

# Minimum Latency Training Strategies for Streaming Sequence-to-Sequence ASR

**Hirofumi Inaguma**<sup>1,2</sup>, Yashesh Gaur<sup>2</sup>, Liang Lu<sup>2</sup>, Jinyu Li<sup>2</sup>, Yifan Gong<sup>2</sup>

Kyoto University, Kyoto, Japan<sup>1</sup>

Microsoft Speech and Language Group, Redmond, WA, USA<sup>2</sup>

IEEE ICASSP 2020



# Background: End-to-end ASR

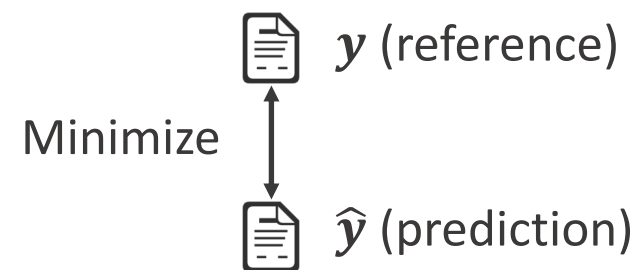
- Input sequence (speech):  $\mathbf{x} = (x_1, \dots, x_T)$
- Output sequence (transcription):  $\mathbf{y} = (y_1, \dots, y_L)$

## Time-synchronous model ( $|\mathbf{x}| = |\hat{\mathbf{y}}|$ )

- Connectionist temporal classification (CTC) [Graves et al., 2006]
- RNN-Transducer (RNN-T) [Graves et al., 2013]
- Recurrent neural aligner (RNA) [Sak et al., 2017]

## Label-synchronous model ( $|\mathbf{x}| \neq |\hat{\mathbf{y}}|$ )

- Attention-based sequence-to-sequence (S2S) [Bahdanau et al., 2016]
- Transformer [Vaswani et al., 2017]



Low accuracy  
Streaming: easy

High accuracy  
Streaming: difficult

# Streaming attention-based S2S ASR

## Neural Transducer [Jailty et al., 2015]


- Perform attention mechanism for a fixed size of block

## Hard monotonic attention [Raffel et al., 2017]

- Learn to detect token boundaries via stochastic binary decision
- Extension: **Monotonic chunkwise attention (MoChA)** [Chiu et al., 2018]

