

---

Slides by Xin Wang  
National Institute of Informatics

---

*© 2018, Xin Wang. All rights reserved.*

This work is licensed under the Creative Commons Attribution 3.0 license.  
See <http://creativecommons.org/> for details.



Note: Natural Japanese speech data belonging to ATR Ximera corpus are deleted  
in this public available version



# Fundamental Frequency Modeling for Neural-Network-based Statistical Parametric Speech Synthesis

Xin Wang

ID: 20151706

2018-07-13

# CONTENTS

## Introduction

- Background
- Topic
- Thesis outline

## Issues and methods

## Summary

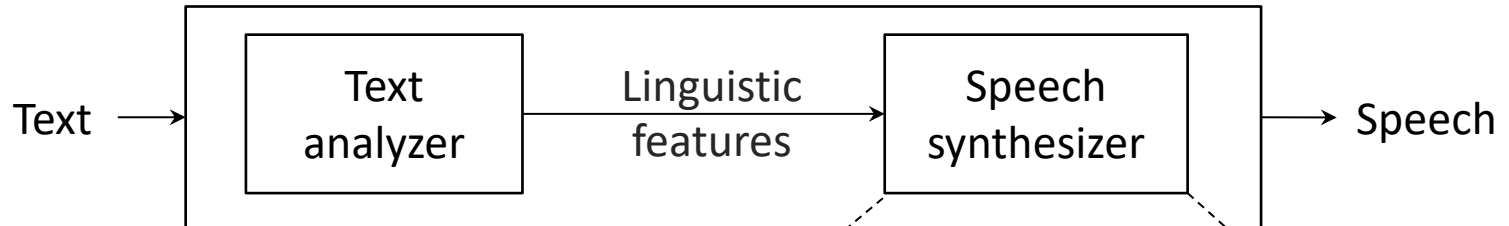


New results / updated explanation

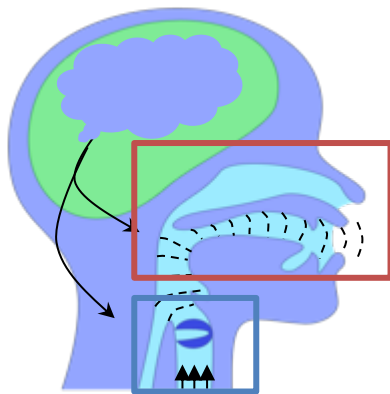
# INTRODUCTION

## Text-to-speech (TTS)

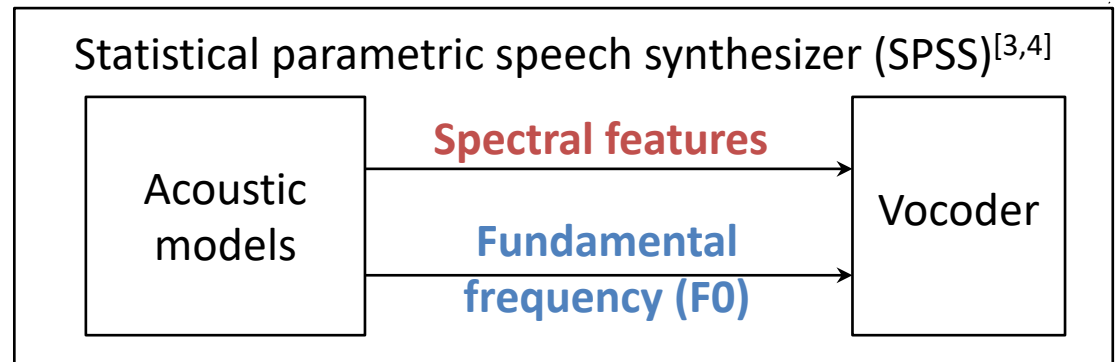
### □ TTS pipeline [1,2]



### □ Synthesizer



(c.f. HTS Slides, by HTS Working Group)



[1] Taylor, P. (2009). Text-to-Speech Synthesis.

[2] Dutoit, T. (1997). An Introduction to Text-to-speech Synthesis.

[3] Tokuda, K., et al., (2013). Speech Synthesis Based on Hidden Markov Models. *Proceedings of the IEEE*, 101(5), 1234–1252.

[4] Zen, H., et al. (2009). Statistical parametric speech synthesis. *Speech Communication*, 51, 1039–1064.

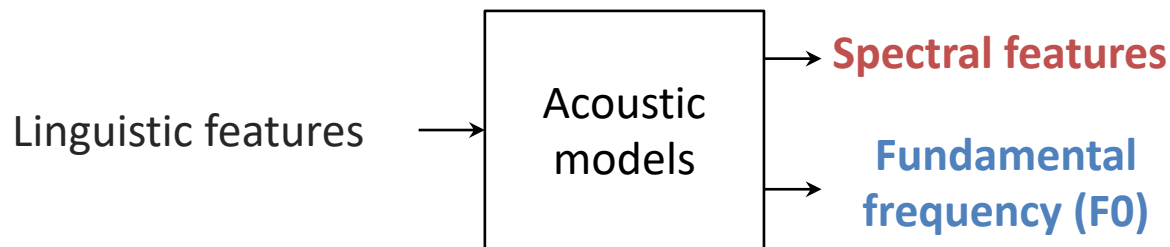
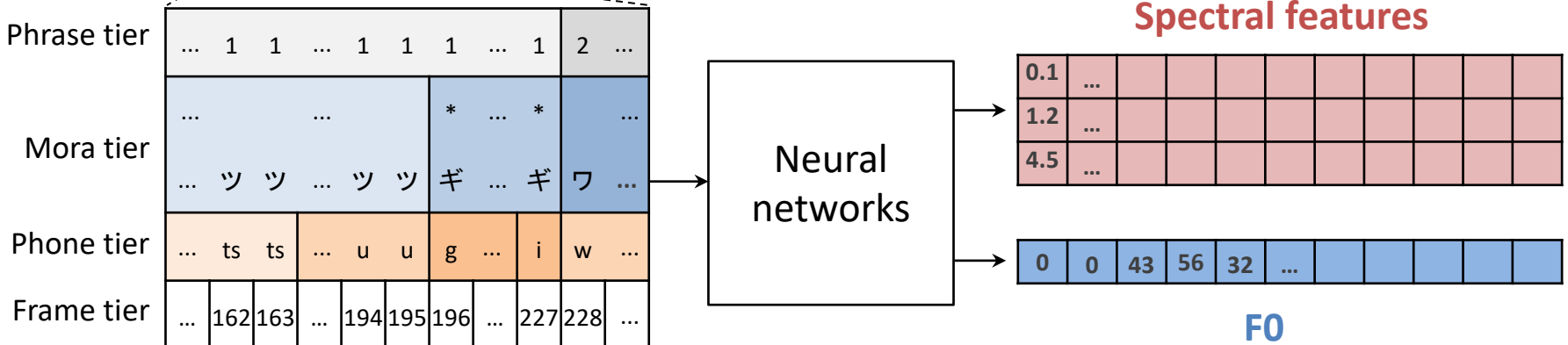


# INTRODUCTION

## Text-to-speech (TTS)

- Neural-network-based acoustic models [5,6,7]

次は新金岡、新金岡です。



[5] H. Zen, A. Senior, and M. Schuster. Statistical parametric speech synthesis using deep neural networks. In Proc. ICASSP, pages 7962–7966, 2013.

[6] Z. H. Ling, et al. Deep learning for acoustic modeling in parametric speech generation: A systematic review of existing techniques and future trends. IEEE Signal Processing Magazine, 32(3):35–52, 2015.

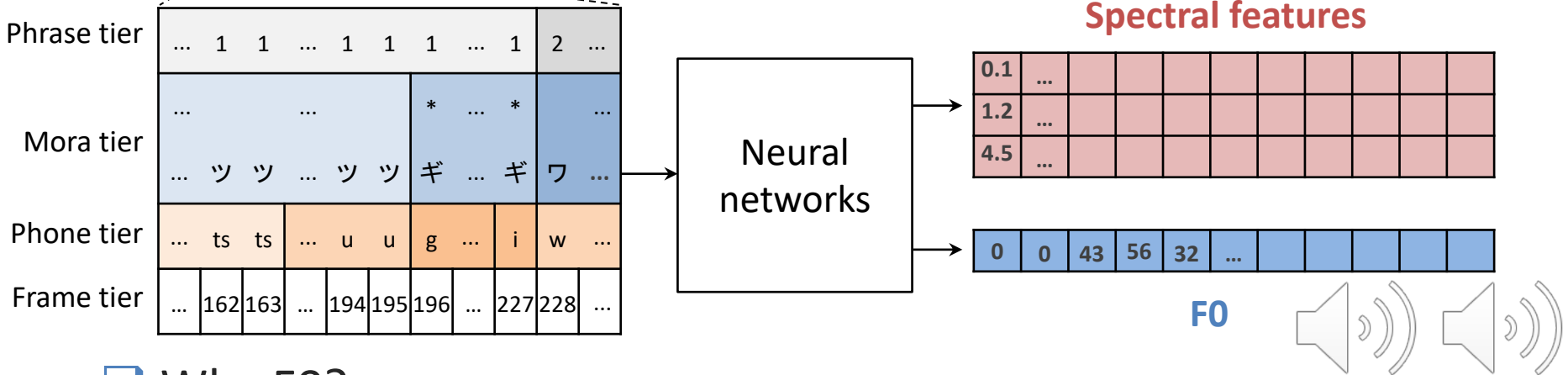
[7] Y. Fan, Y. Qian, F. Xie, and F. K. Soong. TTS synthesis with bidirectional LSTM based recurrent neural networks. In Proc. Interspeech, pages 1964–1968, 2014.

# INTRODUCTION

## Topic

### □ Neural F0 modeling for TTS

次は新金岡、新金岡です。



### □ Why F0?

Speaker A: Who made the marmalade.

Speaker B: Mari<sup>a</sup>n<sup>a</sup>na made the marmalade.

Speaker A: Bob made the marmalade.

Speaker B: (No,) Mari<sup>a</sup>n<sup>a</sup>na made the marmalade.

Speaker B: Marianna made the mar<sup>a</sup>malade.

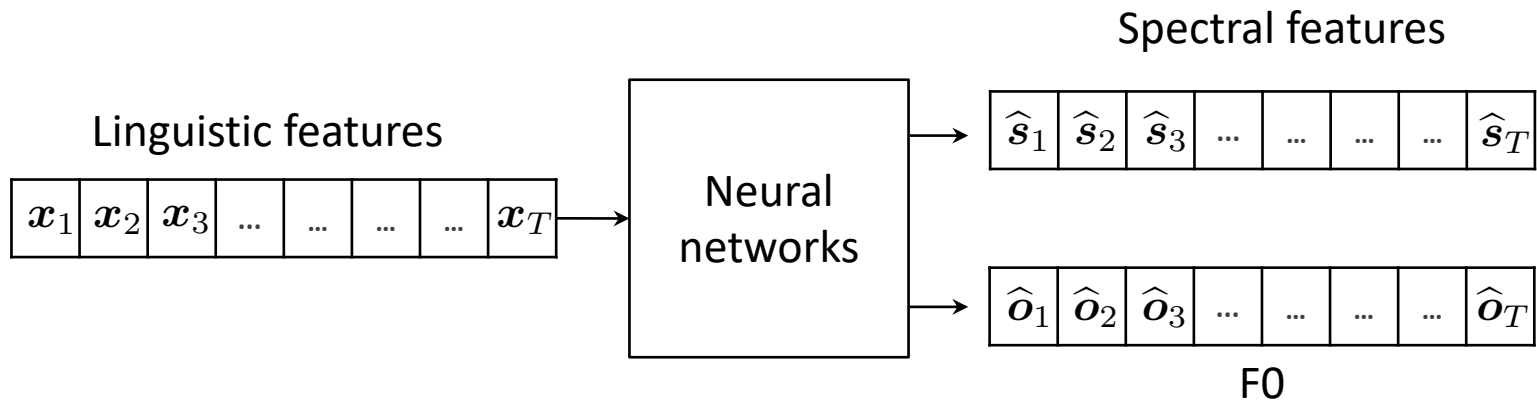
Speaker B: Marianna made the mar<sup>a</sup>malade.

Speaker B: Marianna made the mar<sup>a</sup>malade.

# INTRODUCTION

## Topic

- Issues to be addressed



$$p(\mathbf{o}_{1:T}, \mathbf{s}_{1:T} | \mathbf{x}_{1:T}; \Theta) = \prod_{t=1}^T \mathcal{N}([\mathbf{o}_t, \mathbf{s}_t]; \text{Network}_{\Theta}(\mathbf{x}_{1:T}, t), \beta \mathbf{I})$$

Issue 3: efficient enough?

Issue 1: joint modeling?

Issue 2: temporal dependency?

F0 features [9]

Linguistic features [10]

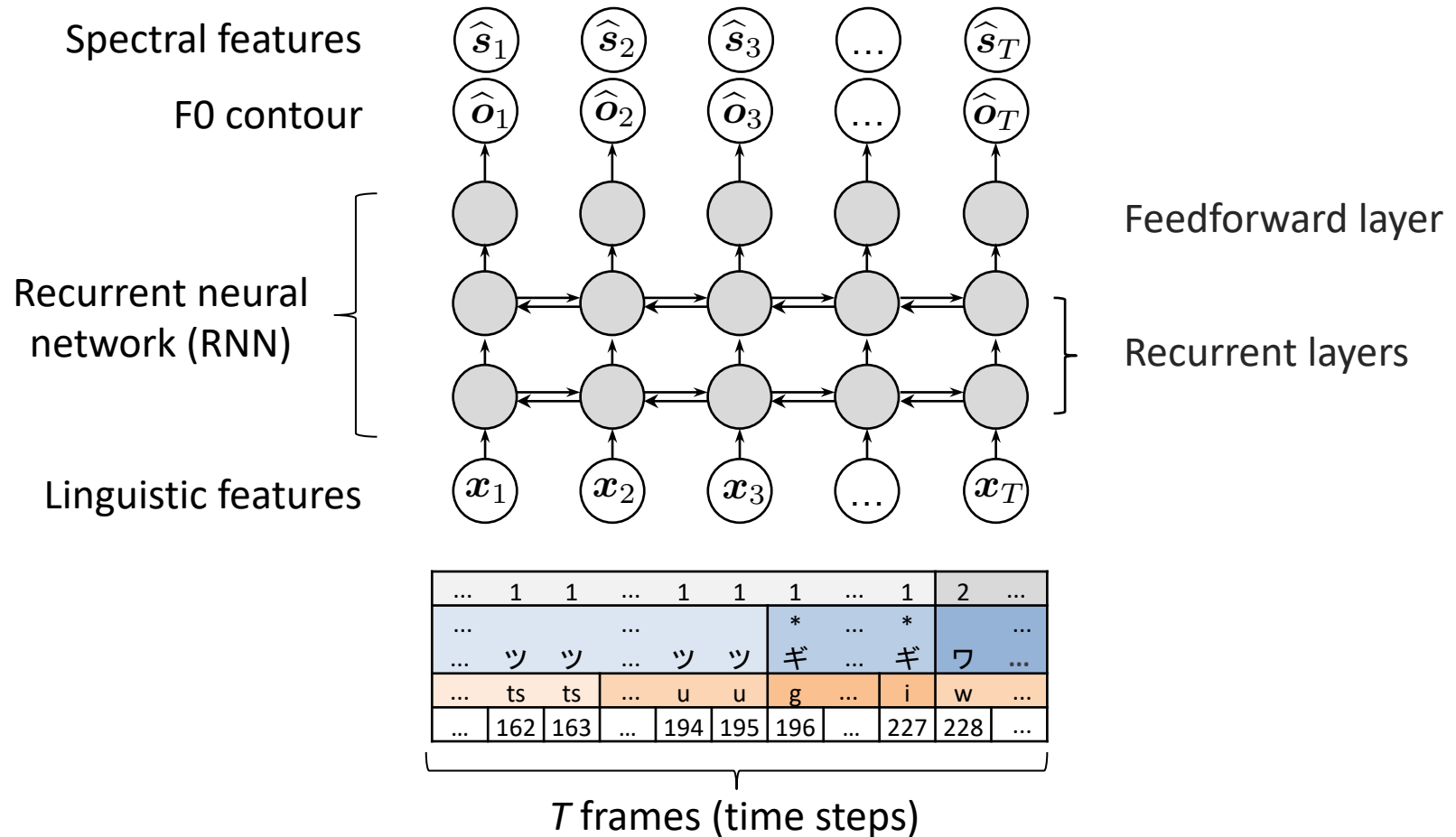
[9] M. S. Ribeiro. Suprasegmental representations for the modeling of fundamental frequency in statistical parametric speech synthesis. PhD thesis, The University of Edinburgh, 2018.

[10] J. Hirschberg. Pitch accent in context predicting intonational prominence from text. *Artificial Intelligence*, 63(1):305–340, 1993.

# INTRODUCTION

## Thesis outline

- Conventional approaches <sup>[7]</sup> (Table 3.1 in thesis)



[7] Y. Fan, Y. Qian, F. Xie, and F. K. Soong. TTS synthesis with bidirectional LSTM based recurrent neural networks. In Proc. Interspeech, pages 1964–1968, 2014.

# INTRODUCTION

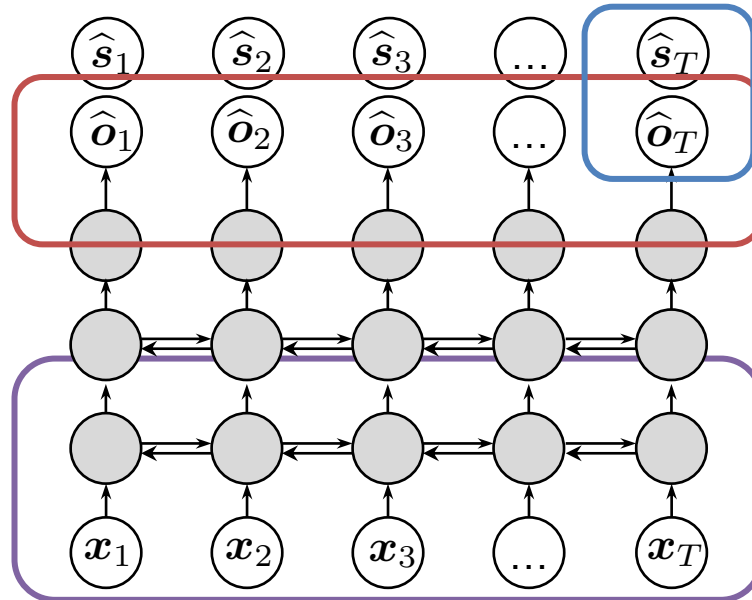
## Thesis outline

- Three issues

Issue 2:  
Temporal dependency?

Issue 3:  
Frame-by-frame  
processing is efficient?

Issue 1:  
Joint modeling?



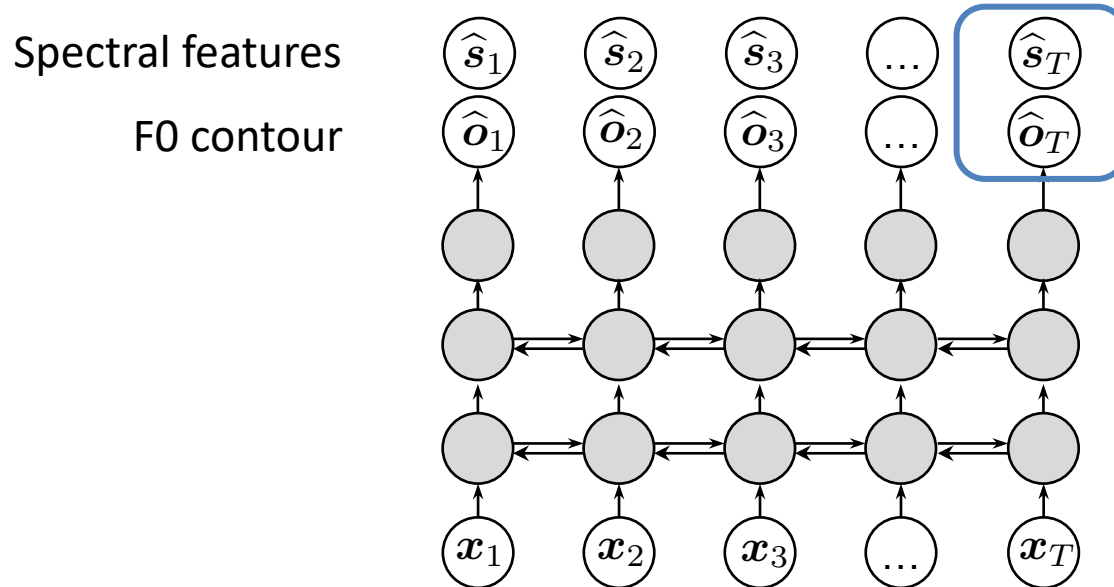
...	1	1	...	1	1	1	...	1	2	...
...	...				*	...	*	...		
...	ツ	ツ	...	ツ	ツ	ギ	...	ギ	ワ	...
...	ts	ts	...	u	u	g	...	i	w	...
...	162	163	...	194	195	196	...	227	228	...

T frames

# INTRODUCTION

## Thesis outline (Chapter 4)

- On issue 1: Joint modeling of F0 and spectral features?



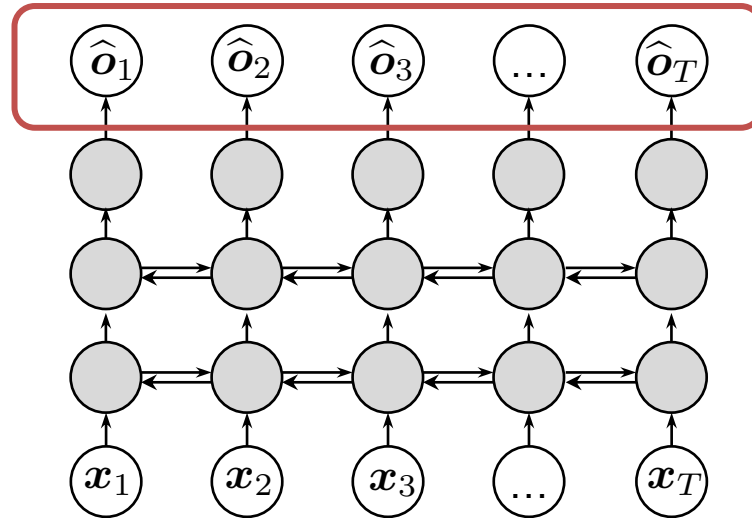
- Investigation using highway networks
  - Spectral features are prioritized
  - Different input/hidden features for F0 and spectral
- × Sub-optimal for F0 modeling
- ✓ Only F0 as target

} Novel analysis

# INTRODUCTION

## Thesis outline (Chapter 5-6)

- On issue 2: Temporal dependency?

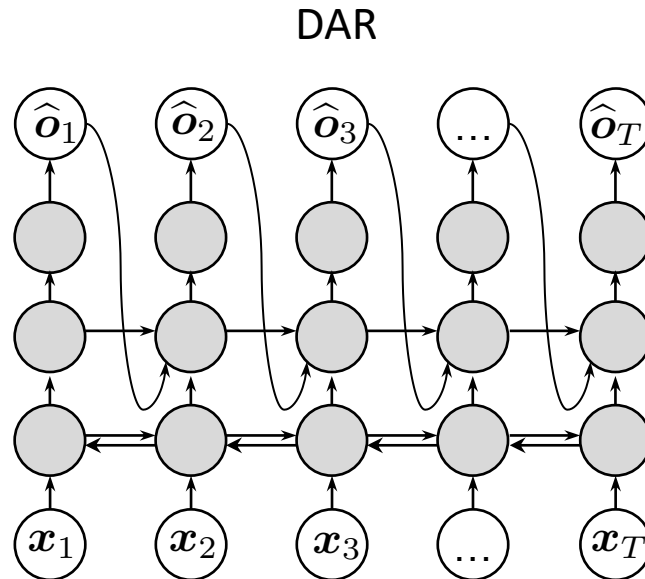


- Evidence from random sampling
- × RNN ignores temporal dependency

# INTRODUCTION

## Thesis outline (Chapter 5-6)

- On issue 2: Temporal dependency?



- ✓ Shallow autoregressive model (SAR)
  - Short-term dependency
- ✓ Deep autoregressive model (DAR)
  - Longer dependency & best results & random sampling

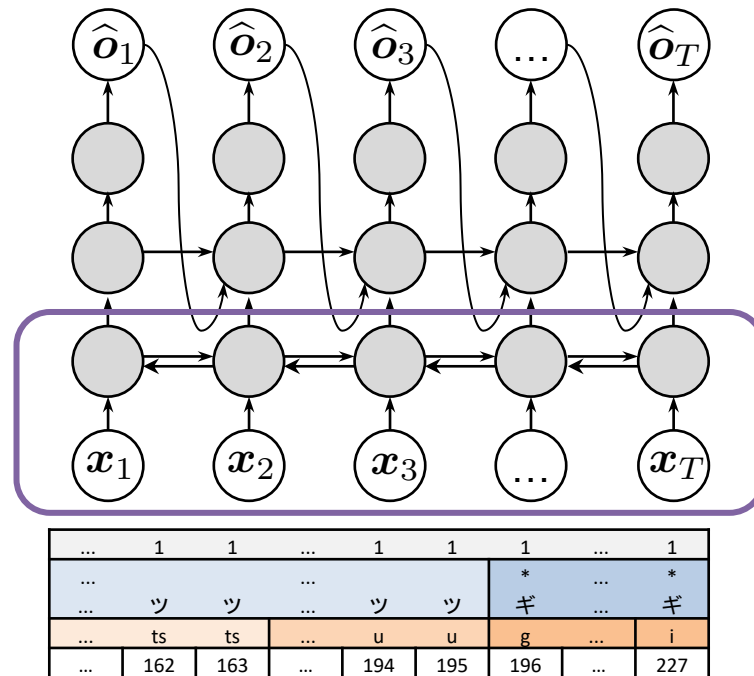
Novel models & interpretations



# INTRODUCTION

## Thesis outline (Chapter 7)

- On issue 3: Frame-by-frame processing?

