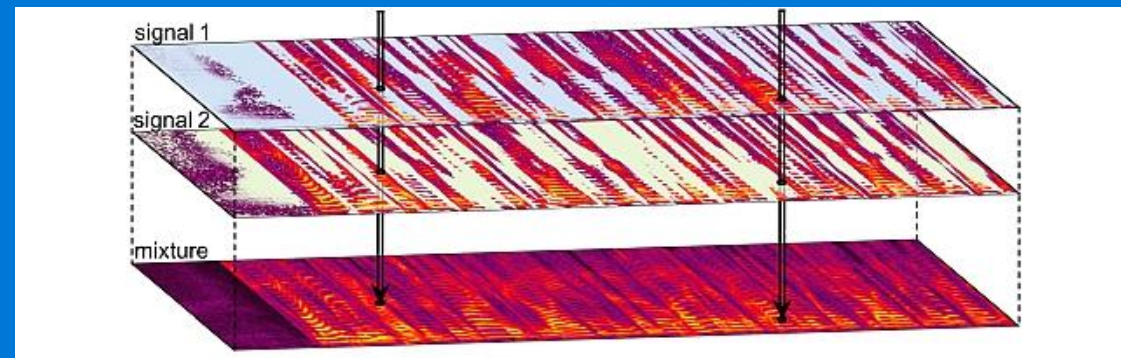


# Sequence Modeling

in Unsupervised Single-channel Overlapped Speech Recognition

Zhehuai Chen and Jasha Droppo

chenzhehuai@sjtu.edu.cn  
jdroppo@microsoft.com



SJTU SPEECH LAB  
上海交通大学智能语音实验室

# Outline

- Introduction
  - Cocktail party problem
  - PIT-TS framework and discriminative training
- Proposed methods
  - Temporal Correlation Modeling
  - Integrating Language Model
- Experiments
- Conclusion

# Introduction

- Cocktail-party problem

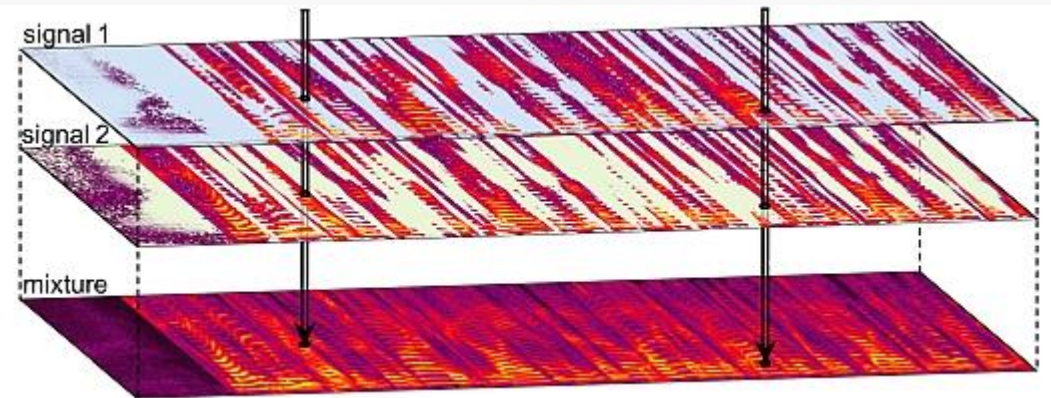


**N=2**

$$P(\mathbf{L}_{u1}, \dots, \mathbf{L}_{uN} | \mathbf{O}_u^{(m)})$$



$$\mathbf{O}_u^{(m)} = \sum_{n=1}^N \mathbf{O}_{un}^{(r)}$$



Assignment error:

e.g. ch-a: how oh you  
ch-b: are no

Cross talk error:

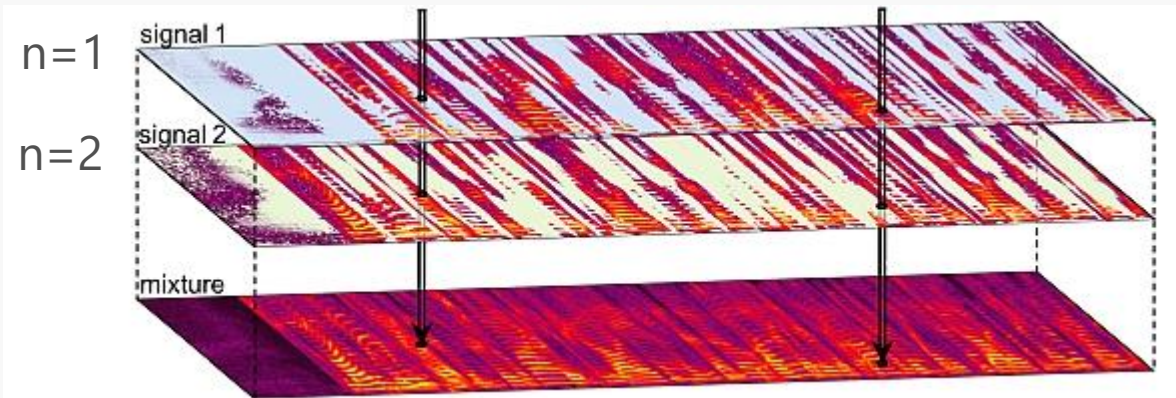
e.g. ch-a: how are you  
ch-b: oh are no

Label assignment problem

$$P(\mathbf{L}_{u1}, \dots, \mathbf{L}_{uN} | \mathbf{O}_u^{(m)}) \approx \prod_{n=1}^N P(\mathbf{L}_{un}^{(r)} | \mathbf{O}_u^{(m)}) \quad (2)$$