## Triggered attention [Moritz et al., 2018]

- Perform global attention over encoder memories trun

- 
- **Good results**
  - **Efficient training**

## Adaptive computation steps (ACS) [Li et al., 2018]

- Learn how many tokens to generate with encoder outputs

## Continuous Integrate-and-Fire (CIF) [Dong et al., 2019]

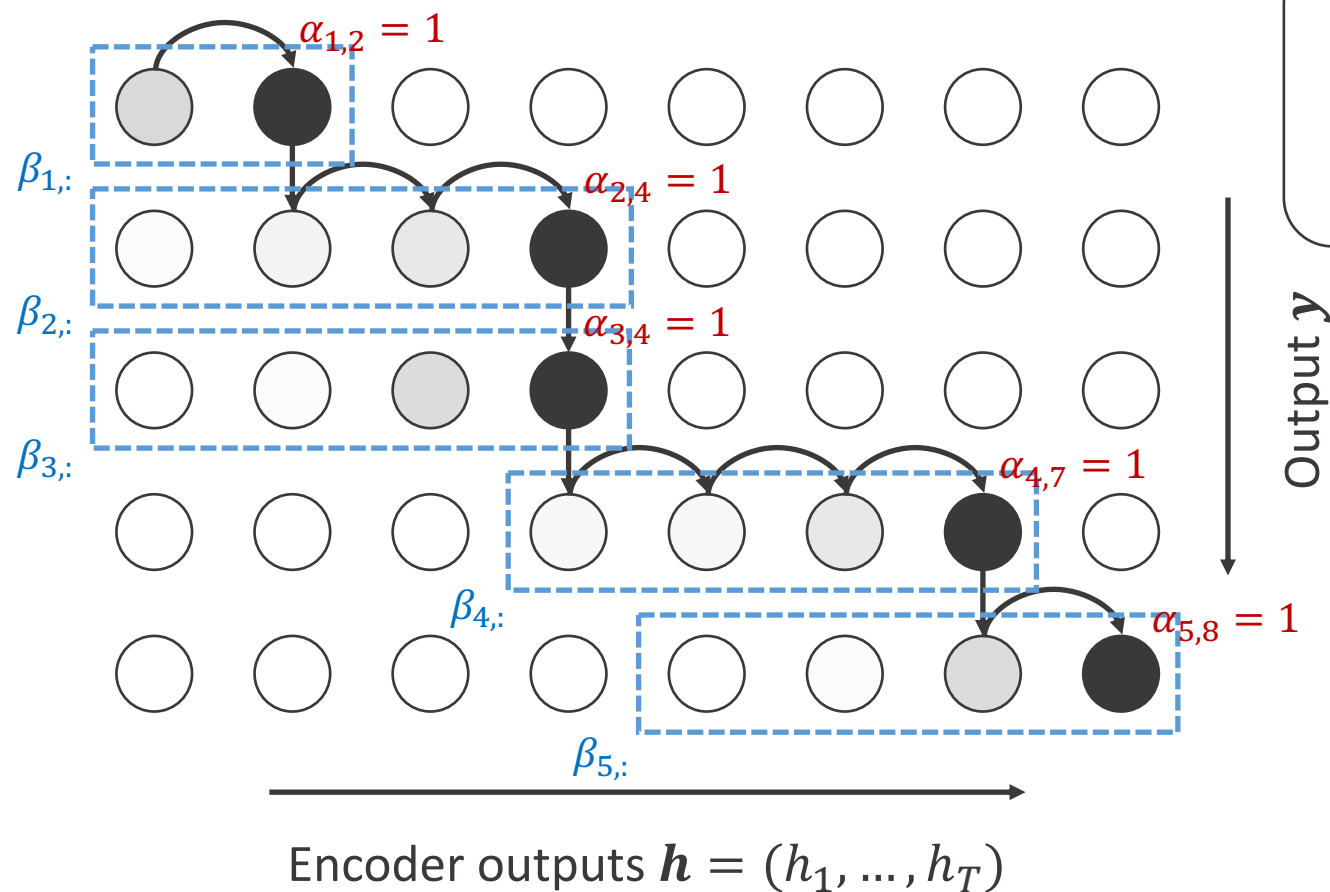
- Fine-grained version of ACS

### And more...

- Windowing approaches
- Reinforcement learning

# MoChA (test time)

e.g.,  $w = 4$  (chunk size: 4)