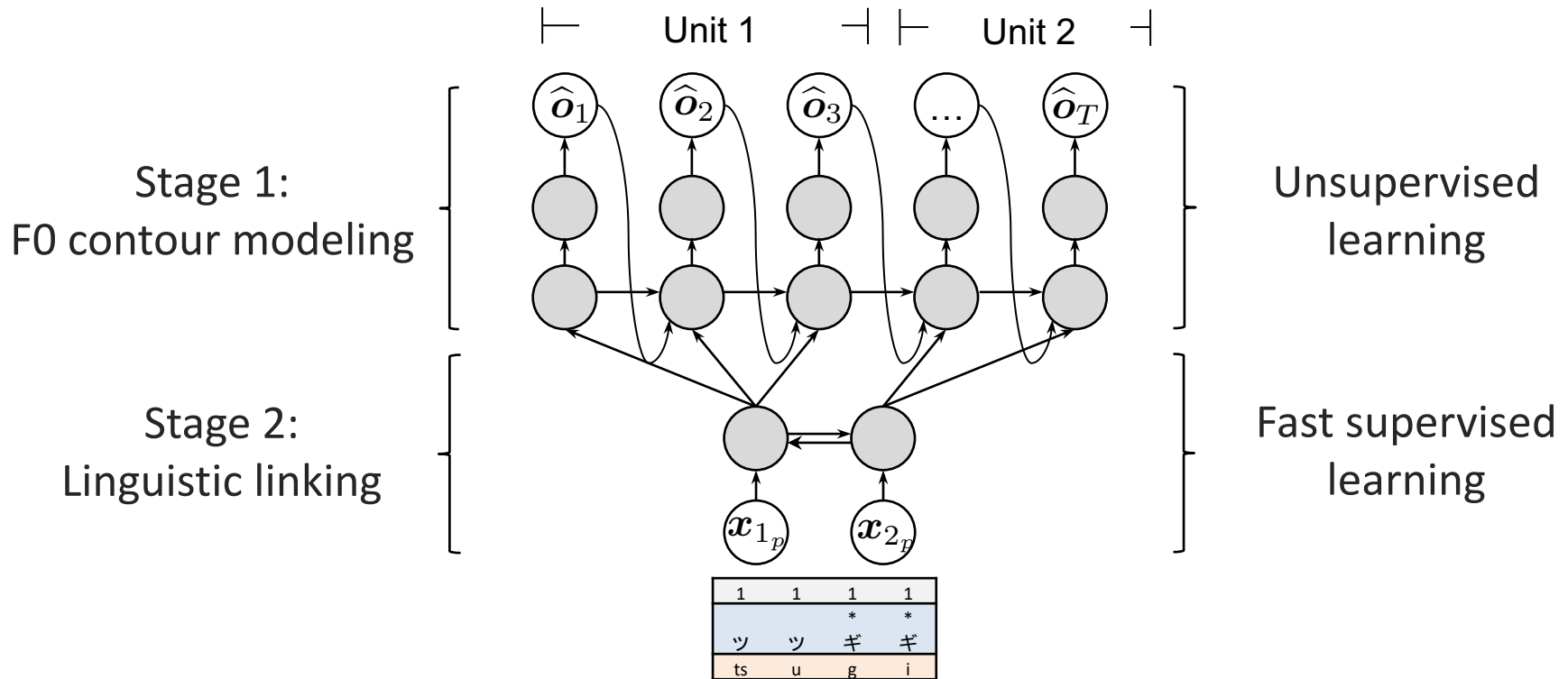


× Inefficient

# INTRODUCTION

## Thesis outline (Chapter 7)

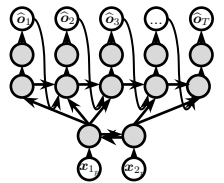
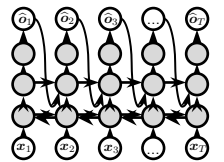
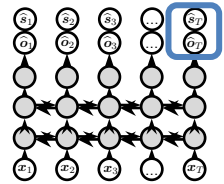
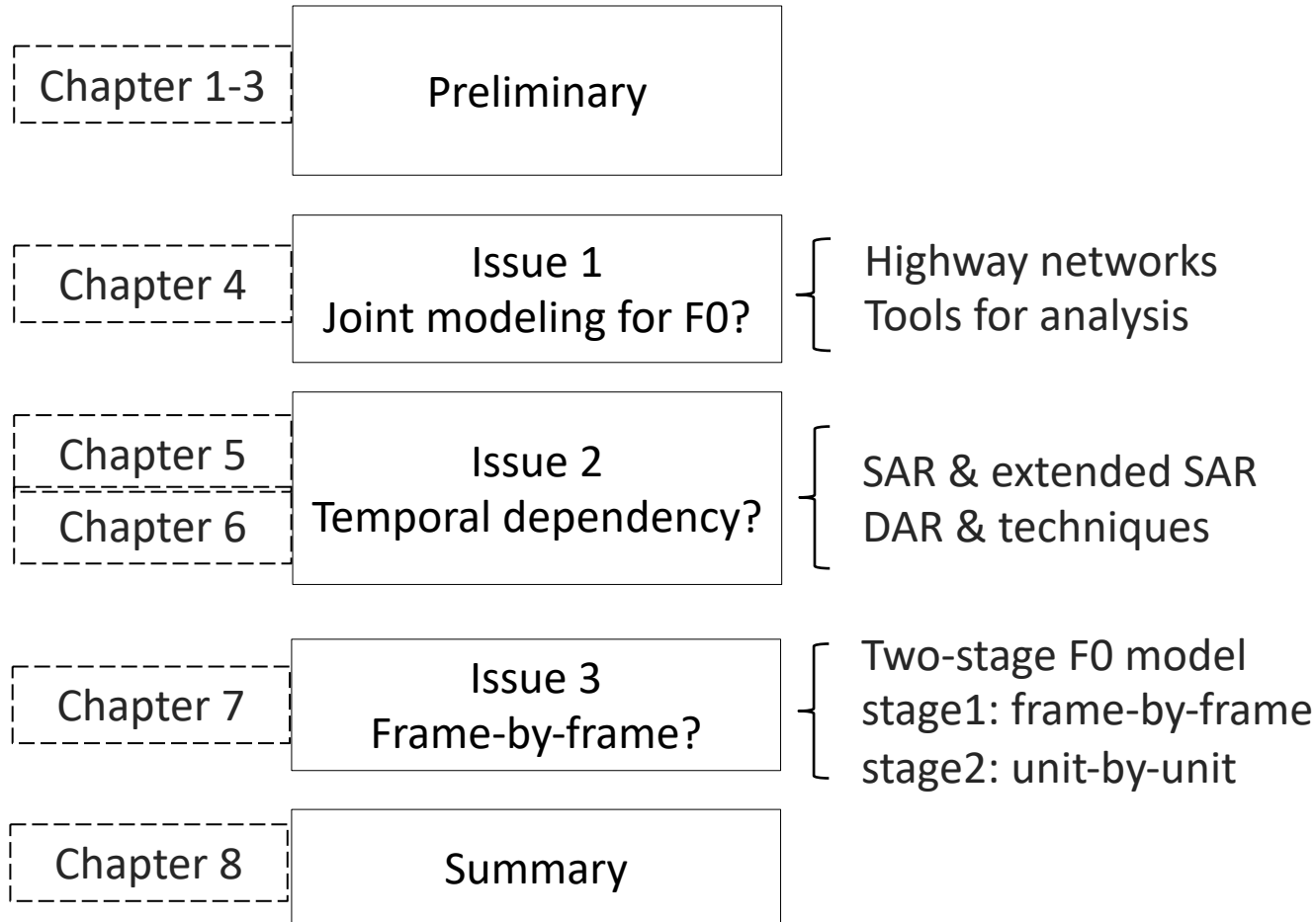
- On issue 3: Frame-by-frame processing?



- ✓ Two-stage model: efficient & interpretable & multi-level Novel model

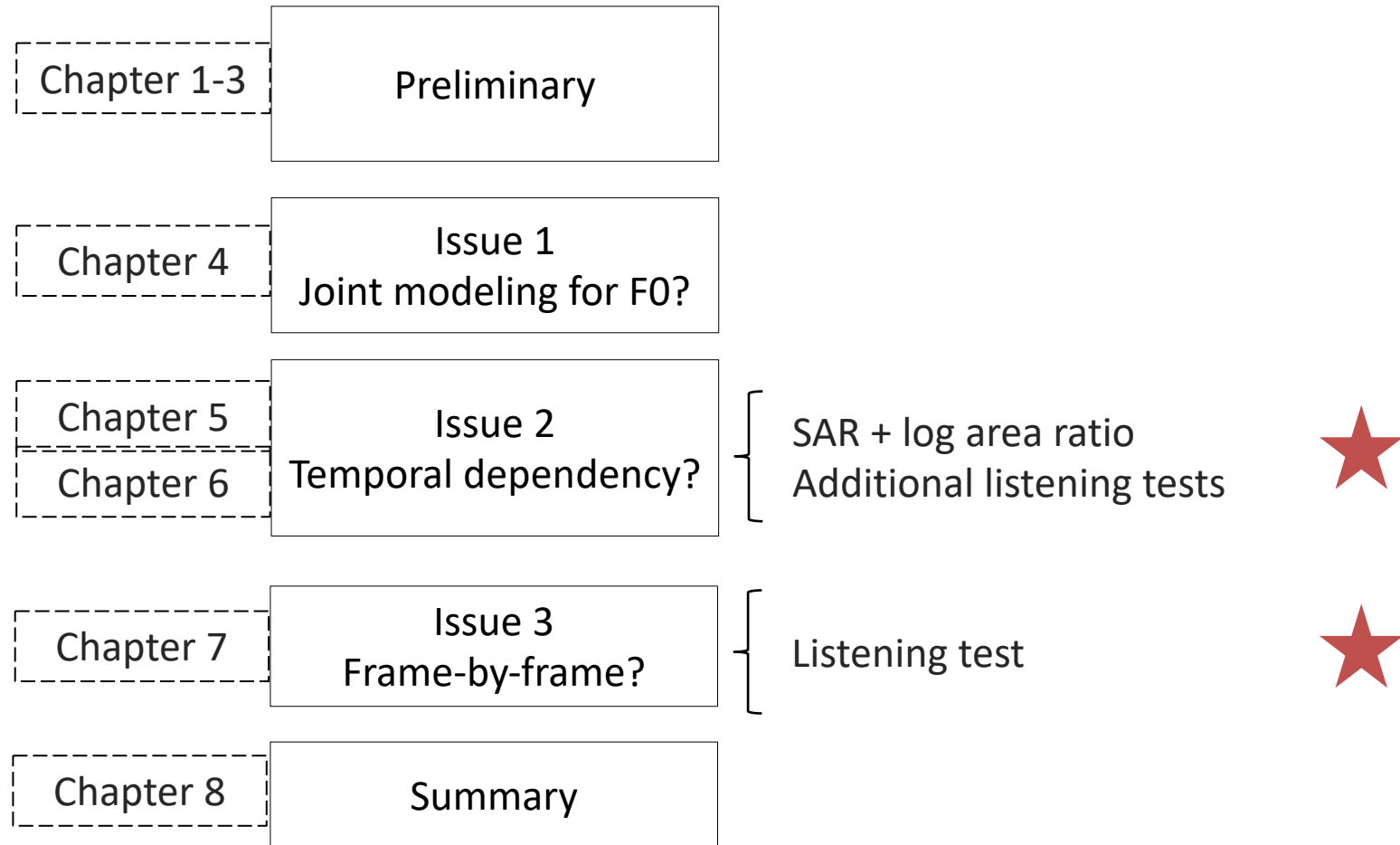
# INTRODUCTION

## Thesis outline



# INTRODUCTION

## Updated results



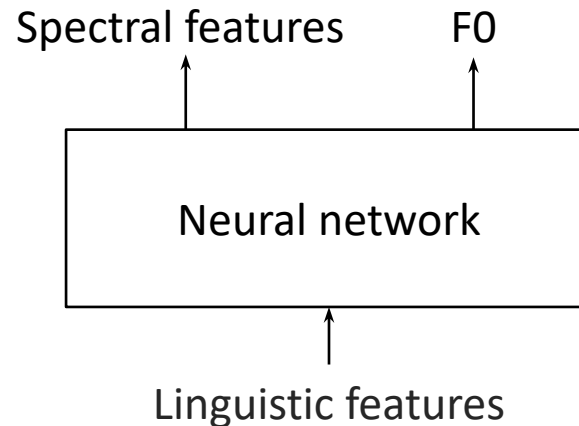
# CONTENTS

- Introduction
- Issue 1: joint modeling of F0 and spectral features
- Issues and methods
- Summary

# ISSUE 1: JOINT LEARNING FOR F0?

## Motivation

- Common approach [5, 7, 11]



- Joint (multi-task) learning
  - ? Beneficial for both targets
  - ? Sharing hidden features



True or not?  
More evidence?

- Empirical results against joint learning [5, 12]

[5] H. Zen, A. Senior, and M. Schuster. Statistical parametric speech synthesis using deep neural networks. In Proc. ICASSP, pages 7962–7966, 2013.

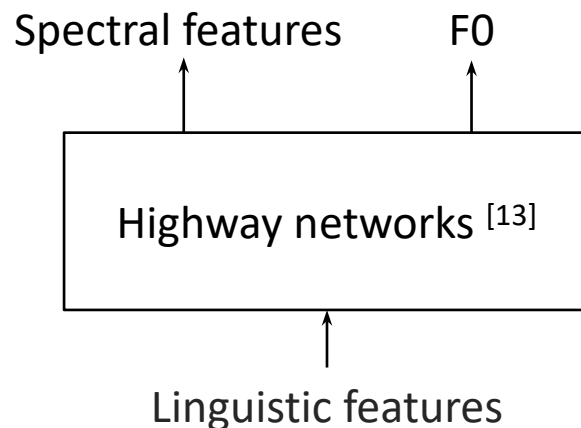
[7] Y. Fan, Y. Qian, F. Xie, and F. K. Soong. TTS synthesis with bidirectional LSTM based recurrent neural networks. In Proc. Interspeech, pages 1964–1968, 2014.

[11] H. Zen and A. Senior. Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis. In Proc. ICASSP, pages 3844–3848, 2014.

[12] S. Kang and H. Meng. Statistical parametric speech synthesis using weighted multi-distribution deep belief network. In Proc. Interspeech, pages 1959–1963, 2014.

# ISSUE 1: JOINT LEARNING FOR F0?

## Method



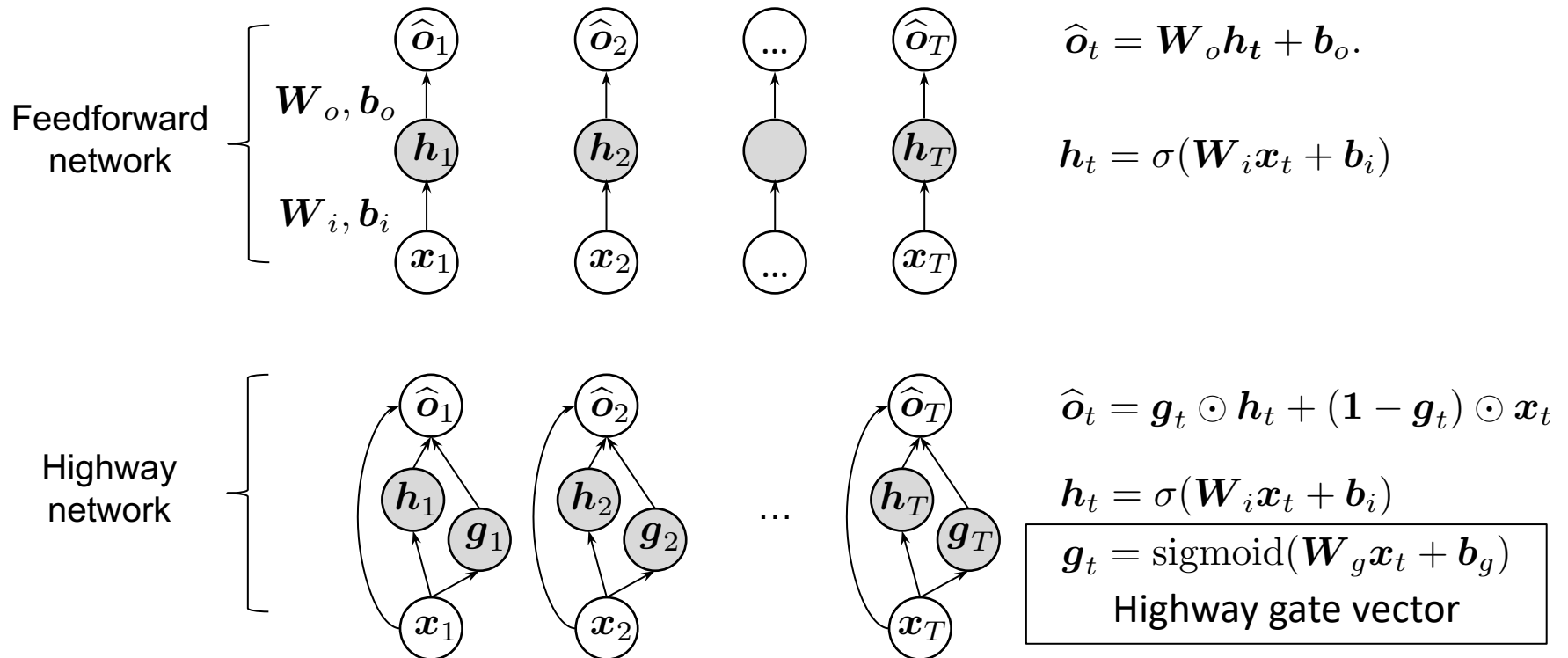
- Joint (multi-task) learning
  - ? Beneficial for both targets ←
  - ? Sharing hidden features
- Model and tools:
  - Highway network [13]
  - Histogram & sensitivity tools

[13] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. In Proc. Deep Learning Workshop, 2015.

# ISSUE 1: JOINT LEARNING FOR F0?

## Method

### □ Definition of highway network

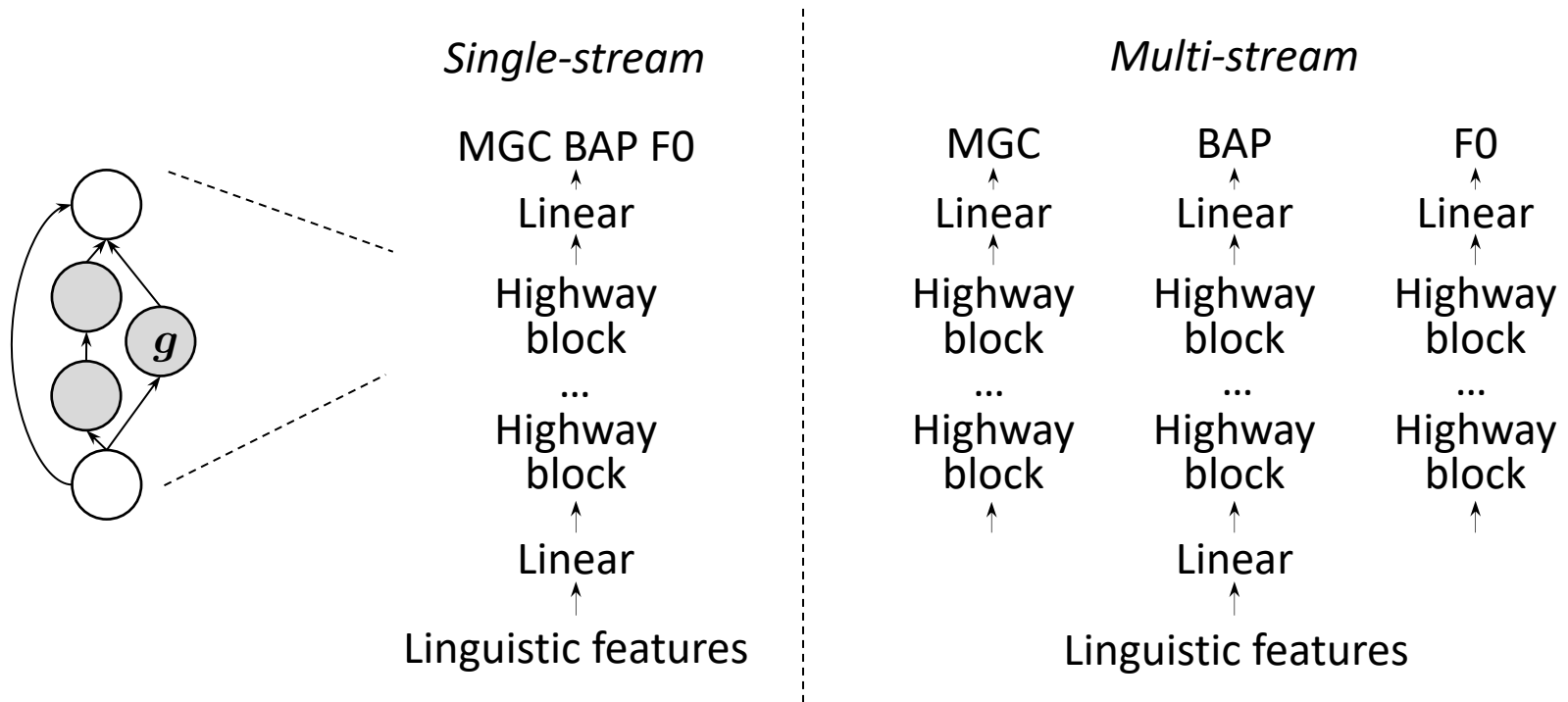




# ISSUE 1: JOINT LEARNING FOR F0?

## Method

- Highway network for acoustic modeling



- Spectral features { MGC: Mel-generalized cepstral (MGC) coefficients [14]  
BAP: band aperiodicity coefficients [15, 16]

[14] K. Tokuda, T. Kobayashi, T. Masuko, and S. Imai. Mel-generalized cepstral analysis a unified approach. In Proc. ICSLP, pages 1043–1046, 1994.

[15] H. Kawahara, J. Estill, and O. Fujimura. Aperiodicity extraction and control using mixed mode excitation and group delay manipulation for a high quality speech analysis, modification and synthesis system straight. In Second International Workshop on Models and Analysis of Vocal Emissions for Biomedical Applications, 2001.

[16] H. Zen and T. Toda. An overview of Nitech HMM-based speech synthesis system for Blizzard Challenge 2005. In Proc. Interspeech, pages 93–96, 2005.

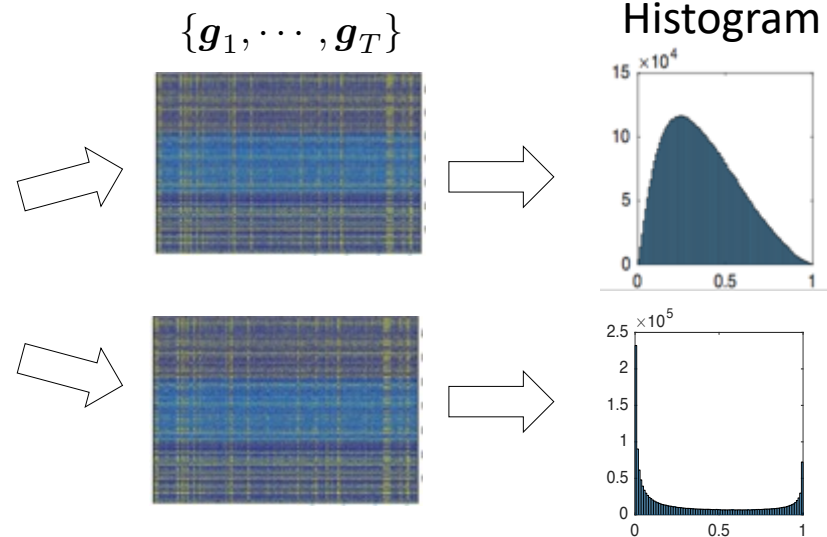
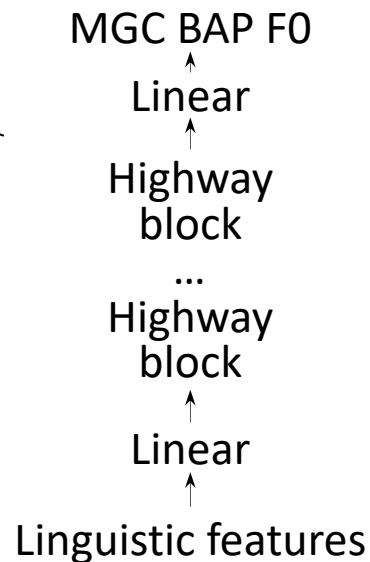
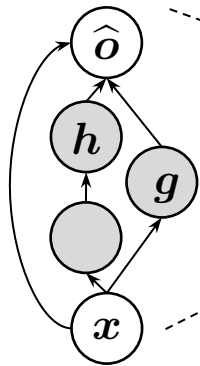
# ISSUE 1: JOINT LEARNING FOR F0?

