

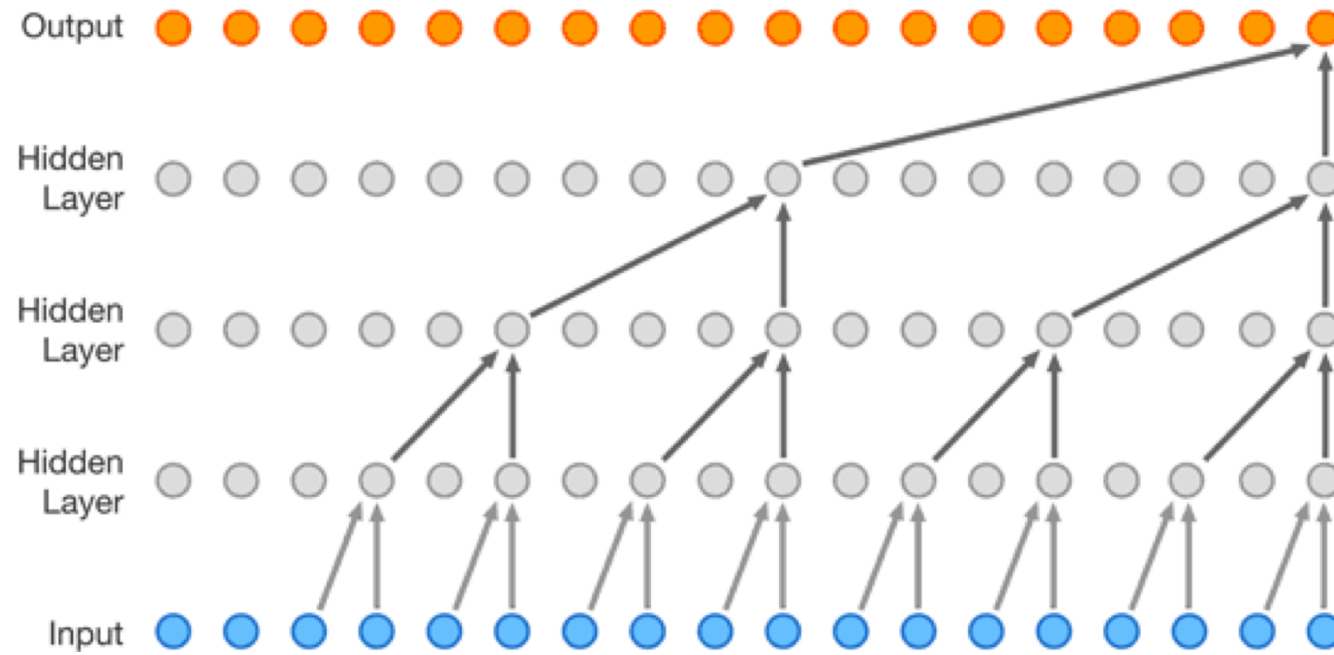
ICML 2019

FloWaveNet: A Generative Flow for Raw Audio

Sungwon Kim¹, Sang-gil Lee¹, Jongyoon Song¹, Jaehyeon Kim², Sungron Yoon^{1,3}

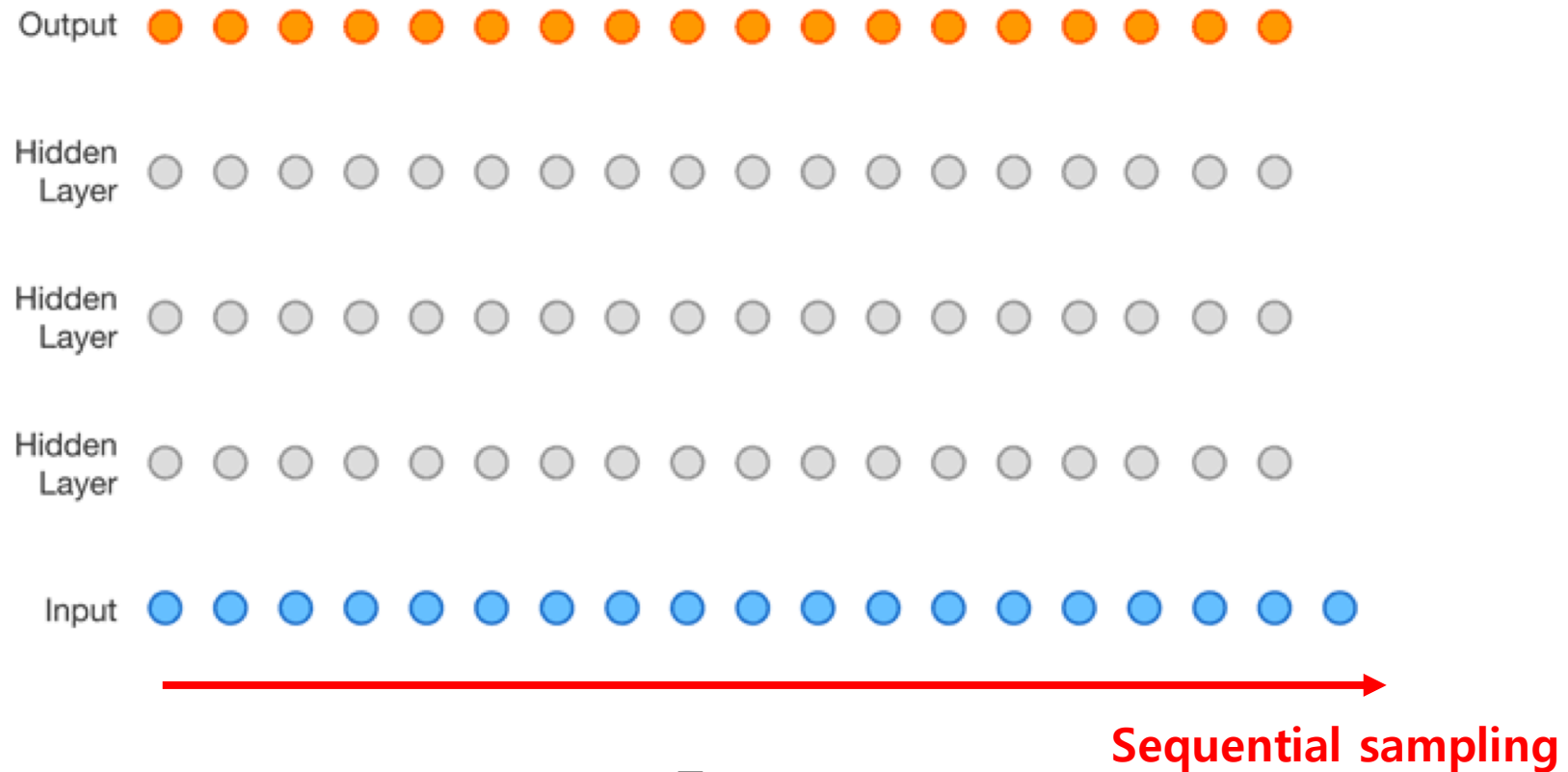