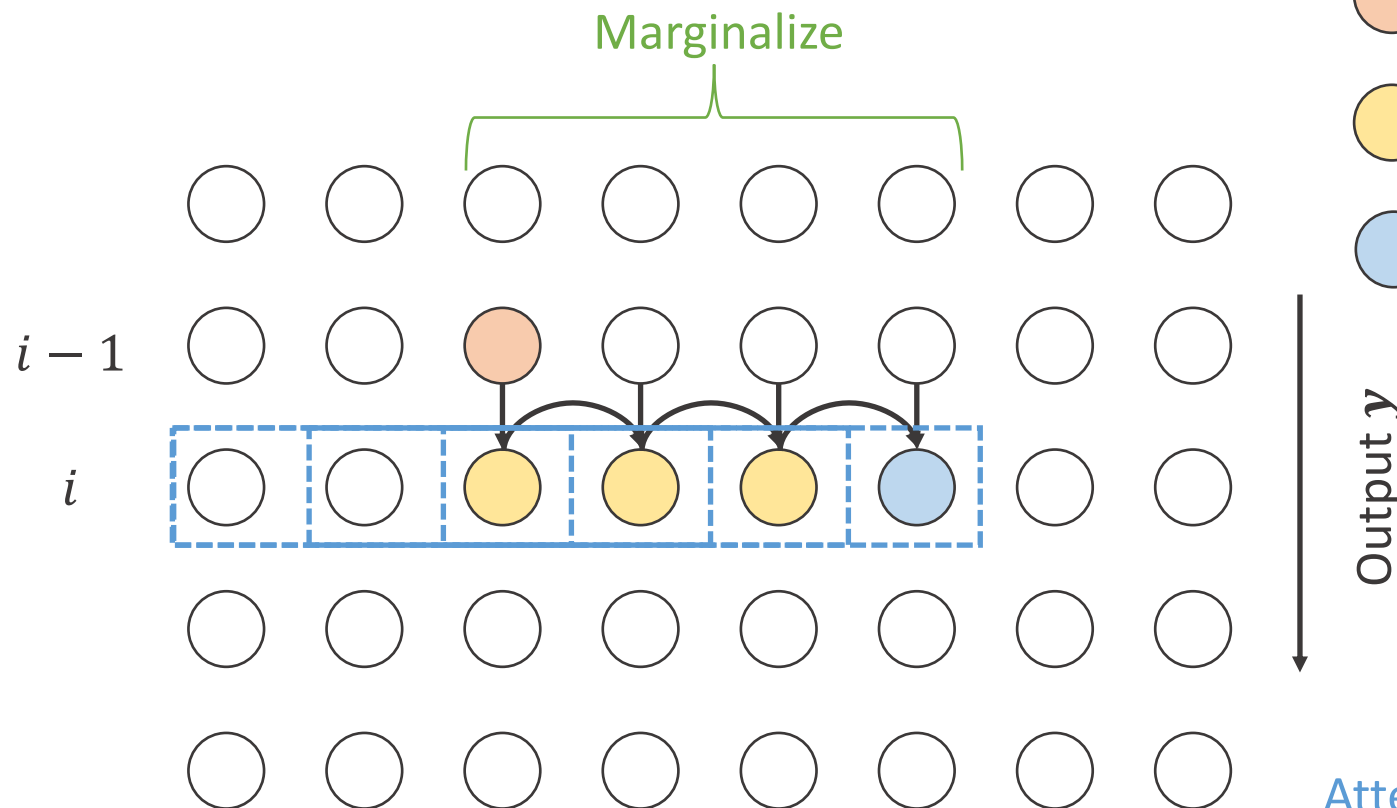





● : Attend ( $\alpha_{i,j} = 1$ )  
○ : Not attend ( $\alpha_{i,j} = 0$ )

Not differentiable

1. **Monotonic attention**: whether to attend or not
2. **Chunkwise attention**: soft attention over a small window

