# Tutorial on end-to-end text-to-speech synthesis

## Part 1 – Neural waveform modeling

Xin WANG

National Institute of Informatics, Japan

2019-01-27

contact: wangxin@nii.ac.jp
we welcome critical comments, suggestions, and discussion

# SELF-INTRODUCTION

- 国立情報学研究所
- 山岸研究室
- 特任研究員

ワン　シン
王　　鑫

- PhD (2015-2018)　総研大・国立情報学研究所
- M.A (2012-2015)　中国科学技術大学

- http://tonywangx.github.io

## XW Research Blog

Research　　Code/Scripts　　Slides

### Introduction

I'm Xin Wang, a student from Yamagishi Lab, National Institute of Informatics, Japan.
If you have any comment and question, please send email to wangxin ~a-t~ nii ~dot~ ac ~dot~ jp.

#### Contents
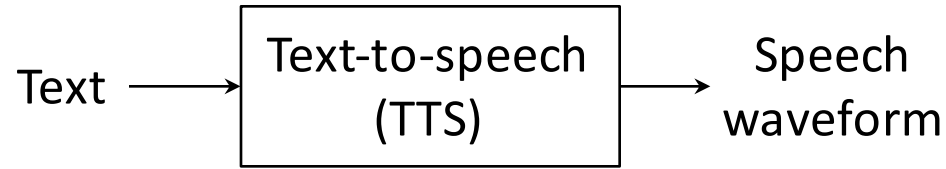- PhD Thesis
- Neural souce-filter waveform model

# CONTENTS

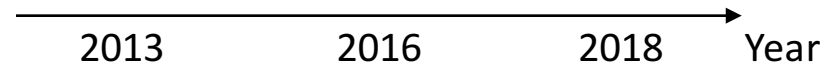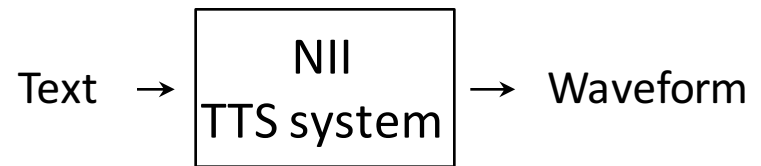- Introduction: text-to-speech synthesis

- Neural waveform models

- Summary & software

# INTRODUCTION

## Text-to-speech synthesis (TTS) [1,2]

Text $\longrightarrow$ | Text-to-speech (TTS) | $\longrightarrow$ Speech waveform

## NII's TTS systems

Text $\rightarrow$ (head icon) $\rightarrow$ Waveform

Text $\rightarrow$ | NII TTS system | $\rightarrow$ Waveform

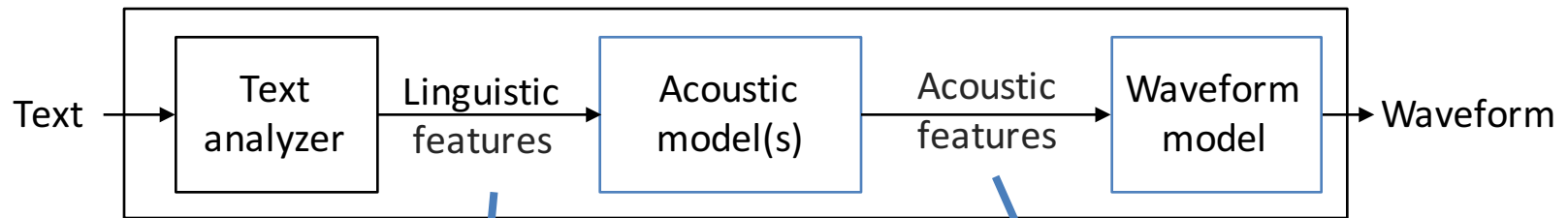2013          2016          2018          Year

[1] P. Taylor. Text-to-Speech Synthesis. Cambridge University Press, 2009.
[2] T. Dutoit. An Introduction to Text-to-speech Synthesis. Kluwer Academic Publishers, Norwell, MA, USA, 1997.
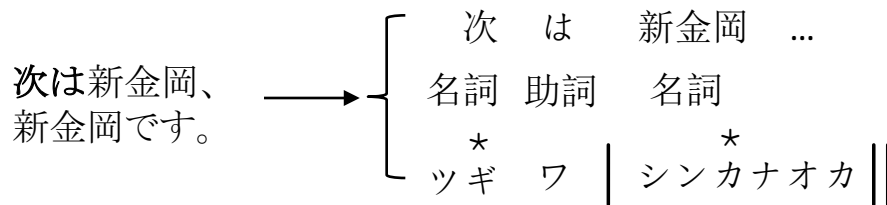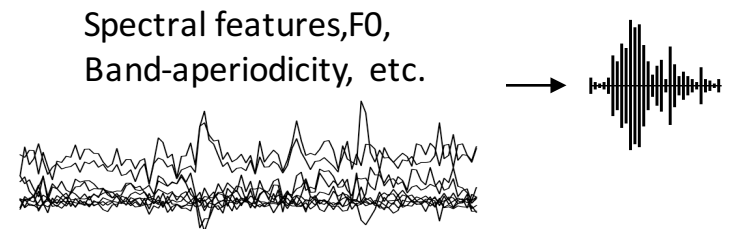
# INTRODUCTION
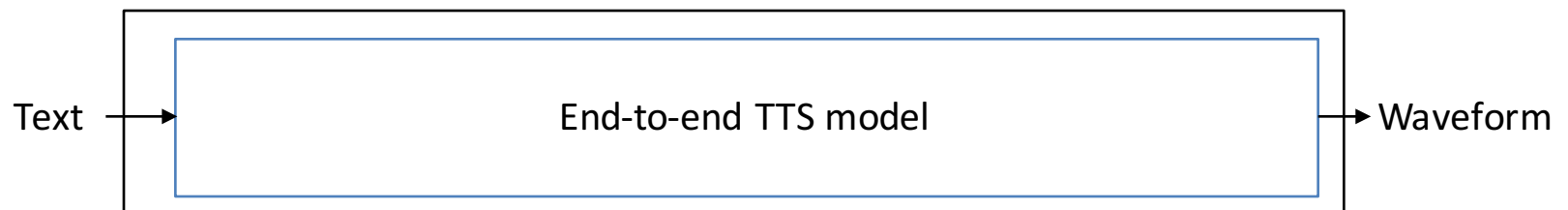
## Text-to-speech synthesis (TTS)

❑ Architectures

Text → **Text analyzer** → Linguistic features → **Acoustic model(s)** → Acoustic features → **Waveform model** → Waveform

**Linguistic features**

次は新金岡、
新金岡です。
→

```
┌  次    は     新金岡   …
│  名詞  助詞   名詞
│        *              *
└  ツギ   ワ  │ シンカナオカ ‖
```
→

**Acoustic features**

Spectral features,F0,
Band-aperiodicity, etc.
→

# INTRODUCTION

## Text-to-speech synthesis (TTS)

❑ Architectures
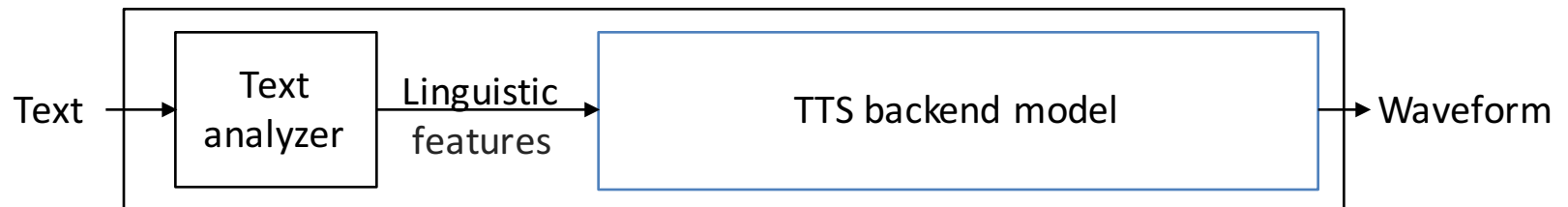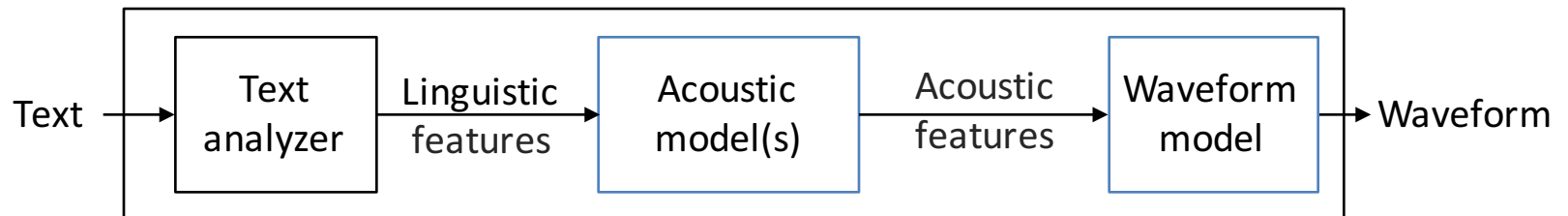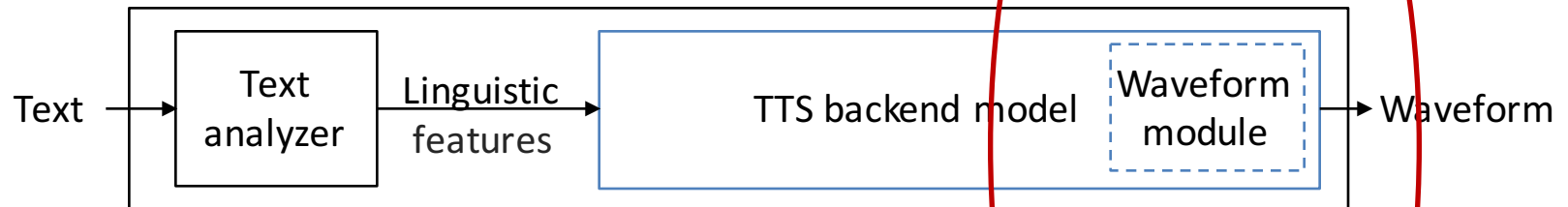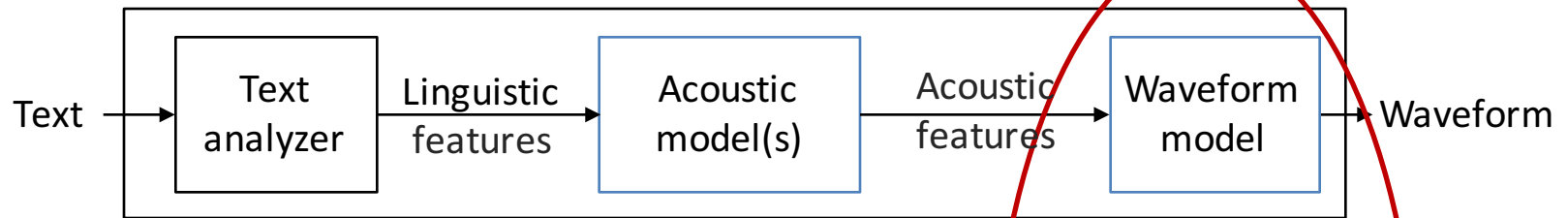
# INTRODUCTION

## Text-to-speech synthesis (TTS)

❑ Architectures

**Tutorial part 1: neural waveform modeling**



**Tutorial part 2: end-to-end TTS**

# CONTENTS

- Introduction: text-to-speech synthesis

- Neural waveform modeling
  - Overview
  - Autoregressive models
  - Normalizing flow
  - STFT-based training criterion

- Summary

# OVERVIEW

## Task definition



Spectrogram

# OVERVIEW

## Task definition

Waveform values

$o_1$ $o_2$ $o_3$ $o_4$ $o_T$

①   ②   ③   ④   ○   ○   ○   ○   ○   ○   ○   ○   ○   …   ○   ○   ○   ①

*T* waveform sampling points

Spectrogram

$c_1$ $c_2$ $c_3$ $c_N$

*N* frames

# OVERVIEW

## Task definition



Waveform values

$o_1 \quad o_2 \quad o_3 \quad o_4 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad o_T$

①　②　③　④　○　○　○　○　○　○　○　○　○　…　○　○　○　①

Neural waveform models

$$p(\boldsymbol{o}_{1:T}|\boldsymbol{c}_{1:N}; \boldsymbol{\Theta})$$

$\boldsymbol{c}_1 \quad \boldsymbol{c}_2 \quad \boldsymbol{c}_3 \qquad\qquad \boldsymbol{c}_N$

# OVERVIEW

## Naïve model

❑ Simple network + maximum-likelihood

- Feedforward network

$$p(\boldsymbol{o}_4|\boldsymbol{c}_{1:N};\boldsymbol{\Theta})$$

Waveform values

Distribution

Hidden layers

Up-sampling

$\boldsymbol{c}_1 \quad \boldsymbol{c}_2 \quad \boldsymbol{c}_3 \qquad \boldsymbol{c}_N$

# OVERVIEW

## Naïve model

❑ Simple network + maximum-likelihood

- RNN

## Naïve model

❑ Simple network + maximum-likelihood

- Dilated CNN [1,2]

Waveform values

Distribution

Hidden layers

Up-sampling

$c_1$  $c_2$  $c_3$    $c_N$

[1] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122, 2015.
[2] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. Lang. Phoneme recognition using time-delay neural networks. IEEE Transactions on Acoustics, Speech, and Signal Processing, 37(3):328–339, 1989.