¹Seoul National University, ²Kakao Corporation,
³ASRI, INMC, Institute of Engineering Research, Seoul National University

WaveNet



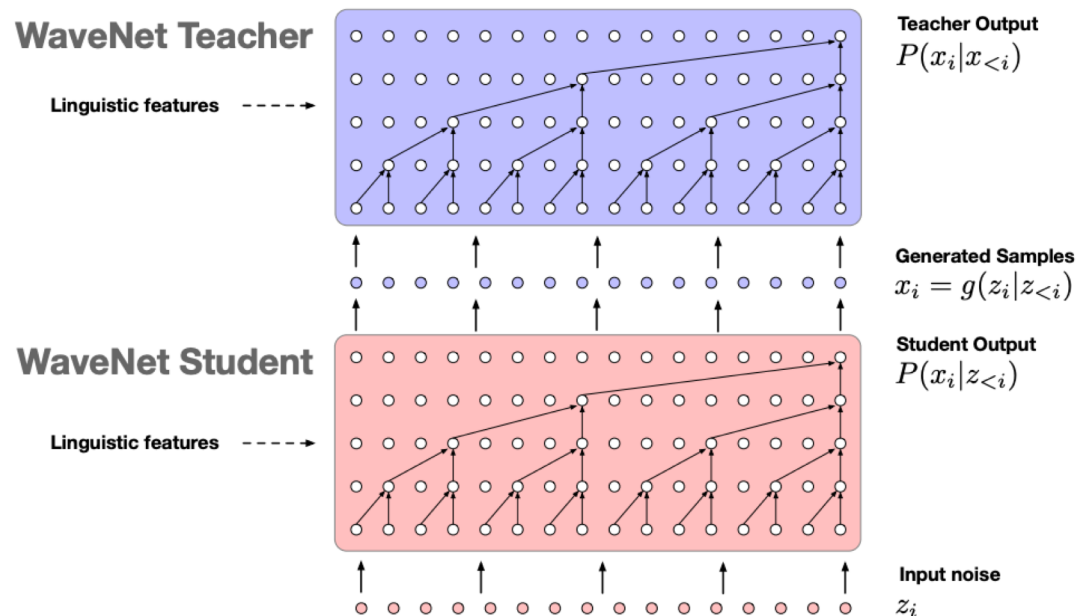
$$\log p_X(x_{1:T}) = \sum_{t=1}^T \log p_X(x_t | x_{<t})$$