## Method

### □ Analysis tools

#### 1. Histogram of gate vectors $g$

$$\hat{o} = g \odot h + (1 - g) \odot x$$



$g \approx 1$  Non-linear transformation  
 $g \approx 0$  No transformation

#### 2. Sensitivity of $g$ to different linguistic features (Sec. 4.3.2)

# ISSUE 1: JOINT LEARNING FOR F0?

## Experiments

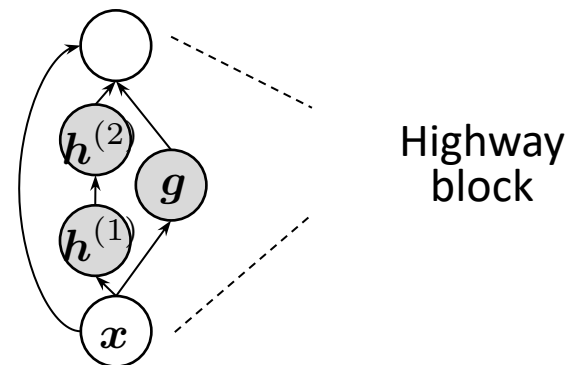
### □ Configuration

- Data: English, 16 hours
- Feature: MGC, BAP, F0 (Interpolated F0 + voicing (U/V))
- Metric: 

{	Root mean square error	(RMSE)
	Correlation coefficients	(CORR)

### □ Three models:

- Single-stream feedforward network
- Single-stream highway network
- Multi-stream highway network



### □ Two tests:

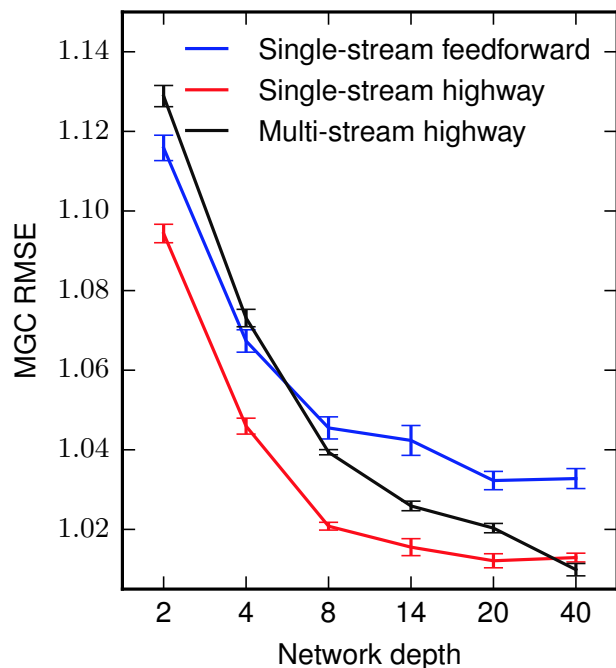
1. Fixed layer width, varying network depth
2. Fixed network depth, varying layer width

# ISSUE 1: JOINT LEARNING FOR F0?

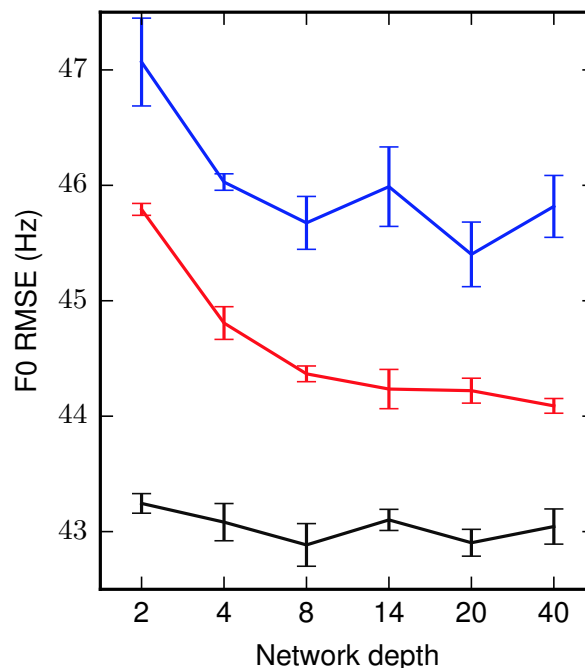
## Experiments

Objective results: increasing network depth

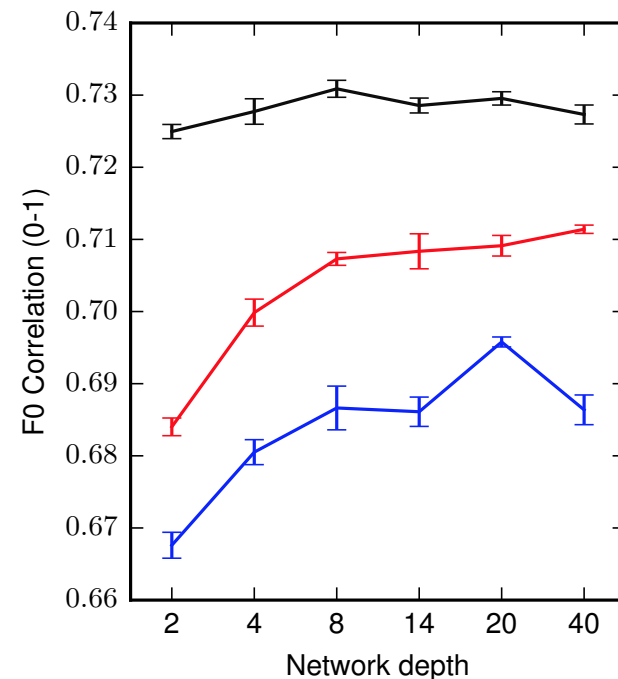
### MGC RMSE



### F0 RMSE



### F0 CORR



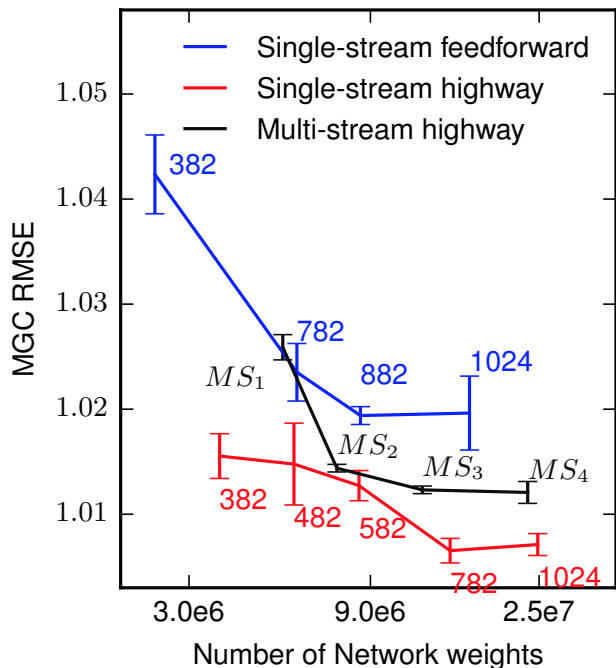
- ❖ Network depth: Number of tanh-based transformation layers
- Single-stream network prioritizes MGC?

# ISSUE 1: JOINT LEARNING FOR F0?

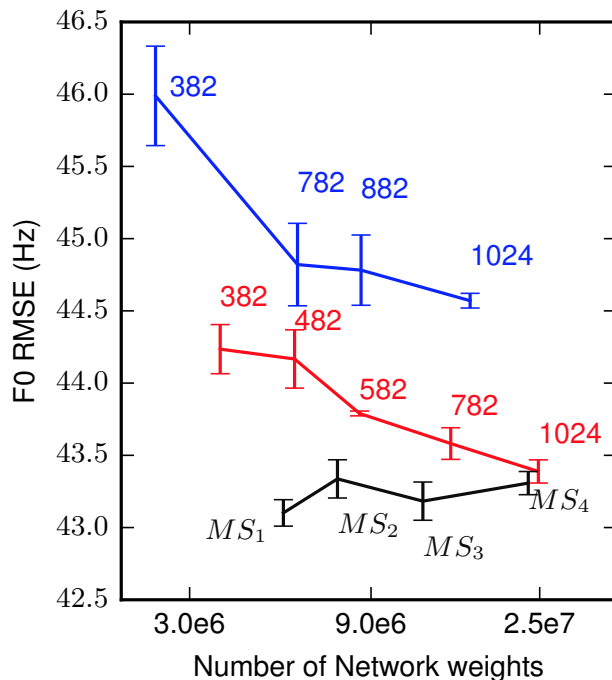
## Experiments

Objective results: increasing width (depth = 14)

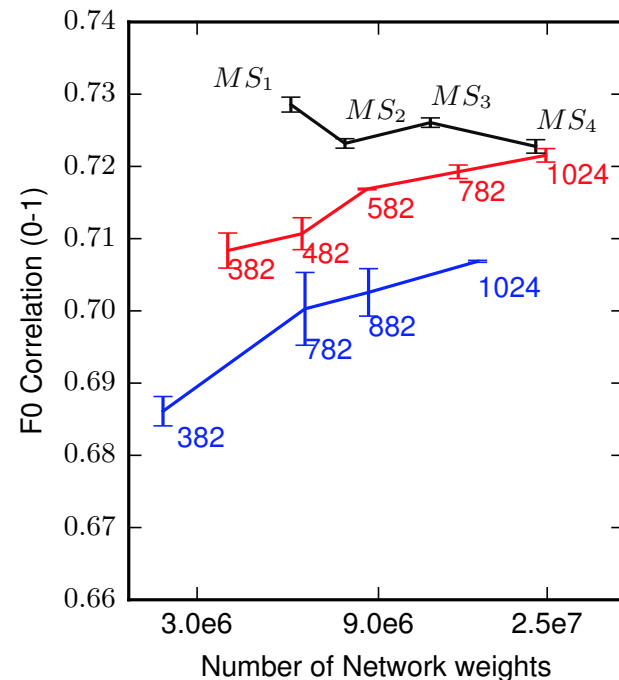
### MGC RMSE



### F0 RMSE



### F0 CORR



- ❖ MS<sub>1</sub> : [MGC 256] – [F0 256]
- ❖ MS<sub>2</sub> : [MGC 382] – [F0 256]
- ❖ MS<sub>3</sub> : [MGC 512] – [F0 382]
- ❖ MS<sub>4</sub> : [MGC 768] – [F0 512]

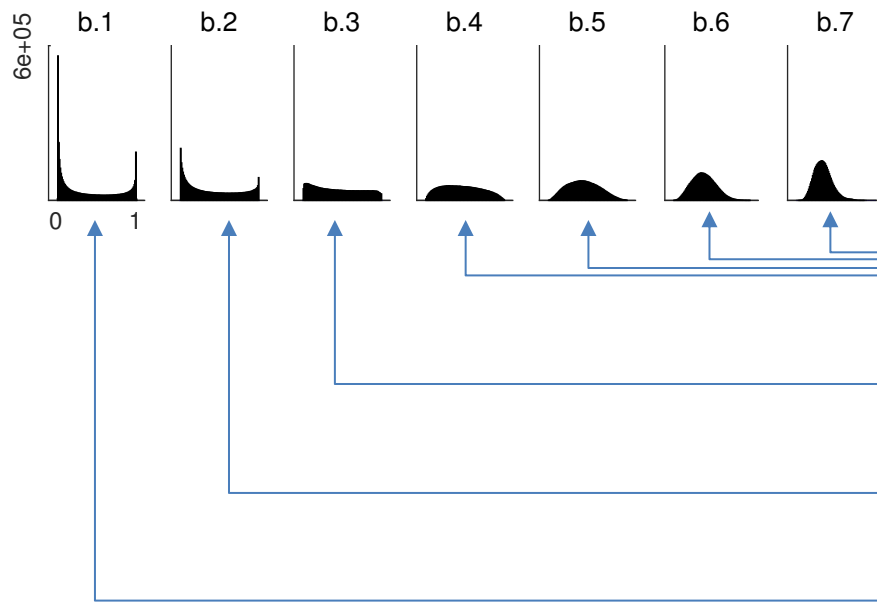
• Single-stream network prioritizes MGC?

# ISSUE 1: JOINT LEARNING FOR F0?

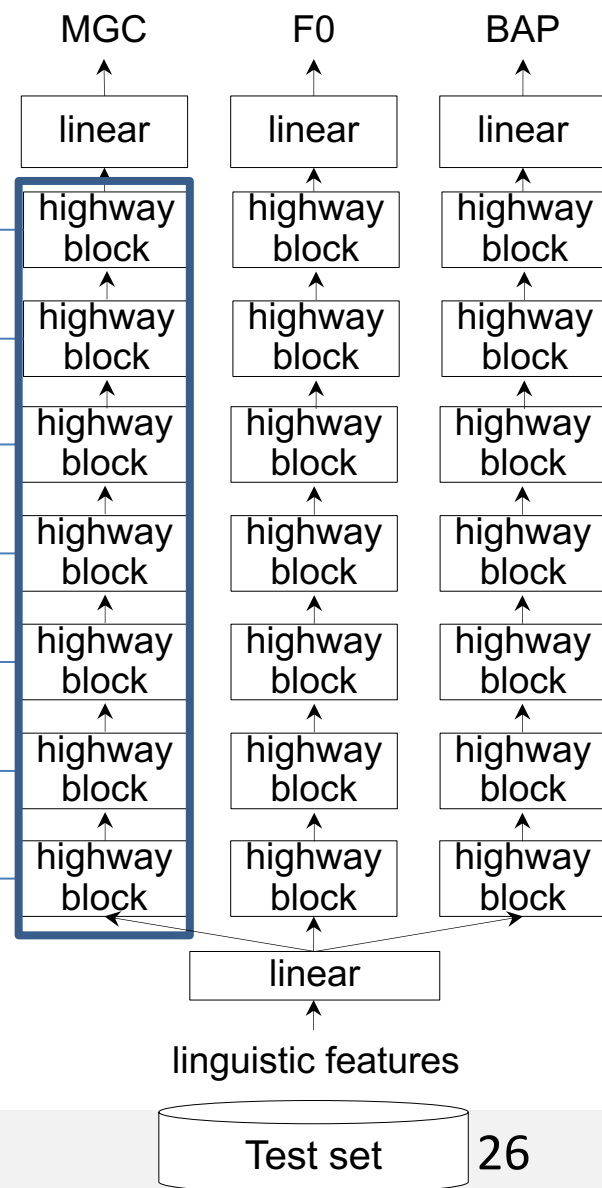
## Experiments

### □ Histogram of $g$

- Multi-stream highway (depth 14)



- $g \approx 1$  Non-linear transformation  
 $g \approx 0$  No transformation

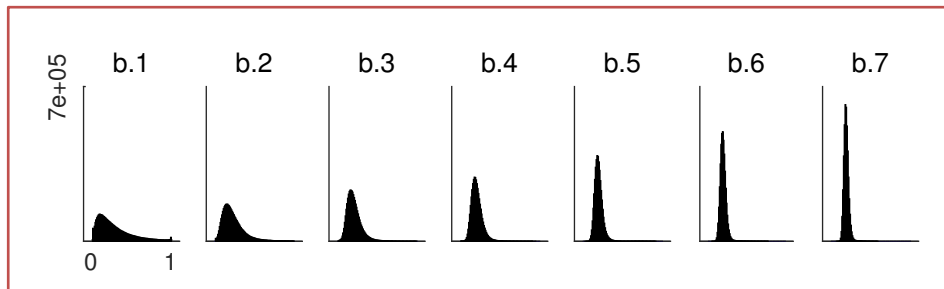
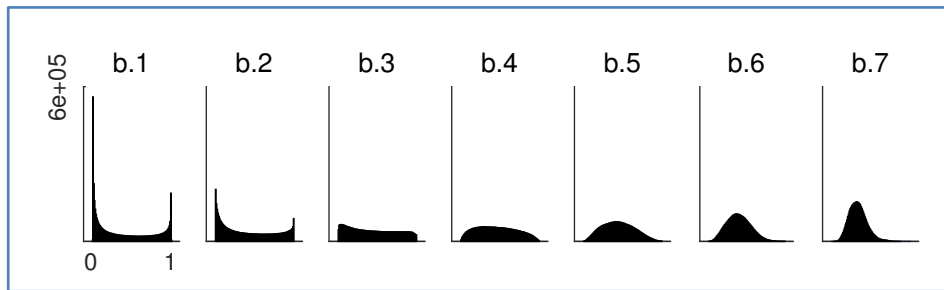


# ISSUE 1: JOINT LEARNING FOR F0?

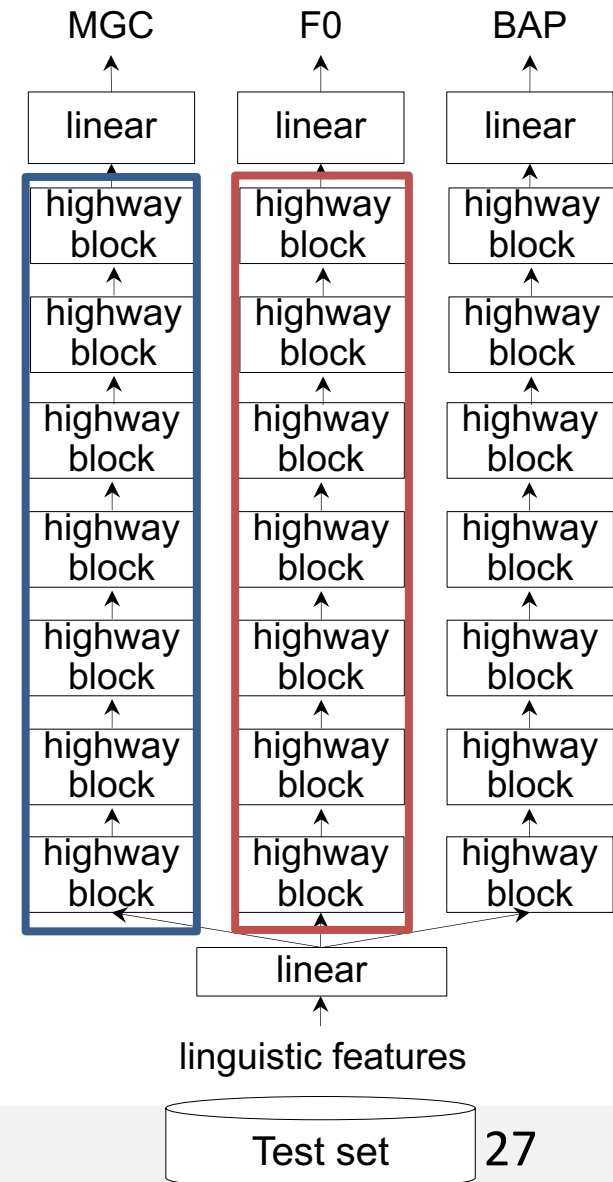
## Experiments

### □ Histogram of $g$

- Multi-stream highway (depth 14, 7 blocks)



- $g \approx 1$  Non-linear transformation
- $g \approx 0$  No transformation
- Different hidden features for MGC and F0

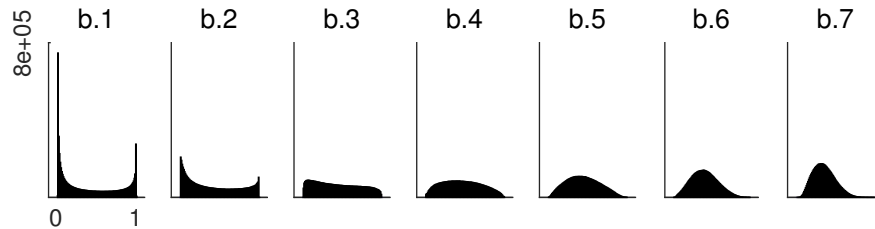


# ISSUE 1: JOINT LEARNING FOR F0?

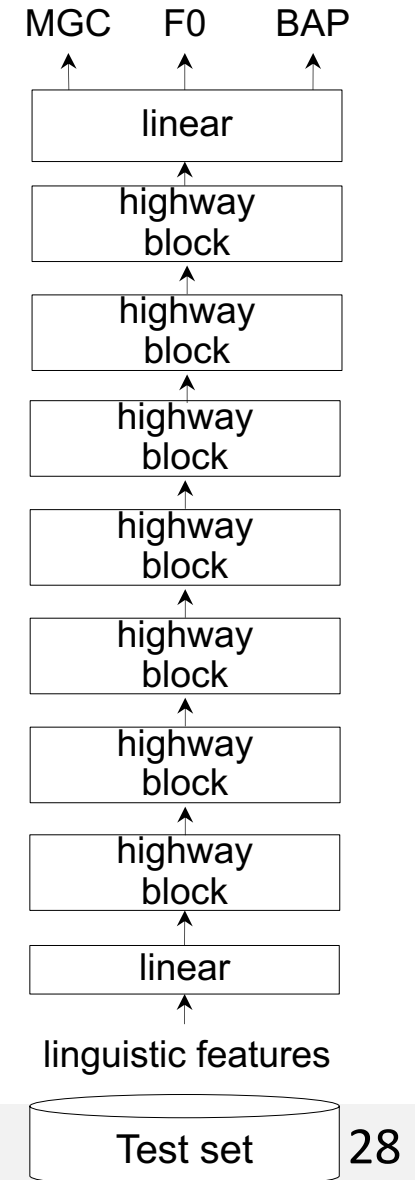
## Experiments

### □ Histogram of $g$

- Single-stream highway (depth 14, 7 blocks)



- Similar to MGC sub-net in multi-stream highway
- Single-stream network prioritizes MGC?





# ISSUE 1: JOINT LEARNING FOR F0?

## Summary

### □ Answer to issue 1

#### Joint (multi-task) learning

- ? Beneficial for both F0 and spectral features
- ? They share hidden features

#### Negative evidence

- Joint learning (single-stream network) prioritizes spectral features
- They use different hidden features
- They use different input features (Sec.4.4.3)
- Results on English and Japanese corpora (Sec.4.5)

NOT for the sake of F0 modeling!

□ Only F0 modeling in the following chapters

□ F0 is useful for MGC modeling? How to do? (slides appendix)



# CONTENTS

- Introduction
- Issue 1: joint modeling of F0 and spectral features
- Issue 2: temporal dependency modeling of F0 contours
- Issues and methods
- Summary

# ISSUE 2: TEMPORAL DEPENDENCY?

## Motivation

- Baseline RNN model [17]

$$\hat{\mathbf{o}}_{1:T} = \{\hat{\mathbf{o}}_1, \dots, \hat{\mathbf{o}}_T\}$$

F0 contour



Recurrent neural network (RNN)



Linguistic features

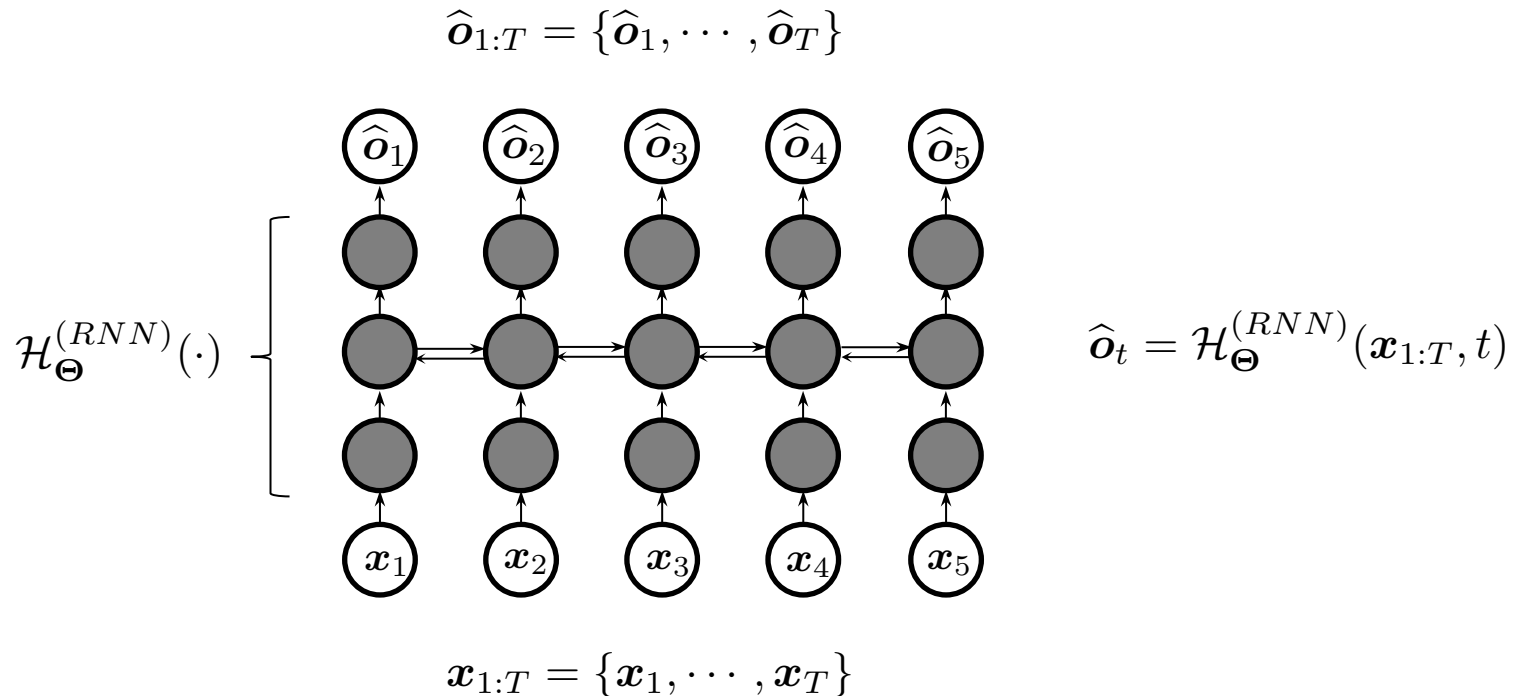
$$\mathbf{x}_{1:T} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$$

- ❖  $T$ : number of frames (time steps)

# ISSUE 2: TEMPORAL DEPENDENCY?

## Motivation

- Baseline RNN model



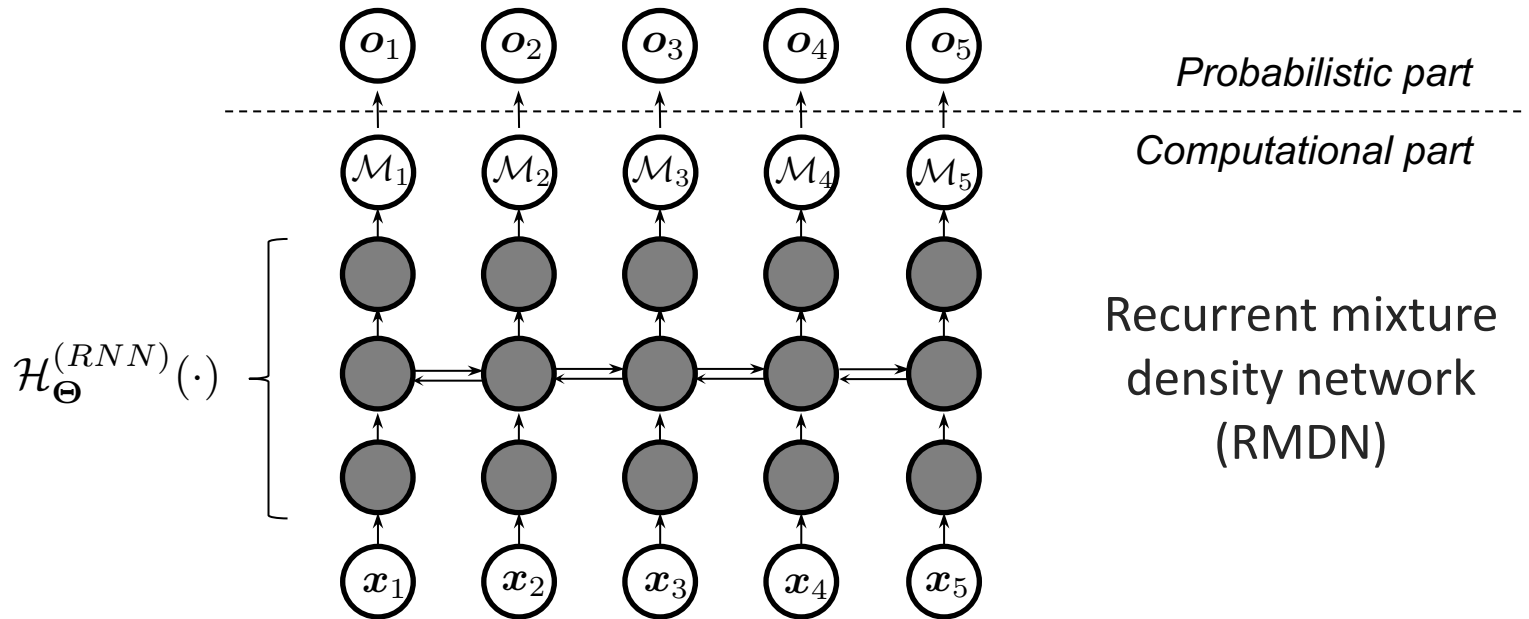
- Learn the correlation between  $o_{t_1}$  and  $o_{t_2}$ ,  $t_2 \neq t_1$



# ISSUE 2: TEMPORAL DEPENDENCY?

## Motivation

### □ Baseline RNN model



$$p(o_{1:T} | x_{1:T}; \Theta) = \prod_{t=1}^T p(o_t | x_{1:T}; \Theta) = \prod_{t=1}^T \mathcal{N}(o_t; \mu_t, \sigma I)$$

$$\mathcal{M}_t = \{\mu_t\}, \quad \text{where} \quad \mu_t = \mathcal{H}_{\Theta}^{(RNN)}(x_{1:T}, t)$$

$$\hat{o}_t = \arg \max_{o_t} p(o_t | x_{1:T}; \Theta^*) = \mu_t \quad \text{Mean-based generation}$$

[18] C. M. Bishop. Mixture Density Networks. Technical report, Aston University, 2004.