Label  
Independence

$$\mathbf{O}_u^{(m)} = \sum_{n=1}^N \mathbf{O}_{un}^{(r)}$$



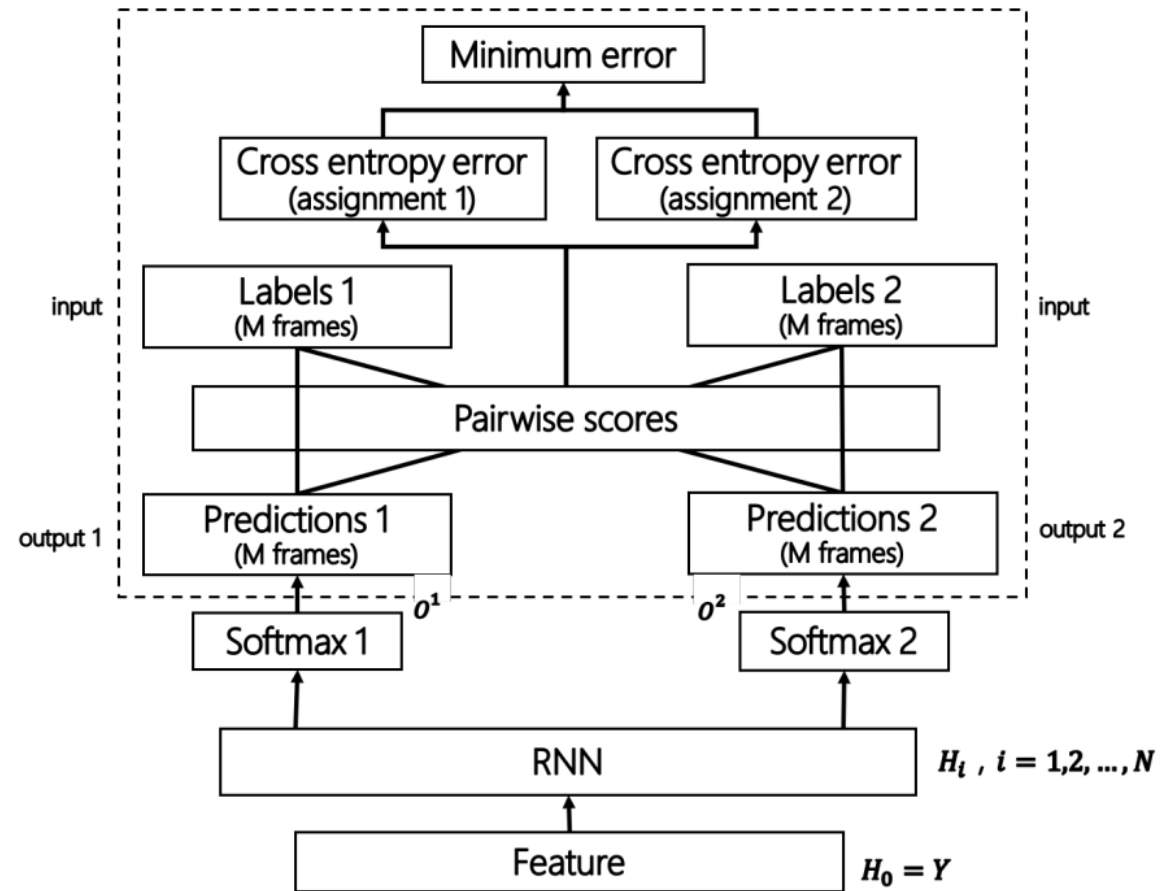
# Permutation Invariant Training for ASR

$$P(\mathbf{L}_{u1}, \dots, \mathbf{L}_{uN} | \mathbf{O}_u^{(m)}) \approx \prod_{n=1}^N P(\mathbf{L}_{un}^{(r)} | \mathbf{O}_u^{(m)}) \quad (2)$$

$$\mathcal{J}_{\text{CE-PIT}} = \sum_u \left[ \min_{s' \in \mathbf{S}} \right] \sum_t \frac{1}{N} \sum_{n \in [1, N]} \text{CE}(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (4)$$

- Disadvantage

- Model Complexity (3 hardest problems)
- Frame CE  $\rightarrow$  Utt. Problem
- No Linguistics



# PIT + Transfer Learning (TS)

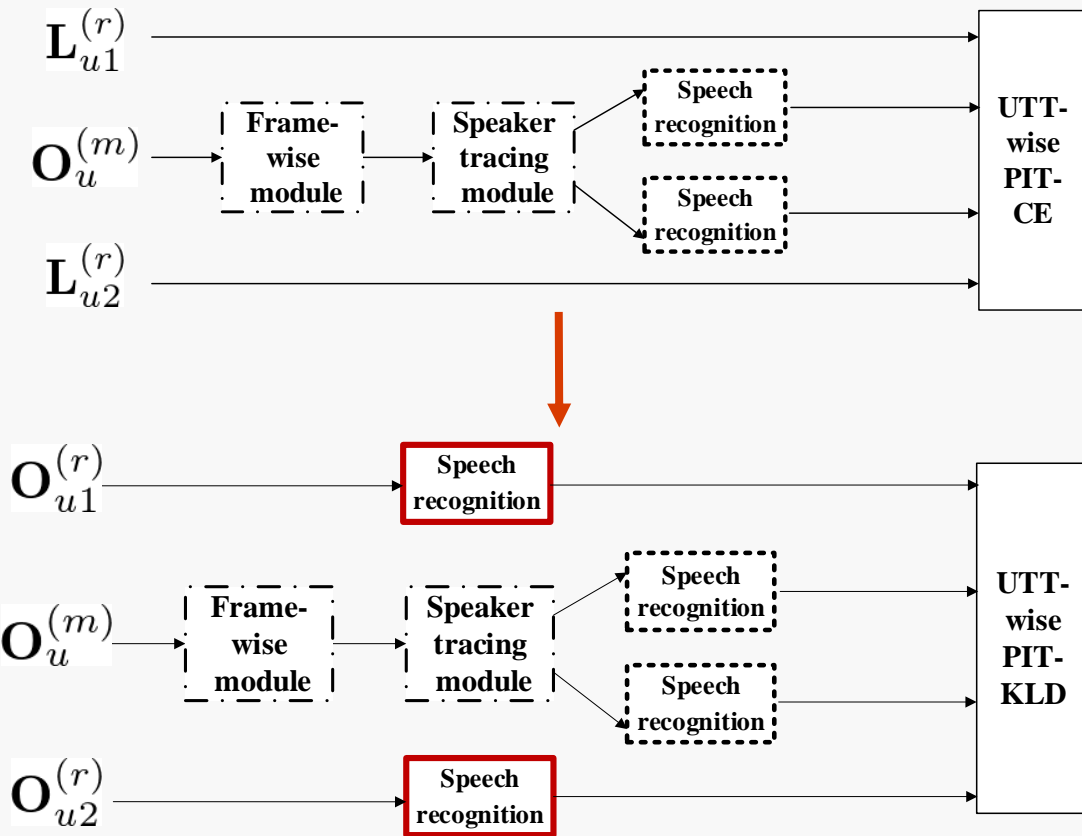
$$\mathcal{J}_{\text{CE-PIT}} = \sum_u \min_{s' \in \mathcal{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} \boxed{CE(l_{utn}^{s'}, l_{utn}^{(r)})} \quad (4)$$

$$\mathcal{J}_{\text{KLD-PIT}} = \sum_u \min_{s' \in \mathcal{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} \boxed{KLD(P(l_{utn}^{(c)} | \mathbf{O}_{un}^{(r)}), P(l_{utn}^{(s')} | \mathbf{O}_u^{(m)}))} \quad (8)$$

Clean infer.