WaveNet



$$\log p_X(x_{1:T}) = \sum_{t=1}^T \log p_X(x_t | x_{<t})$$

Previous parallel speech synthesis models



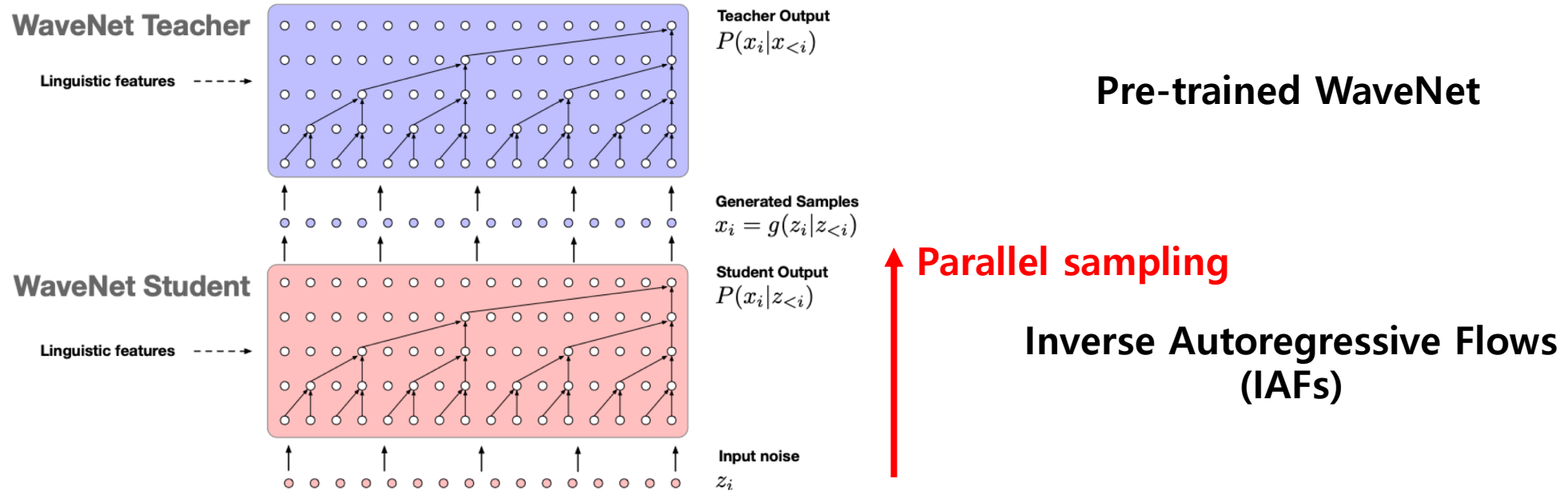
Pre-trained WaveNet

**Inverse Autoregressive Flows
(IAFs)**

Probability Density Distillation

$$KL(P_S(x) || P_T(x))$$

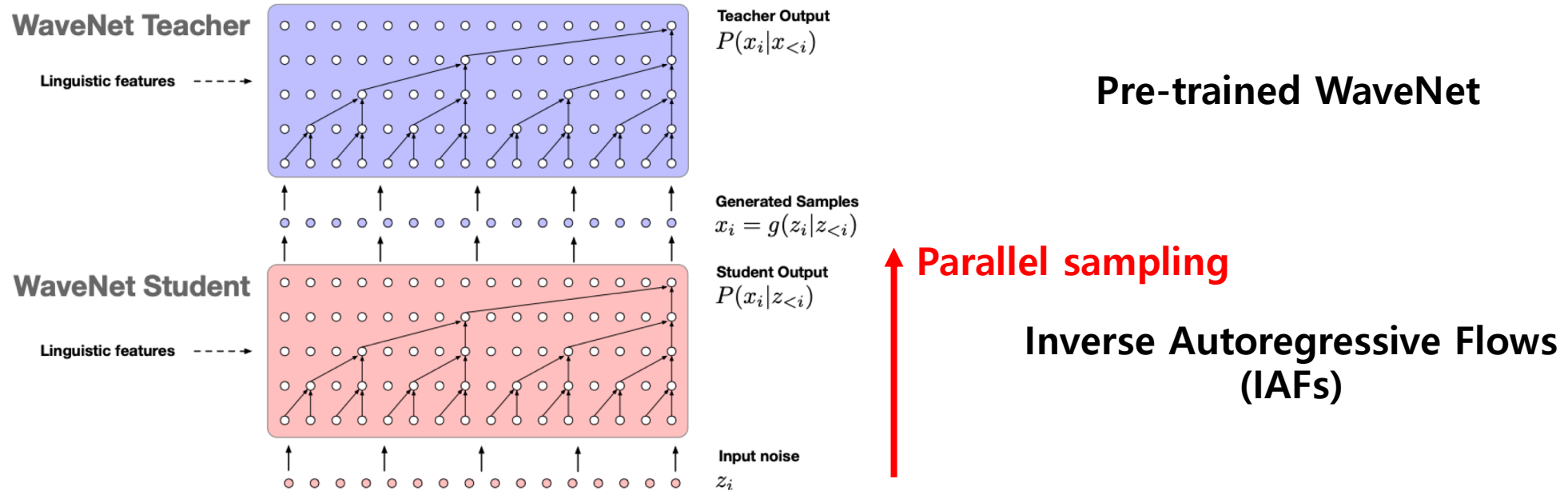
Previous parallel speech synthesis models



Probability Density Distillation

$$KL(P_S(x) || P_T(x))$$

Previous parallel speech synthesis models