[19] C. M. Bishop. Neural networks for pattern recognition. Oxford university press, 1995.

[20] M. Schuster. Better generative models for sequential data problems: Bidirectional recurrent mixture density networks. In Proc. NIPS, pages 589–595, 1999.

# ISSUE 2: TEMPORAL DEPENDENCY?

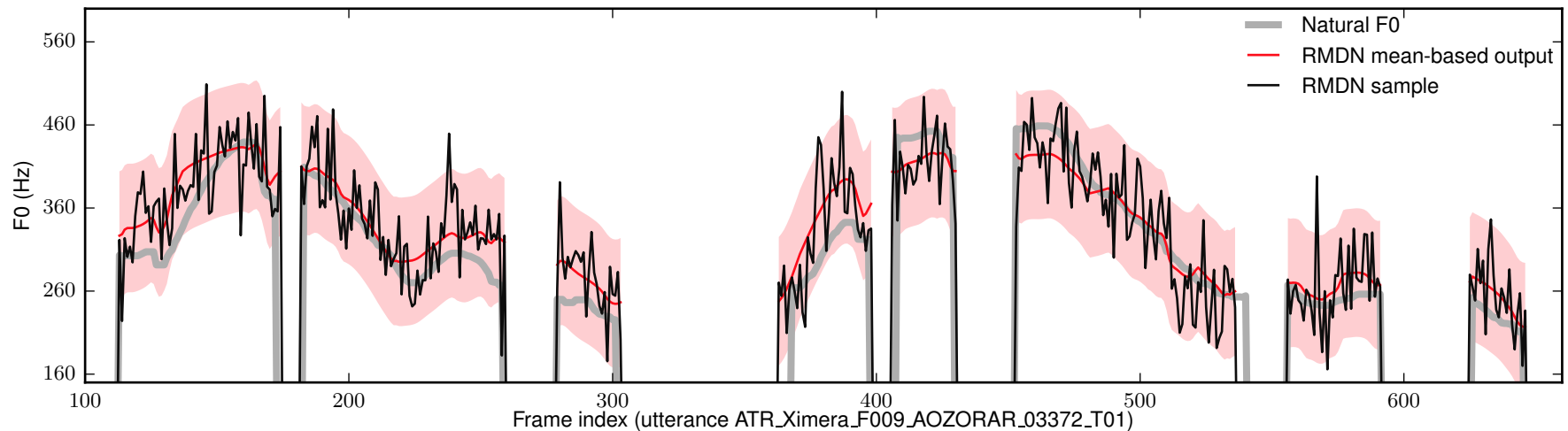
## Motivation

□ Initial answer

Temporal dependency is ignored by RNN/RMDN

$$p(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta) = \prod_{t=1}^T p(o_t | \mathbf{x}_{1:T}; \Theta) = \prod_{t=1}^T \mathcal{N}(o_t; \boldsymbol{\mu}_t, \sigma \mathbf{I})$$

- Evidence from random sampling



# ISSUE 2: TEMPORAL DEPENDENCY?

## Motivation

- Initial answer

Temporal dependency is ignored by RNN/RMDN

- Better models?

RNN/RMDN

$$p(\mathbf{o}_{1:T}) = \prod_{t=1}^T p(\mathbf{o}_t)$$



Autoregressive (AR) [23] models

$$p(\mathbf{o}_{1:T}) = \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{o}_{1:t-1})$$

- Shallow AR models (SAR) & extended SAR
- Deep AR models (DAR)

# CONTENTS

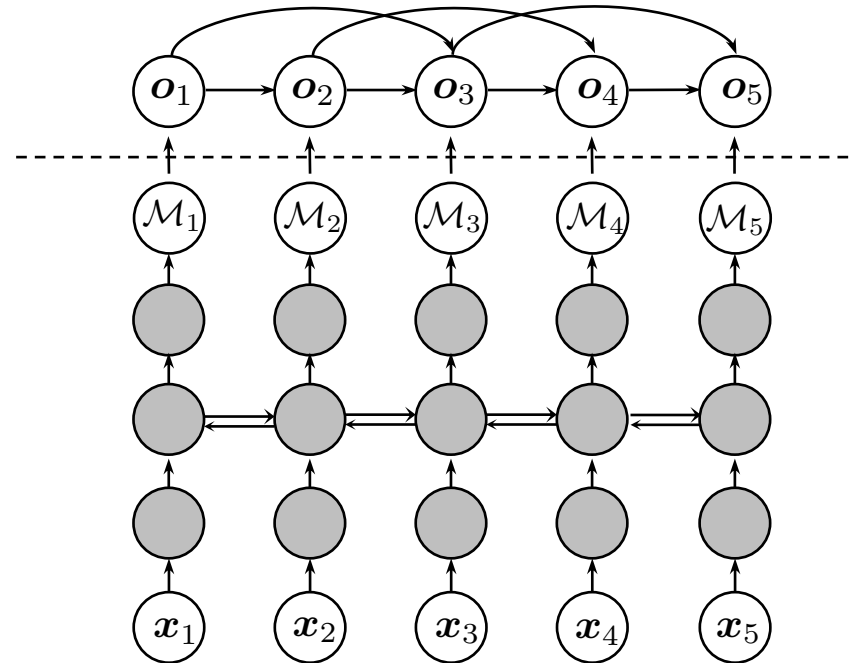
- Introduction
- Issue 1: joint modeling of F0 and spectral features
- Issue 2: temporal dependency modeling of F0 contours
  - SAR and extension
  - DAR
- Issues and methods
- Summary



# ISSUE 2: TEMPORAL DEPENDENCY?

## SAR

### □ Definition

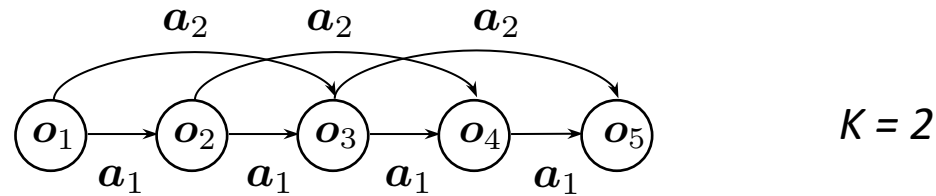


$$p(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{o}_{t-K:t-1}, \mathbf{x}_{1:T})$$

# ISSUE 2: TEMPORAL DEPENDENCY?

## SAR

### □ Definition



$$\begin{aligned} p(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) &= \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{o}_{t-K:t-1}, \mathbf{x}_{1:T}; \Theta, \Psi) \\ &= \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_t + f_{\Psi}(\mathbf{o}_{t-K:t-1}), \boldsymbol{\Sigma}_t) \end{aligned}$$

$$f_{\Psi}(\mathbf{o}_{t-K:t-1}) = \sum_{k=1}^K \mathbf{a}_k \odot \mathbf{o}_{t-k} + \mathbf{b}$$

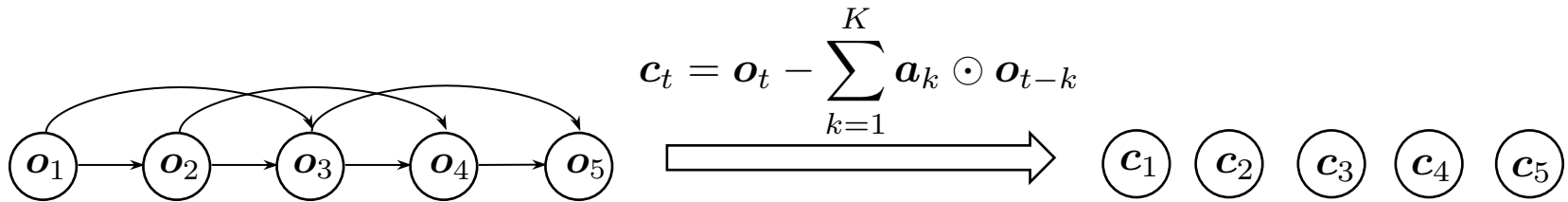
Trainable parameters

- Time invariant  $\Psi = \{\mathbf{a}_1, \dots, \mathbf{a}_K, \mathbf{b}\}$
- $K$ : hyper-parameter

# ISSUE 2: TEMPORAL DEPENDENCY?

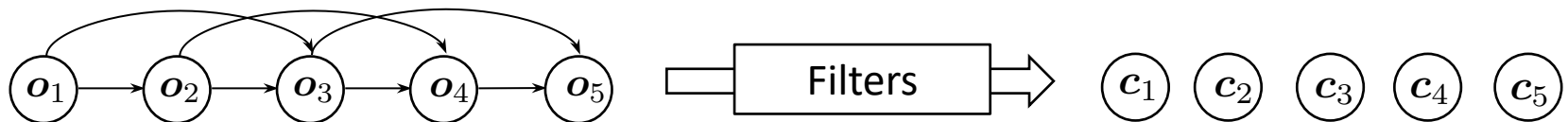
## SAR

### □ Interpretation 1 (Sec. 5.3)



$$\begin{aligned} p(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) &= \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{o}_{t-K:t-1}, \mathbf{x}_{1:T}; \Theta, \Psi) \\ &= \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_t + f_{\Psi}(\mathbf{o}_{t-K:t-1}), \boldsymbol{\Sigma}_t) \\ &= \prod_{t=1}^T p_c(\mathbf{c}_t | \mathbf{x}_{1:T}; \Theta) \end{aligned}$$

### □ Interpretation 2 (Sec. 5.3)

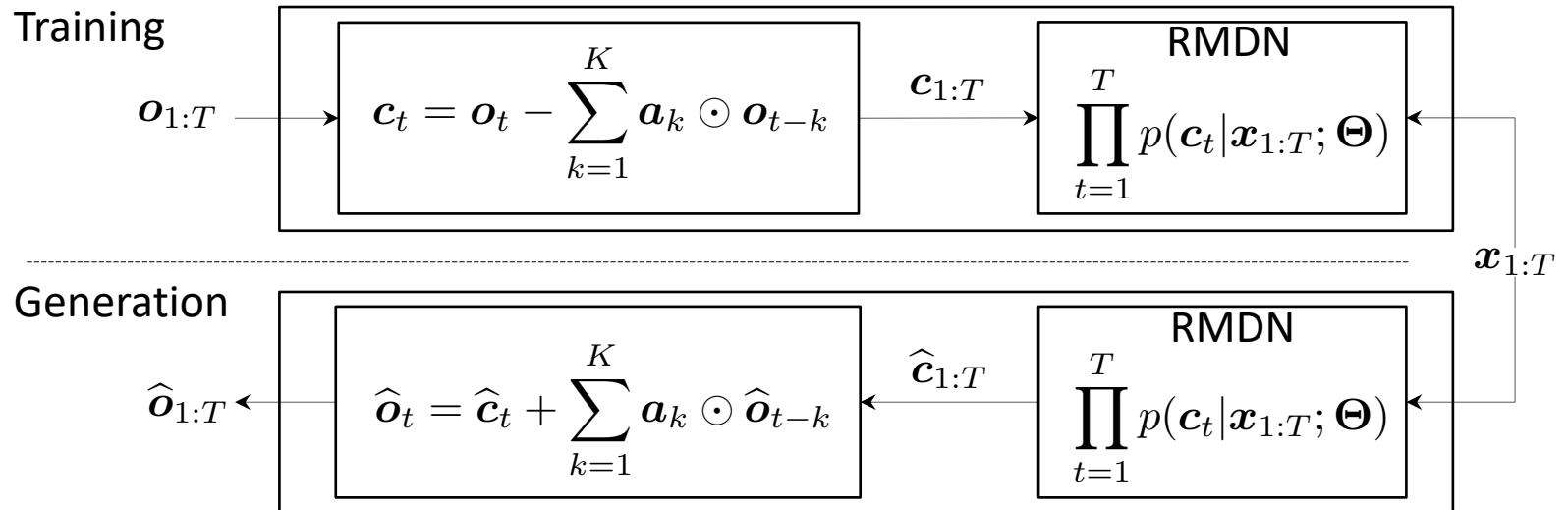


# ISSUE 2: TEMPORAL DEPENDENCY?

## SAR

### □ Interpretation (Sec. 5.3)

- SAR = Linear transformation + RMDN

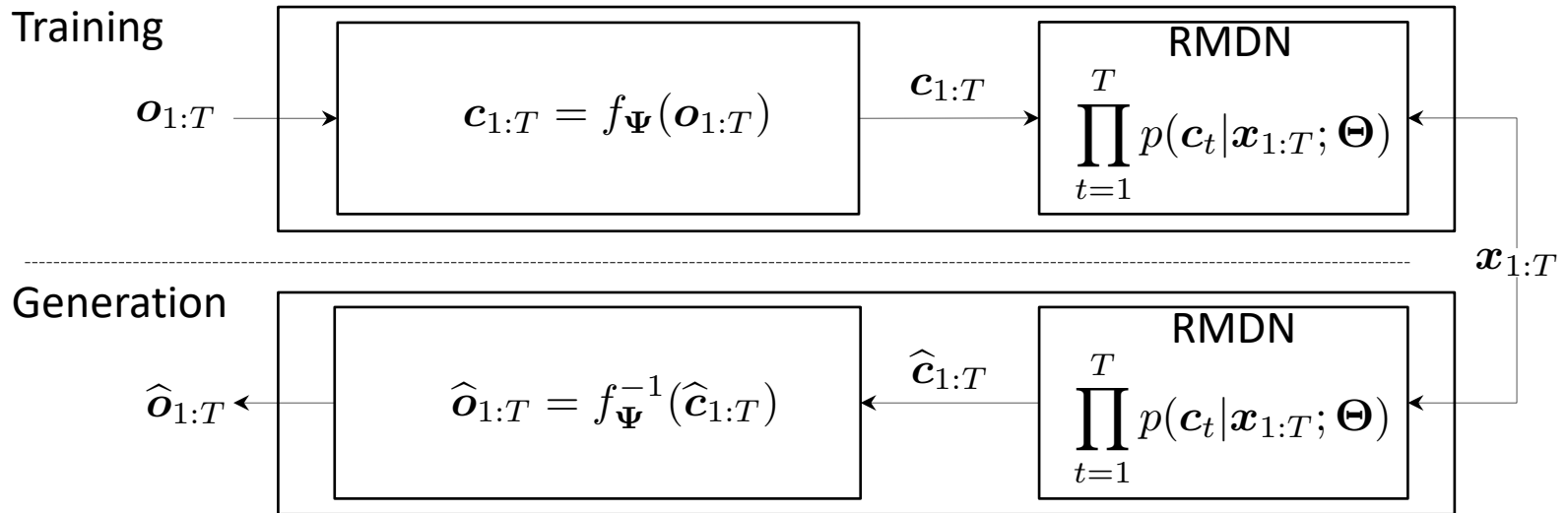


# ISSUE 2: TEMPORAL DEPENDENCY?

## Extended SAR (eSAR)

### □ Motivation

- SAR -> Non-linear transformation + RMDN?



$$p_o(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) = p_c(\mathbf{c}_{1:T} = f_{\Psi}(\mathbf{o}_{1:T}) | \mathbf{x}_{1:T}; \Theta) \left| \det \frac{\partial f_{\Psi}(\mathbf{o}_{1:T})}{\partial \mathbf{o}_{1:T}} \right|$$

- Yes, if  $f_{\Psi}(\mathbf{o}_{1:T})$  is invertible and 'simple'

# ISSUE 2: TEMPORAL DEPENDENCY?

## Extended SAR (eSAR)

### □ Definition

$$p_o(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) = p_c(\mathbf{c}_{1:T} = f_{\Psi}(\mathbf{o}_{1:T}) | \mathbf{x}_{1:T}; \Theta) \left| \det \frac{\partial f_{\Psi}(\mathbf{o}_{1:T})}{\partial \mathbf{o}_{1:T}} \right|$$

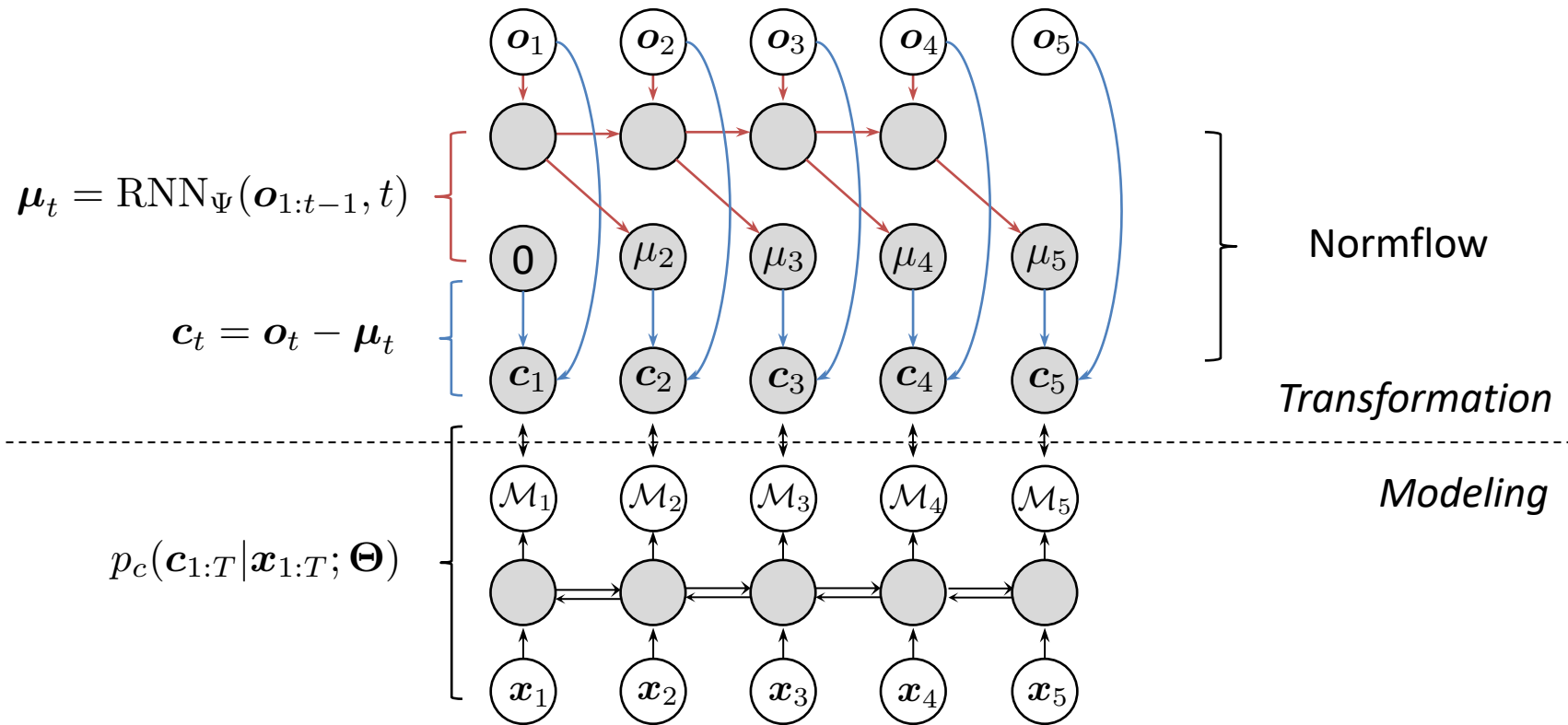
	SAR	eSAR
$\mathbf{c}_{1:T} = f_{\Psi}(\mathbf{o}_{1:T})$	$\mathbf{c}_t = \mathbf{o}_t - \sum_{k=1}^K \mathbf{a}_k \odot \mathbf{o}_{t-k}$	$\mathbf{c}_t = \mathbf{o}_t - \text{RNN}_{\Psi}(\mathbf{o}_{1:t-1}, t)$
$\hat{\mathbf{o}}_{1:T} = f_{\Psi}^{-1}(\hat{\mathbf{c}}_{1:T})$	$\hat{\mathbf{o}}_t = \hat{\mathbf{c}}_t + \sum_{k=1}^K \mathbf{a}_k \odot \hat{\mathbf{o}}_{t-k}$	$\hat{\mathbf{o}}_t = \hat{\mathbf{c}}_t + \text{RNN}_{\Psi}(\hat{\mathbf{o}}_{1:t-1}, t)$
$\det \frac{\partial f_{\Psi}(\mathbf{o}_{1:T})}{\partial \mathbf{o}_{1:T}}$	$\det \frac{\partial f_{\Psi}(\mathbf{o}_{1:T})}{\partial \mathbf{o}_{1:T}} = 1$	$\det \frac{\partial f_{\Psi}(\mathbf{o}_{1:T})}{\partial \mathbf{o}_{1:T}} = 1$

- Volume-preserving [24]

# ISSUE 2: TEMPORAL DEPENDENCY?

## Extended SAR (eSAR)

### Implementation

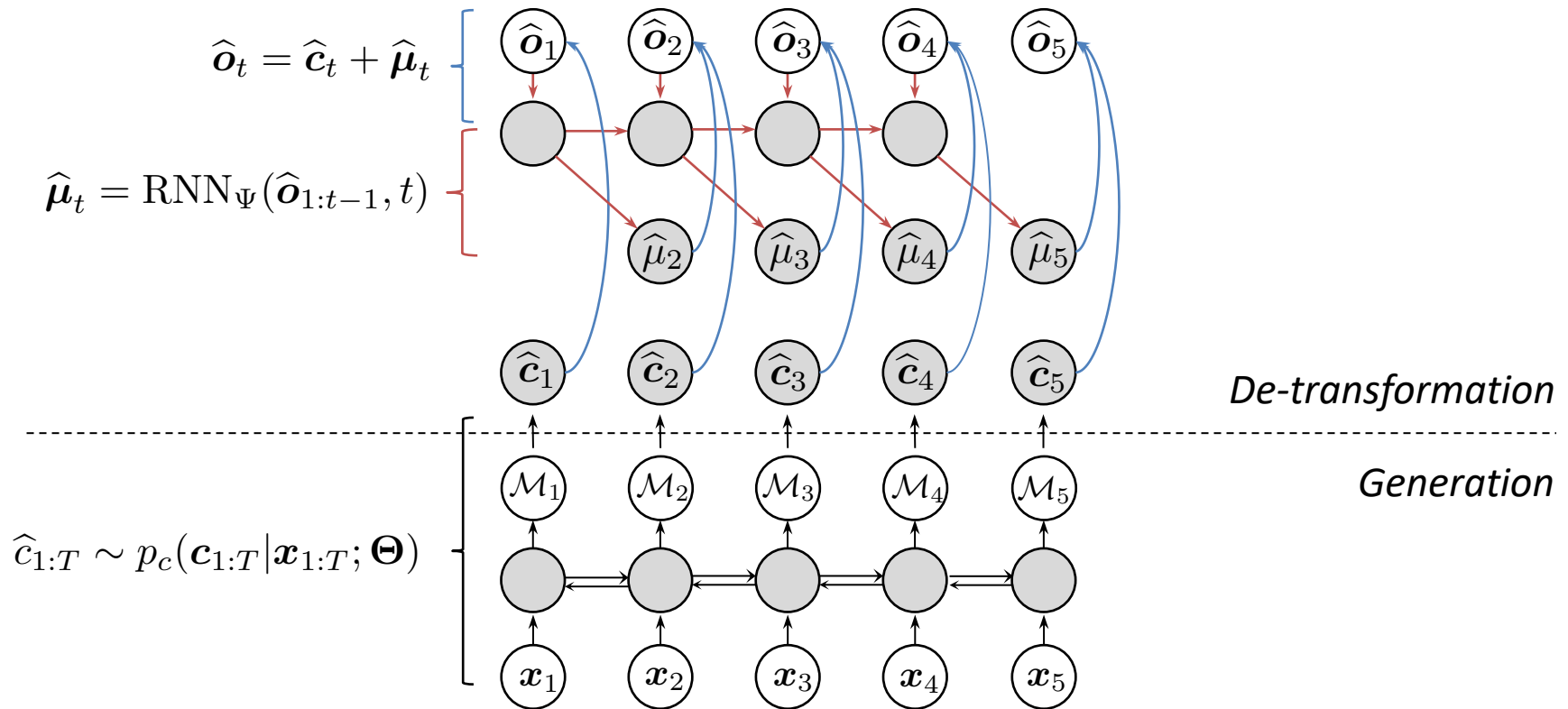


$$p_o(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) = p_c(\mathbf{c}_{1:T} = f_{\Psi}(\mathbf{o}_{1:T}) | \mathbf{x}_{1:T}; \Theta)$$

# ISSUE 2: TEMPORAL DEPENDENCY?

## Extended SAR (eSAR)

### Implementation



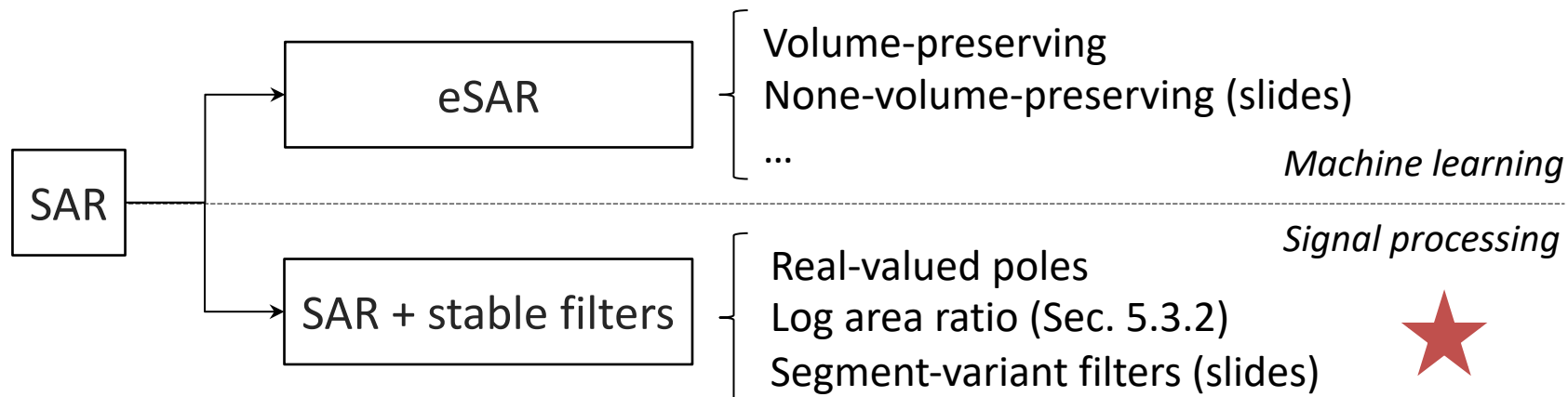
$$p_o(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta, \Psi) = p_c(\mathbf{c}_{1:T} = f_{\Psi}(\mathbf{o}_{1:T}) | \mathbf{x}_{1:T}; \Theta)$$



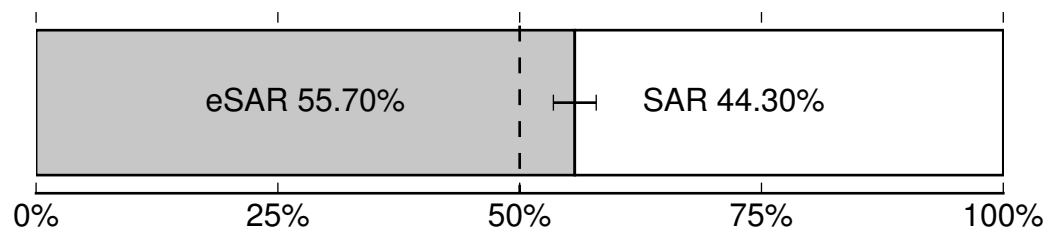
# ISSUE 2: TEMPORAL DEPENDENCY?

