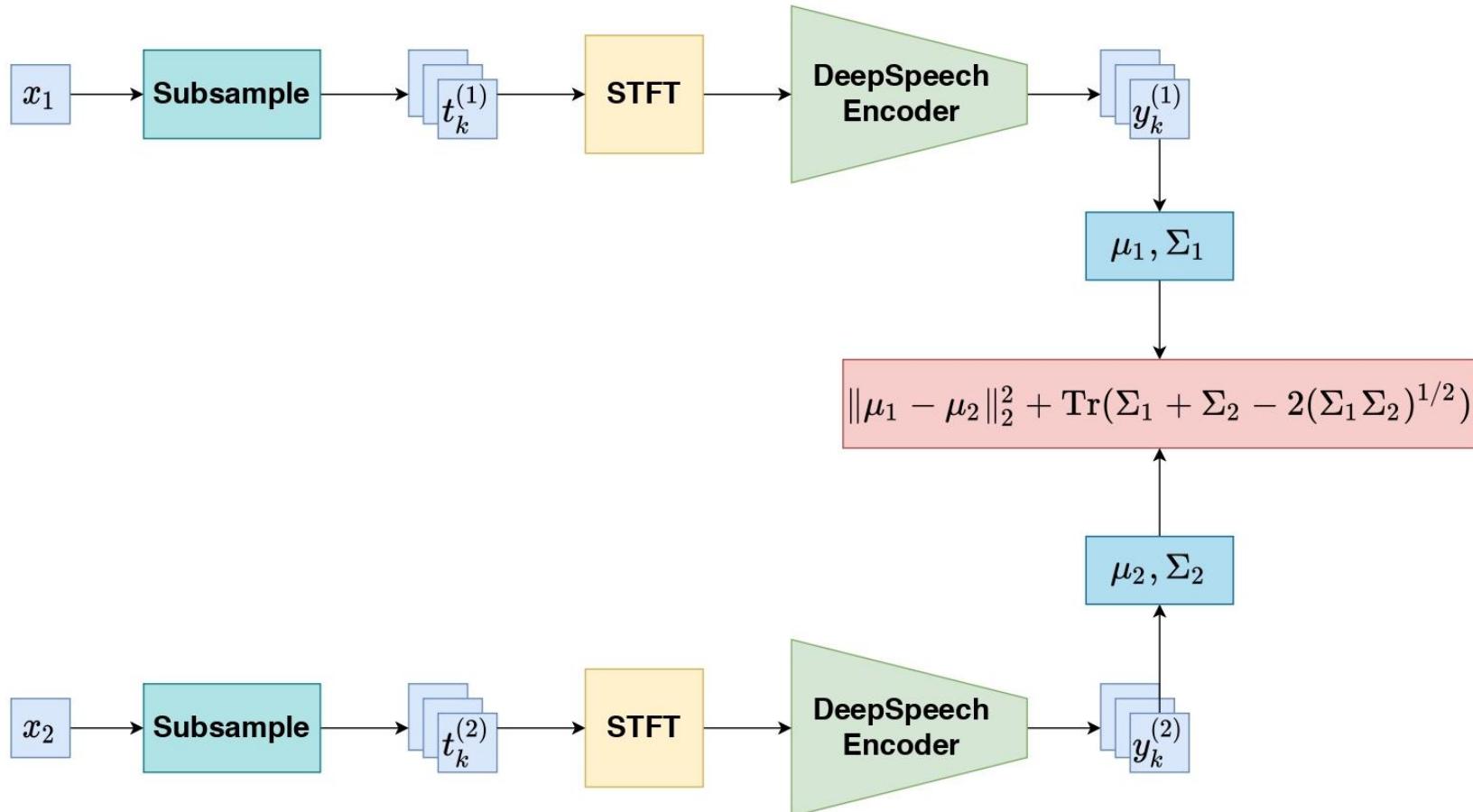
Feature Matching Loss:

$$\mathcal{L}_{\text{FM}}(G, D_k) = \mathbb{E}_{x,s \sim p_{\text{data}}} \left[\sum_{i=1}^T \frac{1}{N_i} \|D_k^{(i)}(x) - D_k^{(i)}(G(s))\|_1 \right]$$