Probability Density Distillation

$$KL(P_S(x) || P_T(x))$$

+

Power Loss
Perceptual Loss
Contrastive Loss
Frame Loss

Our Objectives

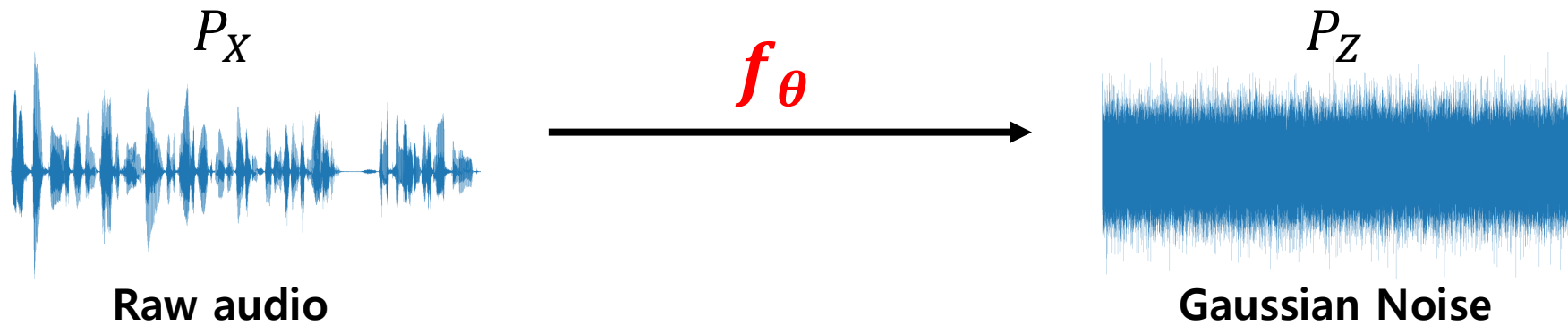
- Simplify the training procedure for parallel sampling
- Maintain the quality of speech samples

Our Objectives

- Simplify the training procedure for parallel sampling
- Maintain the quality of speech samples

Flow-based generative models for raw audio!