# MoChA (training time)



-  : Attend at  $(i - 1)$ -th step
-  : Not attend
-  : Attend at  $i$ -th step

Use expected alignments  
during training  
for backpropagation

Encoder outputs  $\mathbf{h} = (h_1, \dots, h_T)$

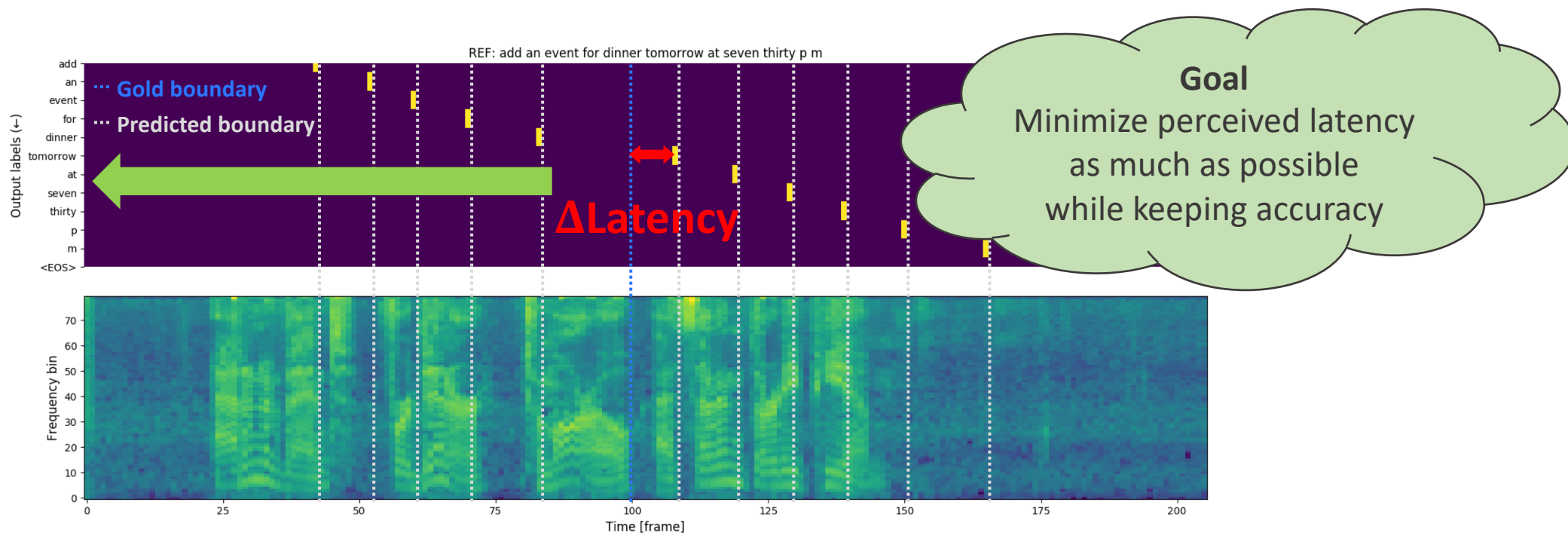
Previous attention

Attend      Not attend

$$\alpha_{i,j} = p_{i,j} \sum_{k=1}^j \left( \alpha_{i-1,k} \prod_{l=k}^{j-1} (1 - p_{i,l}) \right)$$

$$= (1 - p_{i,j-1}) \frac{\alpha_{i,j-1}}{p_{i,j-1}} + \alpha_{i-1,j}$$

# Delayed token generation problem



- Decision boundaries (yellow dots) are delayed from the actual acoustic boundary
  1. Unidirectional encoder (lacking the future information)