PIT model infer.

- Better model convergence
  - Domain adaptation v.s. from scratch





# Linguistics - Multi-outputs Seq. Disc. Training

- Motivation:
  - Both ASR & speaker tracing  $\rightarrow$  sequential
  - Implicit integrating language model
- Formulation:

$$\mathcal{J}_{\text{CE-PIT}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} \boxed{CE(l_{utn}^{(s')}, l_{utn}^{(r)})} \quad (4)$$

- Key challenges:
  - Design the multi-output search space
  - Integrate with label assignment

$$\mathcal{J}_{\text{SEQ-PIT}} = \sum_u \min_{s' \in \mathbf{S}} \frac{1}{N} \sum_{n \in [1, N]} \boxed{\mathcal{J}_{\text{SEQ}}(\mathbf{L}_{un}^{(s')}, \mathbf{L}_{un}^{(r)})} \quad (12)$$

# Proposed methods

- Follow PIT-TS diagram
- Motivation
  - improve sequence modeling & language model
- Method
  - Implicit correlation modeling → explicit
  - Integrate linguistic information

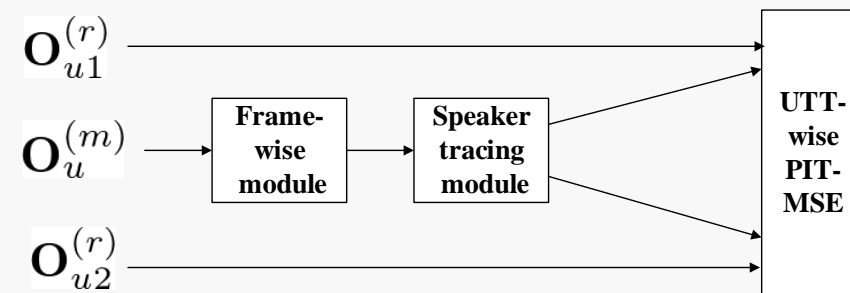


# Acoustics – Temporal Correlation Modeling

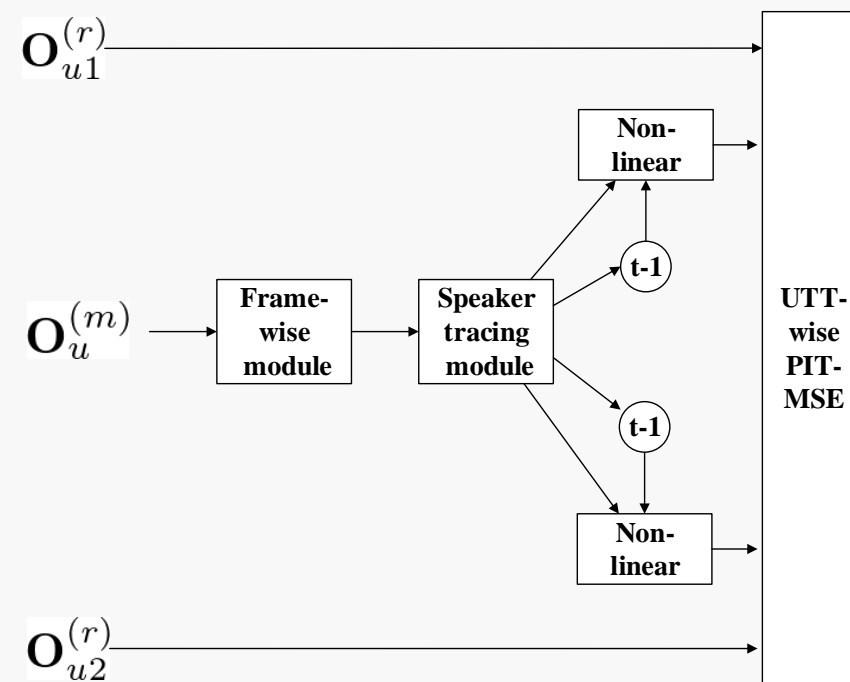
- Motivation
  - Sequential correlation v.s. stream de-correlation
    - the frequency bins between adjacent frames of the same speaker are correlated
- Last inference can improve current inference

Assignment error:

e.g. ch-a: how oh you  
ch-b: are no



(a) Speaker Tracing



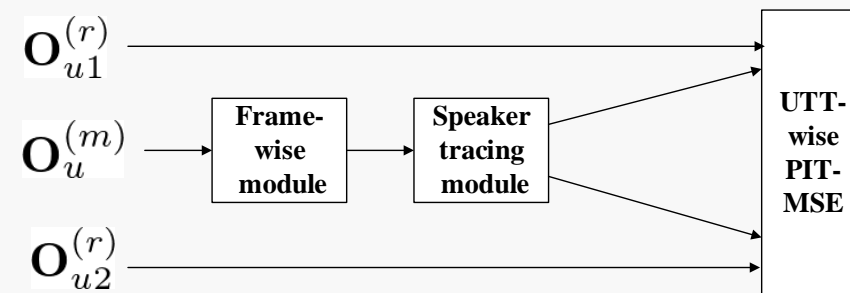
(b) Temporal Correlated Speaker Tracing

# Acoustics – Temporal Correlation Modeling

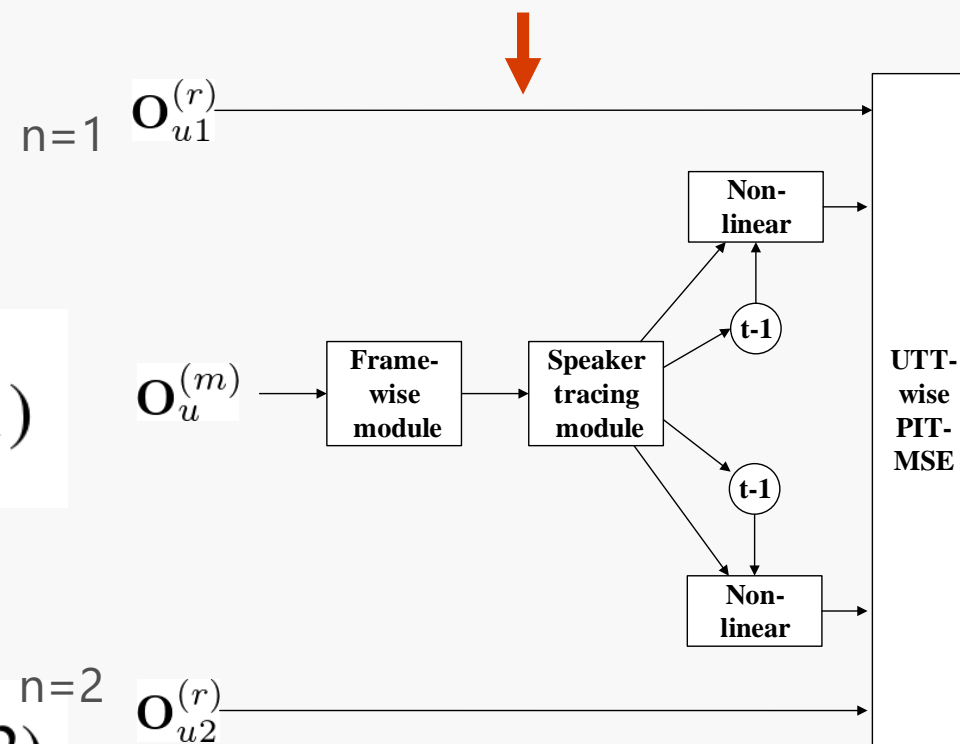
- Motivation
  - **Sequential correlation** v.s. stream de-correlation
  - Last inference can improve current inference
- Sequential labels correlation

$$o_{utn} = \mathcal{F}_{utn}(\mathbf{O}_u^{(m)}) \quad (1)$$

$$\underline{o}_{utn} = \mathcal{F}'_{utn}(\mathbf{O}_u^{(m)}, \underline{o}_{u(t-1)n}) \quad (2)$$



(a) Speaker Tracing



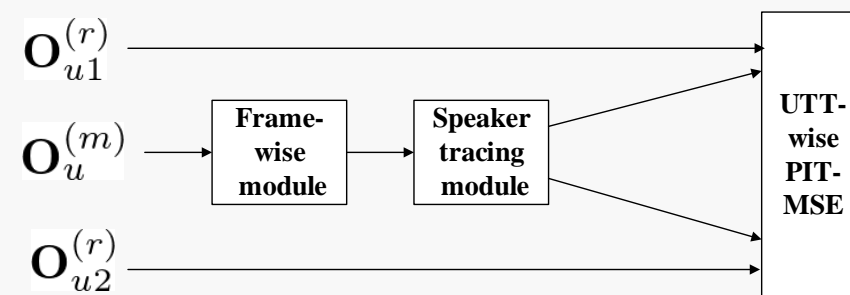
(b) Temporal Correlated Speaker Tracing

# Acoustics – Temporal Correlation Modeling

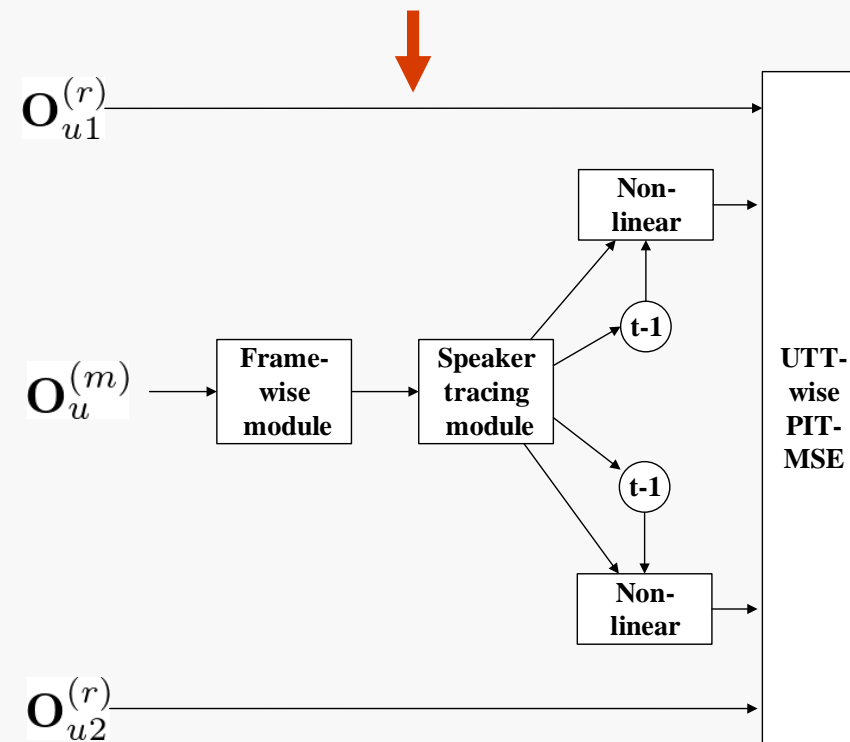
- Motivation
  - Sequential correlation v.s. **stream de-correlation**
  - last inference can improve current inference
- Sequential labels correlation
- alleviates the assignment & cross talk errors

Assignment error:

e.g. ch-a: how oh you  
ch-b: are no



(a) Speaker Tracing



(b) Temporal Correlated Speaker Tracing

# Linguistics – Language Model Integration

- Motivation:
  - Improve **assignment decision** by **combining LM** in training stage
  - Still train a **pure** acoustic model and integrate it with more powerful word level language model in evaluation stage
- Original PIT-CE