# OVERVIEW

## Naïve model

mismatch

Weak temporal model <---------------> strongly correlated data $o_{1:T}$

Waveform values

$o_1$   $o_2$   $o_3$   $o_4$                                    $o_T$

①   ②   ③   ④   ○   ○   ○   ○   ○   ○   ○   ○   ○   ...   ○   ○   ○   ①

Distribution

Hidden layers

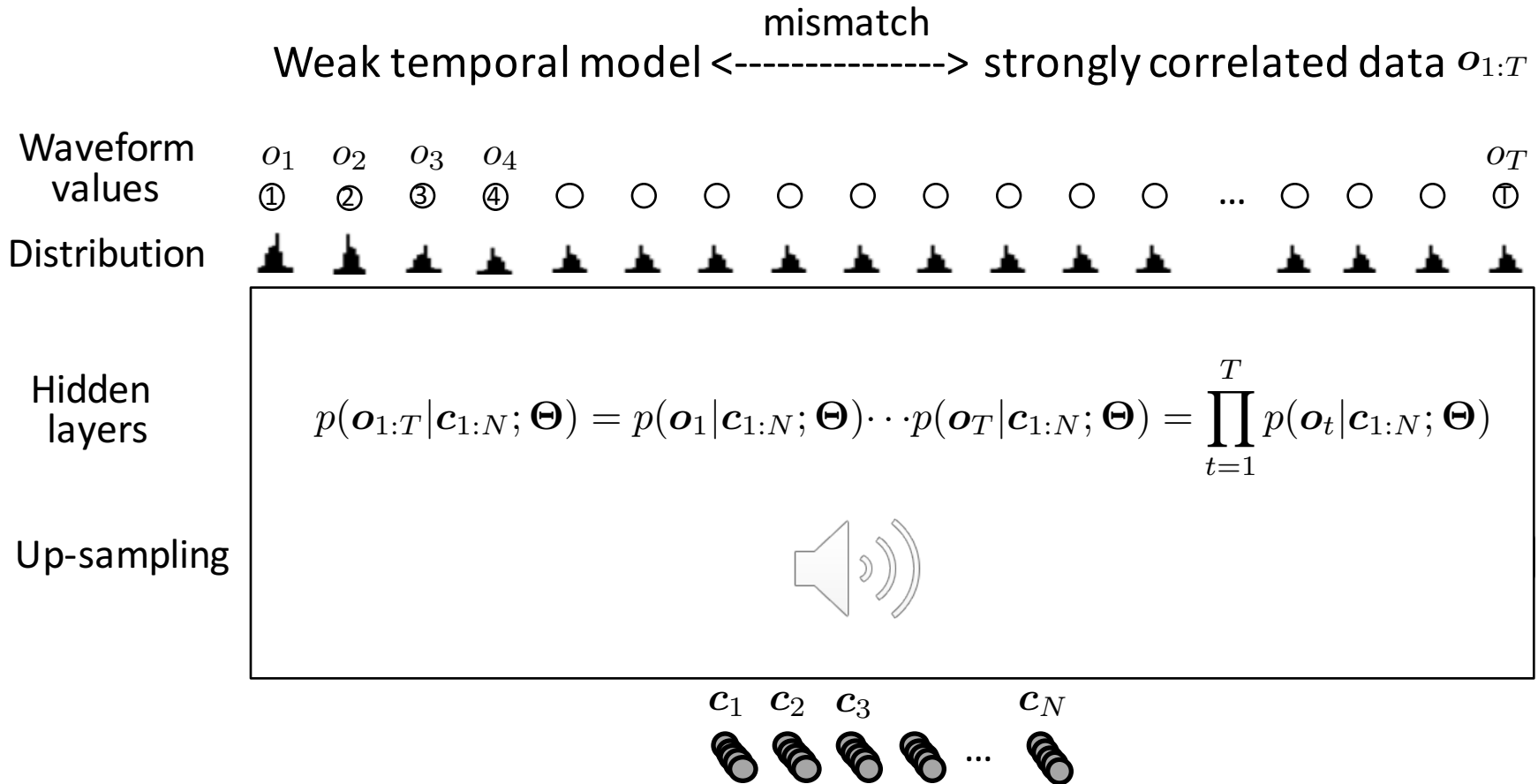$$p(\boldsymbol{o}_{1:T}|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) = p(\boldsymbol{o}_1|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) \cdots p(\boldsymbol{o}_T|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) = \prod_{t=1}^{T} p(\boldsymbol{o}_t|\boldsymbol{c}_{1:N};\boldsymbol{\Theta})$$

Up-sampling

$\boldsymbol{c}_1$   $\boldsymbol{c}_2$   $\boldsymbol{c}_3$       $\boldsymbol{c}_N$

# OVERVIEW

## Towards better models

mismatch

Weak temporal model <---------------> strongly correlated data



Better model    Better criterion    Easier target    Combination

# OVERVIEW



Speech/signal processing

$o_{1:T}$
Knowledge-based Transform
Criterion
Model
$c_{1:N}$

LPCNet   ExcitNet

GlotNet   LP-WaveNet

Subband WaveNet

WaveNet   SampleRNN

Autoregressive models

$o_{1:T}$
Criterion
Model
$c_{1:N}$

WaveRNN   FFTNet

Universal neural vocoder

$o_{1:T}$
Transform
Criterion
Model
$c_{1:N}$

Parallel WaveNet

ClariNet

FloWaveNet

WaveGlow

Normlizing flow / feature transformation

RNN + STFT loss

Neural source-filter model

$o_{1:T}$
Criterion
Model
$c_{1:N}$

Training criterion in frequency domain
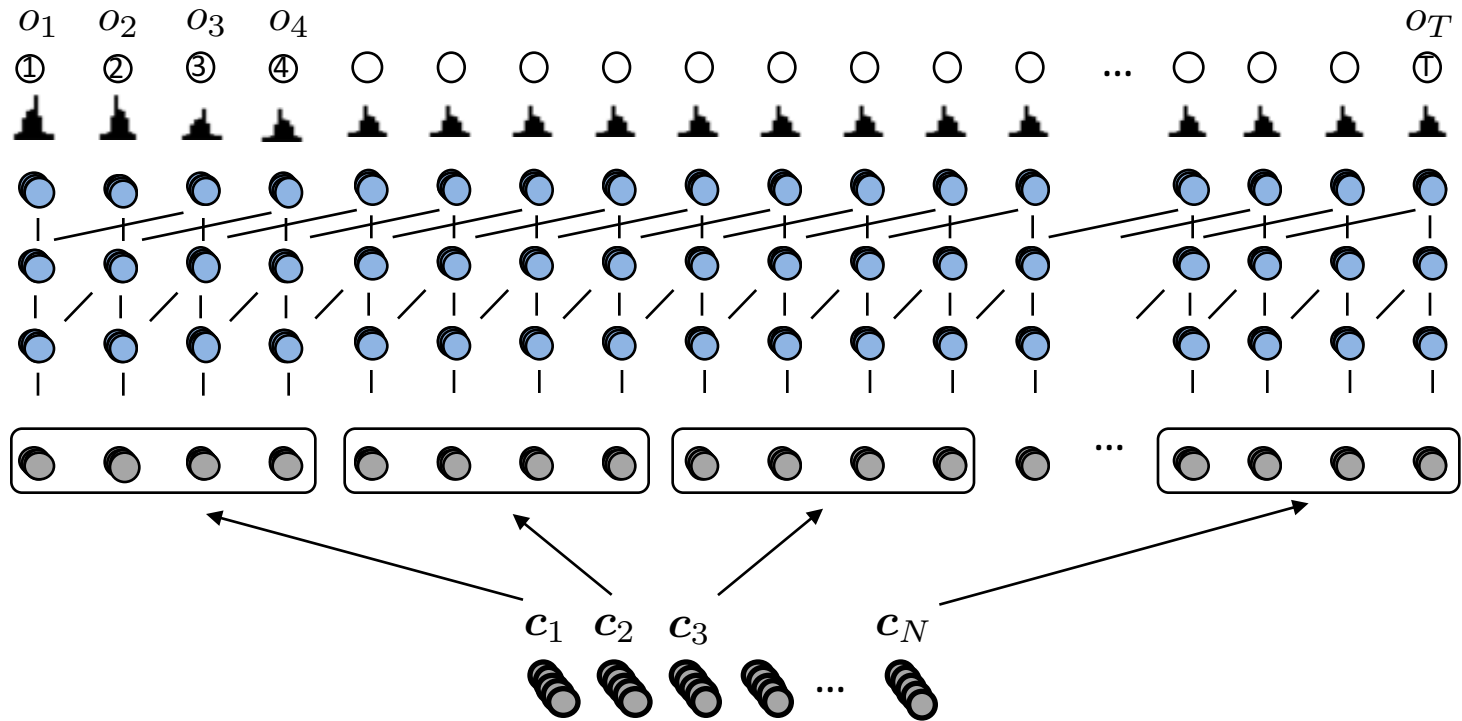
# CONTENTS

- Introduction: text-to-speech synthesis

- Neural waveform models
  - Overview
  - Autoregressive models
  - Normalizing flow
  - STFT-based training criterion

- Summary & software

## Core idea

❑ Naïve model

$$p(\boldsymbol{o}_{1:T}|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) = \prod_{t=1}^{T} p(\boldsymbol{o}_t|\boldsymbol{c}_{1:N};\boldsymbol{\Theta})$$

# PART I: AUTOREGRESSIVE MODELS

## Core idea

❑ Autoregressive (AR) model

$$p(\boldsymbol{o}_{1:T}|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) = \prod_{t=1}^{T} p(\boldsymbol{o}_t|\boldsymbol{o}_{1:t-1}, \boldsymbol{c}_{1:N};\boldsymbol{\Theta})$$

# PART I: AUTOREGRESSIVE MODELS

## Core idea

❑ Autoregressive (AR) model

- Training: use natural waveform for feedback (teacher forcing [1])



[1] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. Neural computation, 1(2):270–280, 1989.
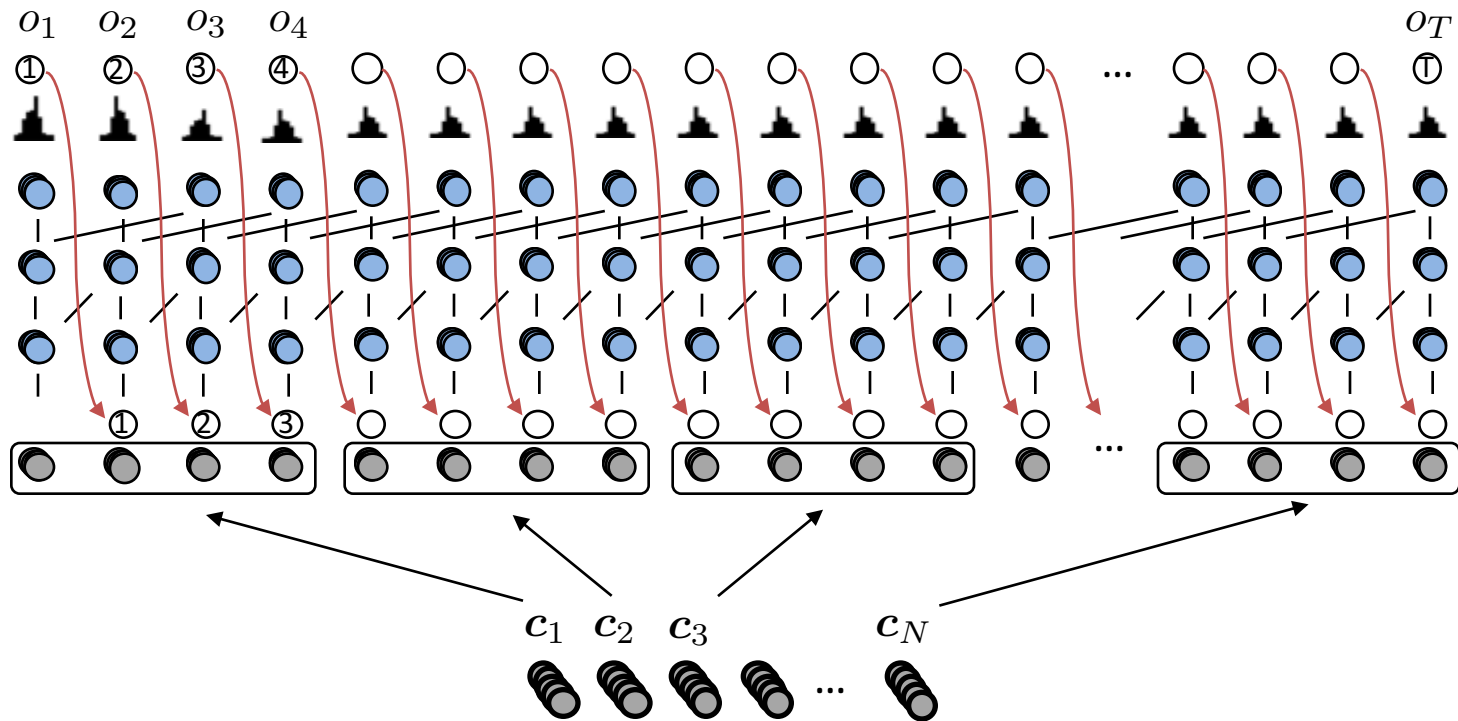
# PART I: AUTOREGRESSIVE MODELS