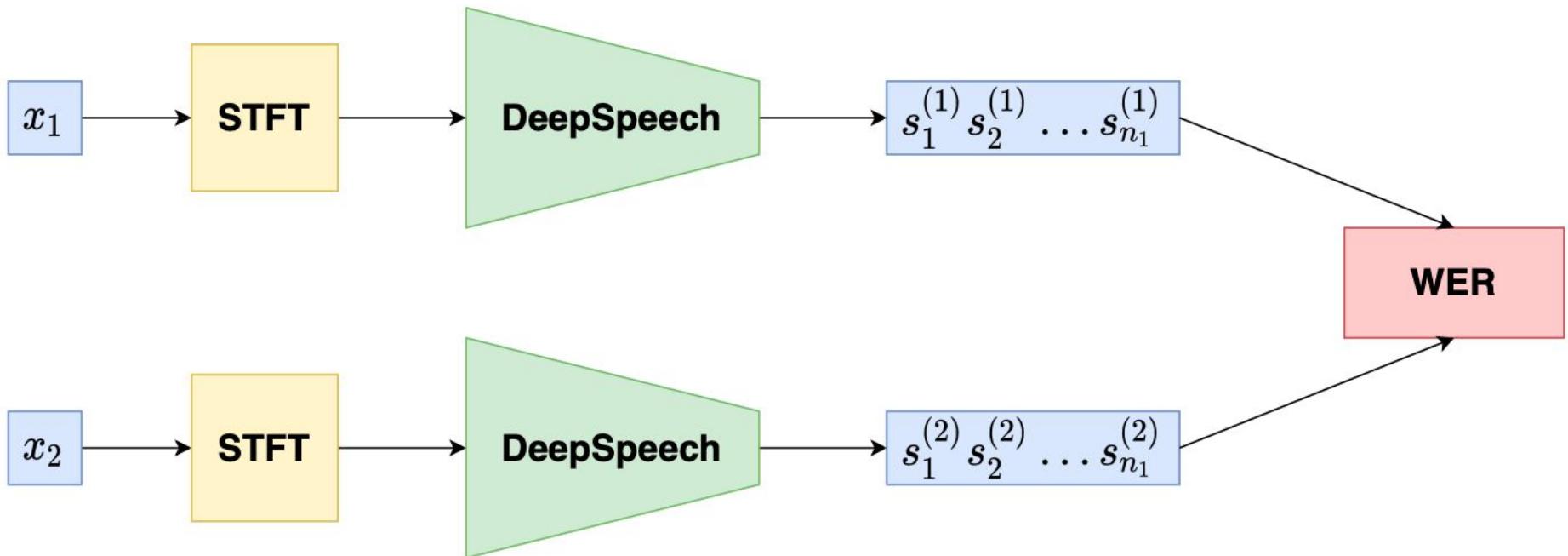
Frechet DeepSpeech Distance

Samsung AI Center - Moscow

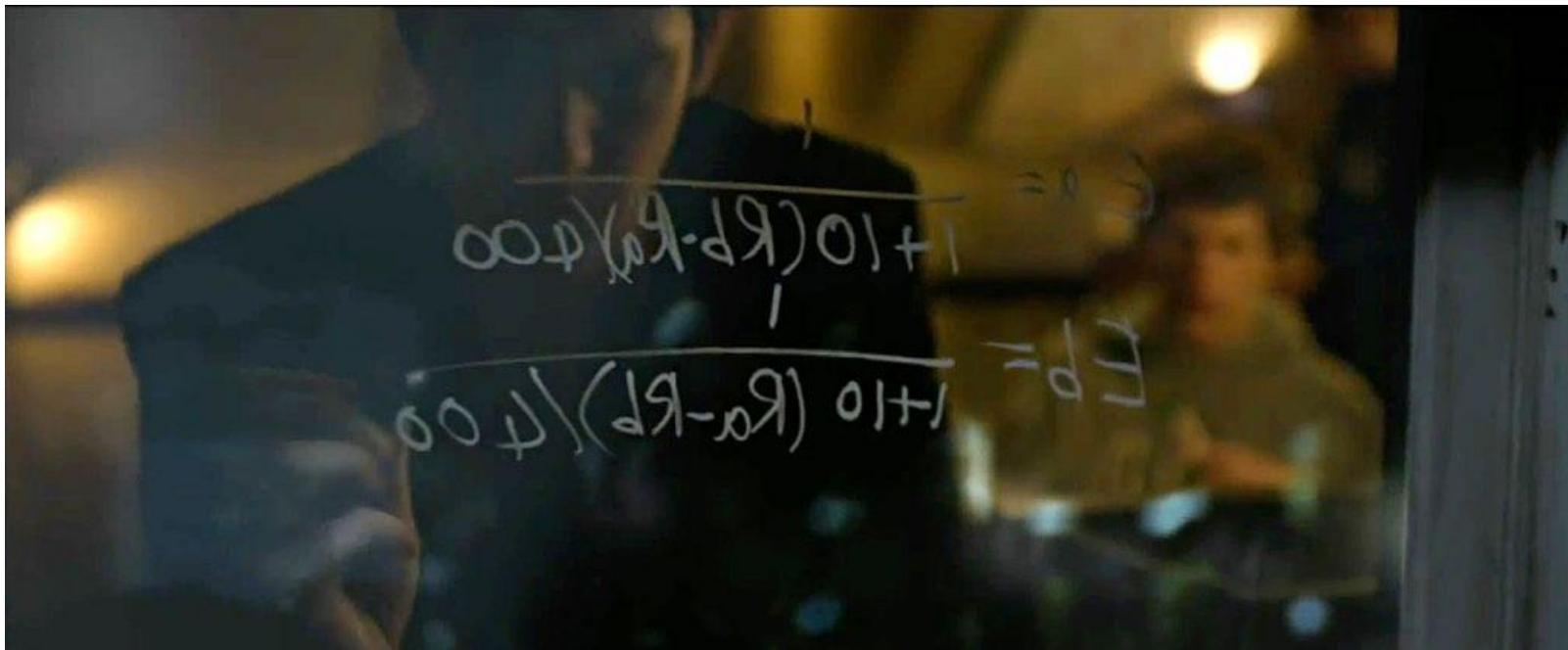


Word Error Rate after ASR

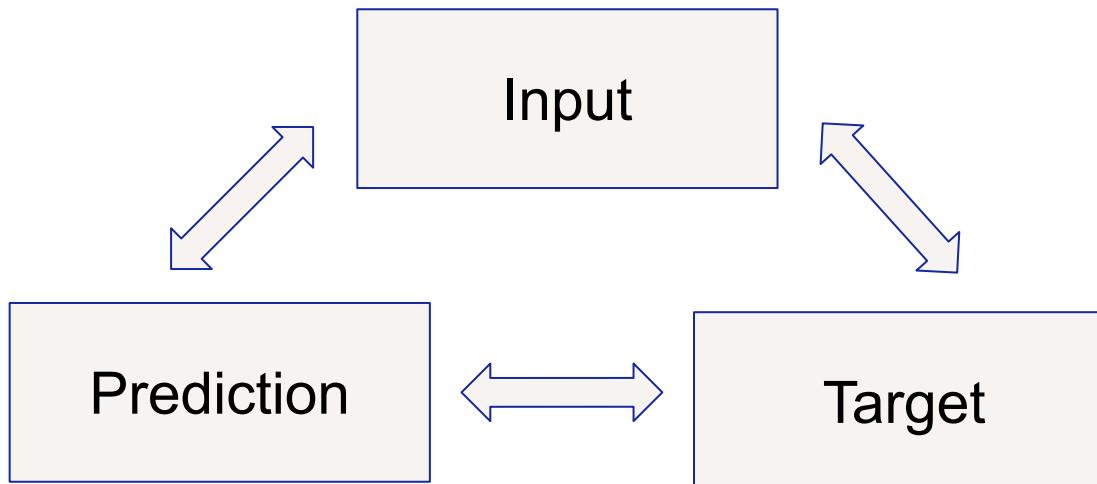
Samsung AI Center - Moscow



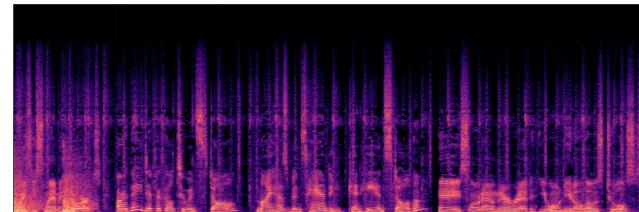
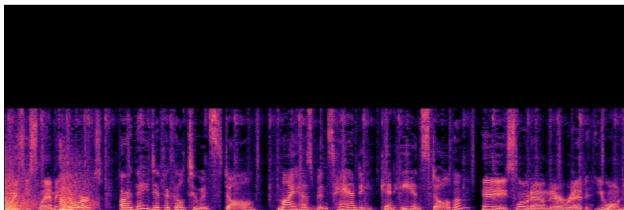
Idea: pairwise comparisons