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (3)$$

# Linguistics – Language Model Integration

- Motivation:
  - Improve **assignment decision** by **combining LM** in training stage
  - Still train a **pure** acoustic model and integrate it with more powerful word level language model in evaluation stage
- Original PIT-CE

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (3)$$

- PIT-MAP:

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)}) / P(l) \cdot P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')})}{P(\mathbf{O}_u^{(m)})}$$

Discriminative training

$$\approx \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda$$

Proposed method

# Linguistics – Language Model Integration

- Motivation:
  - Improve assignment decision by combining LM in training stage
  - Still train a pure acoustic model and integrate it with more powerful word level language model in evaluation stage

- Original PIT-CE

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (3)$$

$$CE(\cdot) \longrightarrow MAP(\cdot)$$

**PIT-trained AM**

- Proposed:

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda \quad (4)$$

# Linguistics – Language Model Integration

- Motivation:
  - Improve assignment decision by combining LM in training stage
  - Still train a pure acoustic model and integrate it with more powerful word level language model in evaluation stage

- Original PIT-CE

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (3)$$

$$CE(\cdot) \longrightarrow MAP(\cdot)$$

**Senone level NNLM**

- Proposed:

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot \left( P(l_{\underline{utn}}^{(r)} | \mathbf{L}_{u(t-1)\underline{n}}^{(s')}) \right)^\lambda \quad (4)$$



# Linguistics – Language Model Integration

- Motivation:
  - Improve assignment decision by combining LM in training stage
  - Still train a pure acoustic model and integrate it with more powerful word level language model in evaluation stage

- Original PIT-CE

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathbf{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)})$$

- Proposed:

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)}) / P(l) \cdot P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')})}{P(\mathbf{O}_u^{(m)})}$$

$$\approx \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda$$

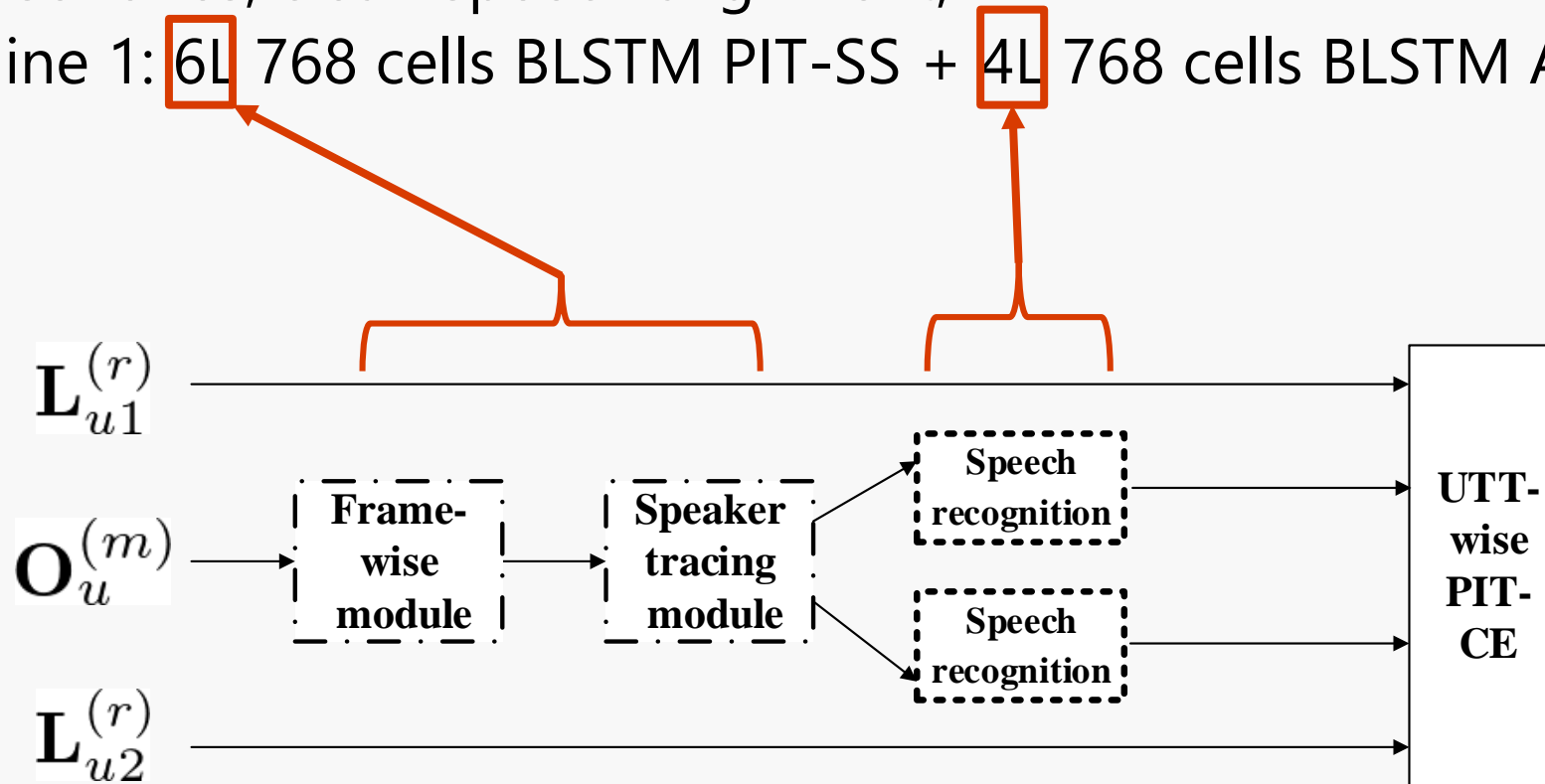
system	Ass.	Opt.
PIT	CE	CE
<b>Proposed</b>	<b>MAP</b>	<b>CE</b>
Disc. Train	MAP	MAP

**Discriminative training**

**Proposed method**

# Experiments

- Setup and baselines:
  - Artificial overlapped SWBD 300 → 150 (→ 50); hub5e-swbd 1831 → 915 utts
  - 9000 senones; clean speech alignment;
  - Baseline 1: 6L 768 cells BLSTM PIT-SS + 4L 768 cells BLSTM ASR