## Core idea

❑ Autoregressive (AR) model

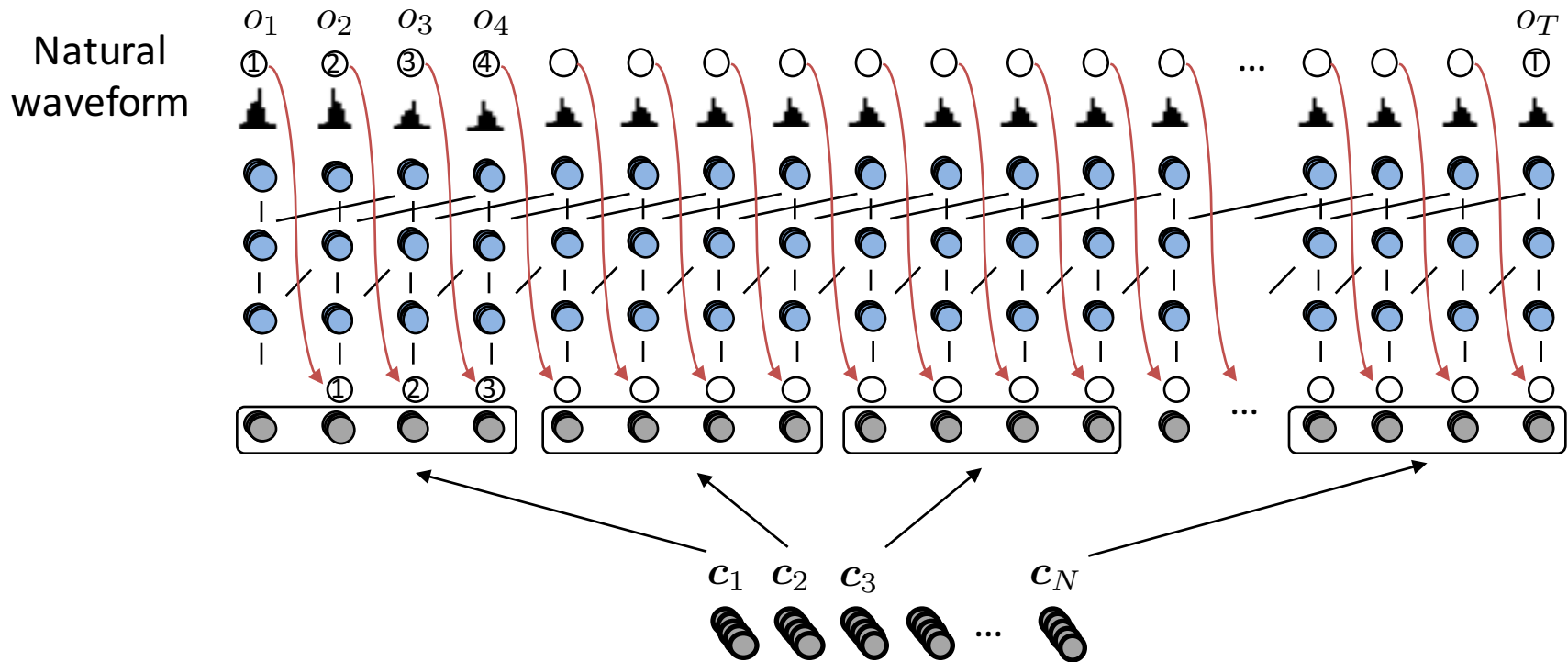- Training: use natural waveform for feedback (teacher forcing [1])



[1] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. Neural computation, 1(2):270–280, 1989.

# PART I: AUTOREGRESSIVE MODELS

## Core idea

❑ Autoregressive (AR) model

- Training: use natural waveform for feedback (teacher forcing [1])

- Generation: $p(\widehat{\boldsymbol{o}}_{1:T}|\boldsymbol{c}_{1:N}; \boldsymbol{\Theta}) = \prod_{t=1}^{T} p(\widehat{\boldsymbol{o}}_t|\widehat{\boldsymbol{o}}_{t-P:t-1}, \boldsymbol{c}_{1:N}; \boldsymbol{\Theta})$
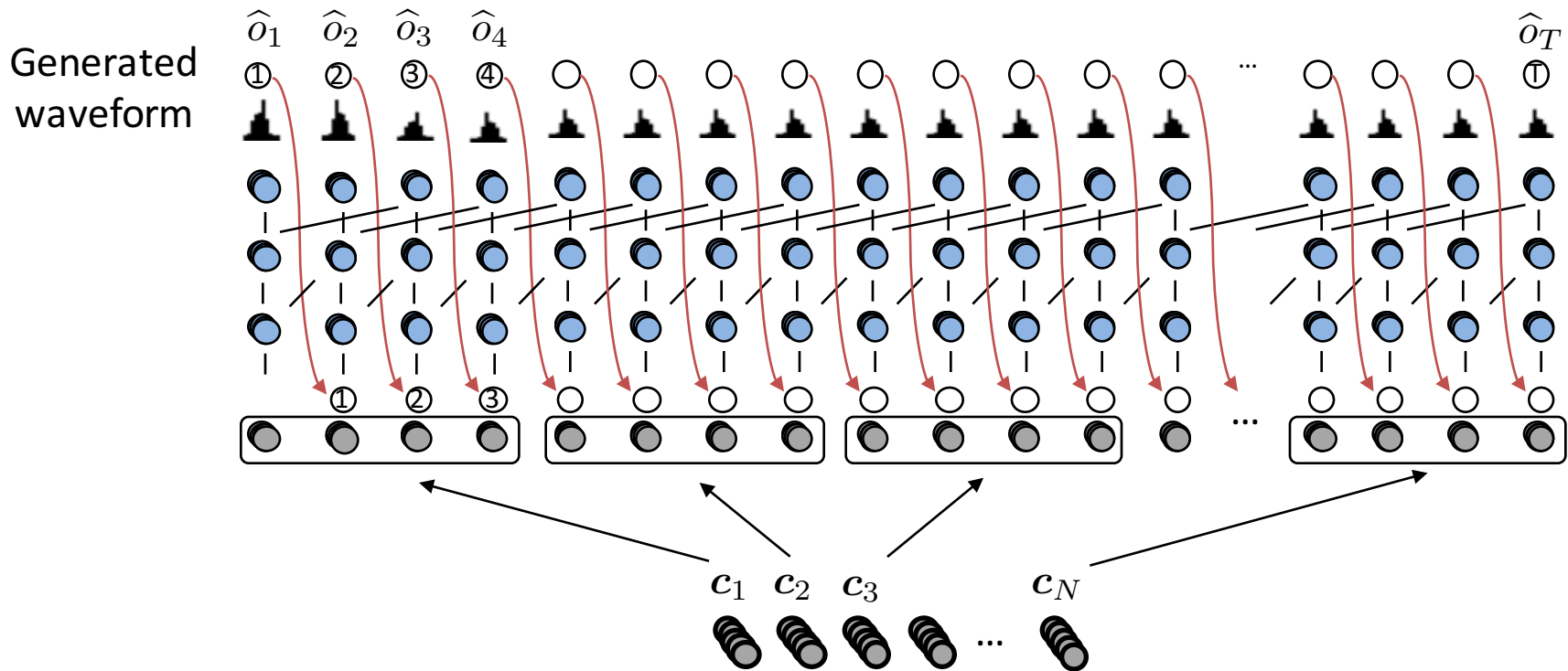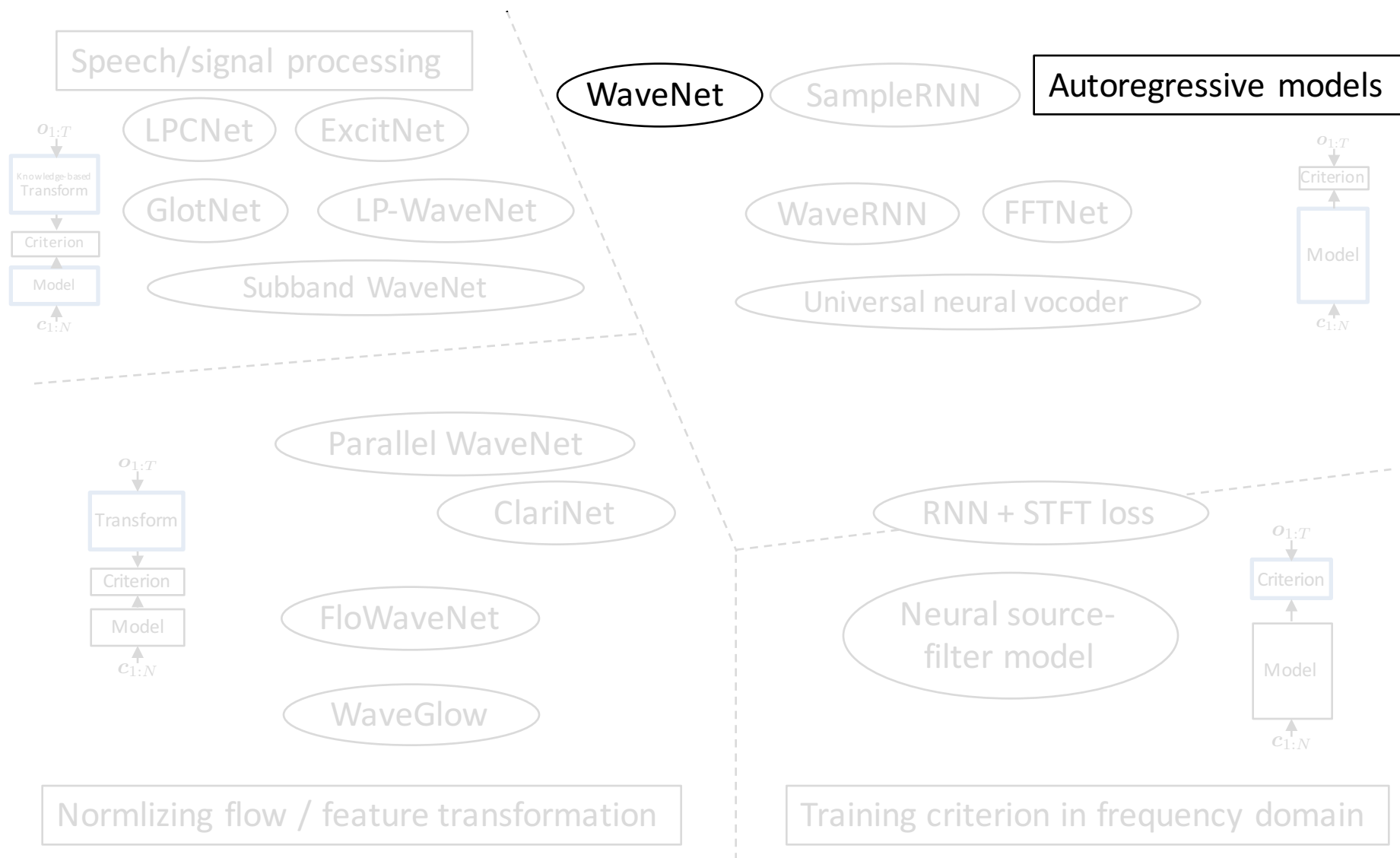
[1] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. Neural computation, 1(2):270–280, 1989.
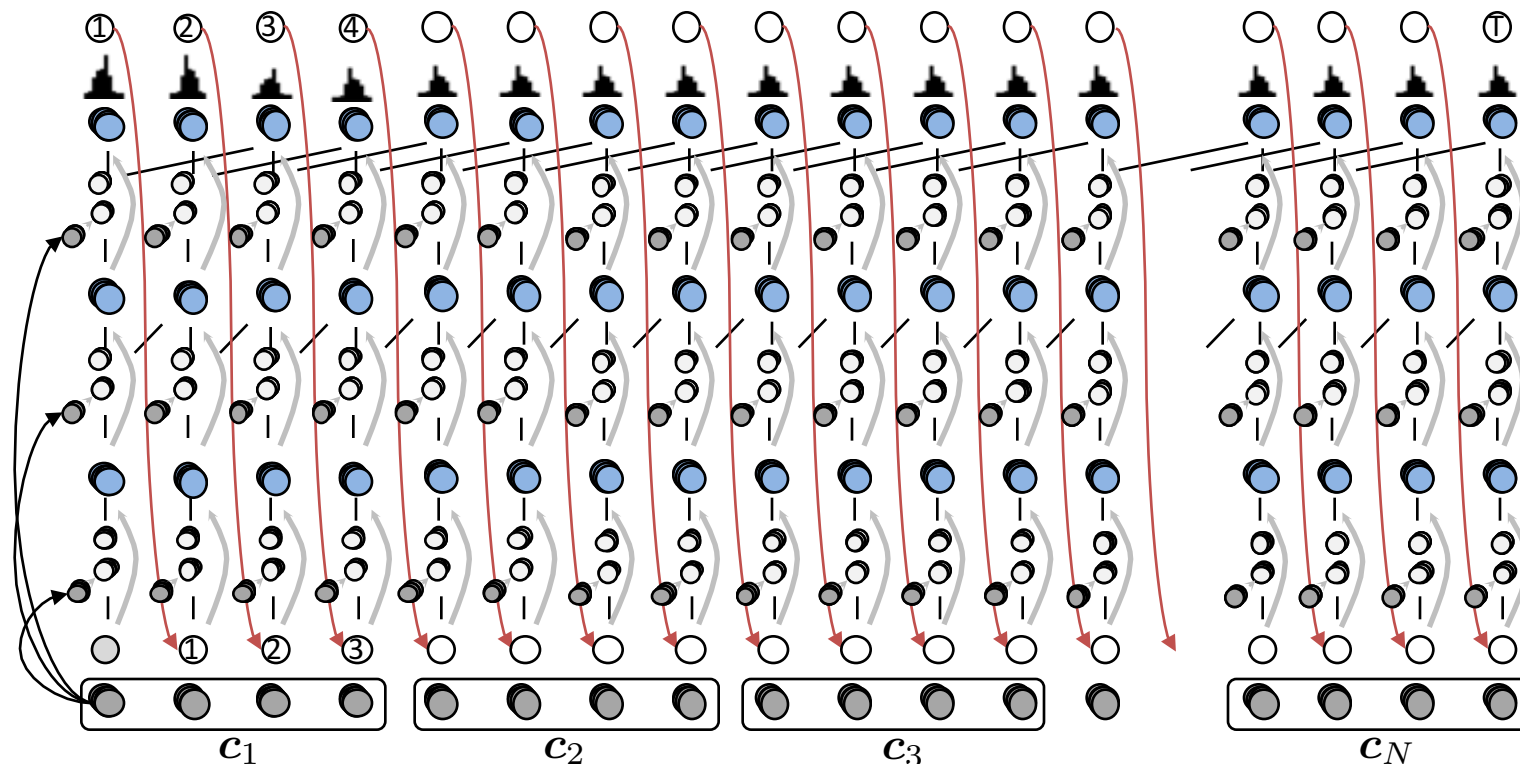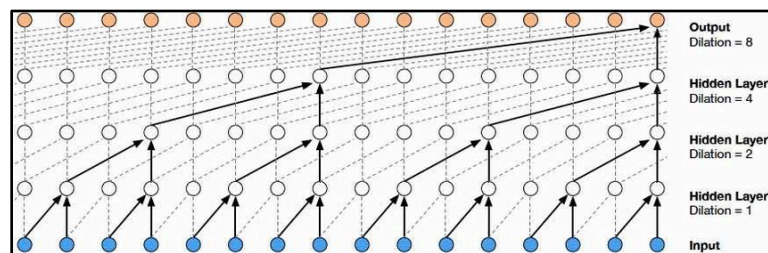
# PART I: AUTOREGRESSIVE MODELS