  2. Sequence-level criterion (utilizing as many future frames as possible to maximize the log-likelihood)
- This leads to increasing user perceived latency
  - Similar behaviors have been reported in CTC [sak et al., 2015] and RNN-T [Li et al., 2019]

# Evaluation metric: latency

- Definition: difference between time-index of a predicted boundary and that of the gold boundary

***Corpus-level latency*** (averaged per token)

$$\Delta_{\text{corpus}} = \frac{1}{\sum_{k=1}^N |\mathbf{y}^k|} \sum_{k=1}^N \sum_{i=1}^{|\mathbf{y}^k|} (\widehat{b}_i^k - b_i^k)$$

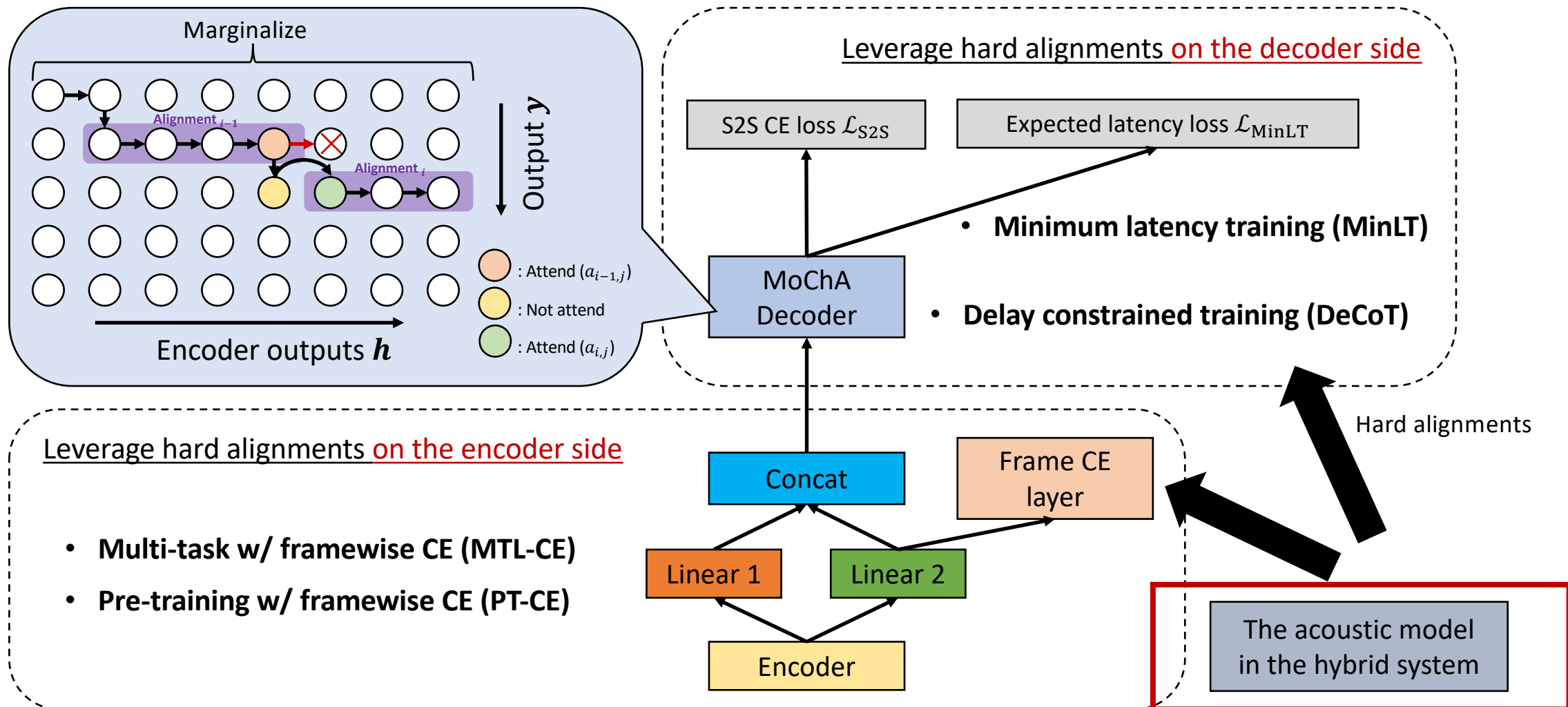
***Utterance-level latency*** (averaged per utterance)