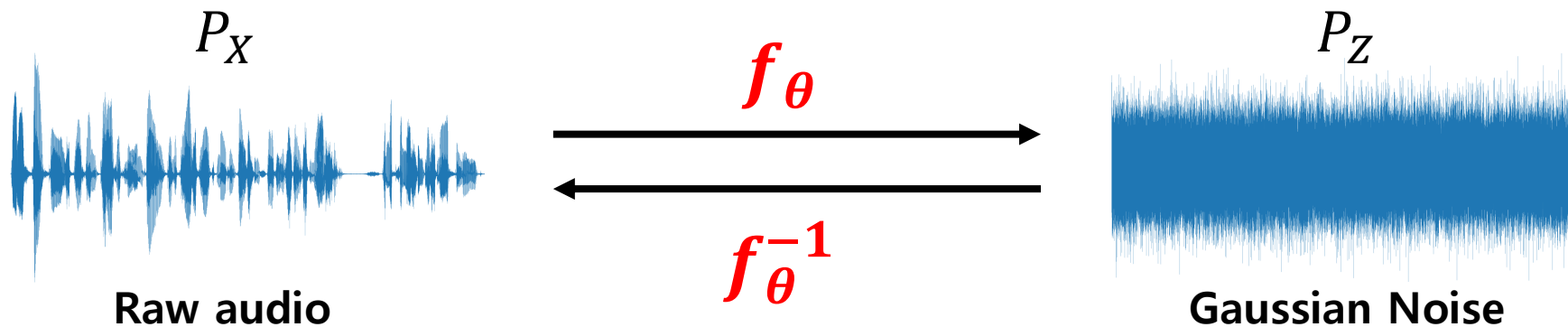
FloWaveNet



Training phase

$$\log p_X(x_{1:T}) = \log p_Z(f_\theta(x_{1:T})) + \log \det \left(\frac{\partial f_\theta(x)}{\partial x} \right)$$

FloWaveNet



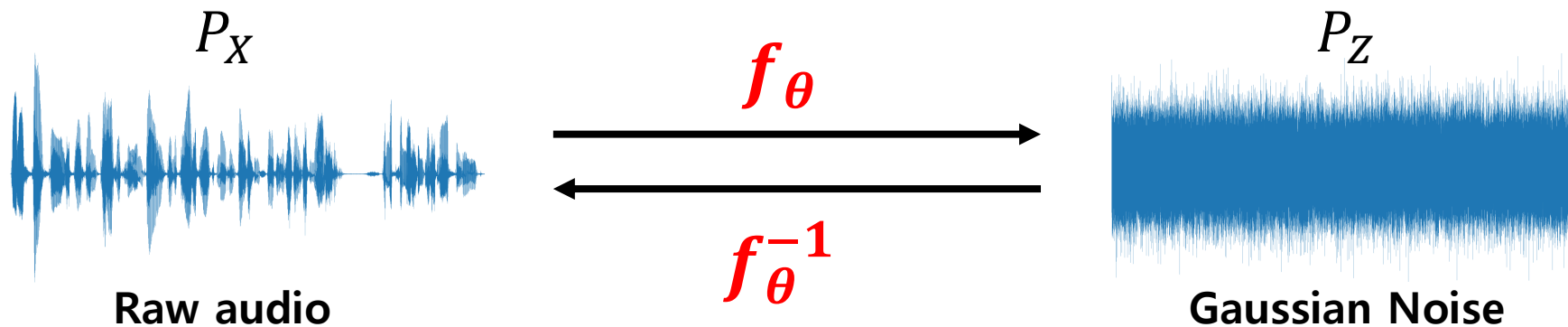
Training phase

$$\log p_X(x_{1:T}) = \log p_Z(f_\theta(x_{1:T})) + \log \det \left(\frac{\partial f_\theta(x)}{\partial x} \right)$$