Speech/signal processing

$o_{1:T}$

Knowledge-based Transform

Criterion

Model

$c_{1:N}$

LPCNet  ExcitNet

GlotNet  LP-WaveNet

Subband WaveNet

WaveNet  SampleRNN

Autoregressive models

WaveRNN  FFTNet

Universal neural vocoder

$o_{1:T}$

Criterion

Model

$c_{1:N}$

Parallel WaveNet

ClariNet

RNN + STFT loss

$o_{1:T}$

Transform

Criterion

Model

$c_{1:N}$

FloWaveNet

WaveGlow

Neural source-filter model

$o_{1:T}$

Criterion

Model

$c_{1:N}$

Normlizing flow / feature transformation

Training criterion in frequency domain

## WaveNet

❑ Network structure



- Details: http://id.nii.ac.jp/1001/00185683/
  http://tonywangx.github.io/pdfs/wavenet.pdf

A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. WaveNet: A generative model for raw audio. arXiv preprint arXiv:1609.03499, 2016.

# PART I: AUTOREGRESSIVE MODELS

Speech/signal processing

WaveNet    SampleRNN

Autoregressive models

$o_{1:T}$

Knowledge-based Transform

Criterion

Model

$c_{1:N}$

LPCNet    ExcitNet

GlotNet    LP-WaveNet

Subband WaveNet

WaveRNN    FFTNet

Universal neural vocoder

$o_{1:T}$

Criterion

Model

$c_{1:N}$

---

❑ **Issues with WaveNet & Sample RNN**

!    Generation is very very … very slow

❑ **Acceleration?**

- **Engineering-based**: parallelize/simplify computation
- **Knowledge-based**:  speech / signal processing theory
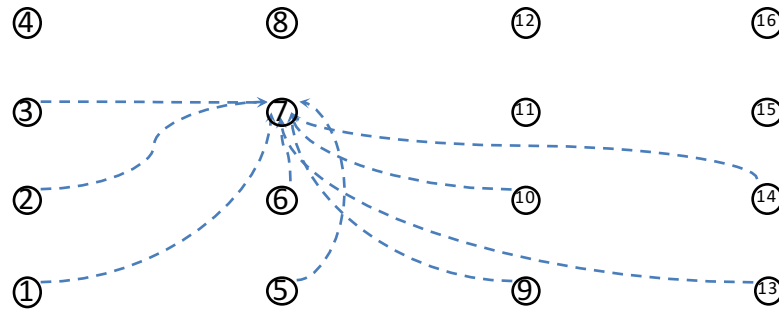
# PART I: AUTOREGRESSIVE MODELS

## WaveRNN

❑ Linear time AR dependency in WaveNet



Waveform values

16 steps to generate

❑ Subscale AR dependency + batch sampling

# PART I: AUTOREGRESSIVE MODELS

## WaveRNN

❑ Linear time AR dependency in WaveNet



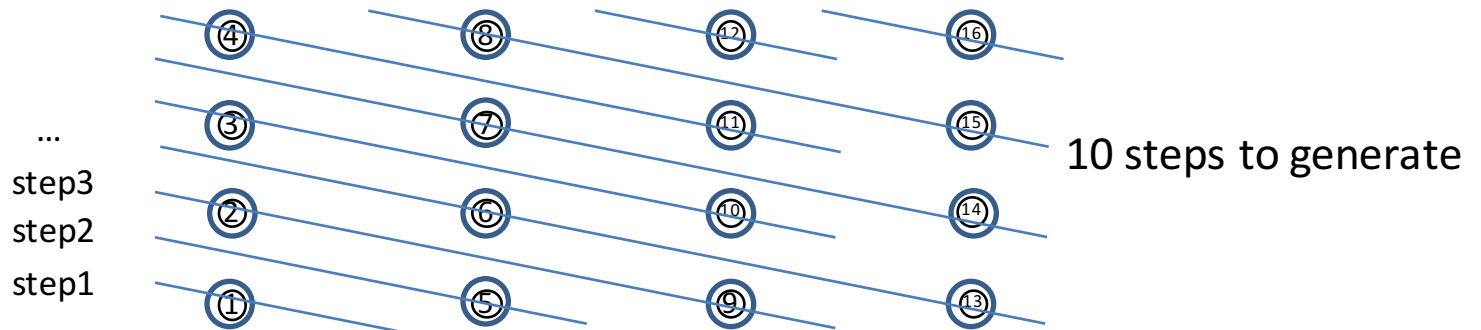16 steps to generate

❑ Subscale AR dependency + batch sampling



10 steps to generate

...
step3
step2
step1
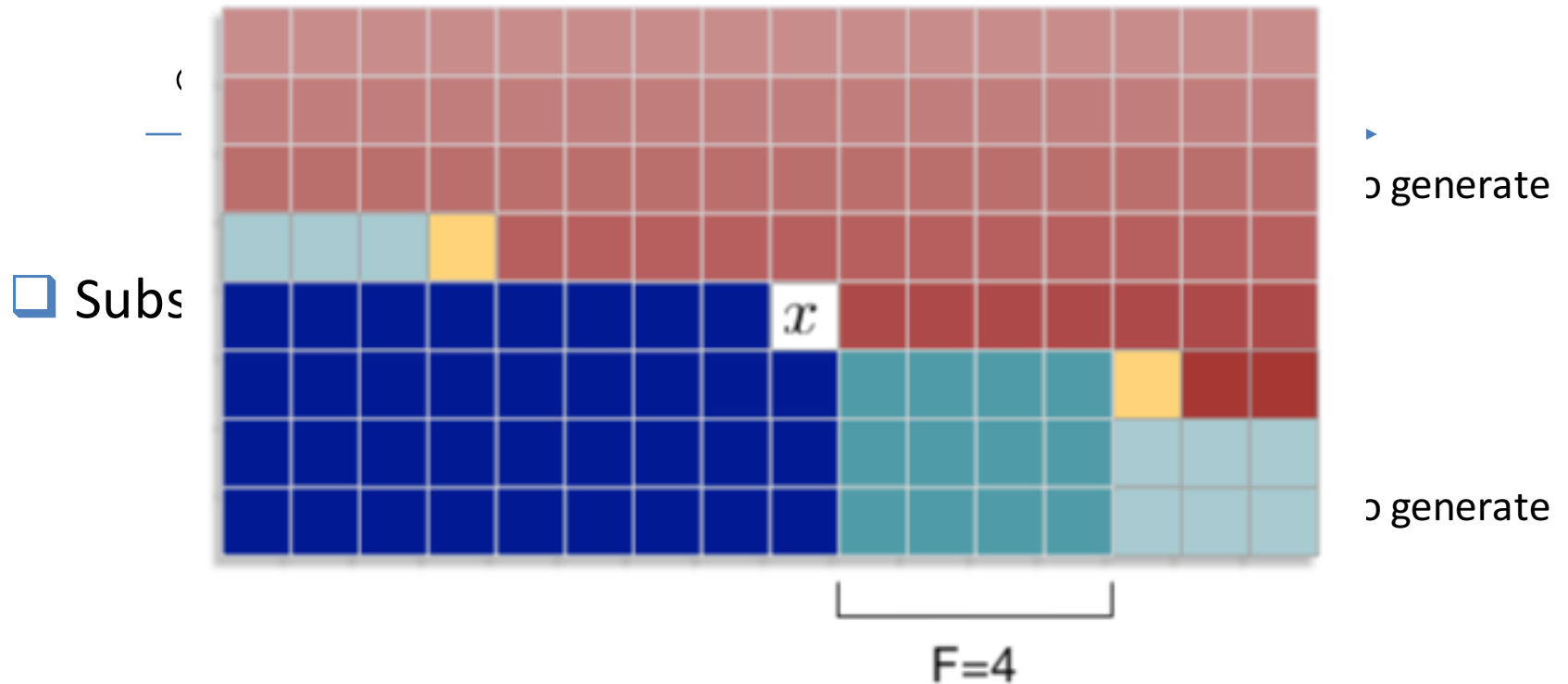
• Generate multiple samples at the same time

## WaveRNN

❑ Linear time AR dependency in WaveNet



- Generate multiple samples at the same time

Speech/signal processing

$o_{1:T}$

Knowledge-based Transform

Criterion

Model

$c_{1:N}$

LPCNet     ExcitNet

GlotNet     LP-WaveNet

Subband WaveNet

WaveNet     SampleRNN

Autoregressive models

$o_{1:T}$

Criterion

Model

$c_{1:N}$

WaveRNN     FFTNet

Universal neural vocoder

❑ AR models

$$p(\boldsymbol{o}_{1:T}|\boldsymbol{c}_{1:N};\boldsymbol{\Theta}) = \prod_{t=1}^{T} p(\boldsymbol{o}_t|\boldsymbol{o}_{1:t-1}, \boldsymbol{c}_{1:N};\boldsymbol{\Theta})$$

!     Generation is still slow

# PART I: AUTOREGRESSIVE MODELS

Speech/signal processing

$o_{1:T}$

Knowledge-based Transform

Criterion

Model

$c_{1:N}$

LPCNet    ExcitNet

GlotNet    LP-WaveNet

Subband WaveNet

WaveNet    SampleRNN

Autoregressive models

$o_{1:T}$

Criterion

Model

$c_{1:N}$

WaveRNN    FFTNet

Universal neural vocoder

$o_{1:T}$

Transform

Criterion

Model

$c_{1:N}$

Non-autoregressive model?

Normlizing flow / feature transformation

# CONTENTS

- Introduction: text-to-speech synthesis

- Neural waveform models
    - Overview
    - Autoregressive models
    - Normalizing flow
    - STFT-based training criterion


- Summary & software

# Part II: Normalizing Flow-Based Models
## General idea

$$p_o(\boldsymbol{o}|\boldsymbol{c};\boldsymbol{\Theta},\boldsymbol{\Psi}) = p_z(\boldsymbol{z} = f_{\boldsymbol{\Psi}}^{-1}(\boldsymbol{o})|\boldsymbol{c};\boldsymbol{\Theta}) \left| \det \frac{\partial f_{\boldsymbol{\Psi}}^{-1}(\boldsymbol{o})}{\partial \boldsymbol{o}} \right|$$



*training*          *generation*

- *o* with strong temporal correlation ---> *z* with weak temporal correlation
- Principle of changing random variable

# PART II: NORMALIZING FLOW-BASED MODELS
## General idea

$$p_o(\boldsymbol{o}|\boldsymbol{c}; \boldsymbol{\Psi}_1, \cdots, \boldsymbol{\Psi}_M) = p_z(\boldsymbol{z} = f_{\boldsymbol{\Psi}_1}^{-1} \circ \cdots \circ f_{\boldsymbol{\Psi}_{M-1}}^{-1} \circ f_{\boldsymbol{\Psi}_M}^{-1}(\boldsymbol{o})|\boldsymbol{c}) \prod_{m=1}^{M} \left| \det \boldsymbol{J}(f_{\boldsymbol{\Psi}_m}^{-1}) \right|$$



*training* | *generation*

- *o* ---> Gaussian noise sequence *z*
- Multiple transformations: normalizing flow

D. Rezende and S. Mohamed. Variational inference with normalizing flows. In Proc. ICML, pages 1530–1538, 2015.

$$p_o(\boldsymbol{o}|\boldsymbol{c}; \boldsymbol{\Psi}_1, \cdots, \boldsymbol{\Psi}_M) = p_z(\boldsymbol{z} = f_{\boldsymbol{\Psi}_1}^{-1} \circ \cdots \circ f_{\boldsymbol{\Psi}_{M-1}}^{-1} \circ f_{\boldsymbol{\Psi}_M}^{-1}(\boldsymbol{o})|\boldsymbol{c}) \prod_{m=1}^{M} \left| \det \boldsymbol{J}(f_{\boldsymbol{\Psi}_m}^{-1}) \right|$$

- Different time complexity $f_{\boldsymbol{\Psi}_m}(\cdot)$, different training strategies

Parallel WaveNet

ClariNet

$f_{\boldsymbol{\Psi}_m}(\cdot)$ in linear time domain

FloWaveNet

WaveGlow

$f_{\boldsymbol{\Psi}_m}(\cdot)$ in subscale time domain

$\boldsymbol{o}_{1:T}$

Transform

Criterion

Model

$\boldsymbol{c}_{1:N}$

A. v. d. Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. v. d. Driessche, E. Lockhart, L. C. Cobo, F. Stimberg, et al. Parallel WaveNet: Fast high-fidelity speech synthesis. arXiv preprint arXiv:1711.10433, 2017.
W. Ping, K. Peng, and J. Chen. Clarinet: Parallel wave generation in end-to-end text-to-speech. arXiv preprint arXiv:1807.07281, 2018.
S. Kim, S.-g. Lee, J. Song, and S. Yoon. Flowavenet: A generative flow for raw audio. arXiv preprint arXiv:1811.02155, 2018.
R. Prenger, R. Valle, and B. Catanzaro. Waveglow: A flow-based generative network for speech synthesis. arXiv preprint arXiv:1811.00002, 2018.