$$\Delta_{\text{utterance}} = \frac{1}{N} \sum_{k=1}^N \frac{1}{|\mathbf{y}^k|} \sum_{i=1}^{|\mathbf{y}^k|} (\widehat{b}_i^k - b_i^k)$$

- Report (1) average, (2) median, (3) 90-th, and (4) 99-th percentile
- Teacher-forcing when calculating latency to match the sequence lengths

# Proposed methods

Where should we apply alignment information in the model?



Leveraging hard alignments extracted from the hybrid system



# 1. Multi-task learning w/ framewise CE (MTL-CE)

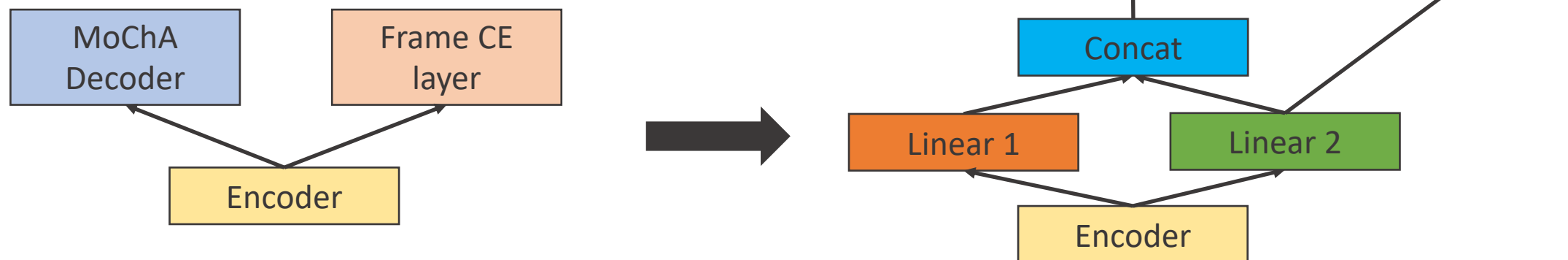
## ◆ Objective function

$$\mathcal{L}_{\text{total}} = (1 - \lambda_{\text{CE}}) \underbrace{\mathcal{L}_{\text{S2S}}(\mathbf{y}|\mathbf{x})}_{\text{MoChA}} + \lambda_{\text{CE}} \underbrace{\mathcal{L}_{\text{CE}}(\mathbf{A}|\mathbf{x})}_{\text{Frame CE}} \quad (0 \leq \lambda_{\text{CE}} \leq 1)$$

- Motivation: align encoder outputs to the true acoustic location

## ◆ Insert linear bottleneck layers

- Inspired by the CTC acoustic model [Yu et al., 2018]