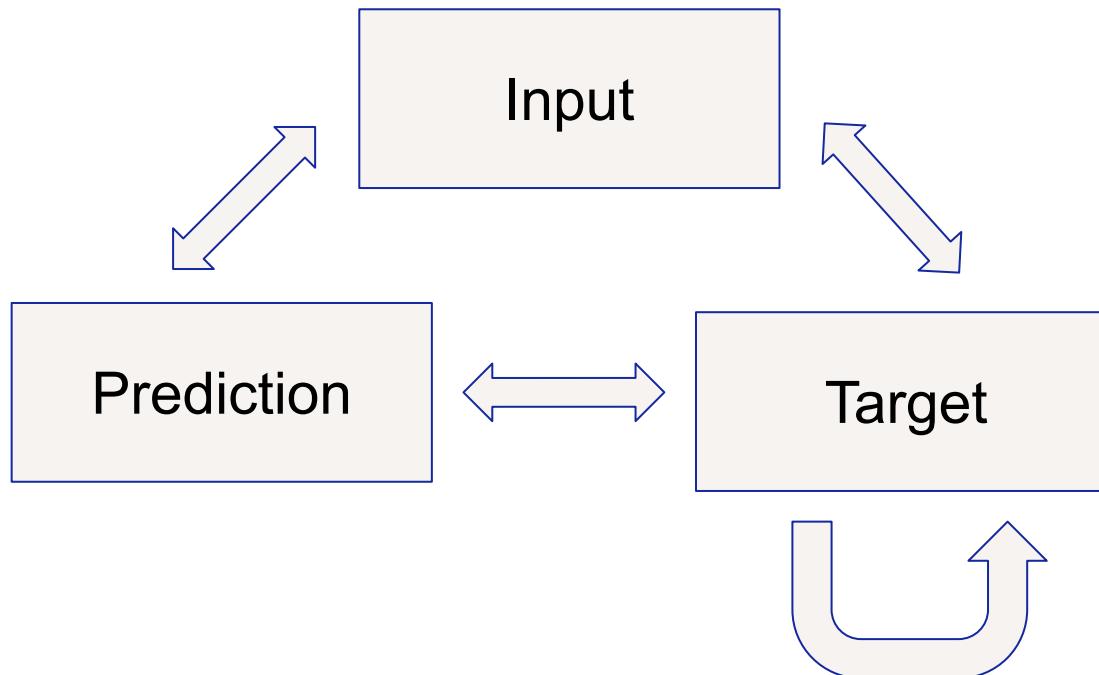
Answers: <, >



Problem



Answers: <, >, =, ≠



Label Studio

<https://labelstud.io/>

Text Classification

To have faith is to trust yourself to the water

Choose text sentiment

Positive^[1] Negative^[2] Neutral^[3]

Entity

Nothing selected

Entities (0)
No Entities added yet

Relations (0)
No Relations added yet

Datasets

LJ Speech

1 speaker,
~ 24 hours

VCTK

110 speakers
~ 44 hours

Mozilla Common VoiceSet

- Allows commercial usage
- Two months of train data (~64Gb)
- 17 train and 10 validation languages
- Highly imbalanced (languages, genders)

LJ Speech

1 speaker,
~ 24 hours

VCTK

110 speakers
~ 44 hours

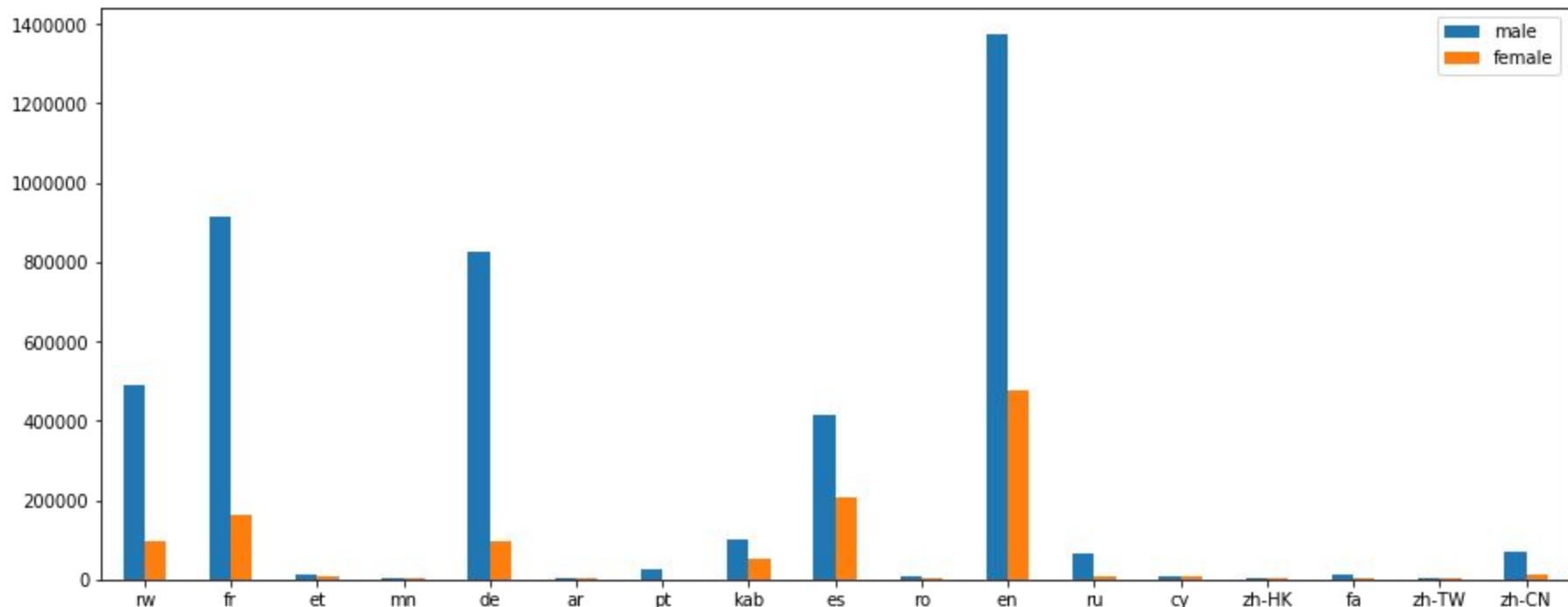
Mozilla Common VoiceSet

- Allows commercial usage
- Two months of train data (~64Gb)
- 17 train and 10 validation languages
- Highly imbalanced (languages, genders)
- Credits to Jules