## ClariNet & parallel WaveNet

❑ Generation process



✓ Parallel computation

✓ Fast generation

\* Initial condition may be implemented differently

## ClariNet & parallel WaveNet

❑ Naïve training process



$$z_t^{(0)} = (z_t^{(1)} - \mu_t)/\sigma_t$$

$\boldsymbol{o}_{1:T}$

$f_{\boldsymbol{\Psi}_M}^{-1}(\cdot)$

$f_{\boldsymbol{\Psi}_{M-1}}^{-1}(\cdot)$

$\cdots$

$\boldsymbol{z}_{1:T}^{(1)}$

$f_{\boldsymbol{\Psi}_1}^{-1}(\cdot)$

$\boldsymbol{z}_{1:T}^{(0)}$

$\mathcal{N}(\boldsymbol{o}; \boldsymbol{I})$

$\boldsymbol{c}_{1:N}$

$\boldsymbol{z}_{1:T}^{(1)}$

$\boldsymbol{\sigma}_{1:T}$
$\boldsymbol{\mu}_{1:T}$

$\text{CNN}_{\boldsymbol{\Psi}_1}(\cdot)$

$\boldsymbol{z}_{1:T}^{(0)}$

! Training time $\sim O(\text{T})$

Training

Generation

(Inverse-AR [1]) Normalizing flow

Original WaveNet

[1] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling. Improved variational inference with inverse autoregressive flow. In Proc. NIPS, pages 4743–4751, 2016.

# PART II: NORMALIZING FLOW-BASED MODELS

## ClariNet & parallel WaveNet

❑ Fast training : knowledge distilling



- Teacher gives $p(\widehat{o}_t | \widehat{\boldsymbol{o}}_{1:t-1}, \boldsymbol{c}_{1:T}, \text{teacher})$
- Student gives $p(\widehat{o}_t | \boldsymbol{z}_{1:T}, \boldsymbol{c}_{1:T}, \text{student})$

Student learns from teacher

Speech/signal processing

$o_{1:T}$

Knowledge-based Transform

Criterion

Model

$c_{1:N}$

LPCNet  ExcitNet

GlotNet  LP-WaveNet

Subband WaveNet

WaveNet  SampleRNN

Autoregressive models

$o_{1:T}$

Criterion

Model

$c_{1:N}$

WaveRNN  FFTNet

Universal neural vocoder

$o_{1:T}$

Transform

Criterion

Model

$c_{1:N}$

Parallel WaveNet

ClariNet

Without teacher WaveNet ?

❑ Parallel WaveNet & ClariNet
  ✓ No AR dependency in generation
  ! Need pre-trained WaveNet

Normlizing flow / feature transformation

## WaveGlow

❑ Why  $f_{\Psi_1}^{-1}(\cdot)$  ~ $O$(T) ?  Dependency in linear time domain



❑ Reduce T? Dependency in subscale time domain

R. Prenger, R. Valle, and B. Catanzaro. Waveglow: A flow-based generative network for speech synthesis. arXiv preprint arXiv:1811.00002, 2018.

# PART II: NORMALIZING FLOW-BASED MODELS

## WaveGlow



Fig. 1: WaveGlow network

R. Prenger, R. Valle, and B. Catanzaro. Waveglow: A flow-based generative network for speech synthesis. arXiv preprint arXiv:1811.00002, 2018.

# PART II: NORMALIZING FLOW-BASED MODELS



Speech/signal processing

$o_{1:T}$ → Knowledge-based Transform → Criterion → Model → $c_{1:N}$

LPCNet   ExcitNet

GlotNet   LP-WaveNet

Subband WaveNet

WaveNet   SampleRNN

Autoregressive models

$o_{1:T}$ → Criterion → Model → $c_{1:N}$

WaveRNN   FFTNet

Universal neural vocoder

$o_{1:T}$ → Transform → Criterion → Model → $c_{1:N}$

Parallel WaveNet

ClariNet

FloWaveNet

WaveGlow

Normlizing flow / feature transformation

?

Simple model without normalizing flow

Neural source-filter model

# CONTENTS

- Introduction: text-to-speech synthesis

- Neural waveform models
  - Overview
  - Autoregressive models
  - Normalizing flow
  - STFT-based training criterion

- Summary & software

## Neural source-filter model

❑ Naïve model and STFT-based criterion

Generated
waveforms



X. Wang, S. Takaki, and J. Yamagishi. Neural source-filter-based waveform model for statistical parametric speech synthesis. arXiv preprint arXiv:1810.11946, 2018.

# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model

❏ Naïve model and STFT-based criterion

X. Wang, S. Takaki, and J. Yamagishi. Neural source-filter-based waveform model for statistical parametric speech synthesis. arXiv preprint arXiv:1810.11946, 2018.

## Neural source-filter model

❑ Naïve model and STFT-based criterion

## Neural source-filter model

❑ Examples



Only harmonics
No formants
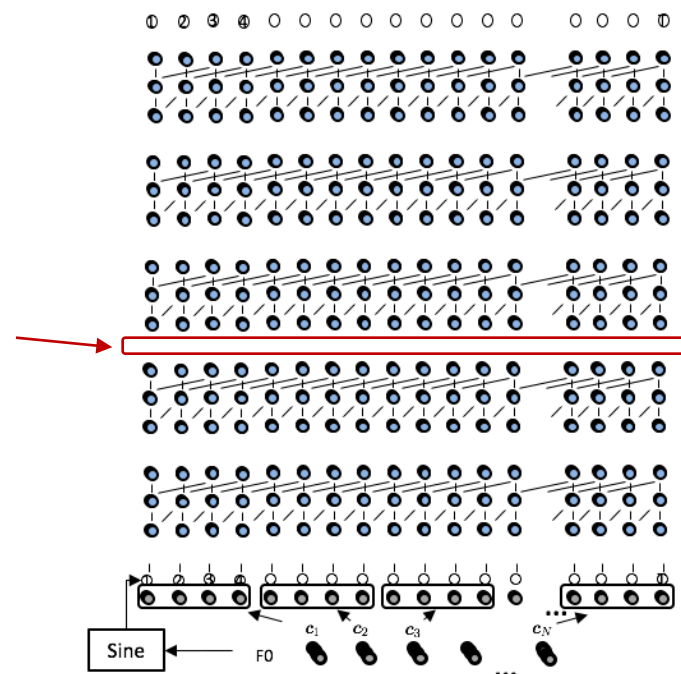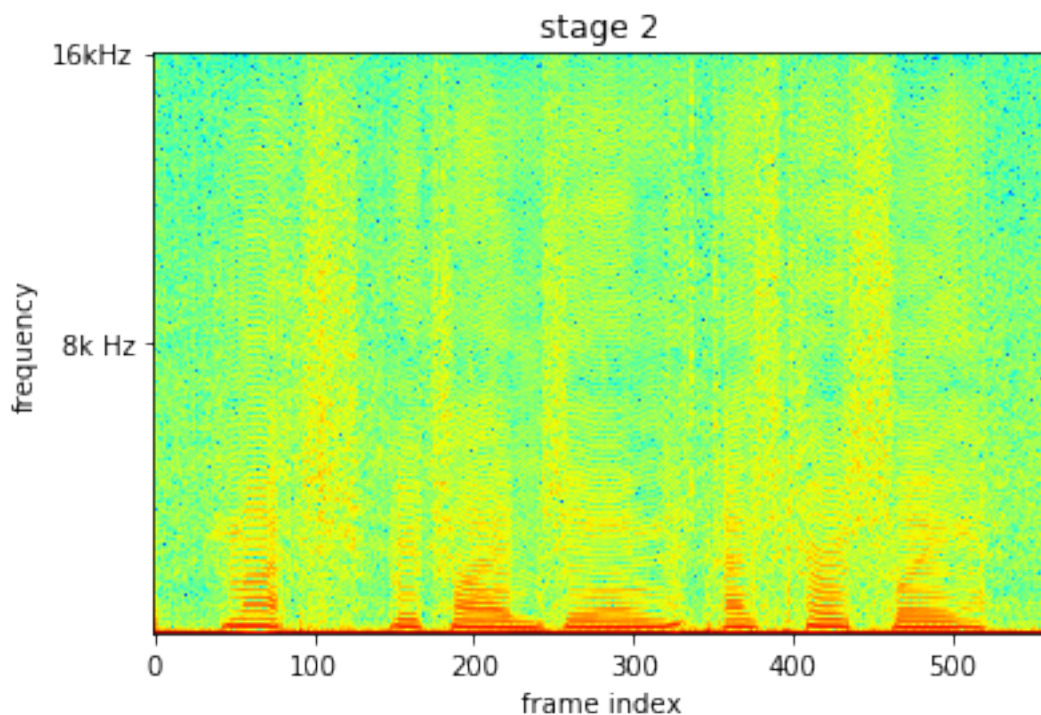
# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model

❑ Examples

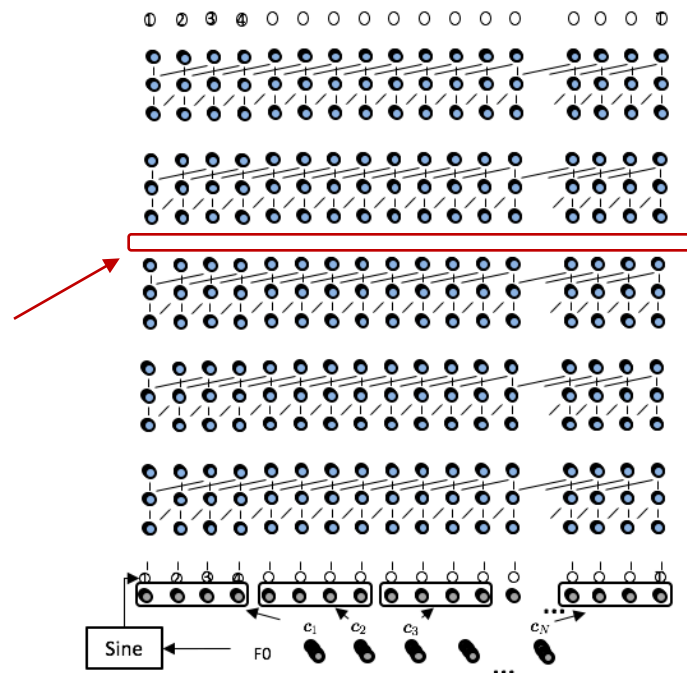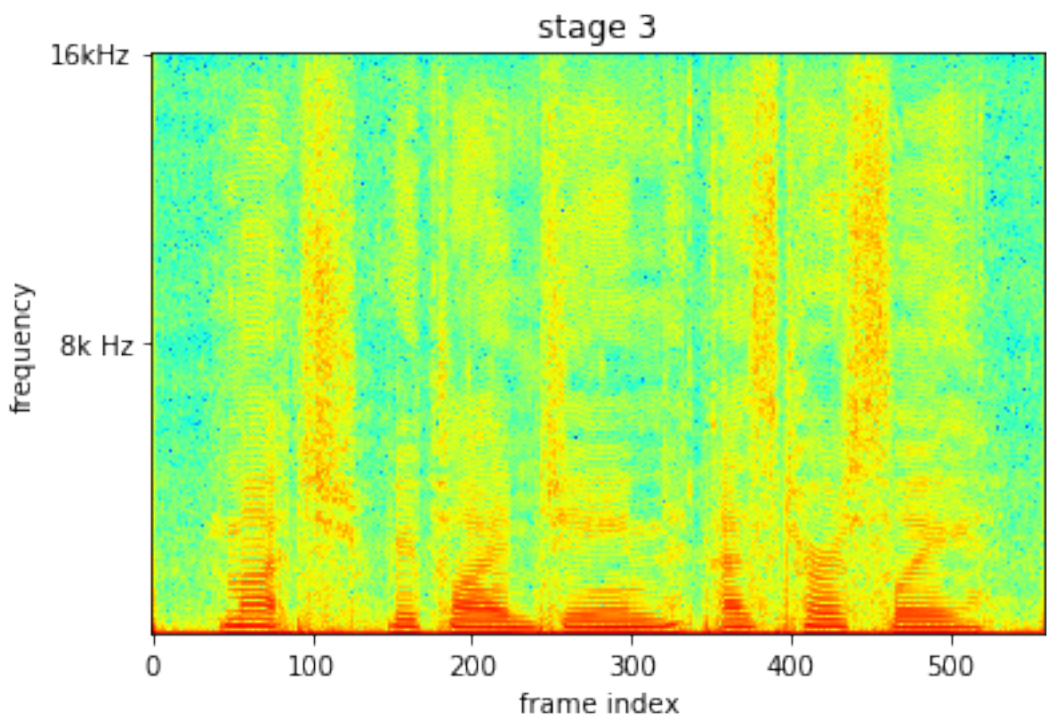# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model