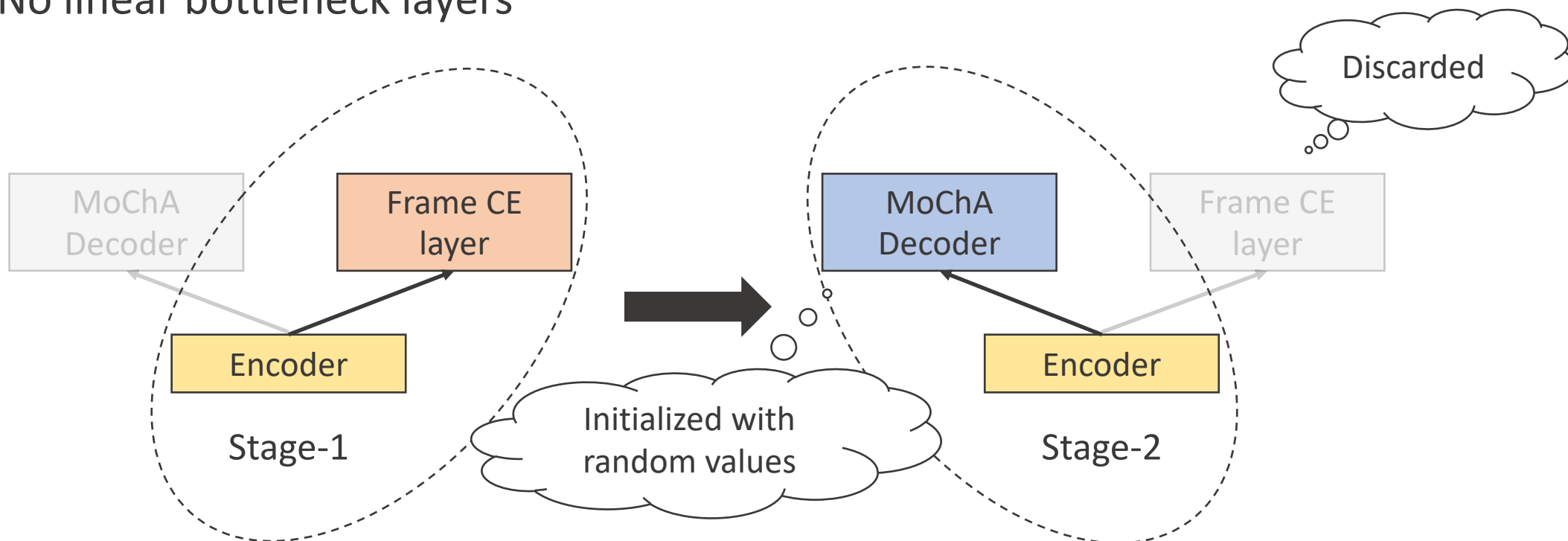


Train both branches  
from scratch

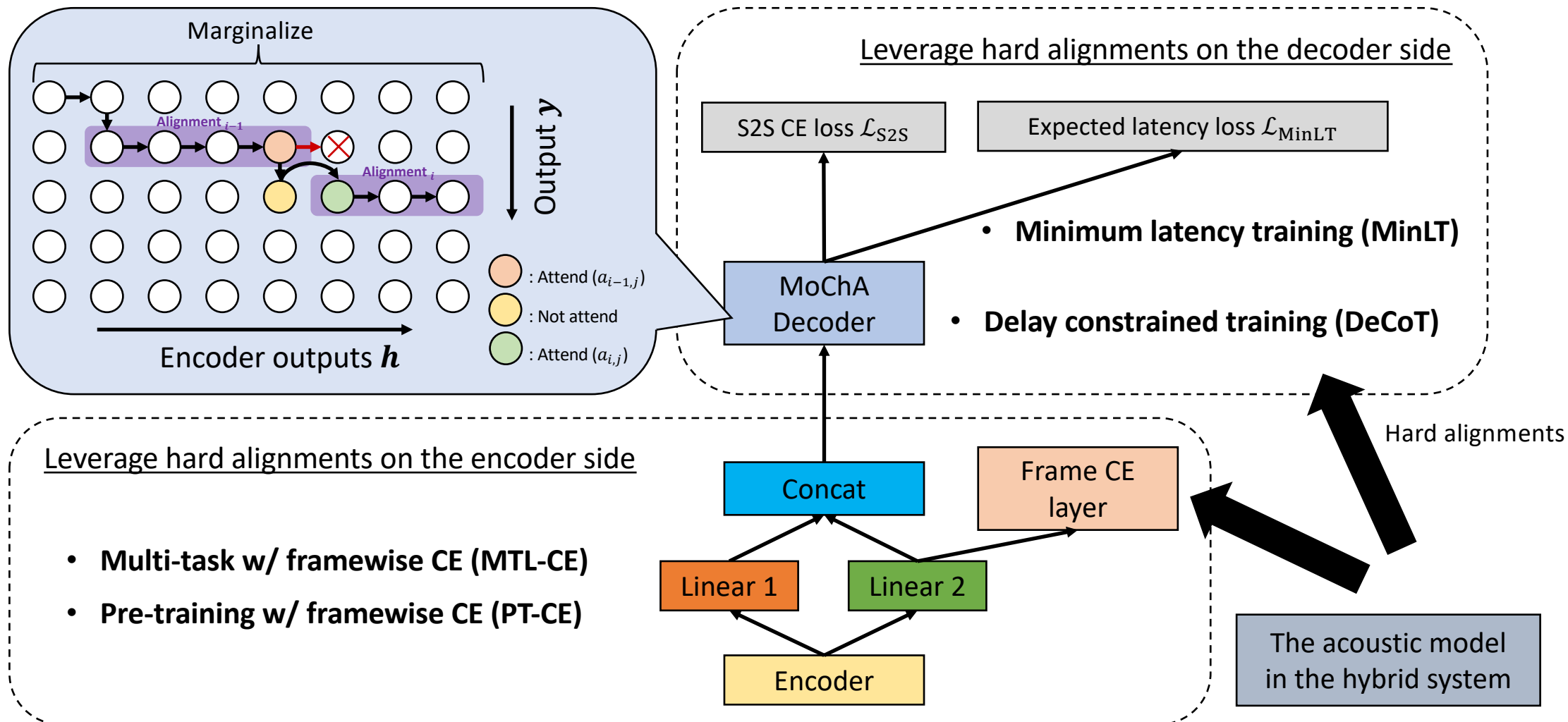
## 2. Pre-training with framewise CE (PT-CE)