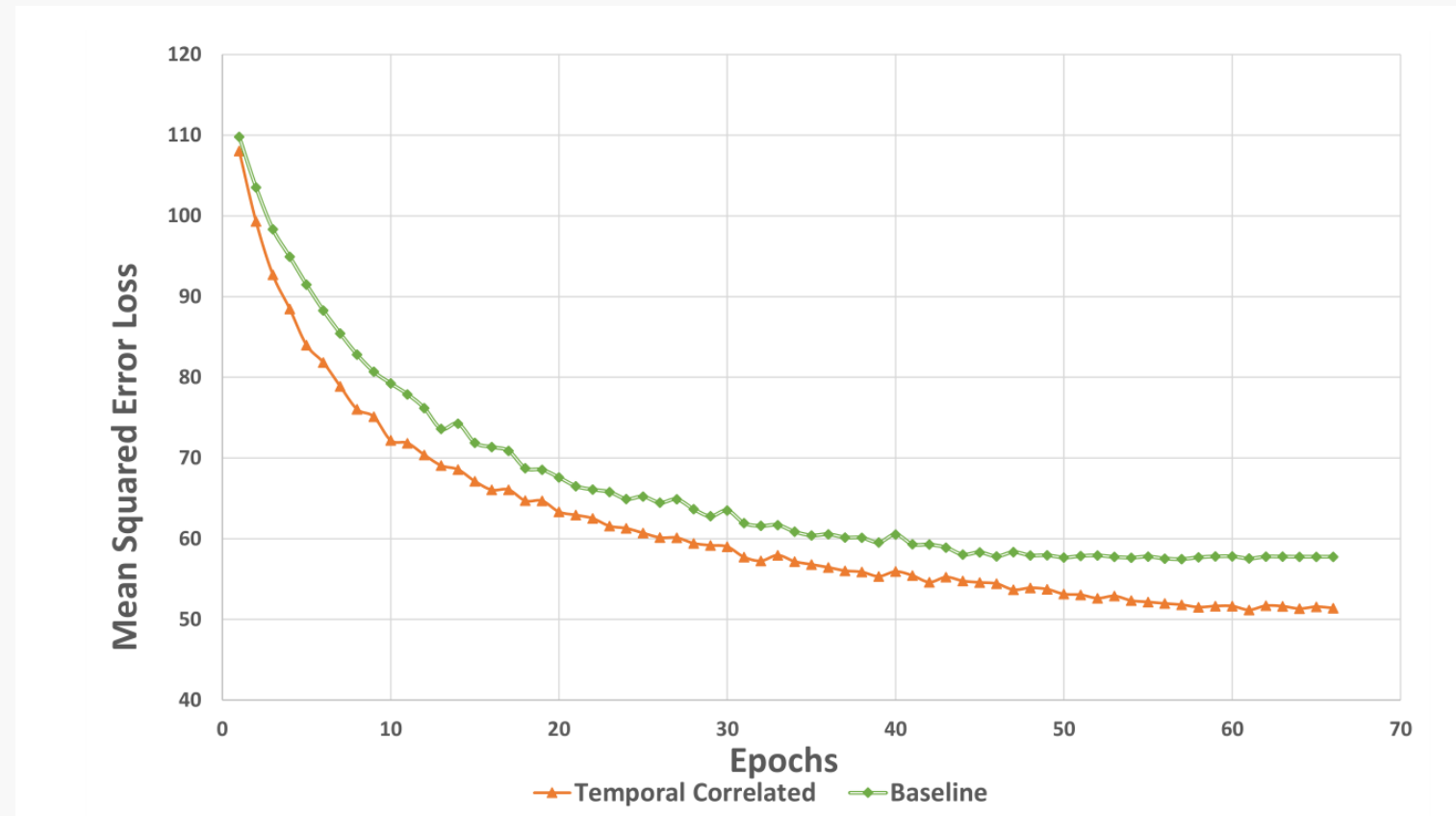
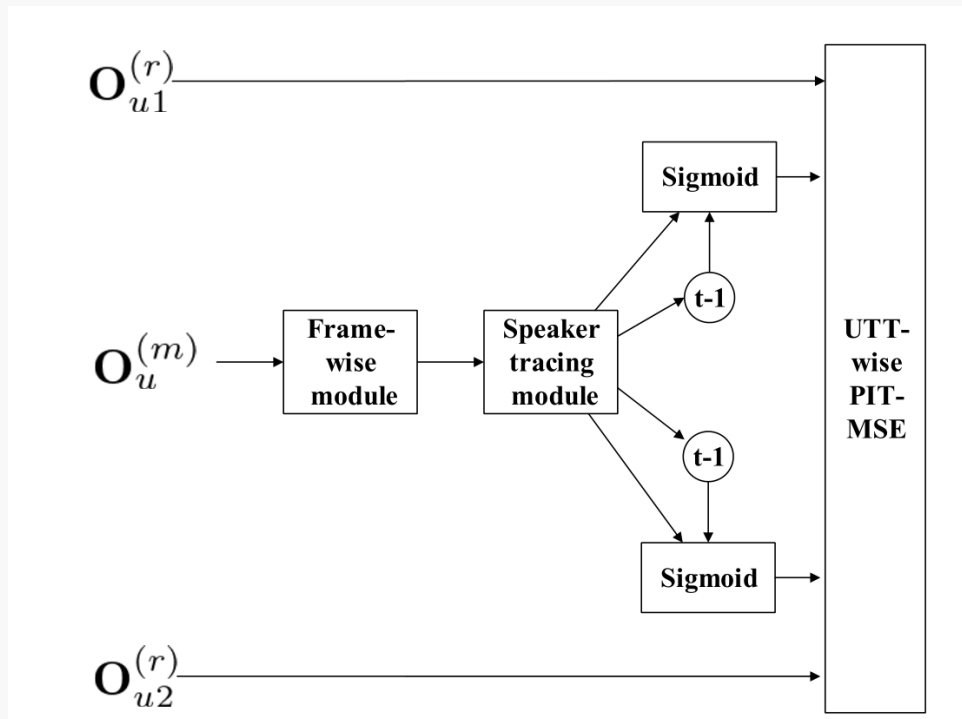
# Experiments