❑ Examples

# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model

❑ Examples

# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model
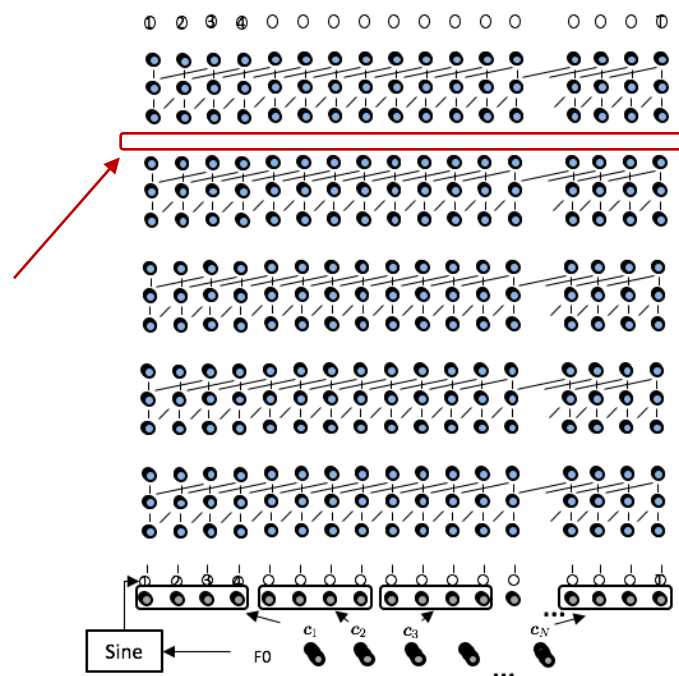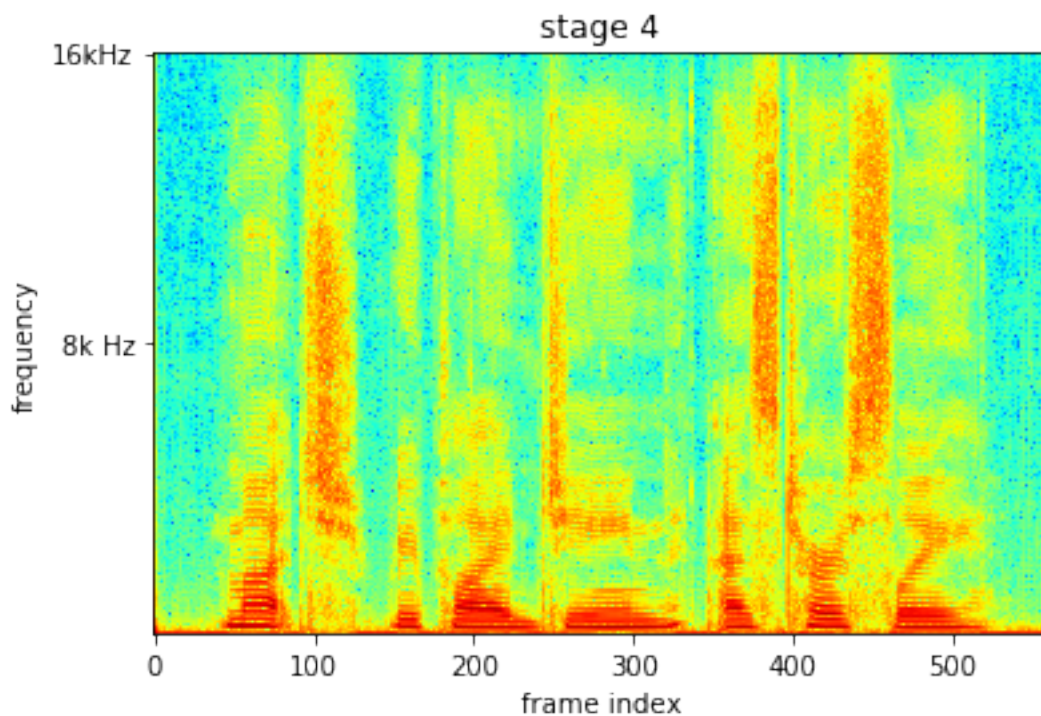
❑ Examples

# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model
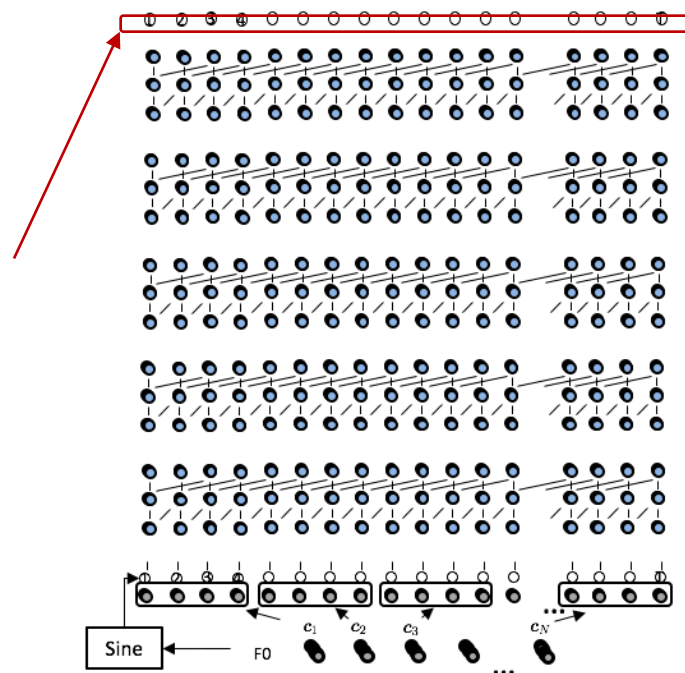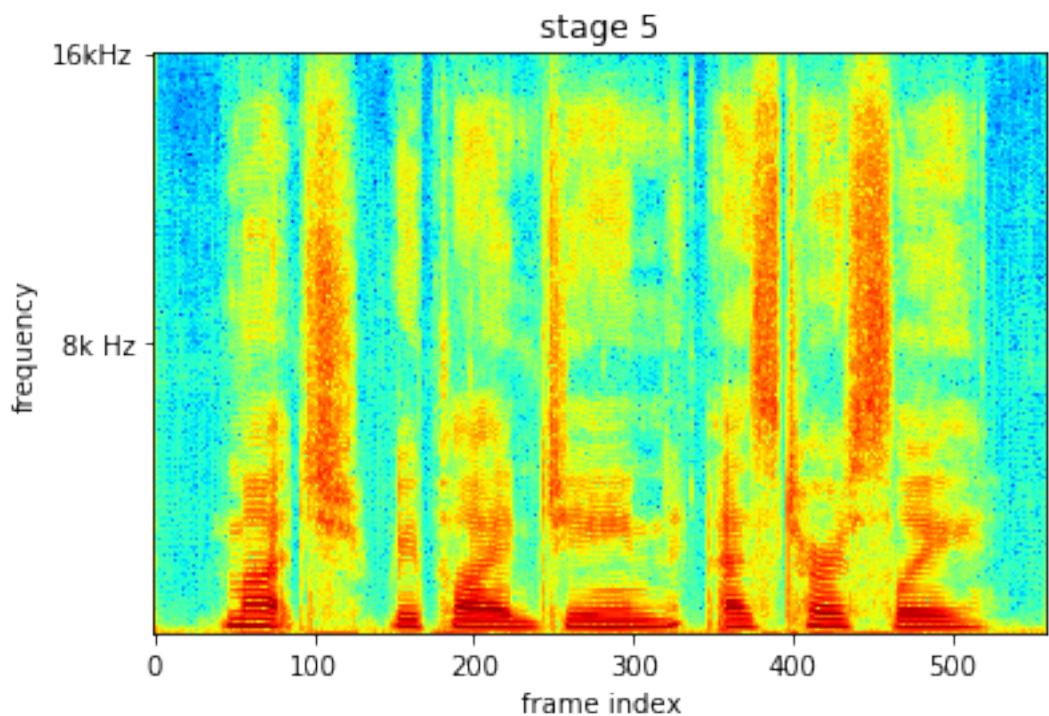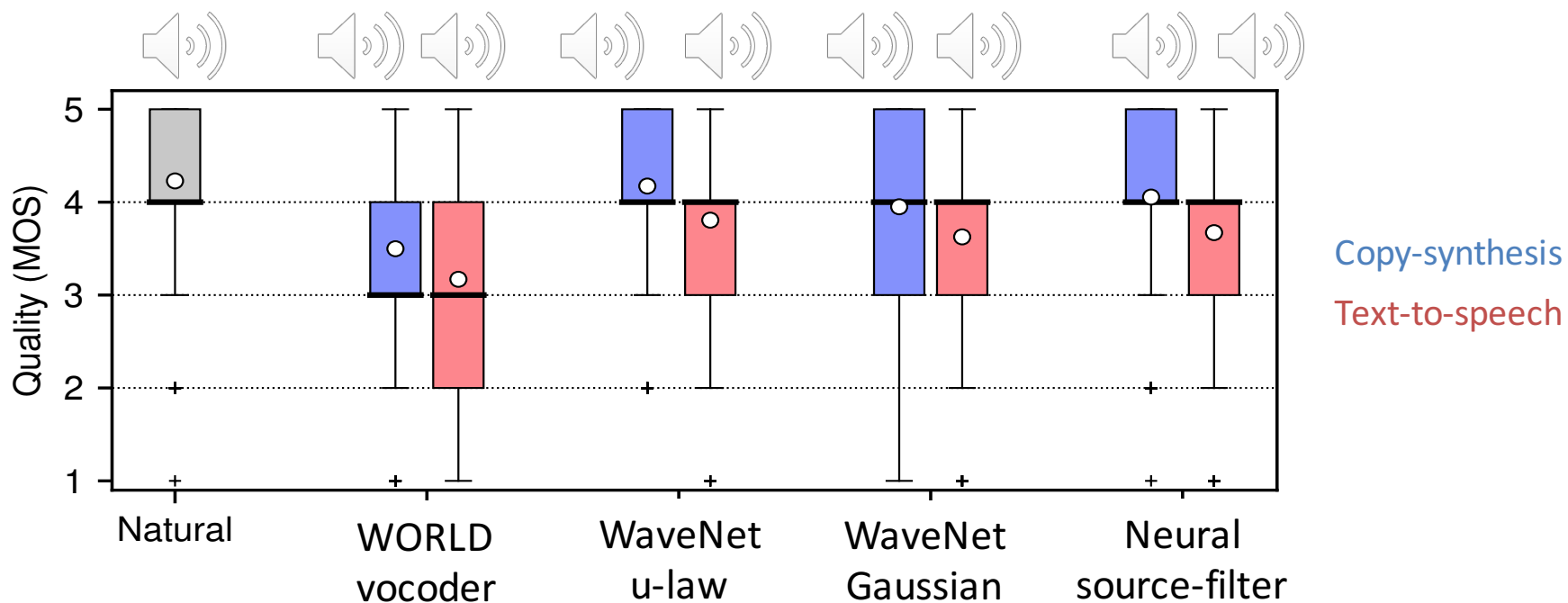
❑ Examples