Sampling phase

$$z = z_{1:T} \sim P_Z(z) = N(0, I), \quad \hat{x} = f_\theta^{-1}(z)$$

FloWaveNet



Training phase

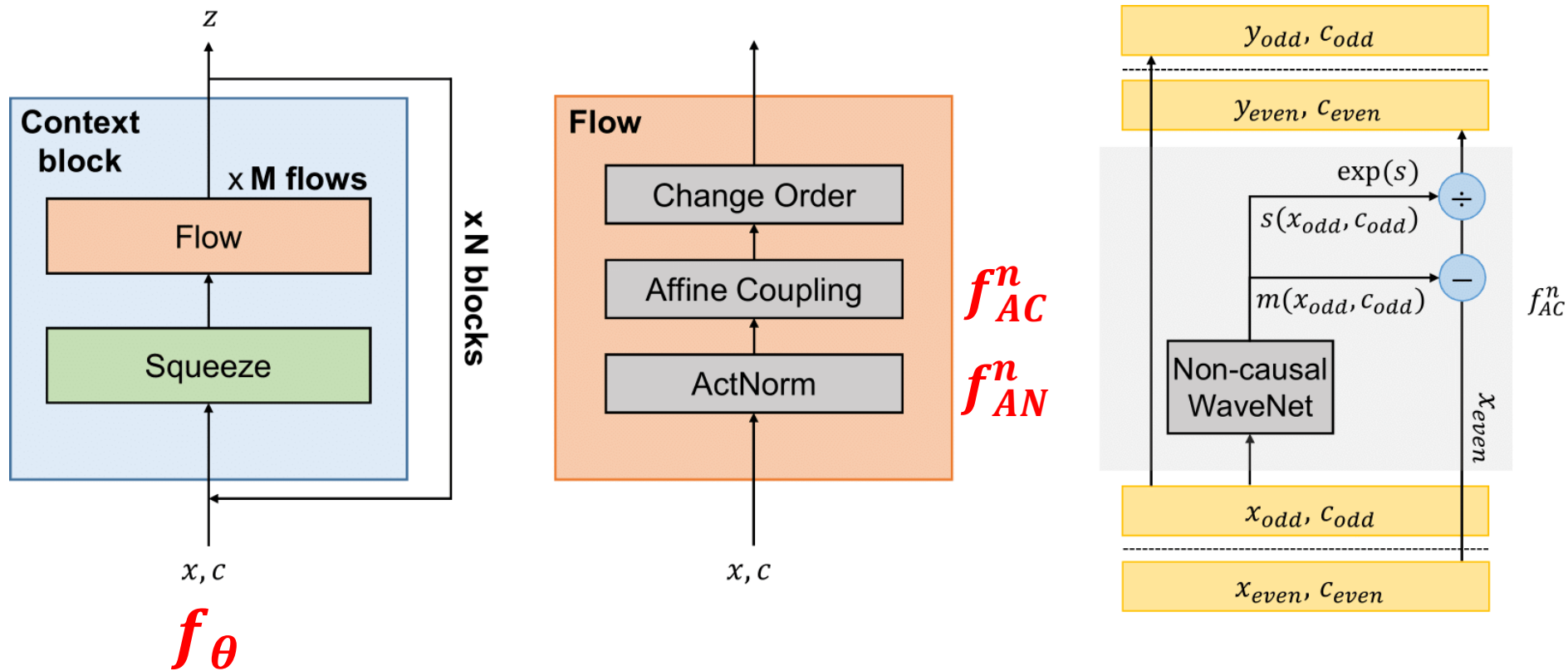
$$\log p_X(x_{1:T}) = \log p_Z(f_\theta(x_{1:T})) + \log \det \left(\frac{\partial f_\theta(x)}{\partial x} \right)$$

Sampling phase

$$z = z_{1:T} \sim P_Z(z) = N(0, I), \quad \hat{x} = f_\theta^{-1}(z)$$

Both the transformation f_θ and f_θ^{-1} are designed to be computed efficiently
→ Efficient training & Parallel sampling

FloWaveNet



$$\log p_X(x_{1:T}) = \log p_Z(f(x_{1:T})) + \sum_n \log \det \left(\frac{\partial (f_{AC}^n \cdot f_{AN}^n)(x)}{\partial x} \right)$$

Mean Opinion Scores



METHODS	5-SCALE MOS	TEST CLL
GROUND TRUTH	4.67 ± 0.076	
MoL WAVE NET	4.30 ± 0.110	4.6546
GAUSSIAN WAVE NET	4.46 ± 0.100	4.6526
GAUSSIAN IAF	3.75 ± 0.159	
FLOWAVE NET	3.95 ± 0.154	4.5457

FloWaveNet \geq Gaussian IAF

Sampling speed

METHODS	ITER/SEC	SAMPLES/SEC
WAVENET	N/A	172
PARALLEL WAVENET	N/A	500K
GAUSSIAN WAVENET	1.329	44
GAUSSIAN IAF	0.636	470K
FLOWAVENET	0.714	420K

FloWaveNet \cong Gaussian IAF \cong Parallel WaveNet >> Autoregressive WaveNet

1000s times faster

Conclusion

- FloWaveNet produces **high quality audio samples** as well as previous distilled models.
- FloWaveNet synthesizes audio samples **in parallel**
 - w/o well pre-trained WaveNet (No distillation!)
 - w/o auxiliary loss terms



Demo page



Code

Poster 6/12 6:30 PM @Pacific Ballroom #2

Thank You!