- Setup and baselines:
  - Artificial overlapped SWBD 300→150 (→50); hub5e-swb 1831 → 915 utts
  - 9000 senones; clean speech alignment;
  - Baseline 1: 6L 768 cells BLSTM PIT-SS + 4L 768 cells BLSTM ASR
  - Baseline 2: + transfer learning (TS, taught by clean teacher)

Neural network	Model	WER
6 BLSTM + 4 BLSTM	PIT-ASR	57.5
	progressive joint training + clean teacher	<b>38.9</b>

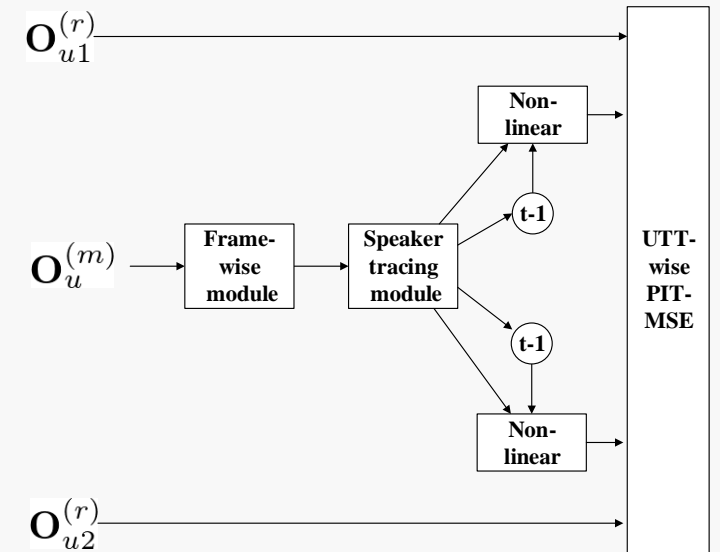
# Experiments – Temporal Correlated

- Baseline: modularization + clean teacher WER=38.9
- Improve in Speaker Tracing:



# Experiments – Temporal Correlated

- Baseline: modularization + clean teacher WER=38.9
- Improve in Speaker Tracing
- WER improve after joint training



Temporal Correlated	# of Sigmoid	WER	Rel. (%)
×	0	38.9	0
	0	37.5	-3.6
✓	1	<b>35.8</b>	<b>-8.0</b>
	2	36.7	-5.7

# Experiments – LM Integration

- Baseline: modularization + clean teacher WER=38.9

$$\mathcal{J}_{\text{U-PIT-CE}} = \sum_u \min_{s' \in \mathcal{S}} \sum_t \frac{1}{N} \sum_{n \in [1, N]} CE(l_{utn}^{(s')}, l_{utn}^{(r)}) \quad (3)$$

$$CE(\cdot) \longrightarrow MAP(\cdot)$$

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda \quad (4)$$

Assign.	Opt.	50 hours		150 hours	
		WER	Rel. (%)	WER	Rel. (%)
CE	CE	38.9	0	32.8	0
MAP	CE	37.3	<b>-4.1</b>	30.9	<b>-5.8</b>

# Experiments – LM Integration

- Baseline: modularization + clean teacher WER=32.8

$$MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) = \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda \quad (4)$$

Assign.	Opt.	50 hours		150 hours	
		WER	Rel. (%)	WER	Rel. (%)
CE	CE	38.9	0	32.8	0
MAP	CE	37.3	<b>-4.1</b>	30.9	<b>-5.8</b>

- with more data, the improvement becomes larger
  - AM becomes stronger
  - Assignment decision is not over-fit to the LM



# Experiments – Compare with disc. training

system	Assign.	Opt.	50 hours	
			WER	Rel. (%)
baseline	CE	CE	38.9	0
LM integration	MAP	CE	37.3	<b>-4.1</b>
LF-DC-bMMI	MAP	<b>MAP</b>	35.6	<b>-8.5</b>

**Discriminative training**

$$\begin{aligned}
 MAP(l_{utn}^{(s')}, l_{utn}^{(r)}) &= \frac{P(\mathbf{O}^{(m)} | \mathbf{L}^{(r)}) \cdot P(\mathbf{L}^{(r)})}{P(\mathbf{O}^{(m)})} \\
 &= \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)}) / P(l) \cdot P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')})}{P(\mathbf{O}_u^{(m)})} \\
 &\approx \frac{P(l_{utn}^{(r)} | \mathbf{O}_u^{(m)})}{P(l)} \cdot (P(l_{utn}^{(r)} | \mathbf{L}_{u(t-1)n}^{(s')}))^\lambda
 \end{aligned}$$

- Differences:

- optimization stage
- NNLM v.s. N-gram in discriminative training
- hardness in modeling  $P(\mathbf{O}_u^{(m)})$

**Proposed method**

# Experiments – Combination

Method	WER	Rel. (%)
baseline	38.9	0
+ Temporal Correlated	35.8	-8.0
+ LM Integration	34.4	-11.5
+ LF-DC-bMMI	31.6	-18.8