## Summary of SAR

- Theoretically appealing



- Performance for F0 modeling (model details later)



- Performance for MGC modeling

- Better than RNN/RMDN [25, 26]

[25] X. Wang, J. Lorenzo-Trueba, S. Takaki, L. Juvola, and J. Yamagishi. A comparison of recent waveform generation and acoustic modeling methods for neural-network-based speech synthesis. In Proc. ICASSP, pages 4804–4808, 2018.

[26] X. Wang, S. Takaki, and J. Yamagishi. An autoregressive recurrent mixture density network for parametric speech synthesis. In Proc. ICASSP, pages 4895–4899, 2017.

# CONTENTS

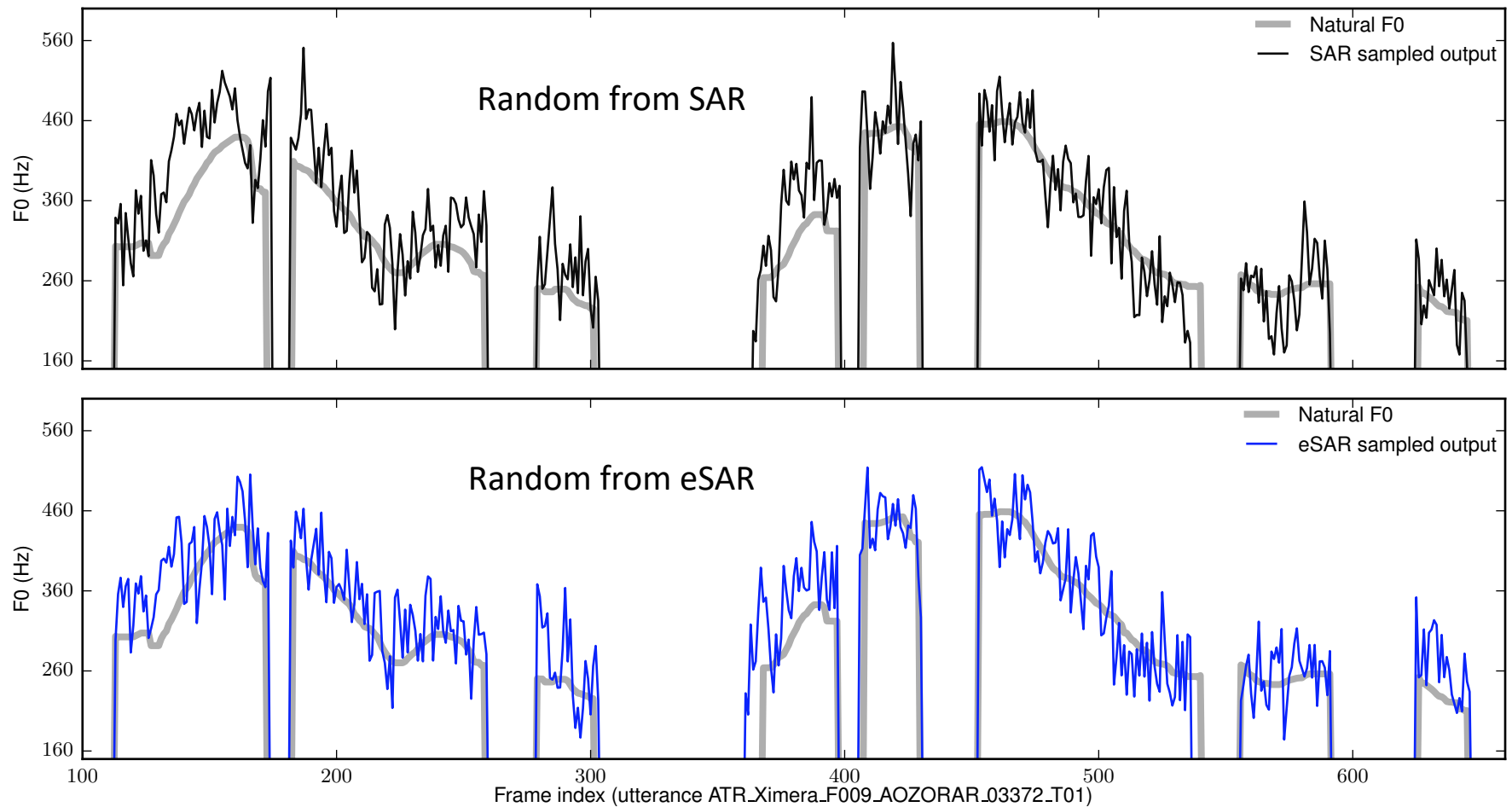
- Introduction
- Issue 1: joint modeling of F0 and spectral features
- Issue 2: temporal dependency modeling of F0 contours
  - SAR and extension
  - DAR
- Issues and methods
- Summary

# ISSUE 2: TEMPORAL DEPENDENCY?

## DAR

### □ Motivation

- Are SAR and eSAR sufficiently good?

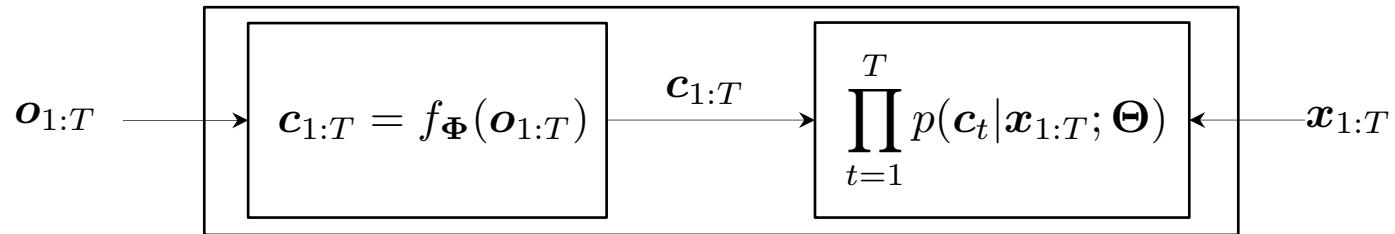


# ISSUE 2: TEMPORAL DEPENDENCY?

## DAR

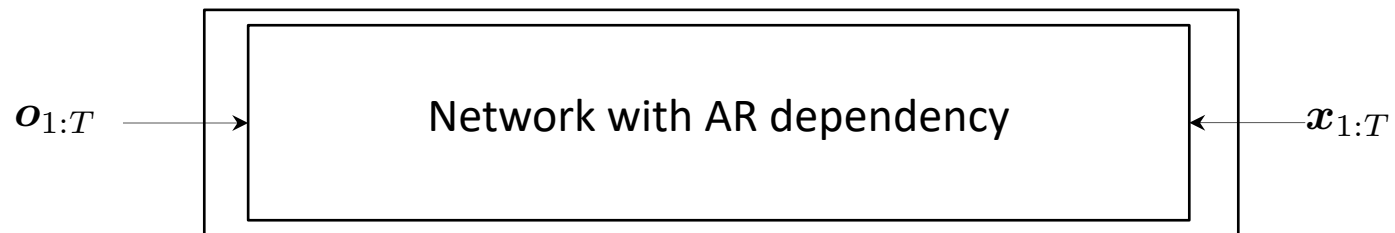
### □ Motivation

- Random sampling on SAR and eSAR?



- SAR: linear  $f_{\Phi}(\cdot)$  (Sec. 6.1)
- eSAR: non-linear but a special form

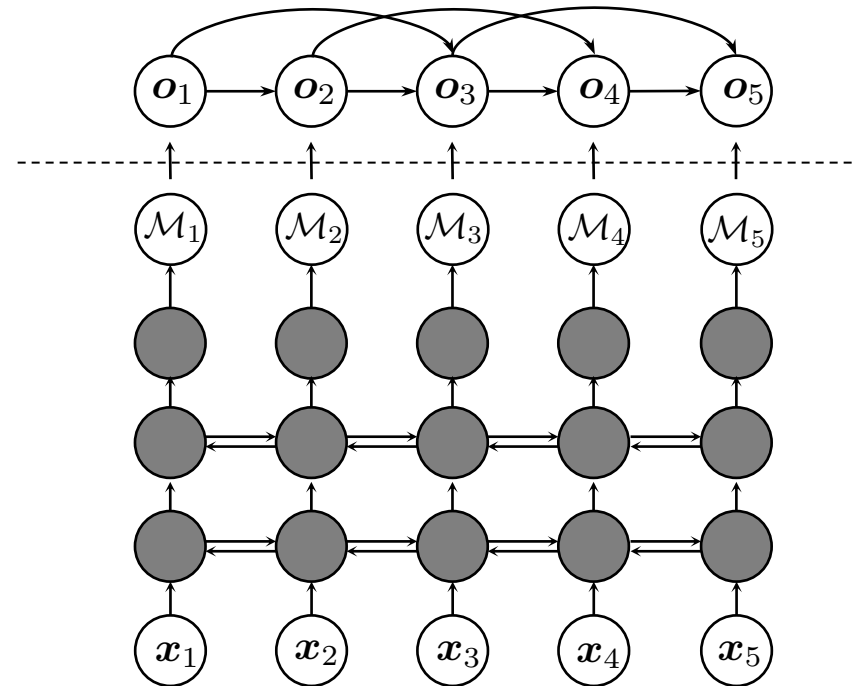
### □ Non-linear and non-invertible AR transformation?



# ISSUE 2: TEMPORAL DEPENDENCY?

## DAR

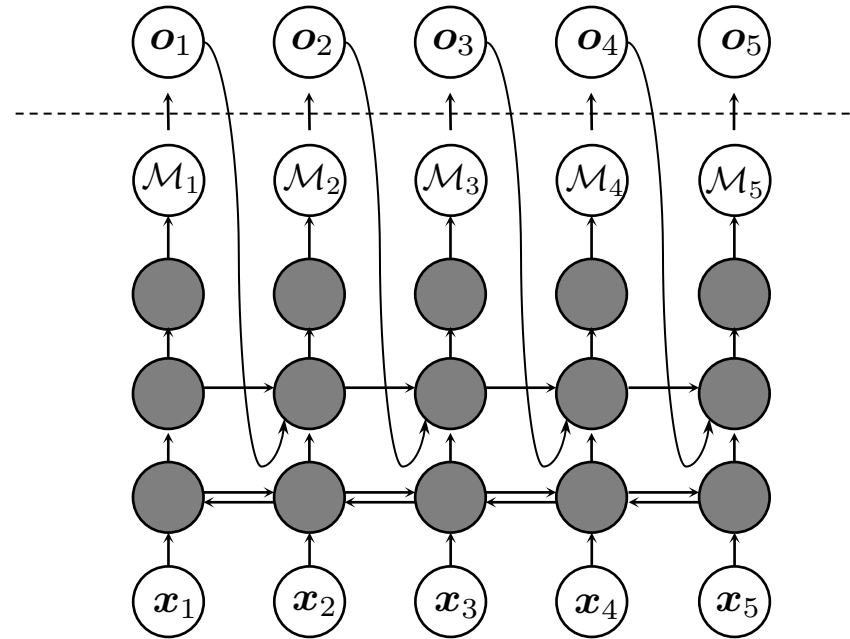
### □ Definition



# ISSUE 2: TEMPORAL DEPENDENCY?

## DAR

### □ Definition



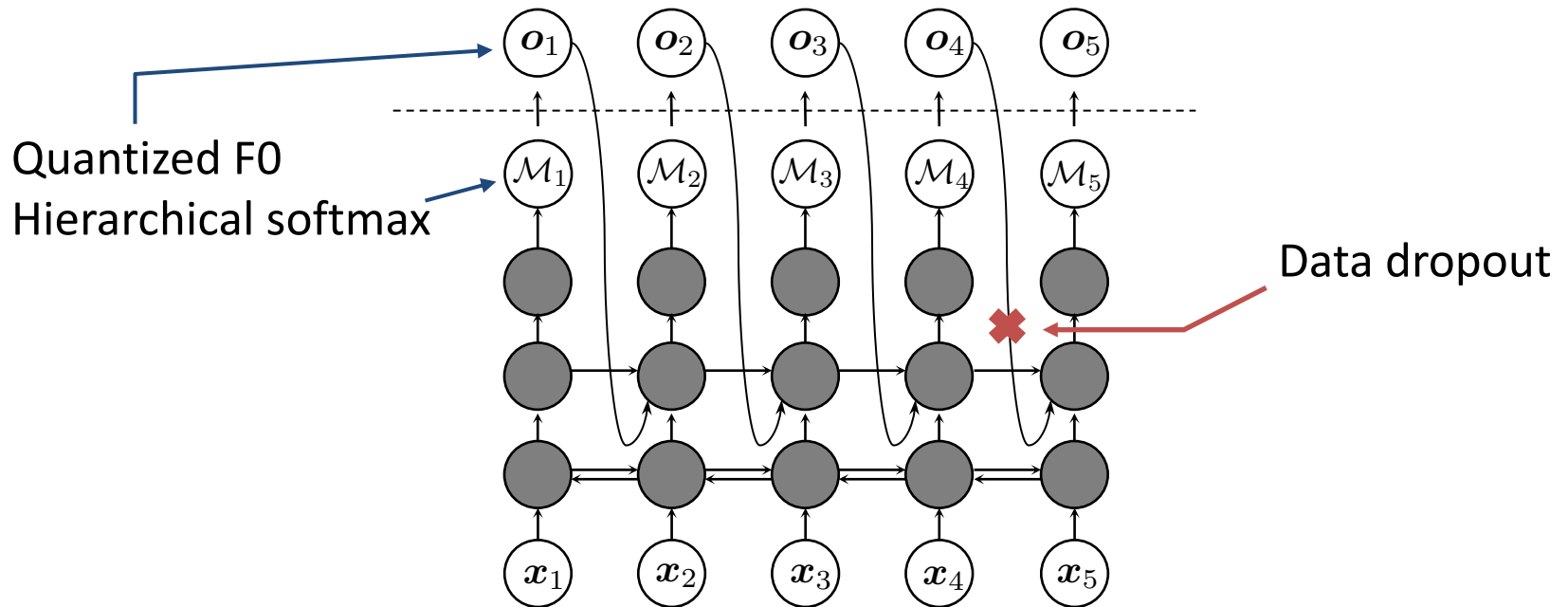
$$p(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta) = \prod_{t=1}^T p(o_t | \mathbf{o}_{1:t-1}, \mathbf{x}_{1:T}; \Theta)$$

- DAR is more general than SAR (Sec 6.2.2 & slides, toy examples)
  1. Longer-time dependency
  2. Non-linear dependency

# ISSUE 2: TEMPORAL DEPENDENCY?

## DAR

### Implementation (Sec. 6.3)



$$P(\mathbf{o}_{1:T} | \mathbf{x}_{1:T}; \Theta) = \prod_{t=1}^T P(o_t | \mathbf{o}_{1:t-1}, \mathbf{x}_{1:T}; \Theta)$$

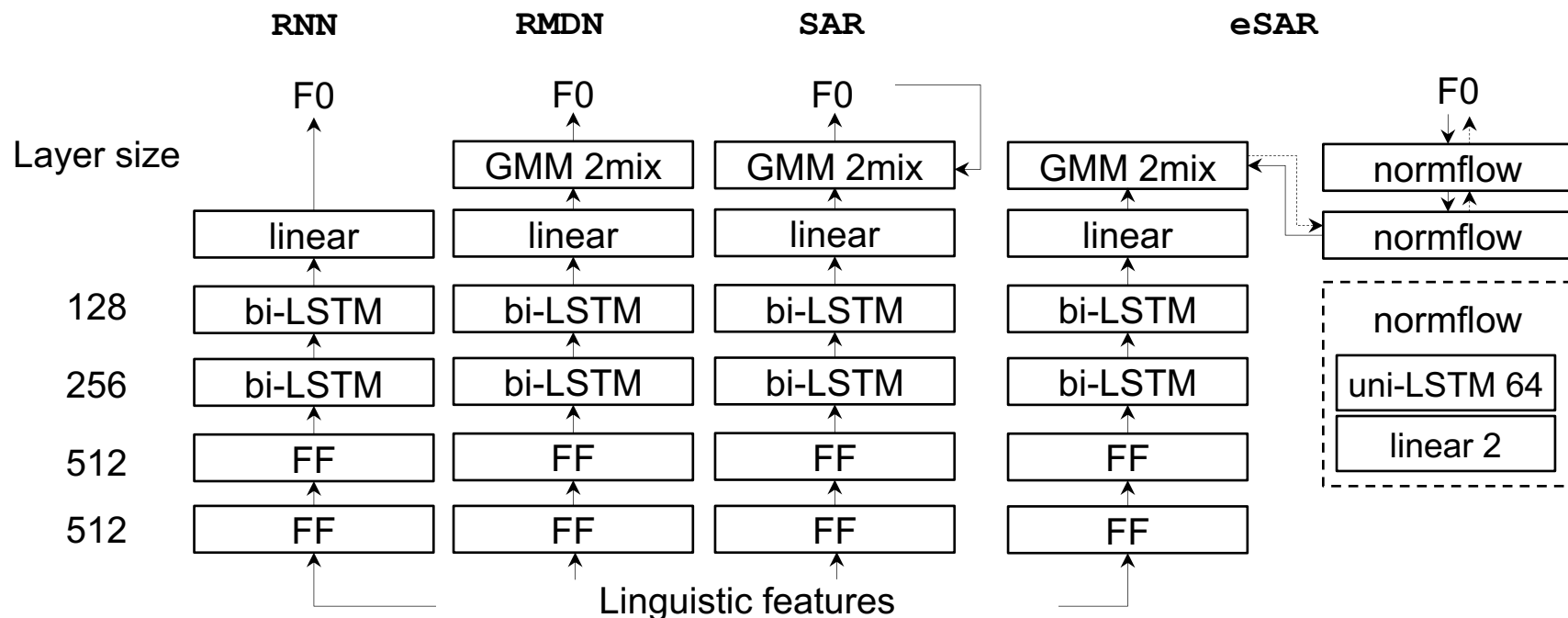
# ISSUE 2: TEMPORAL DEPENDENCY?

## Experiment

### □ Data and features

- Data: Japanese, 48 hours
- Feature: F0 (interpolated F0 value 1 dim + U/V 1 dim)

### □ Models





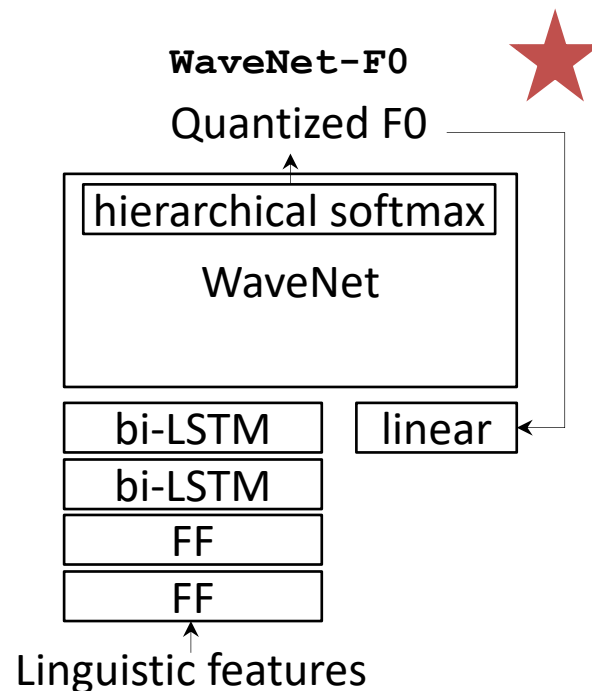
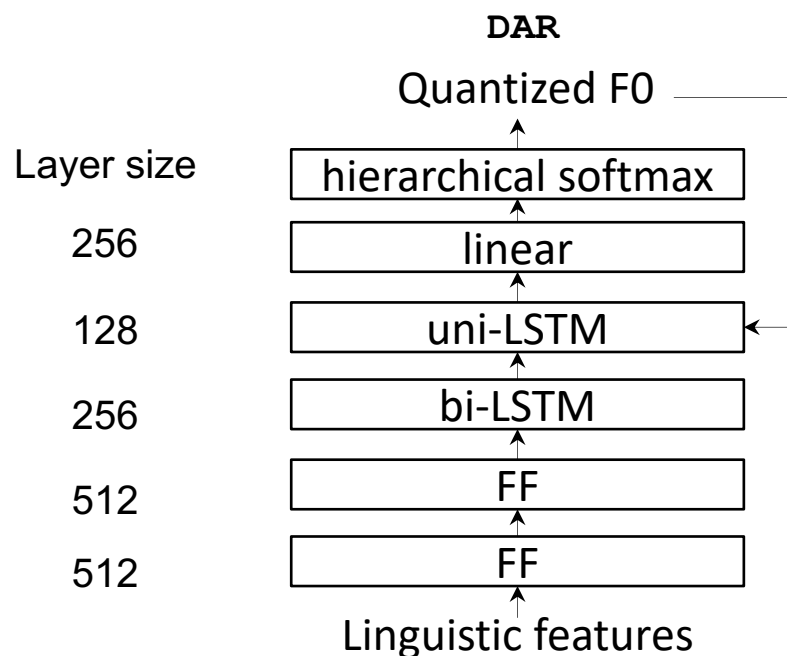
# ISSUE 2: TEMPORAL DEPENDENCY?