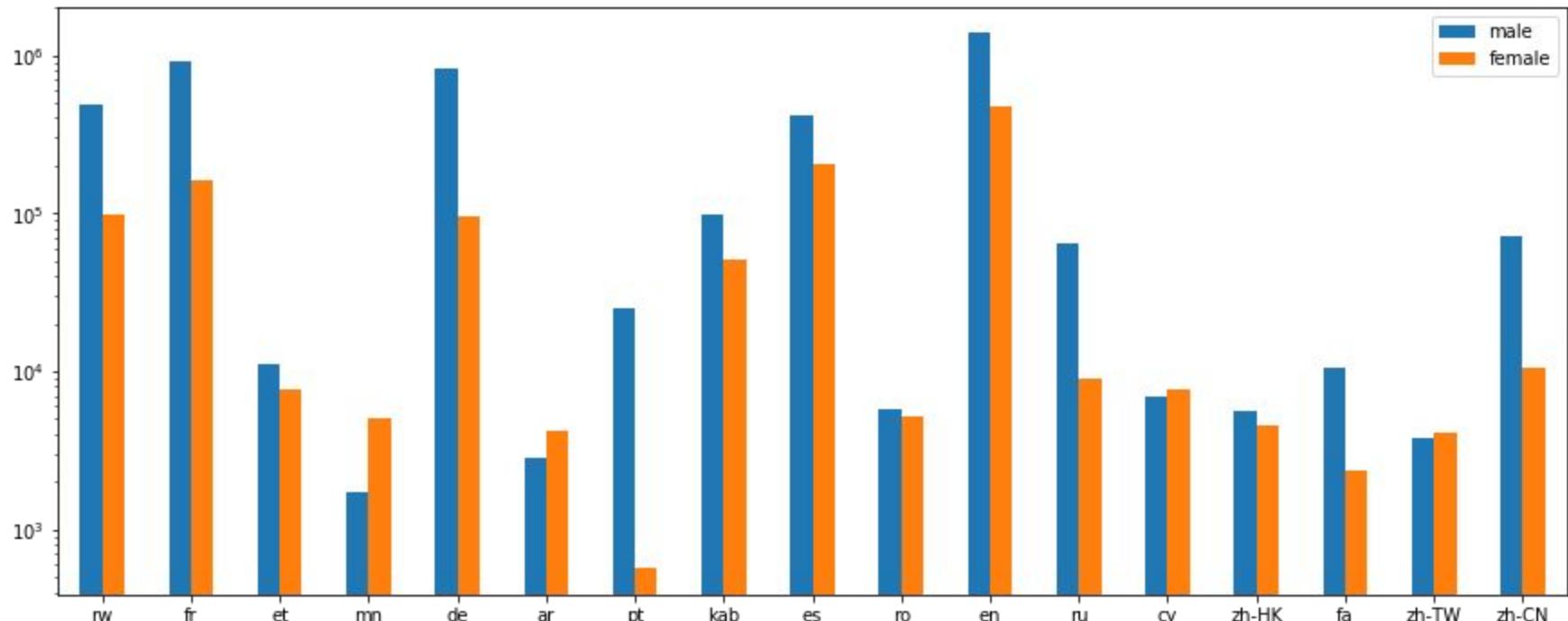


Jules Churkina consulted us while we balanced languages

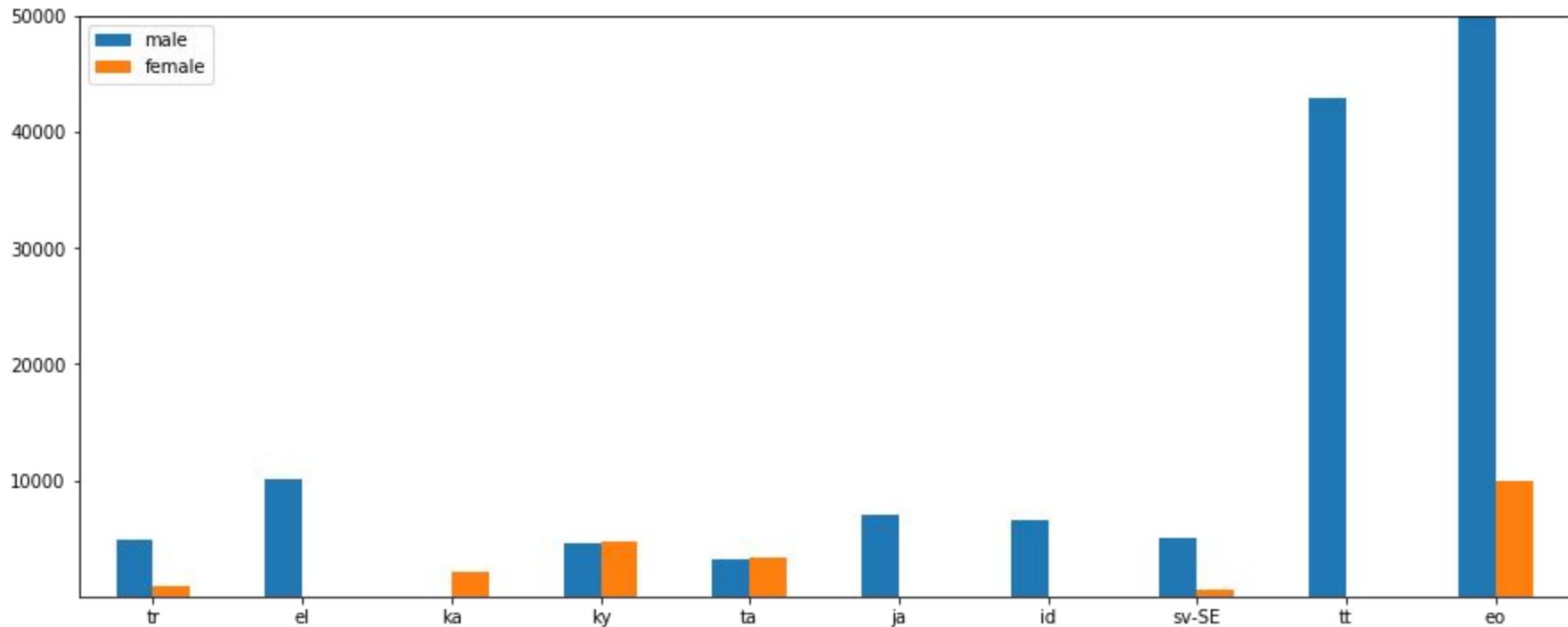
VoiceSet Imbalance: Train Set



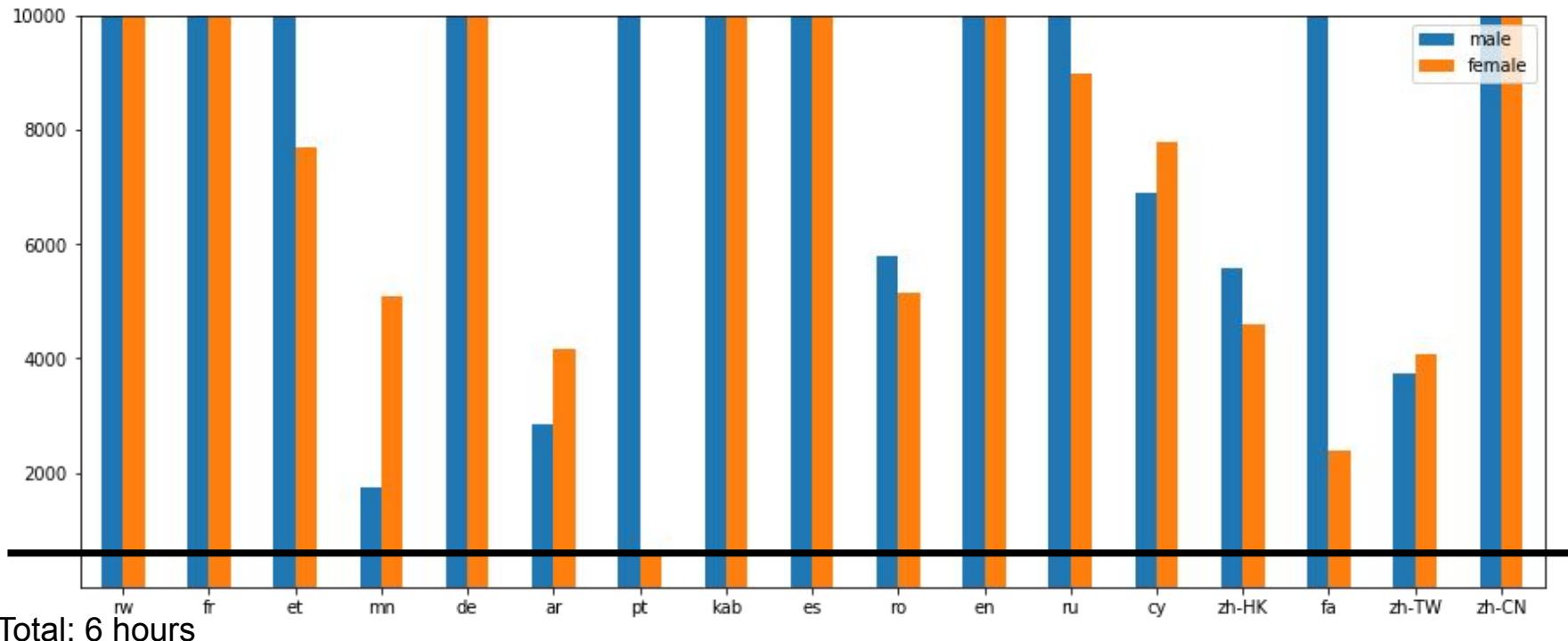
VoiceSet Imbalance: Train Set



VoiceSet Imbalance: Validation Set



VoiceSet Imbalance: Train Set



AudioSet

~ 20K audios

Length of each audio is 10s

Labelling: audio event tags

No raw audio

Downloading is illegal

FreeSound

Each audio is either CC0, CC-BY, Sampling+ or CC-BY-NC

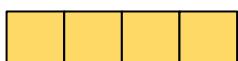
Tags, description and other labels are available

~13Gb of uncompressed filtered wavs

~ 1 day of train set

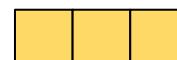
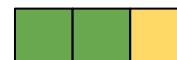
Data Pool

Audios splitted to frames



Batches

All frames of audio
are used



One frame per audio
is used

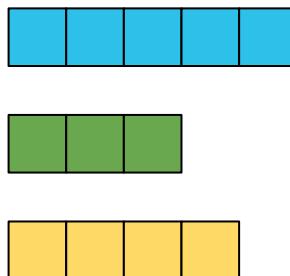


- Slow (most read frames are dropped)
- Uncorrelated batches

- Fast
- Highly correlated batches

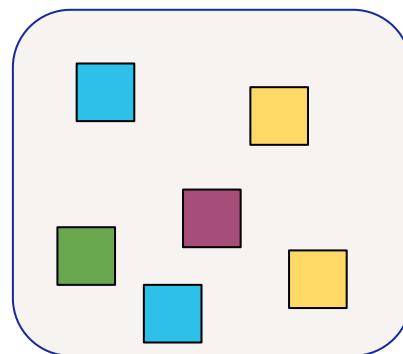
Data Pool

Audios splitted to frames



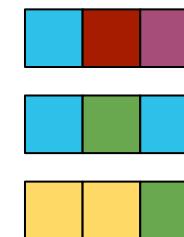
read audios until
required pool size
is reached

Pool



Invariant:
there are at least
K frames in Pool
before sampling

sample batches
randomly from the
pool



- Fast
- Weakly correlated
batches

Databases for Large Datasets

File System

HDF5 (Hierarchical Data Format)

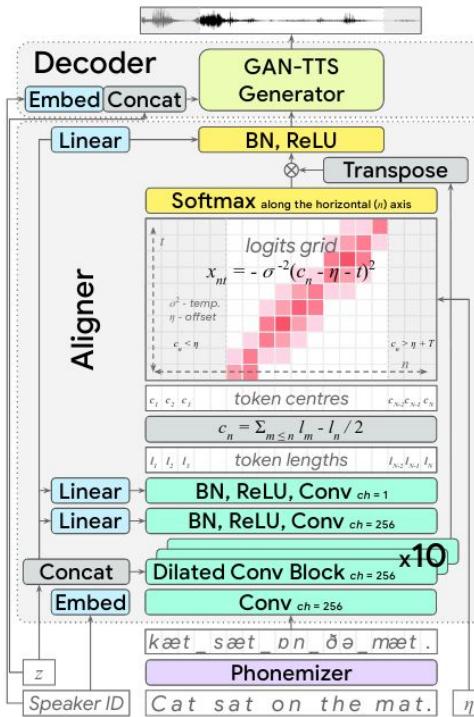
LMDB (Lightning Memory-Mapped Database)

WebDataset (<https://github.com/tmbdev/webdataset/>)



SOTAs

End-to-End Adversarial Text-to-Speech (ICLR 2021 oral)



- Data - Multispeaker (MS) or Single Speaker (SS)