# PART III: STFT-BASED TRAINING CRITERION

## Neural source-filter model

❑ Results
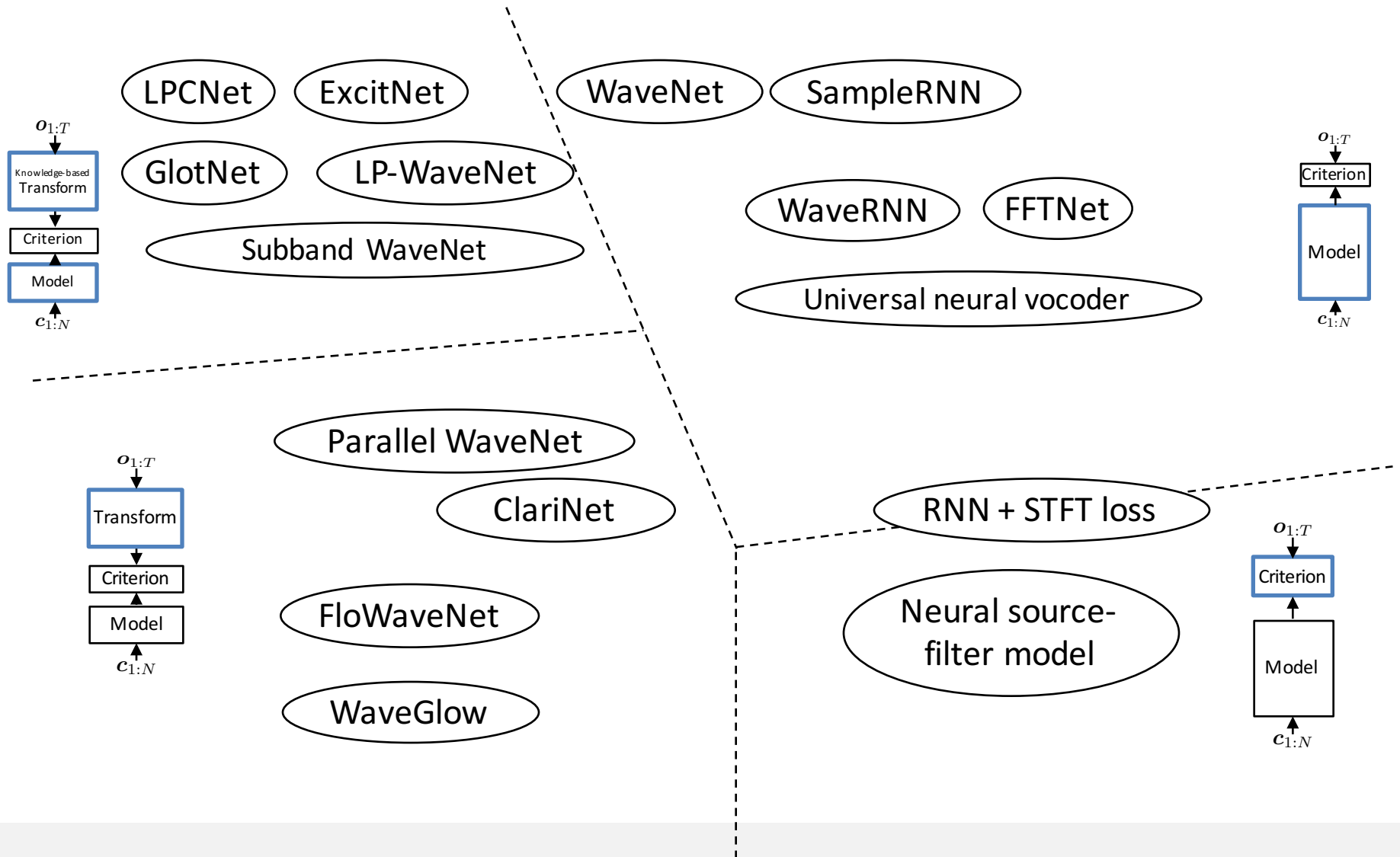
- Faster than WaveNet (at least **100** times)

- Smaller than WaveNet

- Speech quality is similar to WaveNet



Copy-synthesis

Text-to-speech

# CONTENTS

- Introduction: text-to-speech synthesis

- Neural waveform models

- Summary & software

# SUMMARY

LPCNet  ExcitNet

$o_{1:T}$

| Knowledge-based Transform |
| Criterion |
| Model |

$c_{1:N}$

GlotNet  LP-WaveNet

Subband WaveNet

WaveNet  SampleRNN

WaveRNN  FFTNet

Universal neural vocoder

$o_{1:T}$

| Criterion |
| Model |

$c_{1:N}$

Parallel WaveNet

ClariNet

$o_{1:T}$

| Transform |
| Criterion |
| Model |

$c_{1:N}$

FloWaveNet

WaveGlow

RNN + STFT loss

Neural source-filter model

$o_{1:T}$

| Criterion |
| Model |

$c_{1:N}$

56

# SOFTWARE

## NII neural network toolkit



LPCNet  ExcitNet  WaveNet  SampleRNN

GlotNet  LP-WaveNet

WaveRNN  FFTNet

Subband WaveNet

Universal neural vocoder

$o_{1:T}$
Knowledge-based Transform
Criterion
Model
$c_{1:N}$

$o_{1:T}$
Criterion
Model
$c_{1:N}$

Parallel WaveNet

ClariNet

RNN + STFT loss

$o_{1:T}$
Transform
Criterion
Model
$c_{1:N}$

FloWaveNet

Neural source-filter model

WaveGlow

$o_{1:T}$
Criterion
Model
$c_{1:N}$

# SOFTWARE

## NII neural network toolkit

❑ Neural waveform models

1. Toolkit cores: https://github.com/nii-yamagishilab/project-CURRENNT-public.git

2. Toolkit scripts: https://github.com/nii-yamagishilab/project-CURRENNT-scripts

# SOFTWARE

## NII neural network toolkit

❑ Useful slides

- http://tonywangx.github.io/pdfs/CURRENNT_TUTORIAL.pdf



- Other related slides http://tonywangx.github.io/slides.html

# REFERENCE

**WaveNet:**       A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. WaveNet: A generative model for raw audio. arXiv preprint arXiv:1609.03499, 2016.

**SampleRNN:**     S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio. Samplernn: An unconditional end-to-end neural audio generation model. arXiv preprint arXiv:1612.07837, 2016.

**WaveRNN:**       N. Kalchbrenner, E. Elsen, K. Simonyan, et.al. Efficient neural audio synthesis. In J. Dy and A. Krause, editors, Proc. ICML, volume 80 of Proceedings of Machine Learning Research, pages 2410–2419, 10–15 Jul 2018.

**FFTNet:**        Z. Jin, A. Finkelstein, G. J. Mysore, and J. Lu. FFTNet: A real-time speaker-dependent neural vocoder. In Proc. ICASSP, pages 2251–2255. IEEE, 2018.

**Universal vocoder:**  J. Lorenzo-Trueba, T. Drugman, J. Latorre, T. Merritt, B. Putrycz, and R. Barra-Chicote. Robust universal neural vocoding. arXiv preprint arXiv:1811.06292, 2018.

**Subband WaveNet**: T. Okamoto, K. Tachibana, T. Toda, Y. Shiga, and H. Kawai. An investigation of subband wavenet vocoder covering entire audible frequency range with limited acoustic features. In Proc. ICASSP, pages 5654–5658. 2018.

**Parallel WaveNet:**  A. van den Oord, Y. Li, I. Babuschkin, et. al.. Parallel WaveNet: Fast high-fidelity speech synthesis. In Proc. ICML, pages 3918–3926, 2018.

**ClariNet:**      W. Ping, K. Peng, and J. Chen. Clarinet: Parallel wave generation in end-to-end text-to-speech. arXiv preprint arXiv:1807.07281, 2018.

**FlowWaveNet:**   S. Kim, S.-g. Lee, J. Song, and S. Yoon. Flowavenet: A generative flow for raw audio. arXiv preprint arXiv:1811.02155, 2018.

**WaveGlow:**      R. Prenger, R. Valle, and B. Catanzaro. Waveglow: A flow-based generative network for speech synthesis. arXiv preprint arXiv:1811.00002, 2018.

**RNN+STFT:**      S. Takaki, T. Nakashika, X. Wang, and J. Yamagishi. STFT spectral loss for training a neural speech waveform model. In Proc. ICASSP (submitted), 2018.

**NSF:**           X. Wang, S. Takaki, and J. Yamagishi. Neural source-filter-based waveform model for statistical para- metric speech synthesis. arXiv preprint arXiv:1810.11946, 2018.

**LP-WavNet:**     M.-J. Hwang, F. Soong, F. Xie, X. Wang, and H.-G. Kang. Lp-wavenet: Linear prediction-based wavenet speech synthesis. arXiv preprint arXiv:1811.11913, 2018.

**GlotNet:**       L. Juvela, V. Tsiaras, B. Bollepalli, M. Airaksinen, J. Yamagishi, and P. Alku. Speaker-independent raw waveform model for glottal excitation. arXiv preprint arXiv:1804.09593, 2018.

**ExcitNet:**      E. Song, K. Byun, and H.-G. Kang. Excitnet vocoder: A neural excitation model for parametric speech synthesis systems. arXiv preprint arXiv:1811.04769, 2018.

**LPCNet:**        J.-M. Valin and J. Skoglund. Lpcnet: Improving neural speech synthesis through linear prediction. arXiv preprint arXiv:1810.11846, 2018.