## Experiment

### □ Data and features

- Data: Japanese, 48 hours
- Feature: F0 (interpolated F0 value 1 dim + U/V 1 dim)
- Feature: quantized F0 (256 quantization bins)

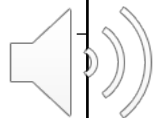
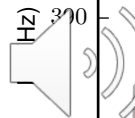
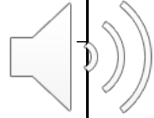
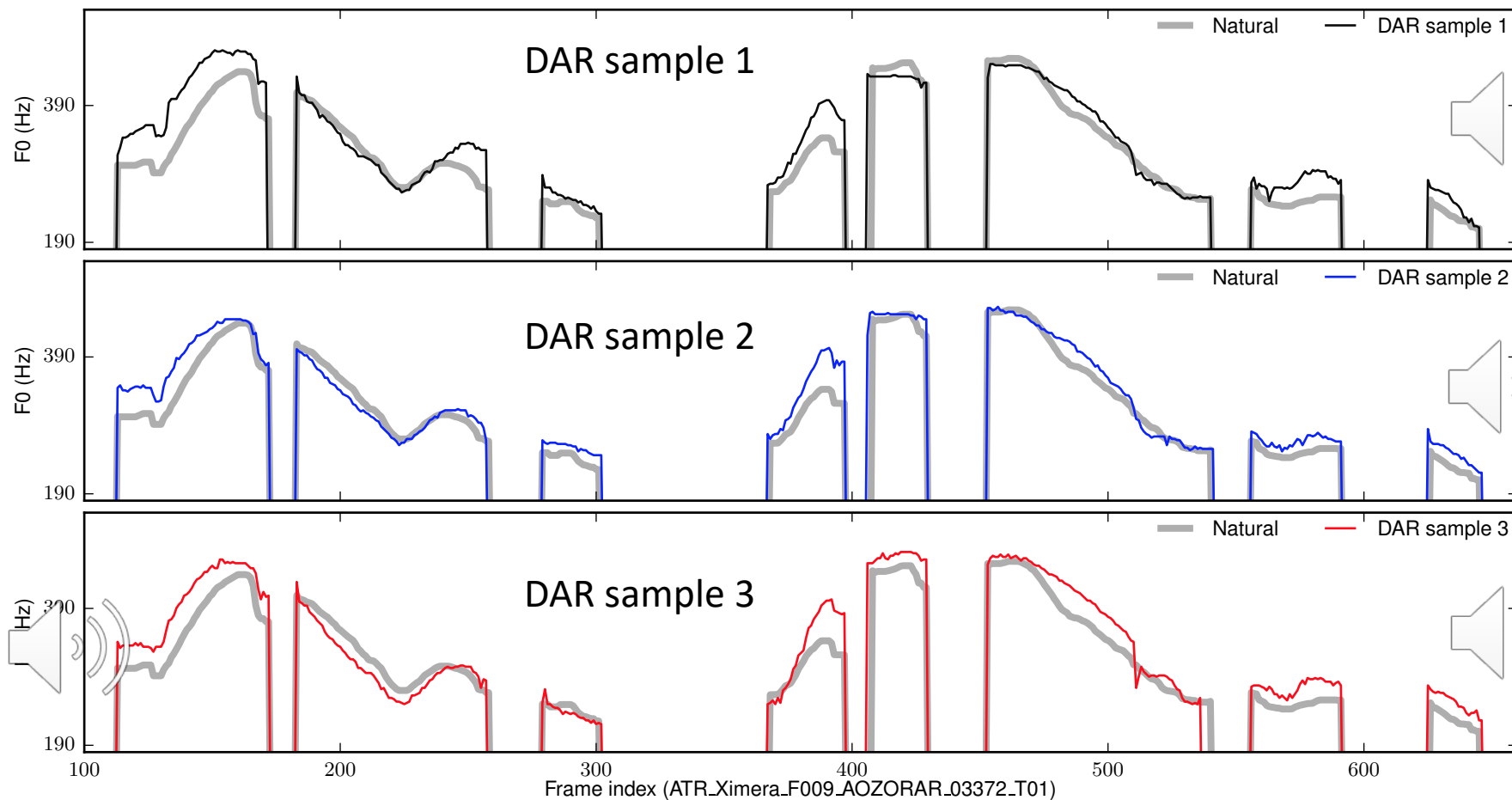
### □ Models



# ISSUE 2: TEMPORAL DEPENDENCY?

## Experiment

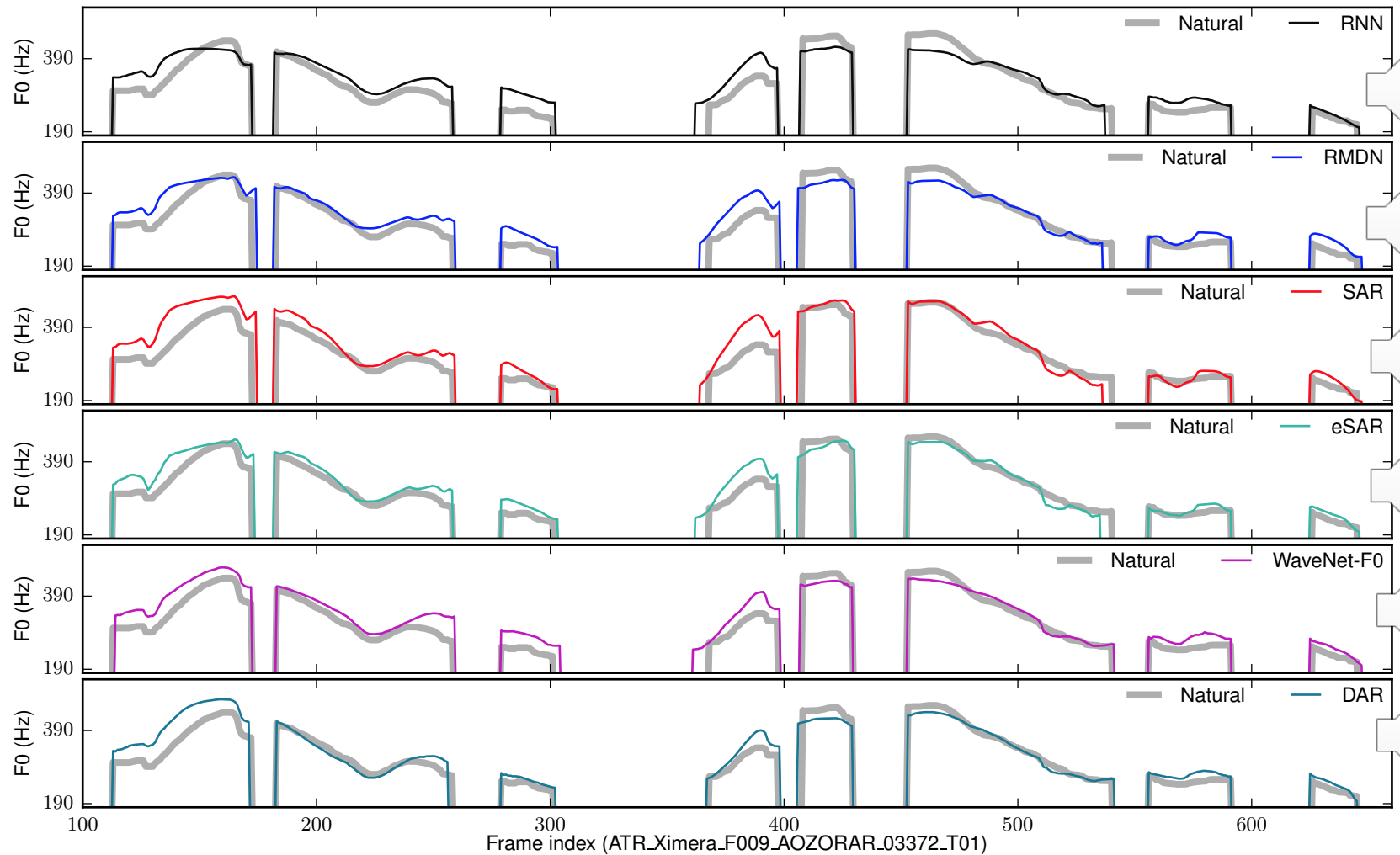
□ Random sampling



# ISSUE 2: TEMPORAL DEPENDENCY?

## Experiment

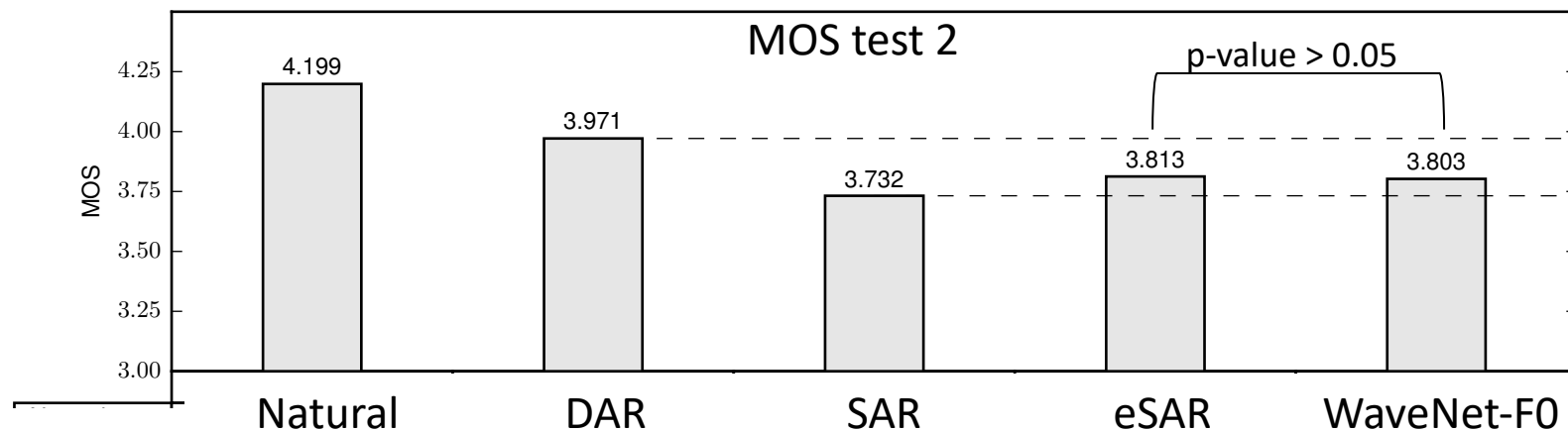
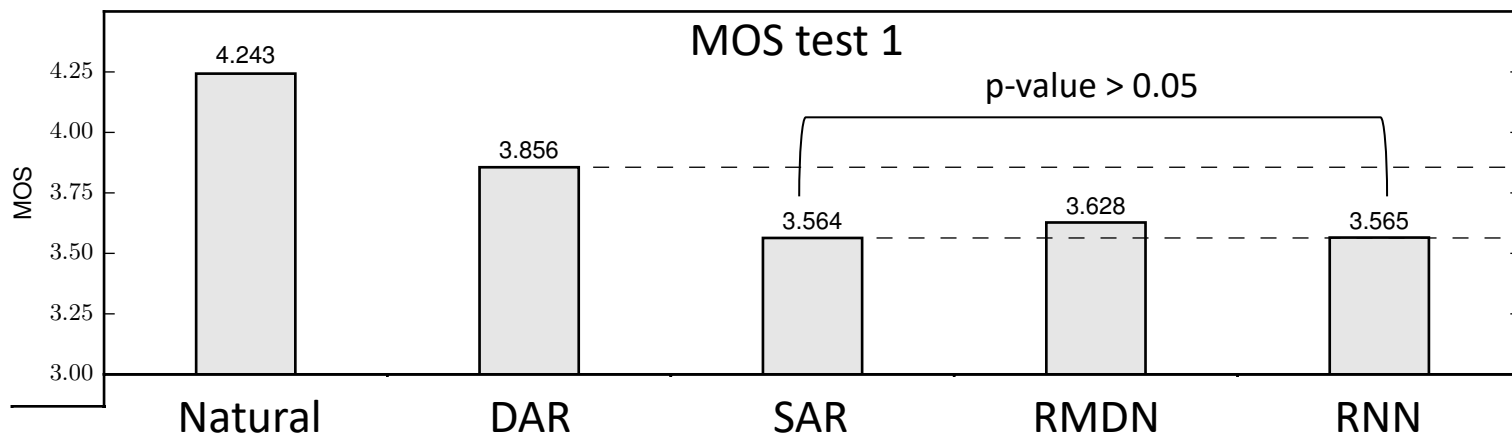
□ Mean-based generation



# ISSUE 2: TEMPORAL DEPENDENCY?

## Experiment

□ Mean-based generation



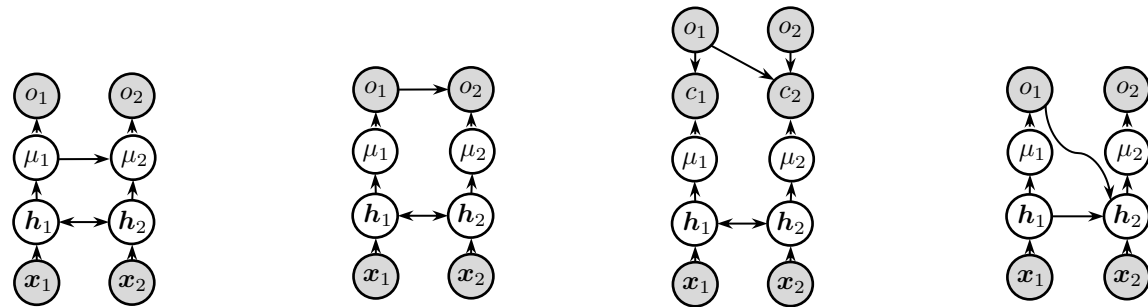
- 500 test utterances, >1000 sets of scores

# ISSUE 2: TEMPORAL DEPENDENCY?

## Summary

- Full answer to issue 2

Temporal dependency is ignored by RNN/RMDN!  
But better model can be defined



	RMDN	SAR	eSAR	DAR
AR linear?	-	Linear	Non-linear (constrained)	Non-linear
AR time span	-	$t - K : t - 1$	$1 : t - 1$	$1 : t - 1$
Tractable?	-	Yes	Somewhat	No
Sampling?	No	No	No	Yes

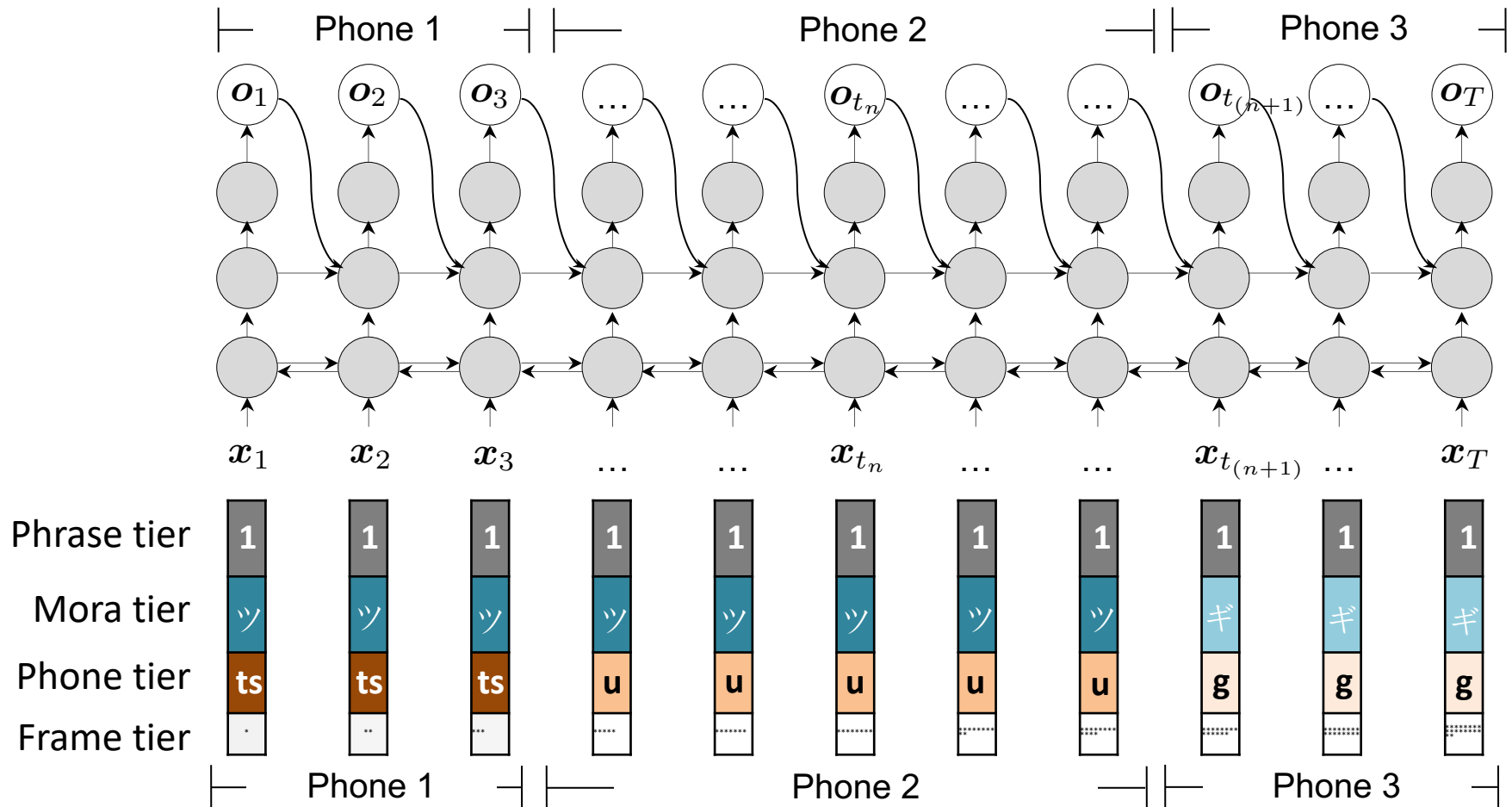
# CONTENTS

- Introduction
- Issue 1: joint modeling of F0 and spectral features
- Issue 2: temporal dependency modeling of F0 contours
- Issue 3: frame-by-frame processing
- Summary

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Motivation

- Inefficient processing



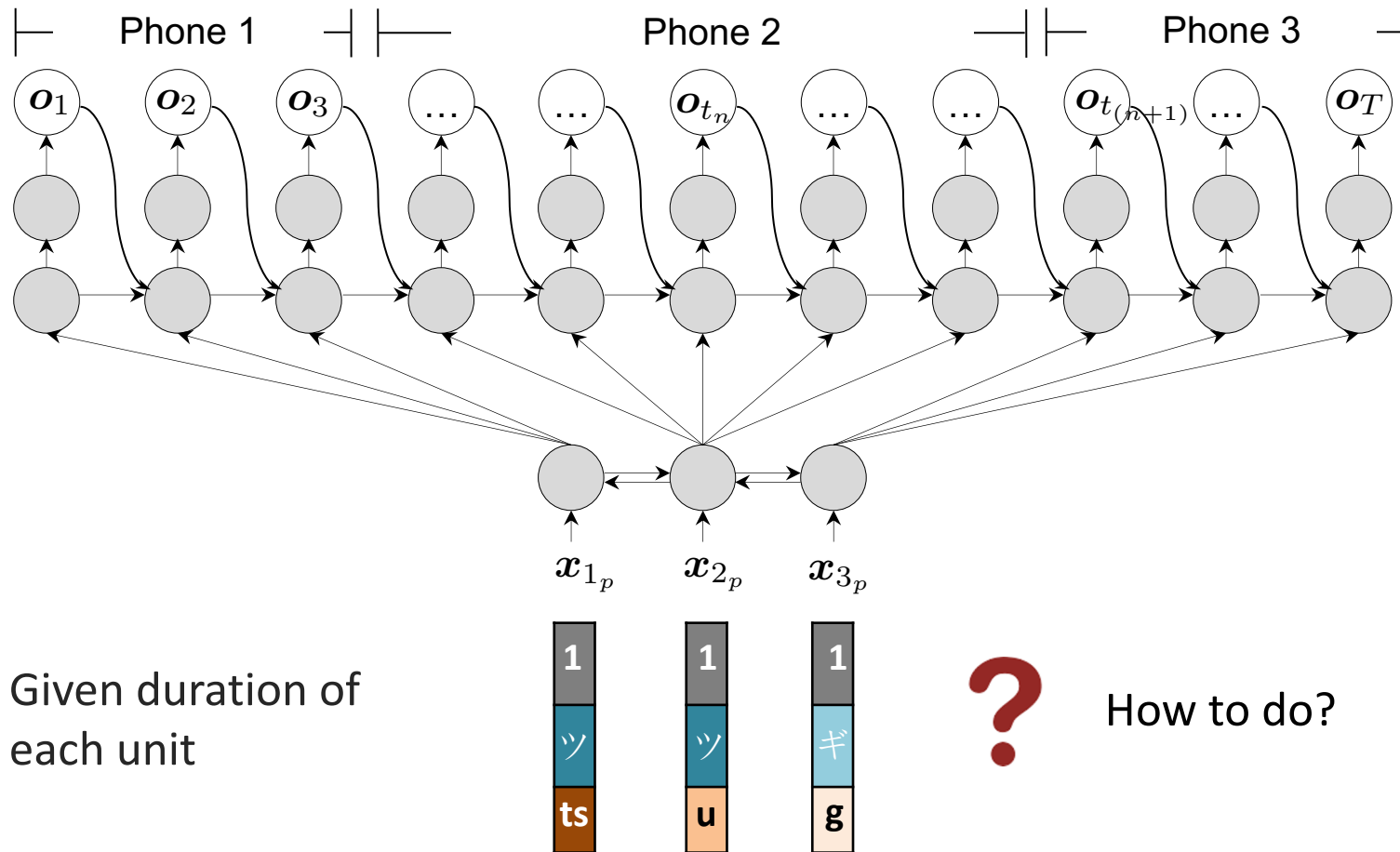




# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Motivation

- More efficient processing?



- Given duration of each unit

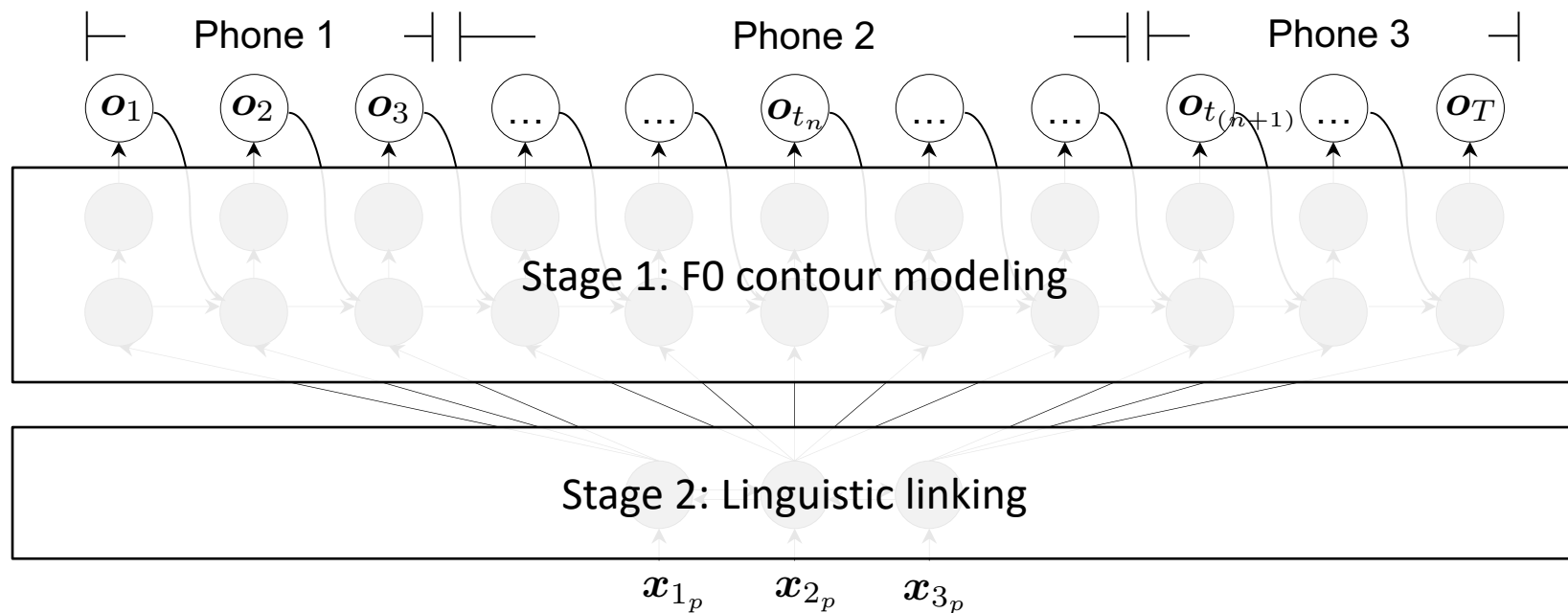


How to do?

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Method

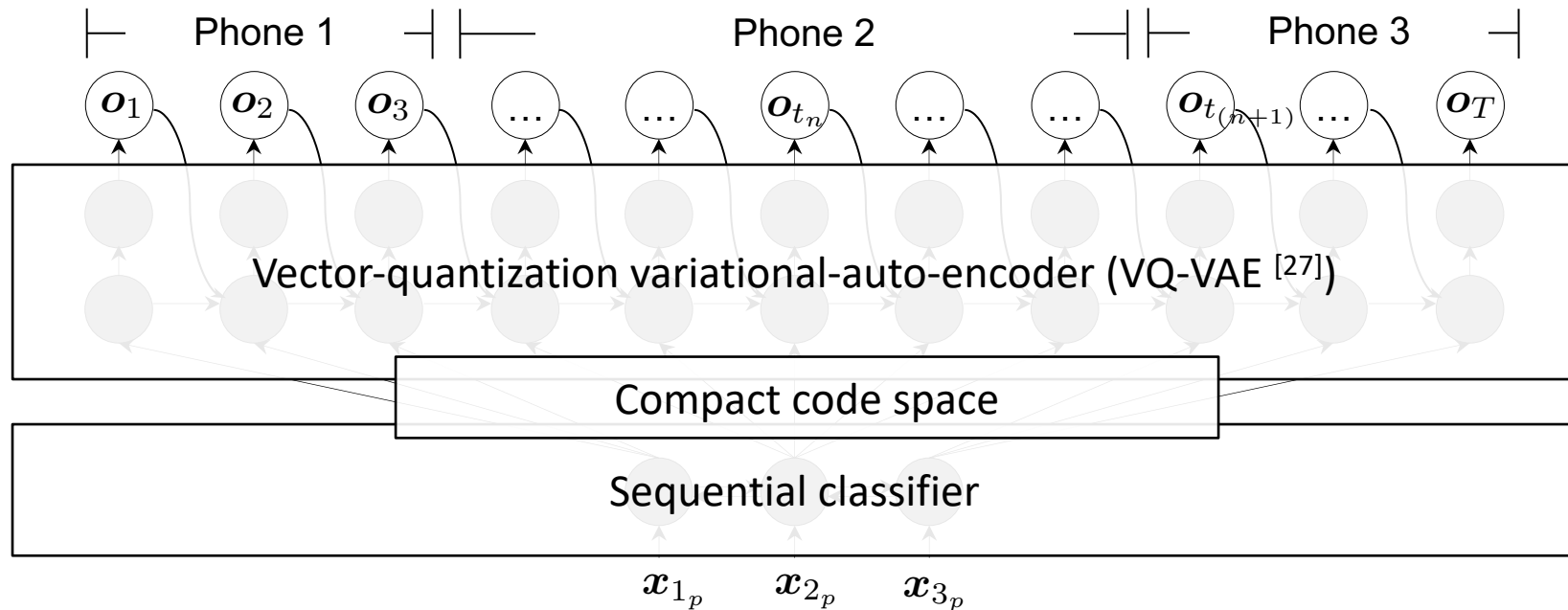
- Two-stage F0 modeling



# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Method

- Two-stage F0 modeling



- Revisit linguistic approaches [28, 29]

[27] A. van den Oord, O. Vinyals, and K. Kavukcuoglu. Neural discrete representation learning. In Proc. NIPS, page to appear, 2017.