- Inputs - characters (Ch) or phonemes (Ph)
- Random Window Discriminators (RWD)
- Mel-spectrogram Discriminator (MSD)

End-to-End Adversarial Text-to-Speech

Samsung AI Center - Moscow

Model	Data	Inputs	RWD	MSD	$\mathcal{L}_{\text{length}}$	$\mathcal{L}_{\text{pred}}$	Align	MOS
Natural Speech				-				4.55 ± 0.075
<i>GAN-TTS</i> (Bińkowski et al., 2020)				-				4.213 ± 0.046
<i>WaveNet</i> (van den Oord et al., 2016)				-				4.41 ± 0.069
<i>Par. WaveNet</i> (van den Oord et al., 2018)				-				4.41 ± 0.078
<i>Tacotron 2</i> (Shen et al., 2018)				-				4.526 ± 0.066
No $\mathcal{L}_{\text{length}}$	MS	Ph	✓	✓	✗	$\mathcal{L}_{\text{pred}}''$	MI	[does not train]
No $\mathcal{L}_{\text{pred}}$	MS	Ph	✓	✓	✓	✗	MI	[does not train]
No Discriminators	MS	Ph	✗	✗	✓	$\mathcal{L}_{\text{pred}}''$	MI	1.407 ± 0.040
No RWDs	MS	Ph	✗	✓	✓	$\mathcal{L}_{\text{pred}}''$	MI	2.526 ± 0.060
No Phonemes	MS	Ch	✓	✓	✓	$\mathcal{L}_{\text{pred}}''$	MI	3.423 ± 0.073
No MelSpecD	MS	Ph	✓	✗	✓	$\mathcal{L}_{\text{pred}}''$	MI	3.525 ± 0.057
No Mon. Int.	MS	Ph	✓	✓	✓	$\mathcal{L}_{\text{pred}}''$	Attn	3.551 ± 0.073
No DTW	MS	Ph	✓	✓	✓	$\mathcal{L}_{\text{pred}}''$	MI	3.559 ± 0.065
Single Speaker	SS	Ph	✓	✓	✓	$\mathcal{L}_{\text{pred}}''$	MI	3.829 ± 0.055
EATS (Ours)	MS	Ph	✓	✓	✓	$\mathcal{L}_{\text{pred}}''$	MI	4.083 ± 0.049

Input text: *In this work, we take on the challenging task of learning to synthesise speech from normalised text or phonemes in an end-to-end manner, resulting in models which operate directly on character or phoneme input sequences and produce raw speech audio outputs.*

Main model:



Single speaker:



No RWDs:



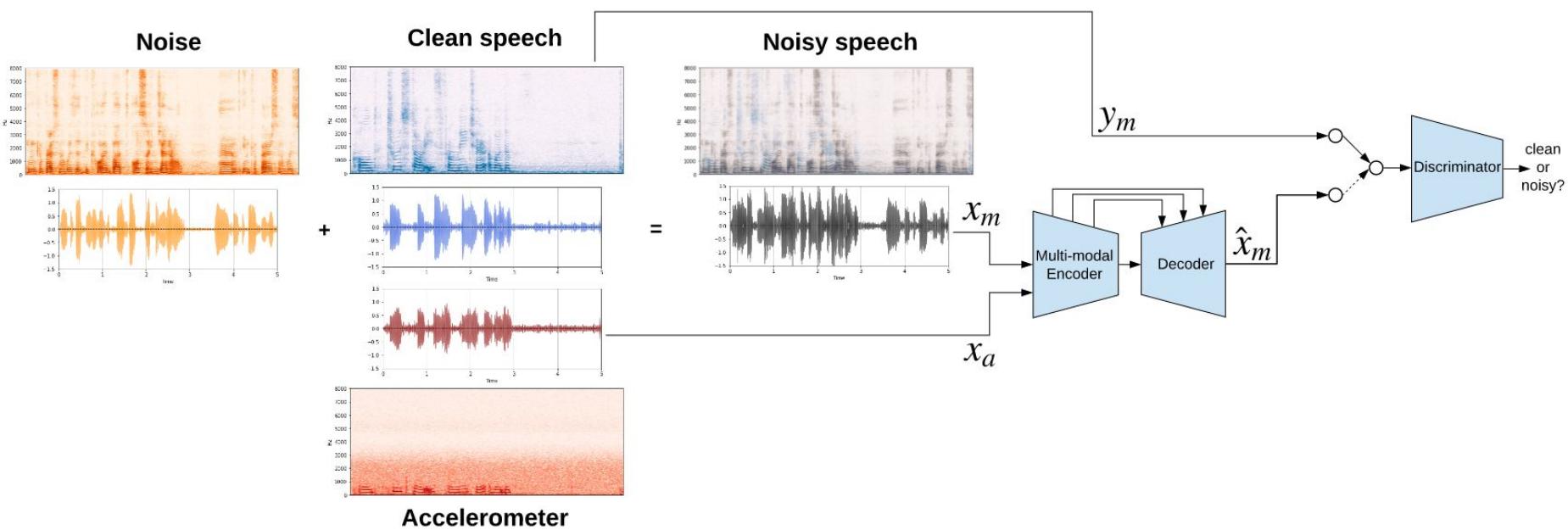
No MelSpecD:



No Discriminators:



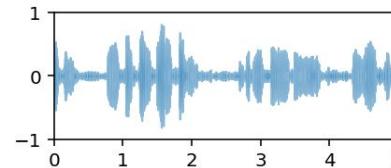
SEANet: A Multi-modal Speech Enhancement Network



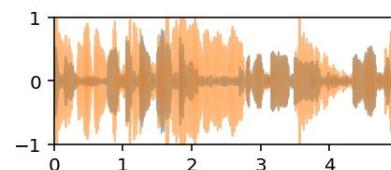
scenario	split	SEANet audio + accel	SEANet audio only
Mixed noise	1	9.9 ± 0.2	8.4 ± 0.2
	2	8.0 ± 0.2	7.9 ± 0.1
	3	8.3 ± 0.1	7.2 ± 0.2
	4	8.8 ± 0.1	8.1 ± 0.1
	5	9.9 ± 0.1	8.4 ± 0.1
	avg.	8.9	8.0
Mixed speech	1	10.1 ± 0.1	-0.9 ± 0.1
	2	8.6 ± 0.1	-0.9 ± 0.1
	3	9.2 ± 0.1	-0.7 ± 0.0
	4	9.0 ± 0.2	-1.0 ± 0.1
	5	11.1 ± 0.2	-0.9 ± 0.1
	avg.	9.6	-0.9

Mixed speech example:

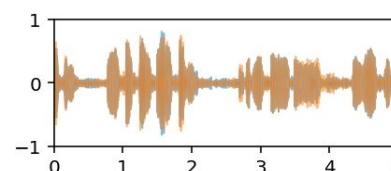
Clean sample:



Noisy sample:



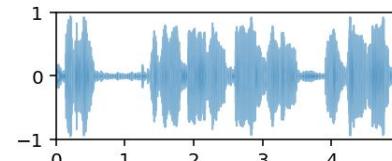
Denoised sample:



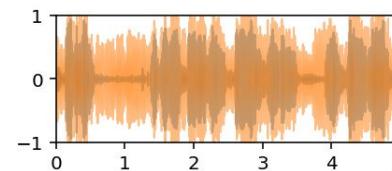
scenario	split	SEANet audio + accel	SEANet audio only
Mixed noise	1	9.9 ± 0.2	8.4 ± 0.2
	2	8.0 ± 0.2	7.9 ± 0.1
	3	8.3 ± 0.1	7.2 ± 0.2
	4	8.8 ± 0.1	8.1 ± 0.1
	5	9.9 ± 0.1	8.4 ± 0.1
	avg.	8.9	8.0
Mixed speech	1	10.1 ± 0.1	-0.9 ± 0.1
	2	8.6 ± 0.1	-0.9 ± 0.1
	3	9.2 ± 0.1	-0.7 ± 0.0
	4	9.0 ± 0.2	-1.0 ± 0.1
	5	11.1 ± 0.2	-0.9 ± 0.1
	avg.	9.6	-0.9

Mixed noise example:

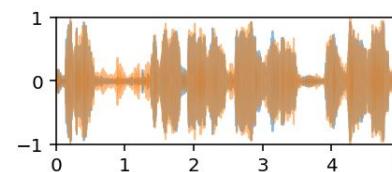
Clean sample:



Noisy sample:

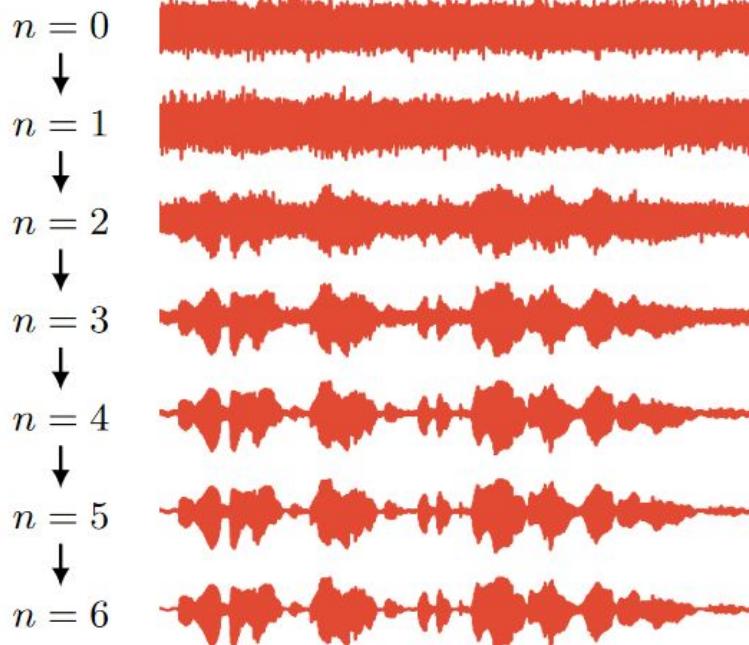


Denoised sample:

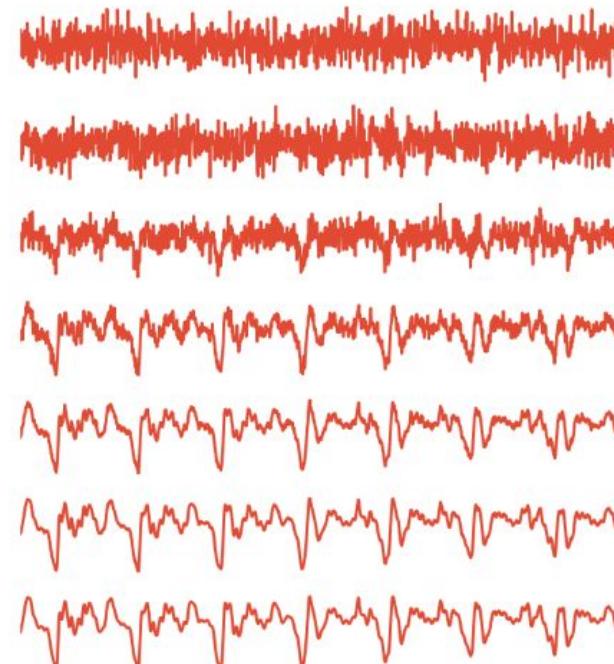


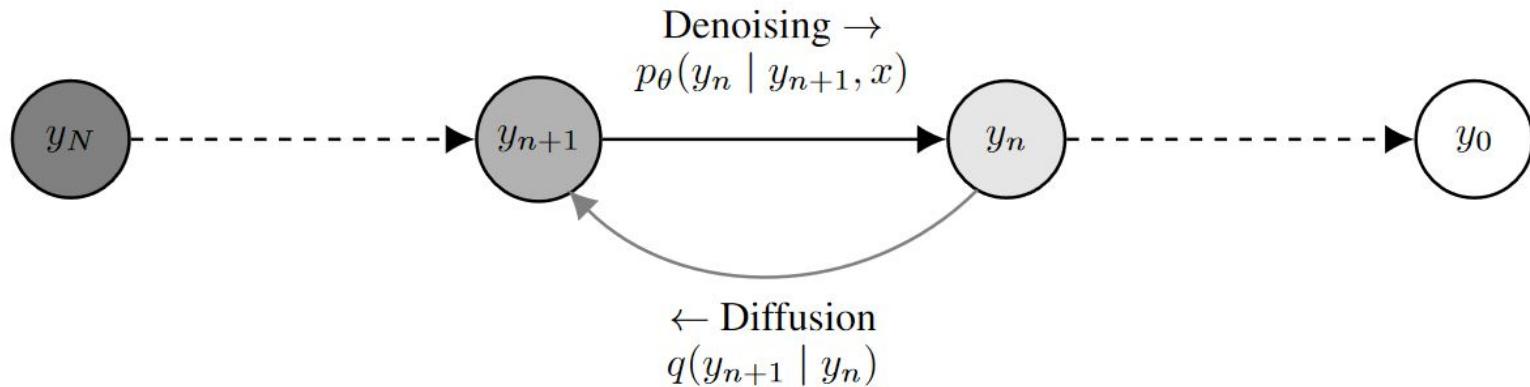
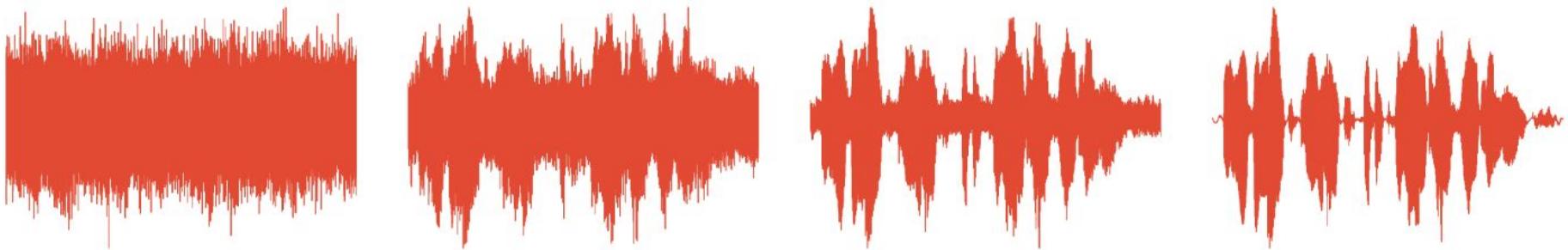
Diffusion Models