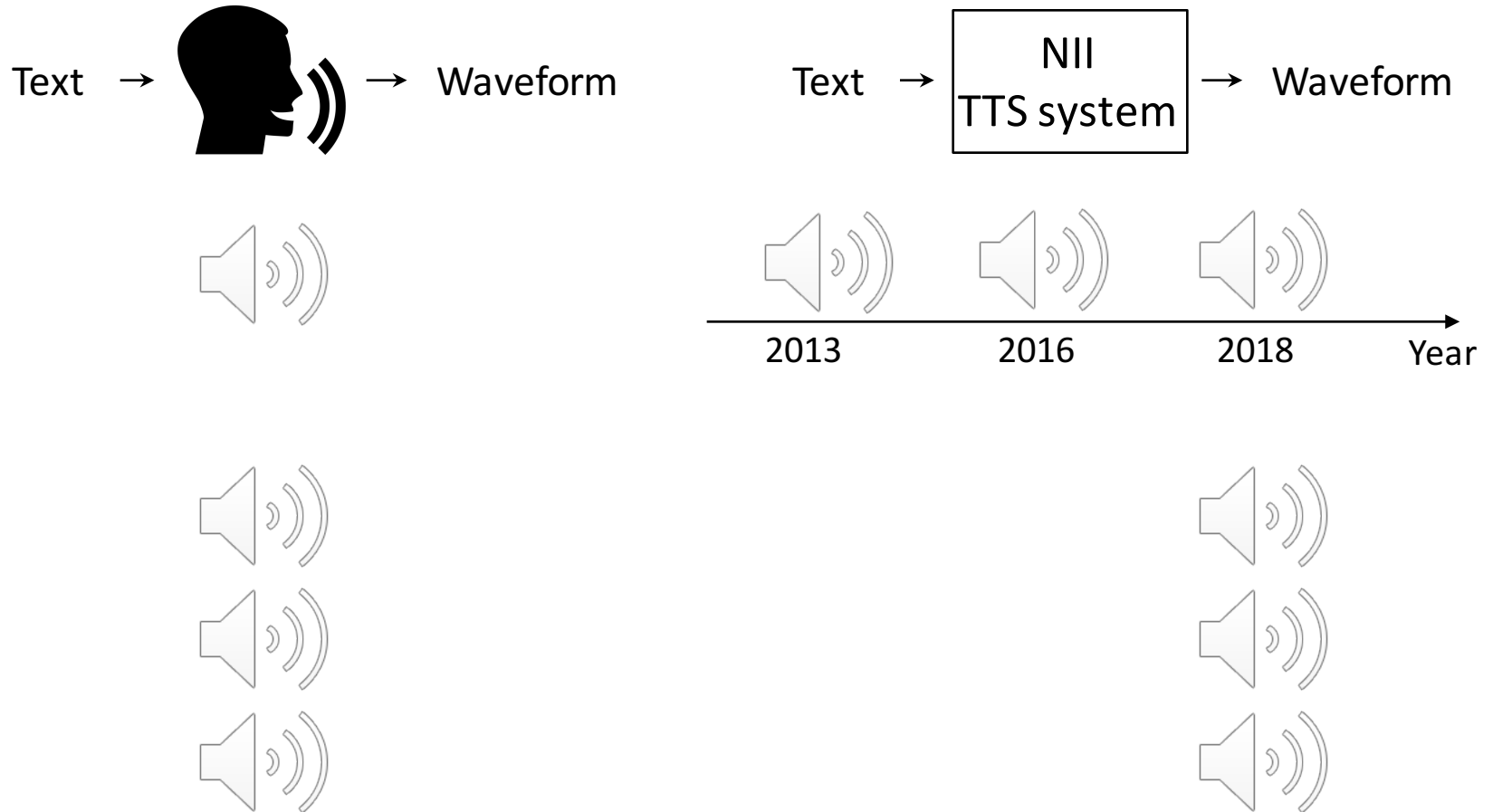# End of Part 1

## Codes, scripts, slides

http://nii-yamagishilab.github.io

# INTRODUCTION

## Text-to-speech (TTS)

❑ Speech samples from NII's TTS system
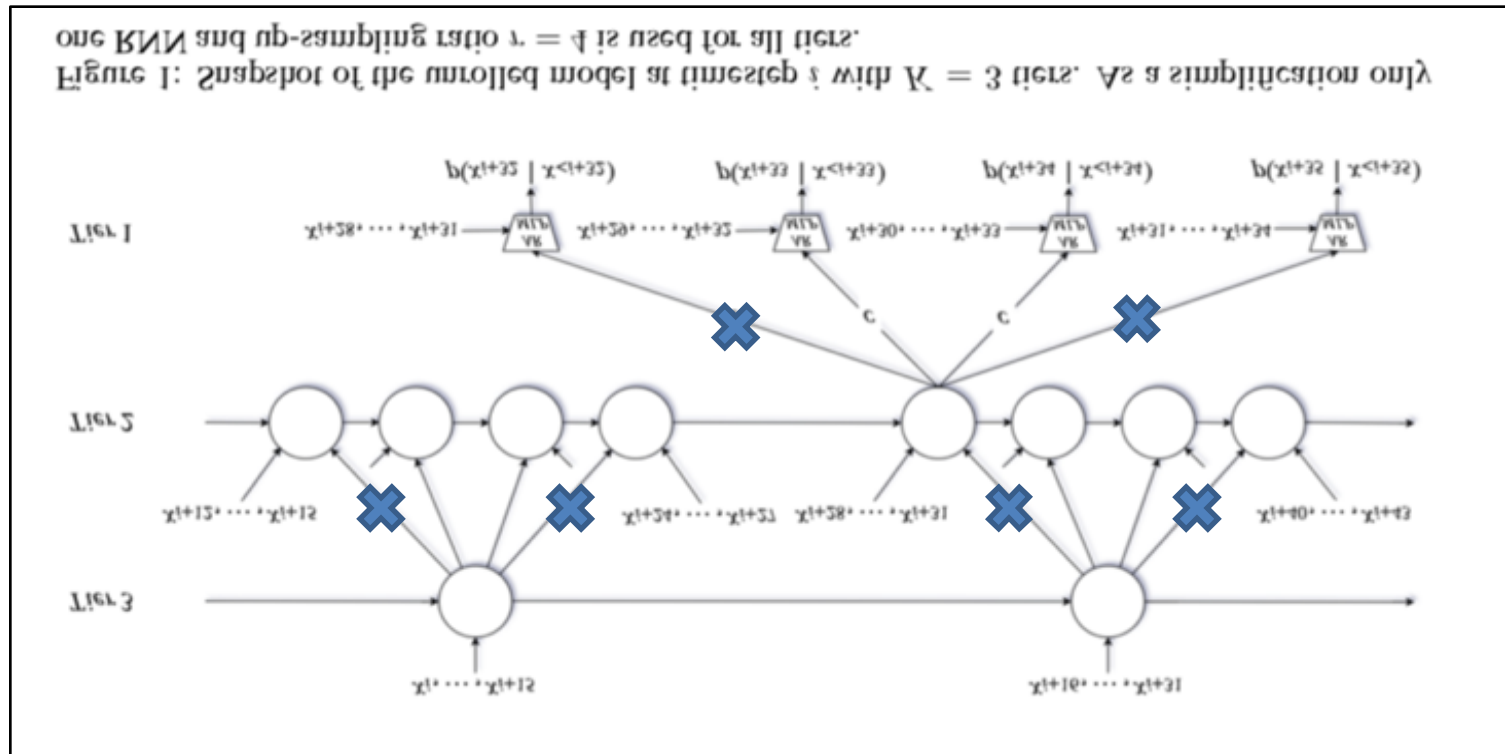
# PART I: AUTOREGRESSIVE MODELS

## SampleRNN

❑ Idea

- Hierarchical / multi-resolution dependency



Figure 1: Snapshot of the unrolled model at timestep $i$ with $K = 3$ tiers. As a simplification only one RNN and up-sampling ratio $r = 4$ is used for all tiers.

S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. Courville, and Y. Bengio. Samplernn: An unconditional end-to-end neural audio generation model. arXiv preprint arXiv:1612.07837, 2016.

## SampleRNN

❑ Example network structure

- Training
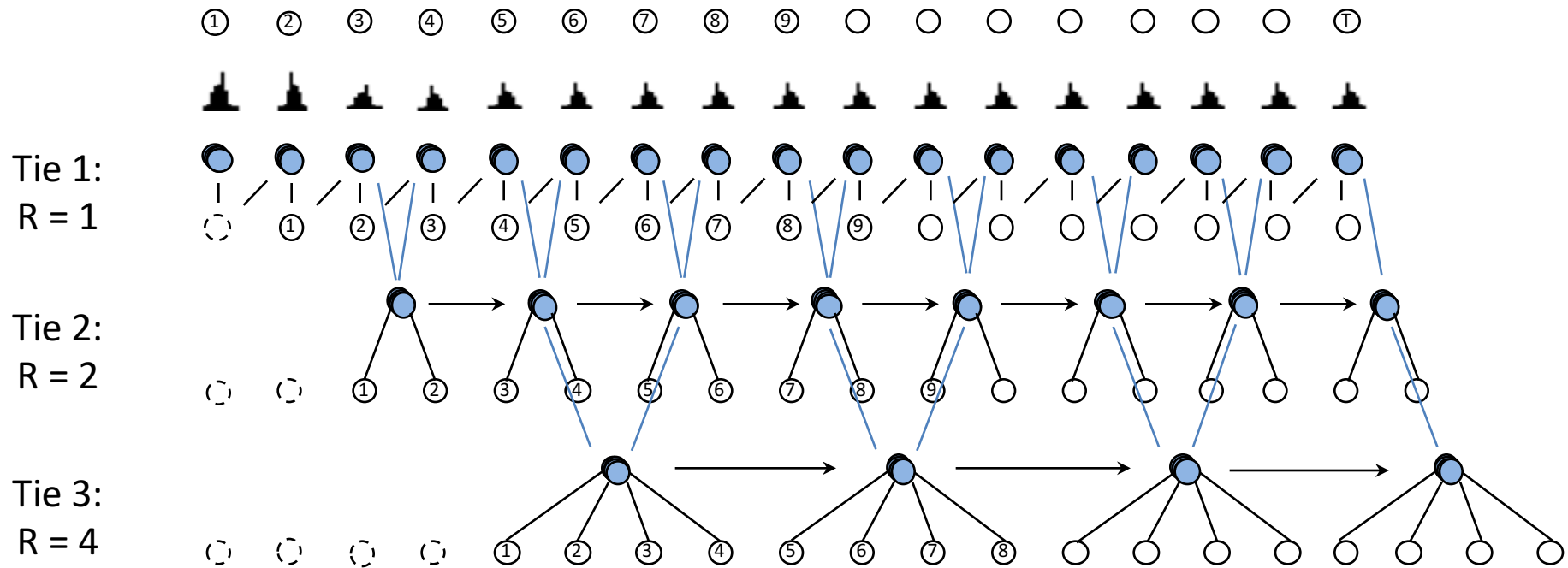


- R: time resolution increased by * 2
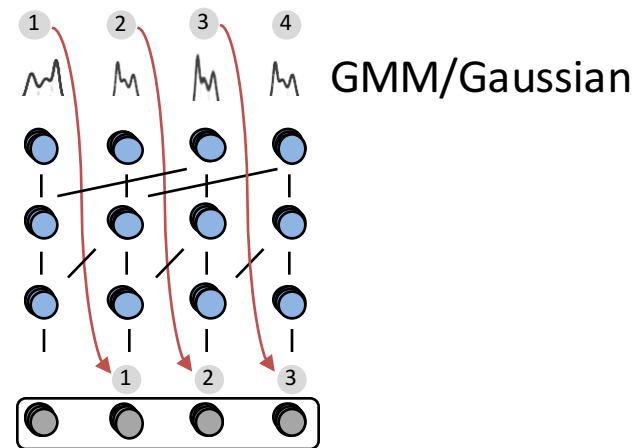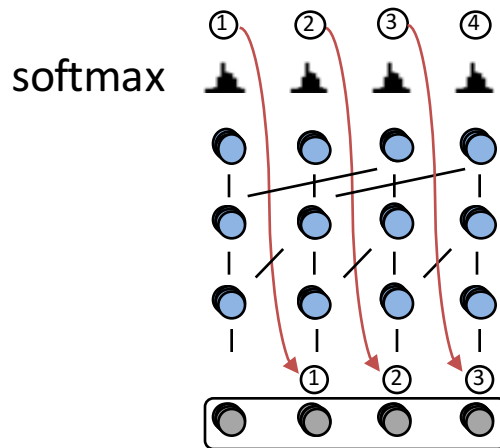
## SampleRNN

❑ Example network structure

- Generation process

# PART I: AUTOREGRESSIVE MODELS

## WaveNet

❑ Variants

- *u*-Law discrete waveform ---> continuous-valued waveform

  ▪ Mixture of logistic distribution [1]

  ▪ GMM / Single-Gaussian [2]



softmax                                      GMM/Gaussian

- Quantization noise shaping [3], related noise shaping method [4]

[1] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. arXiv preprint arXiv:1701.05517, 2017.

[2] W. Ping, K. Peng, and J. Chen. Clarinet: Parallel wave generation in end-to-end text-to-speech. arXiv preprint arXiv:1807.07281, 2018.
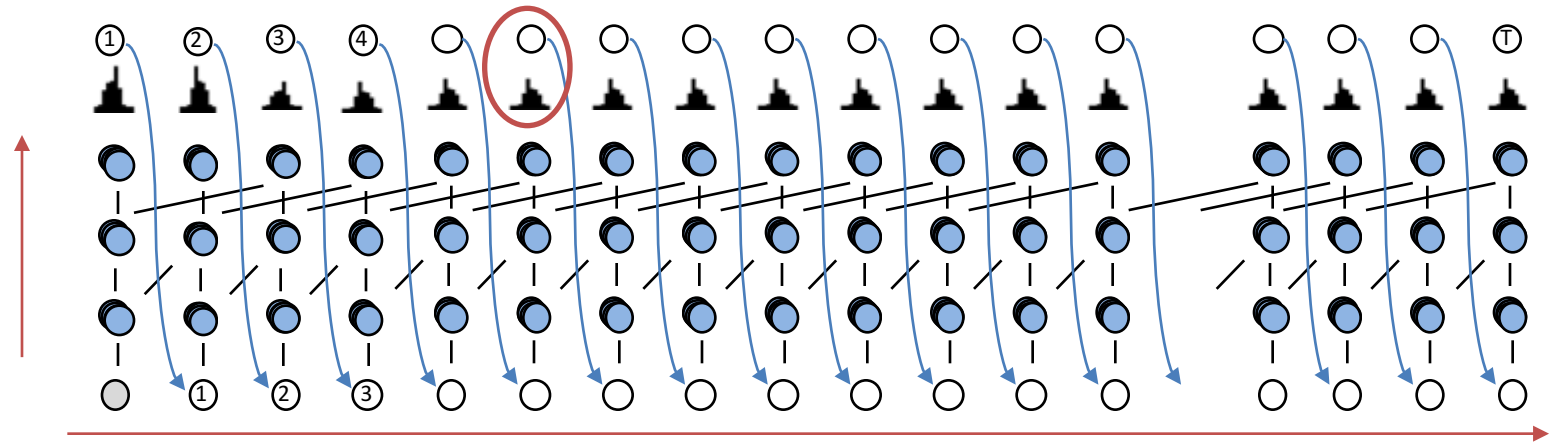
[3] T. Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda. Mel-cepstrum-based quantization noise shaping applied to neural-network-based speech waveform synthesis. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 26(7):1173–1180, 2018.

[4] K. Tachibana, T. Toda, Y. Shiga, and H. Kawai. An investigation of noise shaping with perceptual weighting for WaveNet-based speech generation. In Proc. ICASSP, pages 5664–5668. IEEE, 2018.

# PART I: AUTOREGRESSIVE MODELS
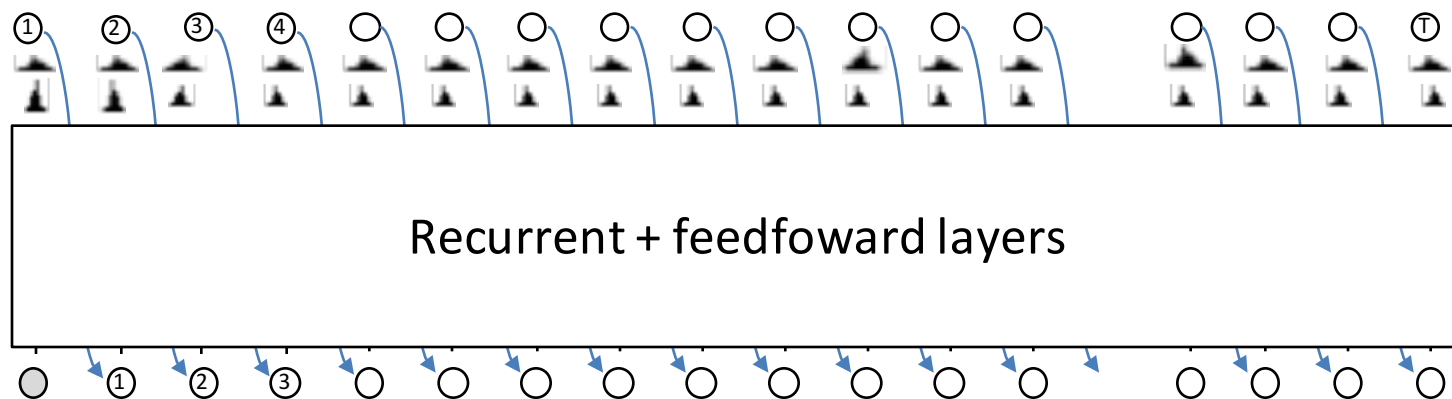
## WaveRNN

❑ WaveNet is inefficient



- Computation cost
    1. Impractical for 16 bit PCM (softmax of size 65536)
    2. Very deep network (50 dilated CNN …)
    3. …
- Time latency
    1. Generation time ~ $O$(waveform_length)

## WaveRNN

❑ WaveRNN strategies



- Computation cost
  1. ~~Impractical for 16 bit PCM (softmax of size 65536)~~    Two-level softmax
  2. ~~Very deep network (50 dilated CNN ...)~~              RNN + Feedforward
  3. ...
- Time latency
  1. ~~Generation time ~ *O*(waveform_length)~~    Subscale dependency + batch

N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu. Efficient neural audio synthesis. In J. Dy and A. Krause, editors, Proc. ICML, volume 80 of Proceedings of Machine Learning Research, pages 2410–2419, Stock- holmsm̈assan, Stockholm Sweden, 10–15 Jul 2018. PMLR.