- Operate in different levels → can be combined

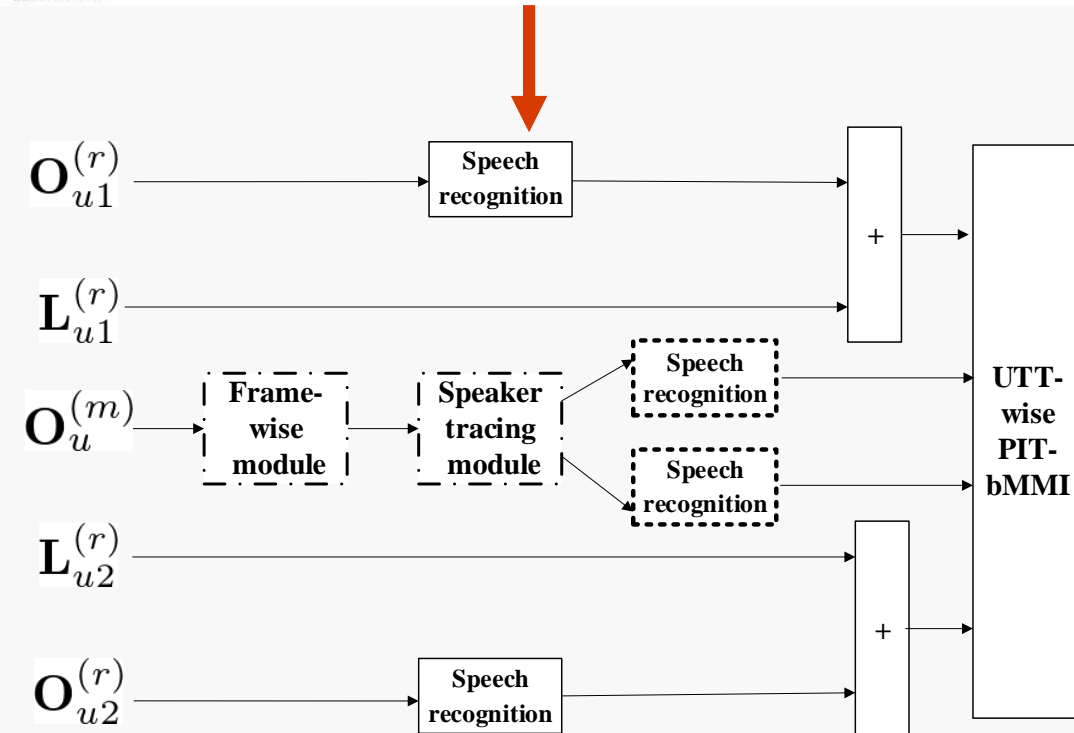
# Experiments – Combination

Method	WER	Rel. (%)
baseline	38.9	0
+ Temporal Correlated	35.8	-8.0
+ LM Integration	34.4	-11.5
+ LF-DC-bMMI	31.6	-18.8
+ MMI clean teacher	35.8	-8.0
+ LF-DC-bMMI	35.2	-9.5

- Operate in different levels → can be combined
- Better than only utilize TS + discriminative training

# Our final system

- Acoustics
  - Modular Initialization 4%
    - CNN 10%
  - Transfer Learning Based Joint Training 20%
  - **Temporal Correlation Modeling 8%**
- Linguistics
  - Multi-outputs Sequence Discriminative Training 8%
  - **Integrating Language Model in Assignment Decision 4%**

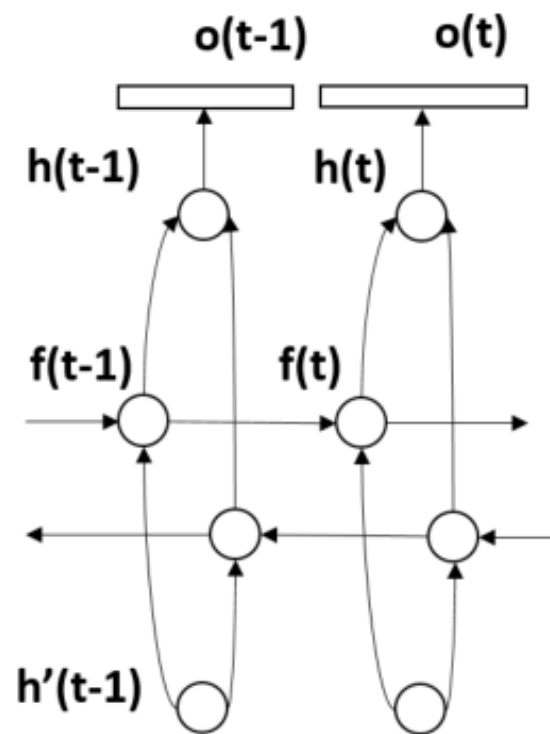


[1] Zhehuai Chen, Jasha Droppo, Sequence Modeling in Unsupervised Single-channel Overlapped Speech Recognition, accepted by IEEE International Conference on Acoustics, Speech and Signal Processing(ICASSP), Calgary, Canada, 2018.

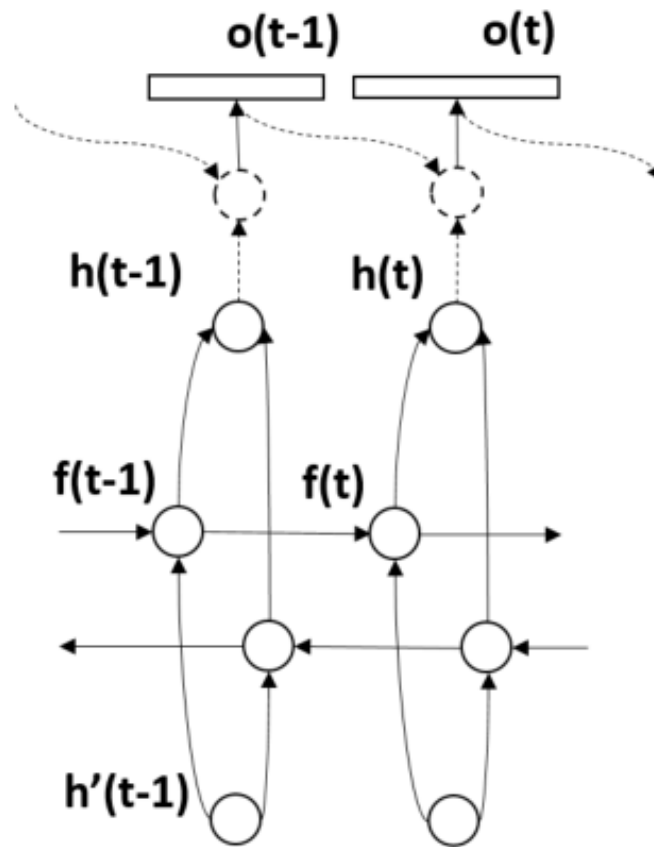
[2] Zhehuai Chen, Jasha Droppo, Jinyu Li, Wayne Xiong, Progressive Joint Modeling in Unsupervised Single-channel Overlapped Speech Recognition. IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 26, no. 1, pp. 184-196, Jan. 2018. doi: 10.1109/TASLP.2017.2765834.

# Backup materials

# Temporal correlation modeling in BLSTM



(a) BLSTM



(b) Temporal Correlated BLSTM

# Experiments – Example 50hrs (F-F)

- Clean ASR (90+WER)
- 1 PIT-CE
- 2 Transf.
- 3 +MMI teacher
- 4 +seq. disc. tr.







# Experiments – Example 150hrs (F-F)

- Clean ASR (90+WER)
- 1 PIT-CE
- 2 Transf.
- 3 +CNN
- 4 +seq. disc. tr.