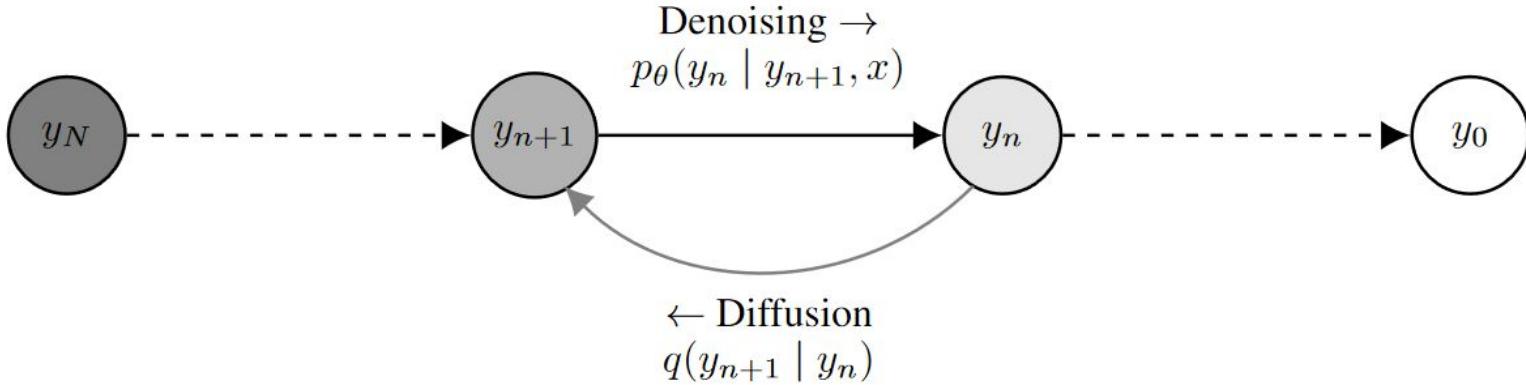
signal



zoomed view of
a 50 ms segment



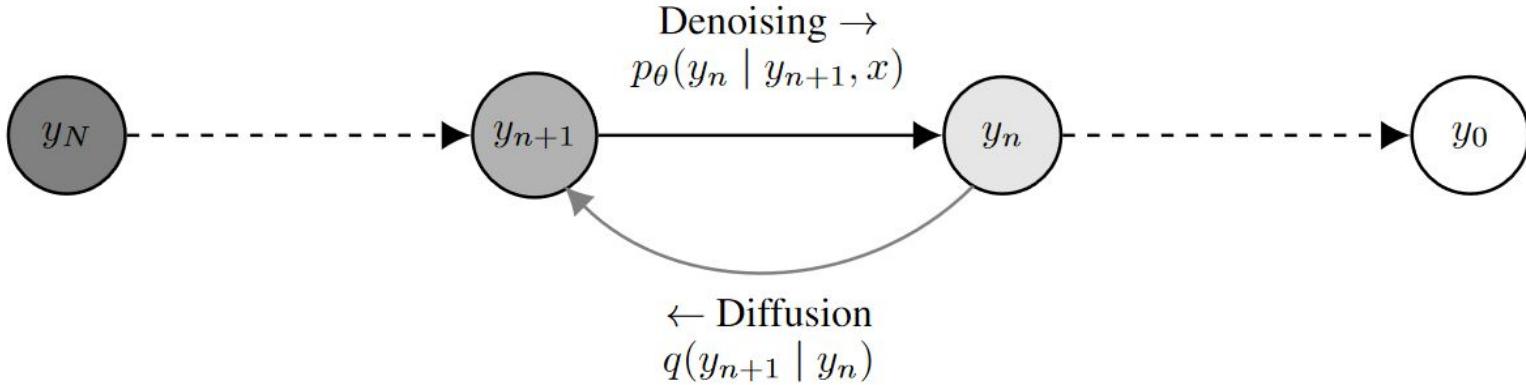




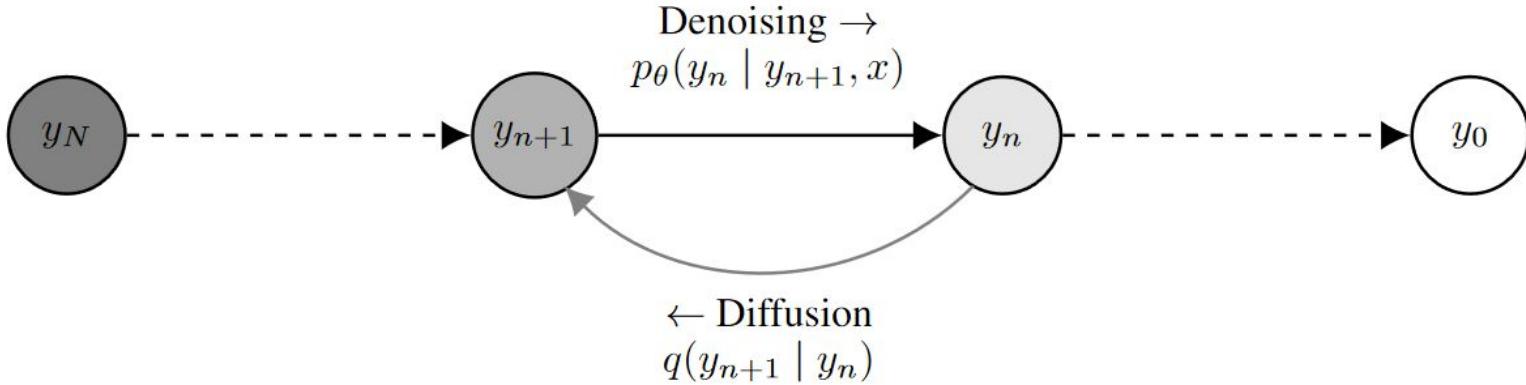
$$q(y_{n+1} | y_n) = \mathcal{N}(y_{n+1} | \sqrt{1 - \beta_n} y_n, \beta_n I)$$

$$p(y_n | y_{n+1}) = \mathcal{N}(y_n | \mu(y_{n+1}, n), \sigma(y_{n+1}, n)^2 I)$$

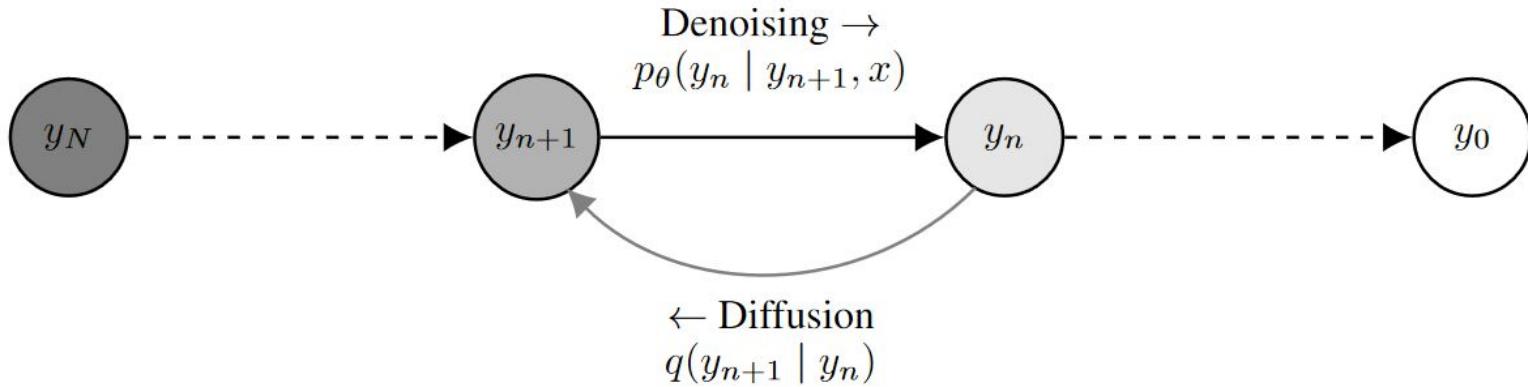
$$p(y_N) = \mathcal{N}(y_N | 0, I)$$



$$p_\theta(y_0) \geq \mathbb{E}_{y_1, \dots, y_N \sim q(y_1, \dots, y_N | y_0)} \log \frac{p_\theta(y_0, y_1, \dots, y_N)}{q(y_1, \dots, y_N | y_0)}$$



$$\mathbb{E}_{y_1, \dots, y_N \sim q(y_1, \dots, y_N | y_0)} \log \frac{p(y_N) p_\theta(y_{N-1} | y_N) \dots p_\theta(y_0 | y_1)}{q(y_1 | y_0) \dots q(y_N | y_{N-1})}$$



$$\mathbb{E}_{y_1, \dots, y_N \sim q(y_1, \dots, y_N \mid y_0)} \left(\sum_{n=0}^{N-1} \log \frac{p_\theta(y_n \mid y_{n+1})}{q(y_{n+1} \mid y_n)} + p(y_N) \right)$$

References (ICLR 2021)

- WaveGrad (Google Brain)
- DiffWave (Baidu Research, NVIDIA)
- For images: Score-based Generative Modeling Through Stochastic Differential Equations (Stanford University, Google Brain)

Thank You