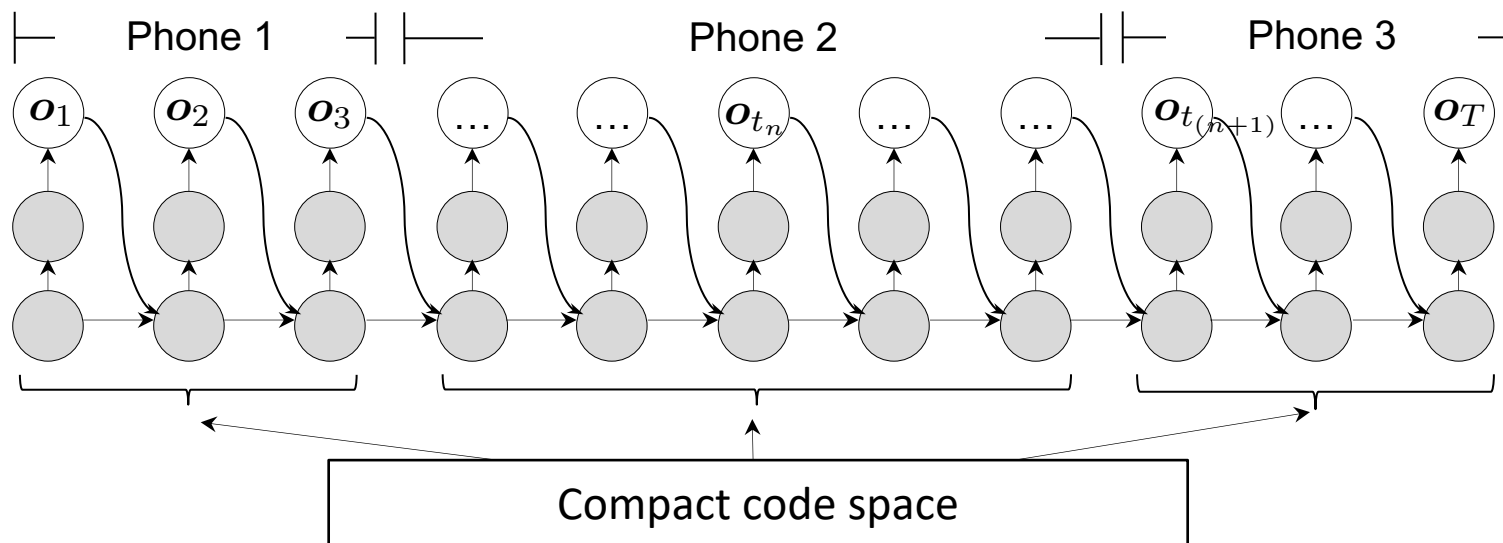
[28] K. E. Dusterhoff, A. W. Black, and P. A. Taylor. Using decision trees within the Tilt intonation model to predict F0 contours. Proc. Eurospeech, pages 1627–1630, 1999.

[29] K. Hirose, K. Sato, Y. Asano, and N. Minematsu. Synthesis of F0 contours using generation process model parameters predicted from unlabeled corpora: Application to emotional speech synthesis. Speech communication, 46(3):385–404, 2005.

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Stage 1: F0 contour modeling

□ VQ-VAE

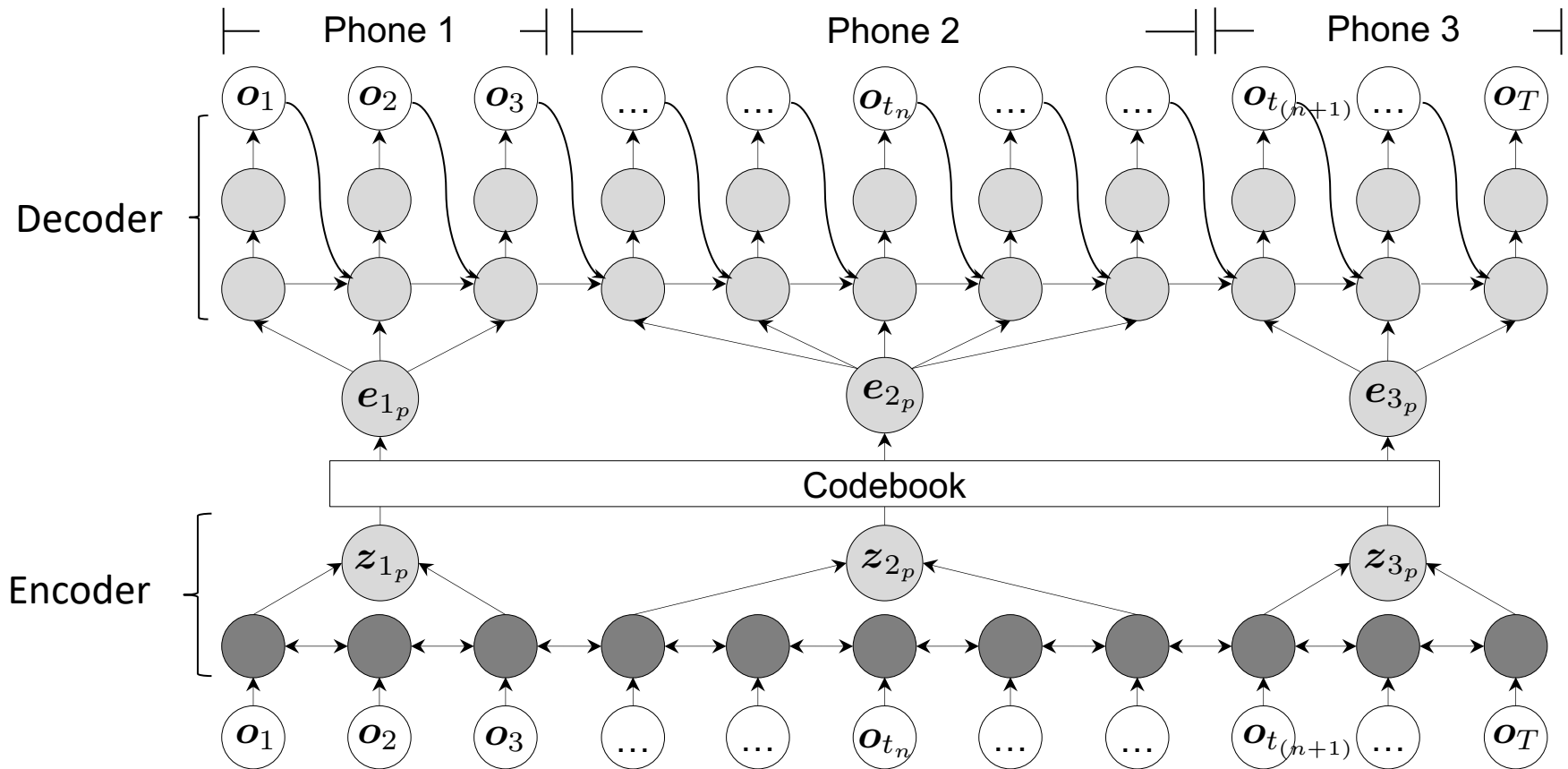


- Unsupervised learning
  - One code for varied-length unit
  - Multiple linguistic levels
- } Goals

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Stage 1: F0 contour modeling

□ VQ-VAE

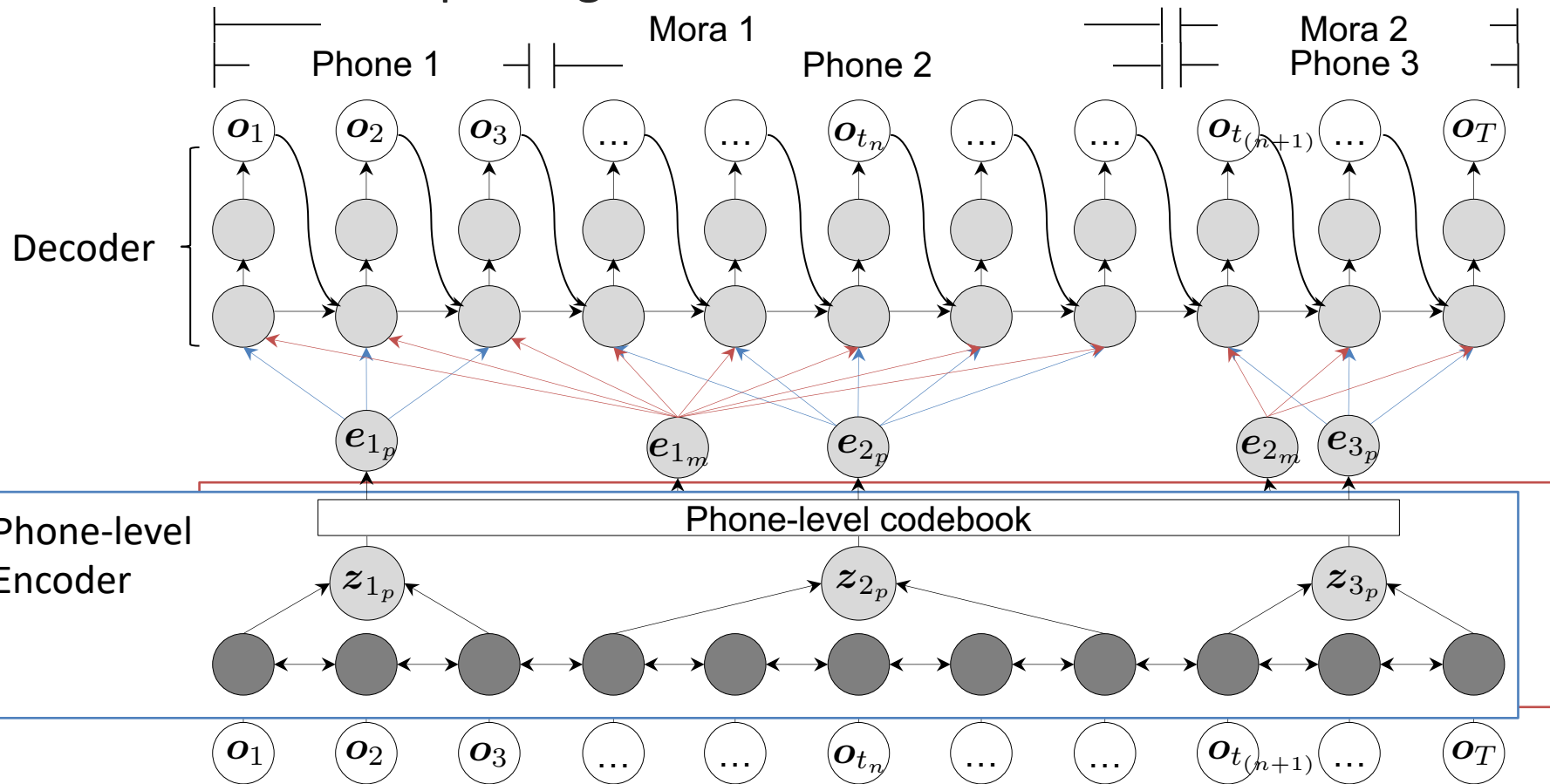


- One code per phone

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Stage 1: F0 contour modeling

□ VQ-VAE multiple linguistic tiers

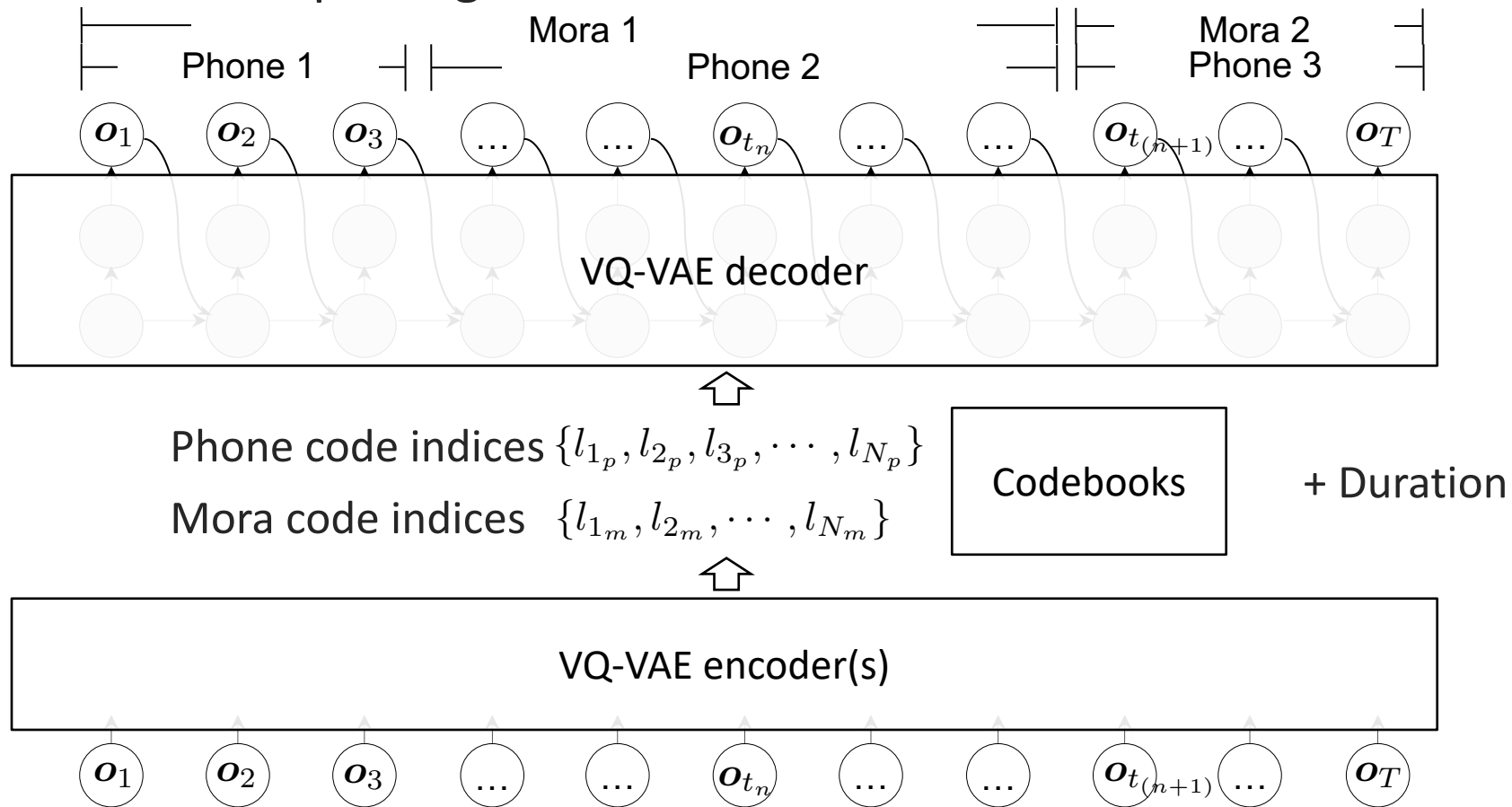


- One code per phone & one code per mora

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

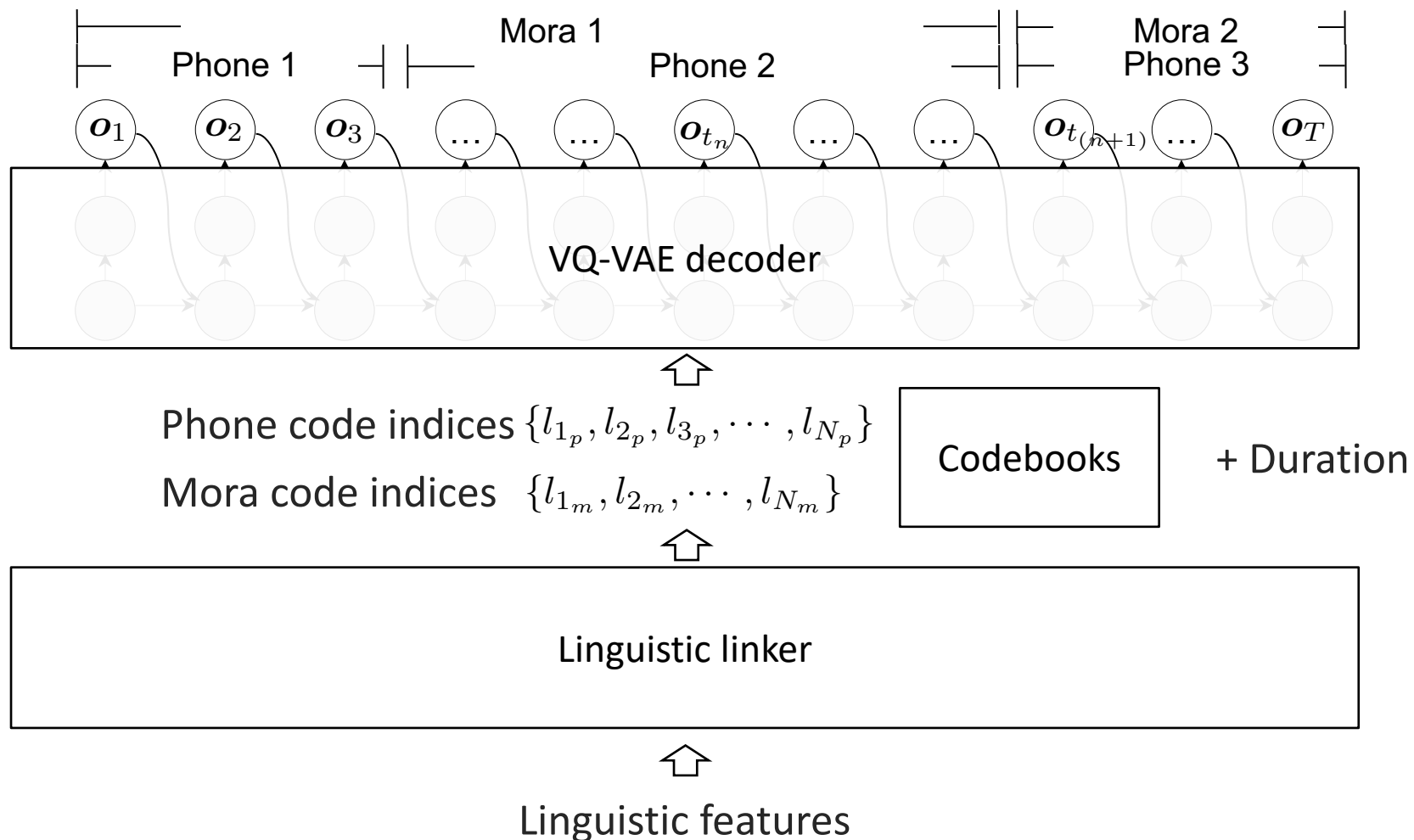
## Stage 1: F0 contour modeling

### □ VQ-VAE multiple linguistic tiers



# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Stage 2: Linguistic linking





# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Stage 2: Linguistic linking

Phone code indices  $\{l_{1_p}, l_{2_p}, l_{3_p}, \dots, l_{N_p}\}$

Mora code indices  $\{l_{1_m}, l_{2_m}, \dots, l_{N_m}\}$



RNN-based sequential classifier

(Sec. 7.2.2)

- Clockwork RNN <sup>[30]</sup>
- Highway <sup>[31]</sup>
- AR & feedback links
- Dropout <sup>[32]</sup>



Linguistic features

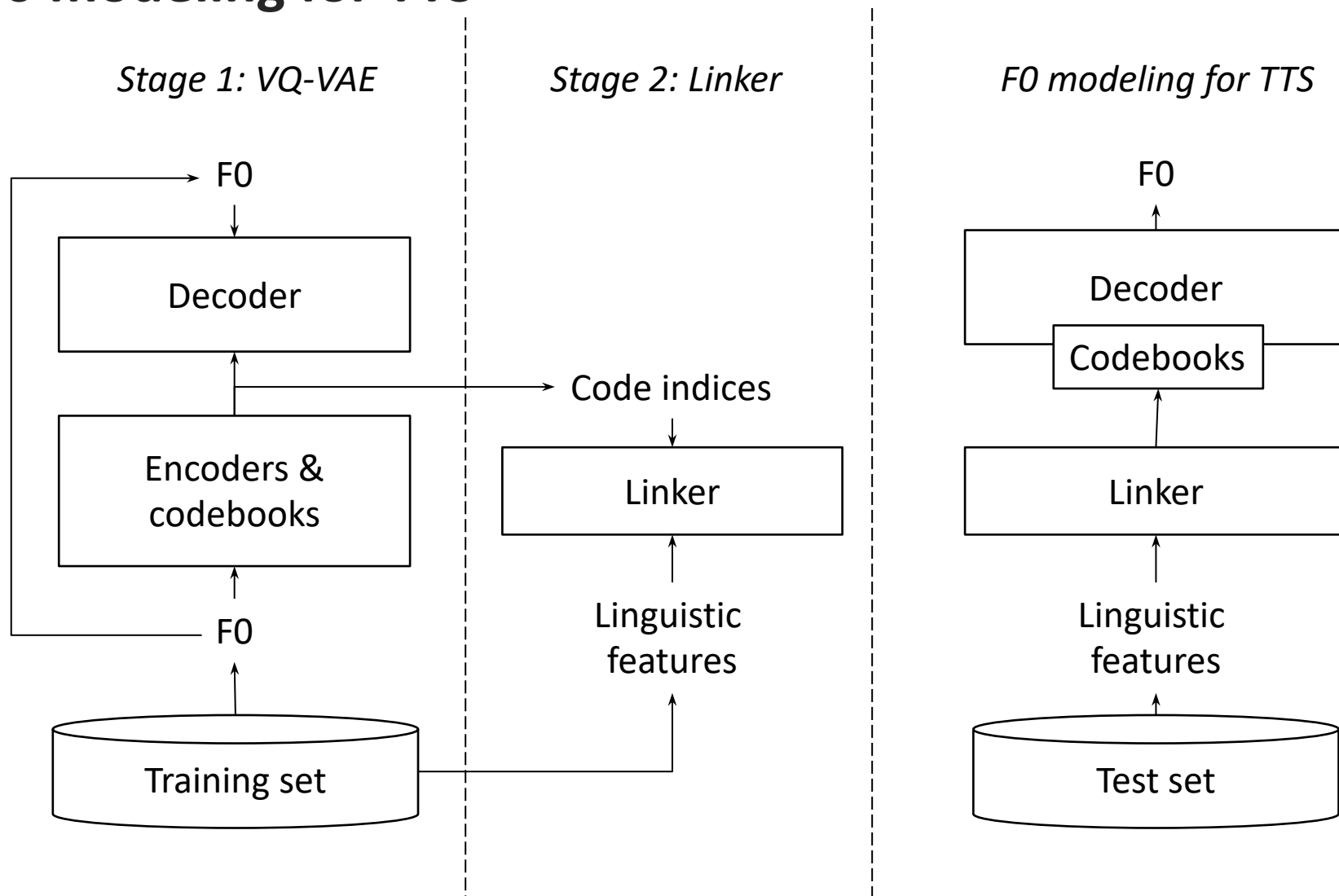
[30] J. Koutnik, K. Greff, F. Gomez, and J. Schmidhuber. A Clockwork RNN. In Proc. ICML, pages 1863–1871, 2014.

[31] K. Greff, R. K. Srivastava, and J. Schmidhuber. Highway and residual networks learn unrolled iterative estimation. In Proc. ICLR, 2017.

[32] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research, 15(1):1929–1958, 2014.

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## F0 modeling for TTS

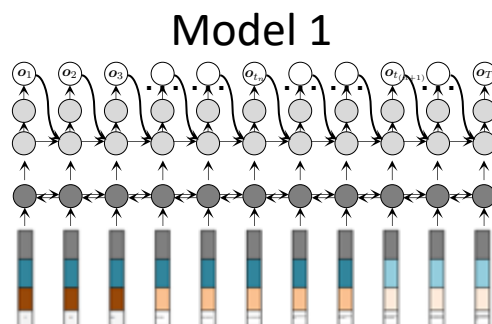


# ISSUE 3: FRAME-BY-FRAME PROCESSING?

Experiments on stage 1 (sec.7.3.2)

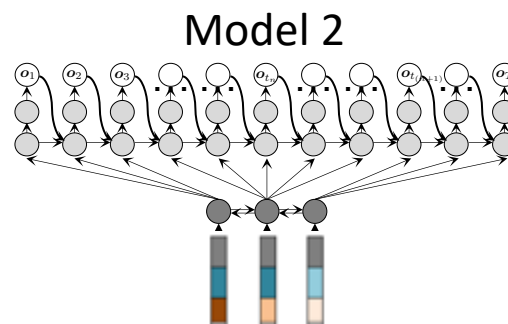
Experiments on whole model

□ Models



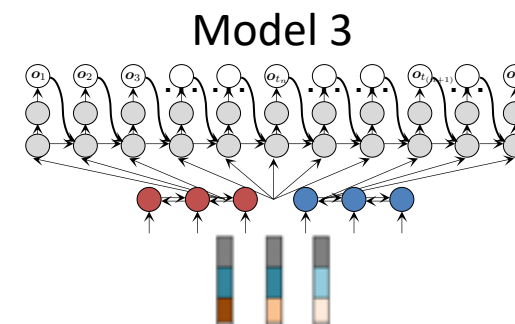
VQ-VAE decoder  
(phone-level)

Frame-by-frame  
Linker



VQ-VAE decoder  
(phone-level)

Phone-by-phone  
Linker



VQ-VAE decoder  
(phone-mora-level)

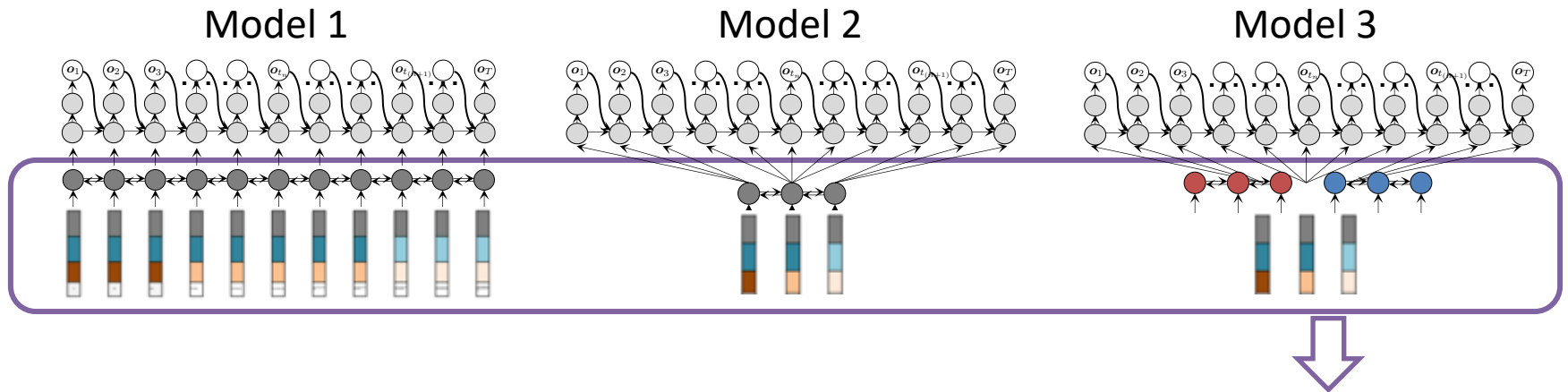
Phone-by-phone  
Linker + Mora lock

- Given natural duration

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Experiments

### Objective results

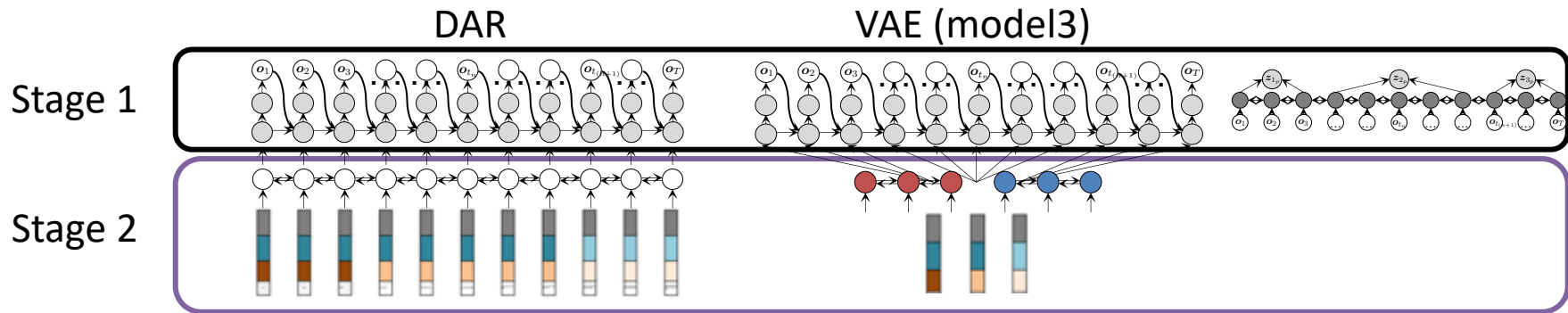


	RMSE	CORR	U/V	Time cost (s/epoch)
Model 1	34.3	0.839	7.96%	1300
Model 2	27.1	0.906	6.36%	54
Model 3	<b>25.5</b>	<b>0.916</b>	<b>4.87%</b>	<b>65</b>

# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## VAE vs DAR

### Objective results



	RMSE	CORR	U/V	Number of parameters (m)			Time cost (s/epoch)		
				Stage 1	Stage 2	Sum	Stage 1	Stage 2	Sum
DAR	28.3	0.903	3.46%	0.36	1.11	1.48	~700	~1300	2000
VAE model3	25.5	0.916	4.87%	0.44	0.67	1.11	1500	65	1565

- VQ-VAE encoder needs time and memory

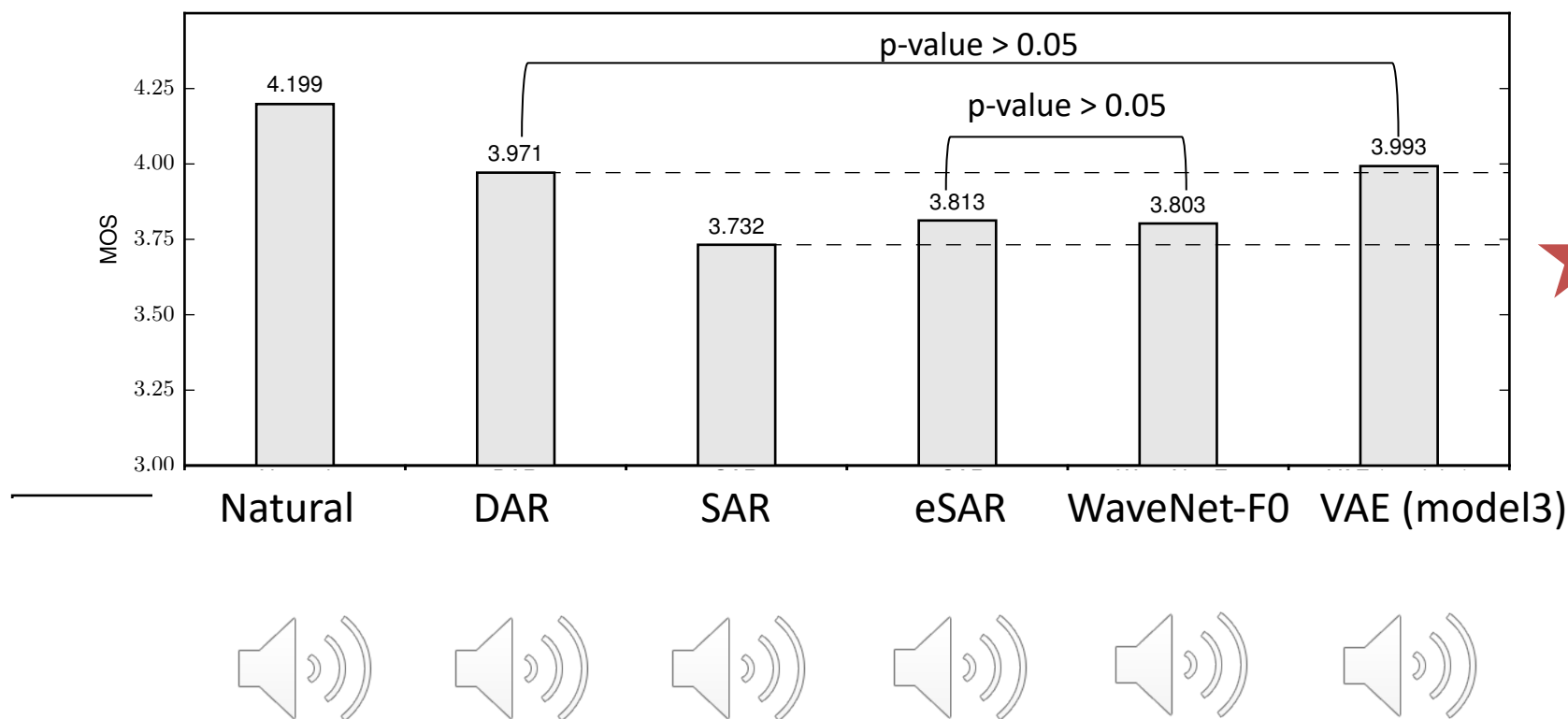


# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## VAE vs DAR

☐ Subjective test

MOS test



- 500 test utterances, >1000 sets of scores

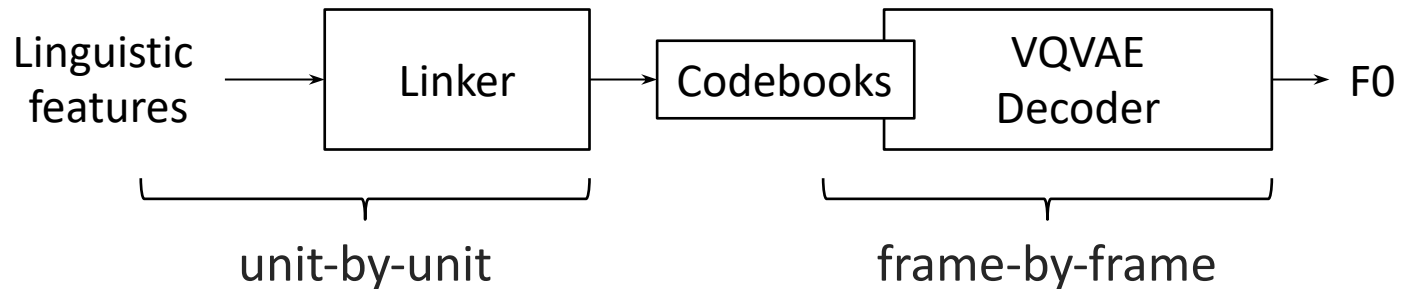
# ISSUE 3: FRAME-BY-FRAME PROCESSING?

## Summary

- Answer to issue 3

It could be more efficient

- How



- Results

- Multiple linguistic tiers
- More efficient than DAR: smaller + faster + F0 CORR > 0.91
- Interpretable latent code spaces (Sec. 7.3.2)
- Random sampling OK (slides)

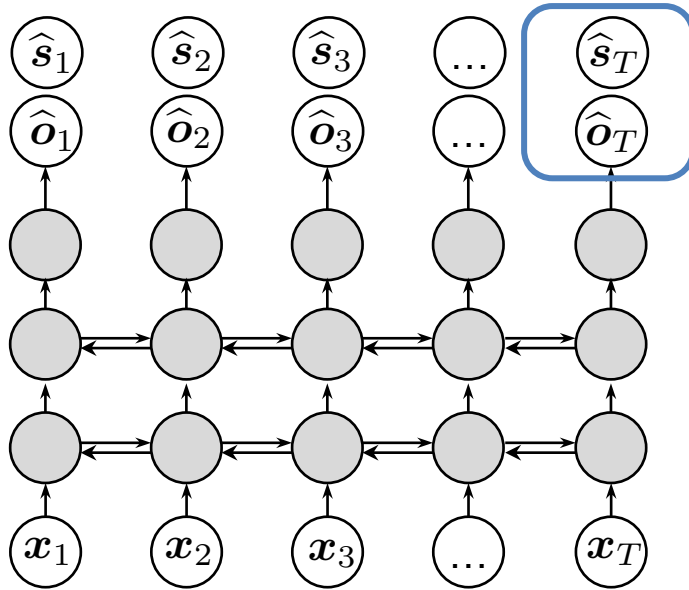
# CONTENTS

- ❑ Introduction
- ❑ Issue 1: joint modeling of F0 and spectral features
- ❑ Issue 2: temporal dependency modeling of F0 contours
- ❑ Issue 3: frame-by-frame processing
- ❑ Conclusion

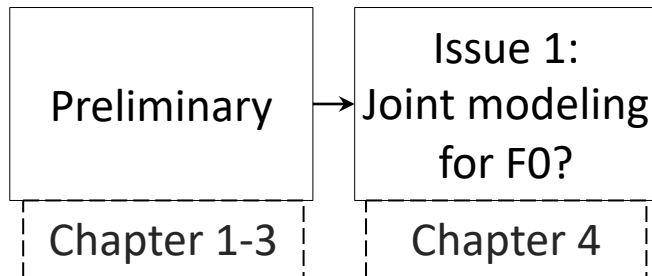


# CONCLUSION

## Summary



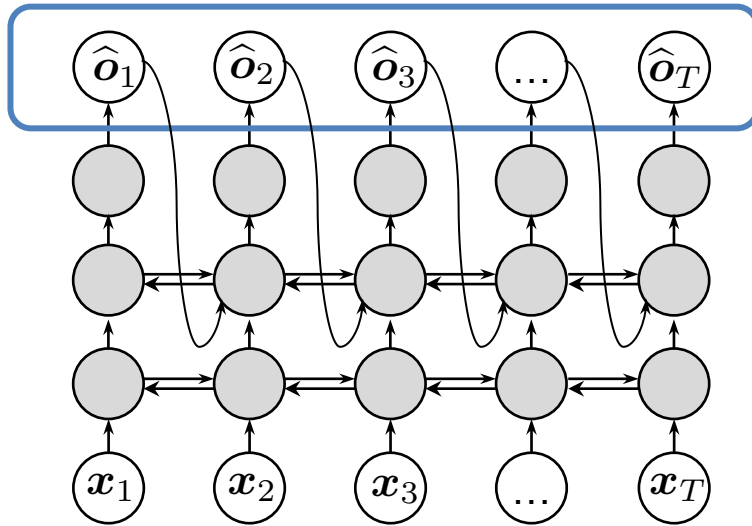
- ? Joint modeling of F0 and spectral?
- × Sub-optimal for F0 modeling
- Methods:
  - Highway networks
  - Histogram + sensitivity analysis
- Results:
  - Spectral features are prioritized
  - Different input/hidden features for F0 and spectral features



# CONCLUSION

## Summary

DAR



? Temporal dependency in RNN/RMDN?

× Ignored by RNN/RMDN

• Models:

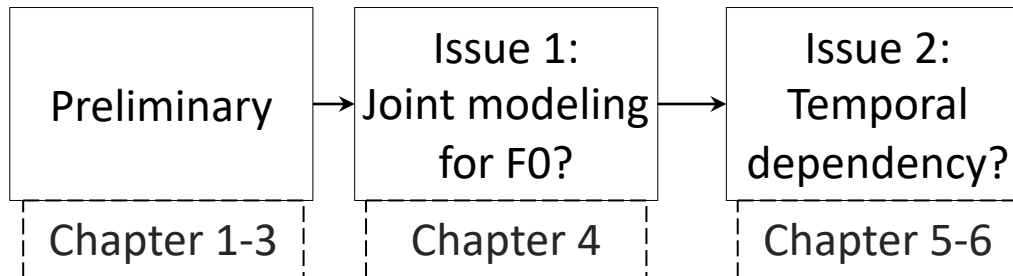
○ SAR : tractable dependency

○ DAR : non-linear + longer dependency

• Results:

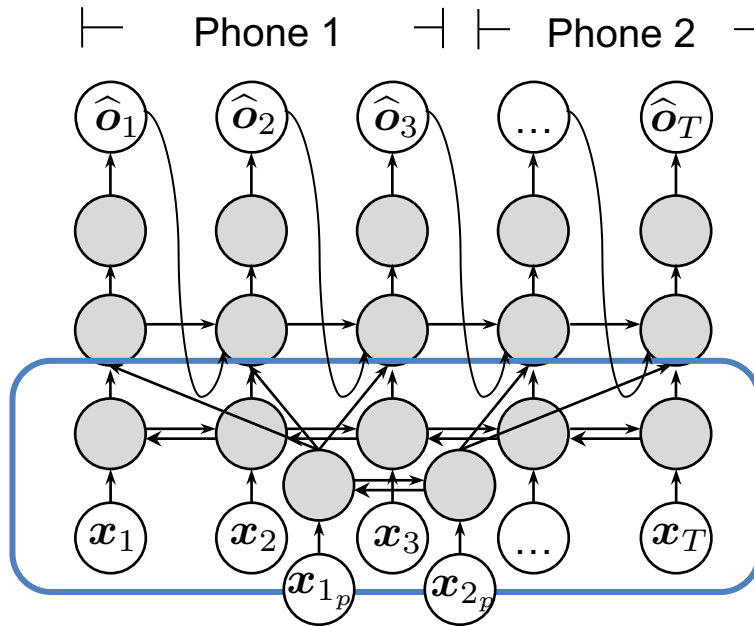
○ DAR: F0 CORR > 0.90, MOS score

○ DAR supports random sampling!



# CONCLUSION

## Summary



? Frame-by-frame processing

× Inefficient

• Two-stage F0 model:

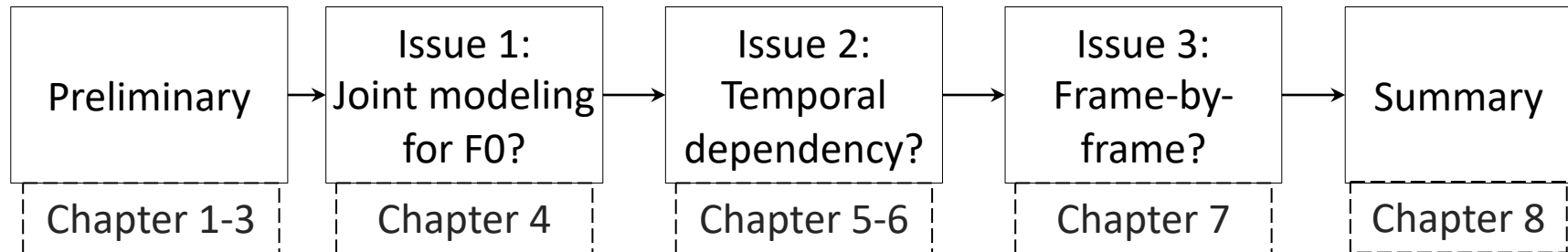
○ F0 contour coding: VQ-VAE + DAR

○ Linguistic linking: unit-by-unit classifier

• Results:

○ F0 CORR > 0.91

○ More efficient (smaller, faster)

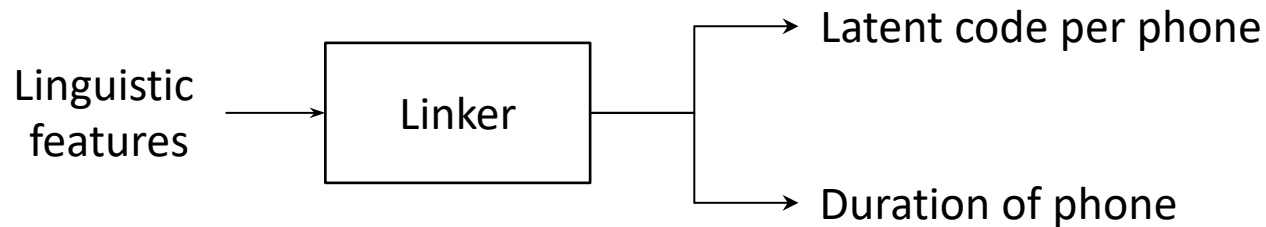


# CONCLUSION

## New topics?

### ❑ Towards complete prosody modeling

- Not only F0 but also duration
- Easy for joint modeling



### ❑ Waveform modeling

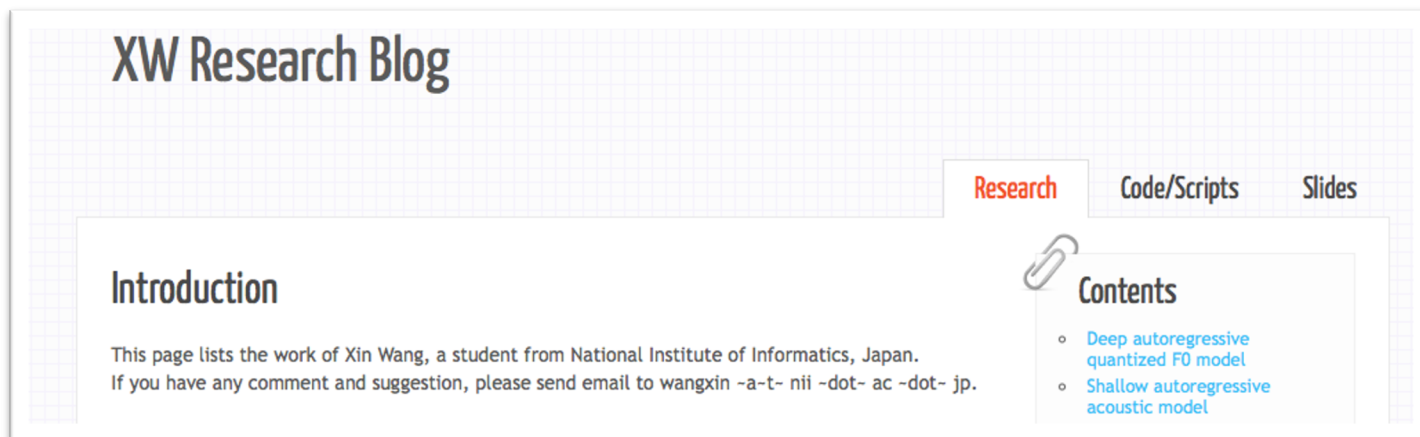
- F0 and waveform are 1 dimensional signals
- Signal processing methods available

SAR + log area ratio + segment-variant filters

# Thank you for your attention

## Q & A

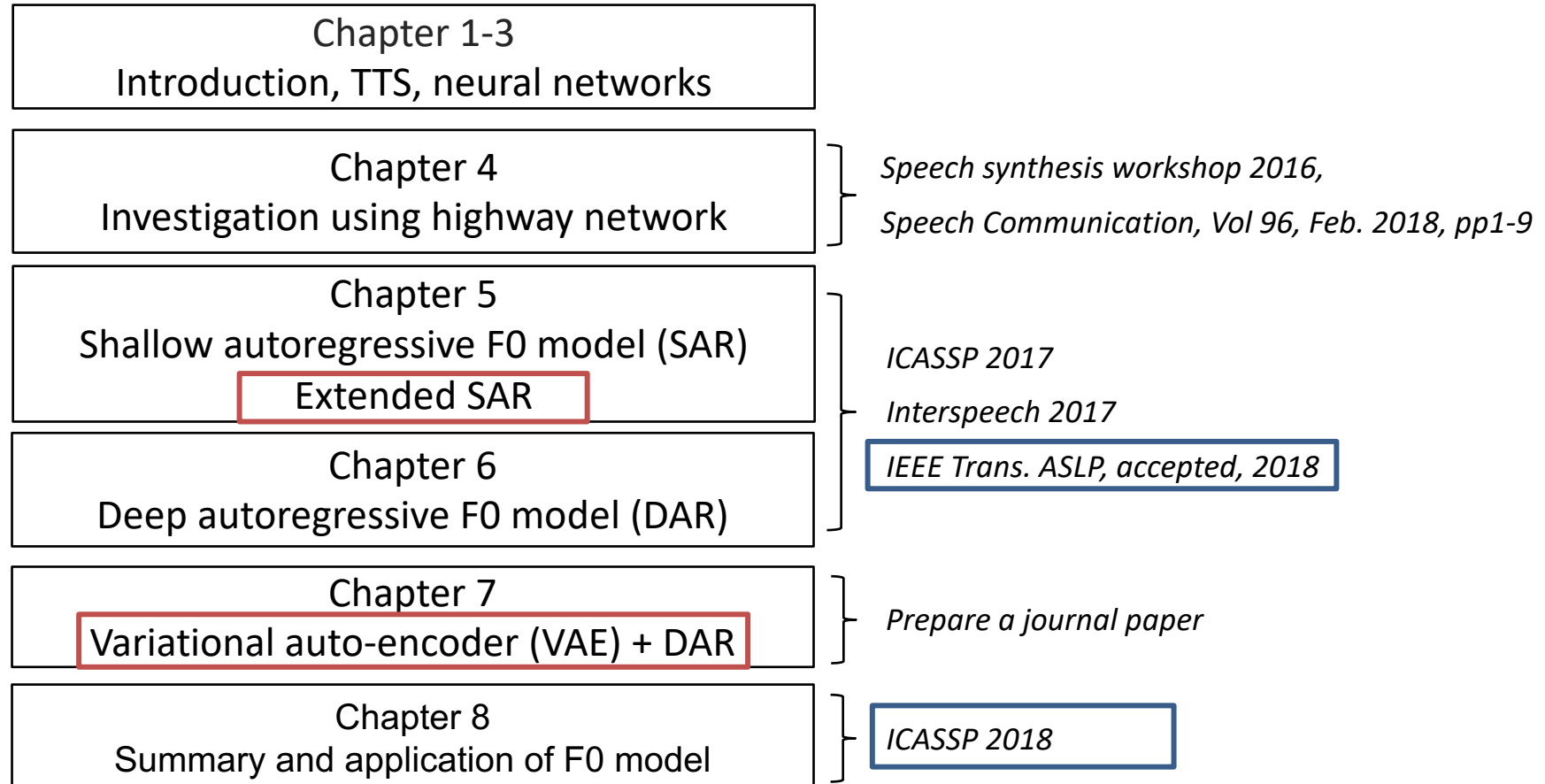
Codes, scripts, slides: [tonywangx.github.io](https://tonywangx.github.io)



The screenshot shows the 'XW Research Blog' website. At the top left is the title 'XW Research Blog'. Below it is a navigation bar with three tabs: 'Research' (highlighted in red), 'Code/Scripts', and 'Slides'. On the left side, there is a section titled 'Introduction' with the text: 'This page lists the work of Xin Wang, a student from National Institute of Informatics, Japan. If you have any comment and suggestion, please send email to wangxin -a-t- nii -dot- ac -dot- jp.' On the right side, there is a 'Contents' section with a paperclip icon, listing two items: 'Deep autoregressive quantized F0 model' and 'Shallow autoregressive acoustic model'.

# THESIS OUTLINE

## □ Neural F0 modeling



✓ Meet basic requirement on publication: 3 journals + 6 conferences

# WORK DURING PH.D.

## □ Work not included in Ph.D. thesis

- Word embedding as input features (IEICE 2018, Interspseech2016)
- Impact of training data size (SSW 2016)

## □ Collaborated work

- DAR using manually-annotated / corrupted data (Interspeech 2018)  
Provide DAR, SAR, and WaveNet-vocoder
- Voice cloning using found data (Odyssey 2018)  
Provide WaveNet-vocoder, and acoustic models
- Cyborg speech: multilingual speech synthesis (ICASSP 2018)  
Provide acoustic models
- Speech synthesis from MFCC (ICASSP 2018)  
Provide DAR on MFCC
- Controllable speech synthesis (Interspeech 2017)  
Provide acoustic models