### ◆ 2-staged training

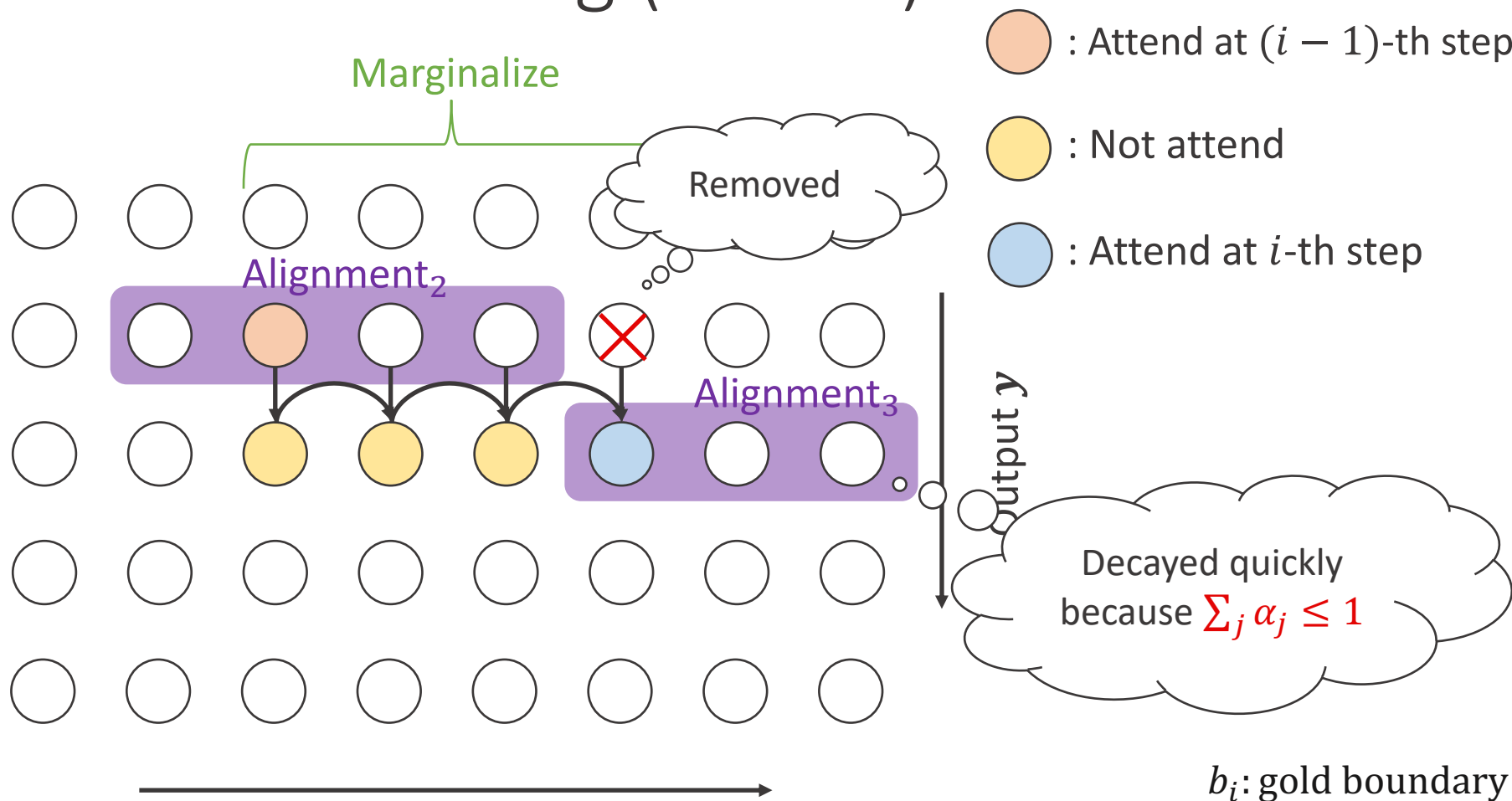
- Motivation
  - Start training from well-aligned encoder representations
  - Do not have to tune the CE weight  $\lambda_{CE}$
- No linear bottleneck layers



# Proposed methods



### 3. Delay constrained training (DeCoT)



Remove inappropriate paths whose boundaries surpass the actual acoustic boundary more than fixed acceptable latency  $\delta$  [frames]

$$\alpha_{i,j} = \begin{cases} p_{i,j} \left( (1 - p_{i,j-1}) \frac{\alpha_{i,j-1}}{p_{i,j-1}} + \alpha_{i-1,j} \right) & (j \leq b_i + \delta) \\ 0 & (\text{otherwise}) \end{cases}$$

### 3. Delay constrained training (DeCoT)

#### Regularization with quantity loss

- Add a regularization term to keep  $\sum_j \alpha_j = 1$
- Originally proposed in CIF [Dong et al., 2019] with a different motivation

$L$ : the number of tokens in the reference

$$\mathcal{L}_{\text{QUA}} = |L - \sum_{i=1}^L \sum_{j=1}^T \alpha_{i,j}|$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{S2S}} + \lambda_{\text{QUA}} \mathcal{L}_{\text{QUA}} \quad (\lambda_{\text{QUA}} \geq 0)$$

## 4. Minimum latency training (MinLT)

### ◆ Objective function

- Directly minimize the expected latency  $\mathcal{L}_{\text{MinLT}}$  by utilizing hard alignments  $A$

Expected boundary

$$\mathcal{L}_{\text{MinLT}} = \frac{1}{|y|} \sum_{i=1}^L \left| \sum_{j=1}^T j \alpha_{i,j} - b_i \right| \quad (b_i: \text{reference boundary for } i\text{-th token})$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{S2S}} + \lambda_{\text{MinLT}} \mathcal{L}_{\text{MinLT}} \quad (\lambda_{\text{MinLT}} \geq 0)$$

- Motivation: reduce latency flexibly
  - DeCoT assumes the fixed latency for each token

### ◆ Related work

- Latency loss has been investigated in simultaneous NMT [Arivazhagan et al., 2019]
- Non-silence frames are not distributed uniformly over the input speech in ASR

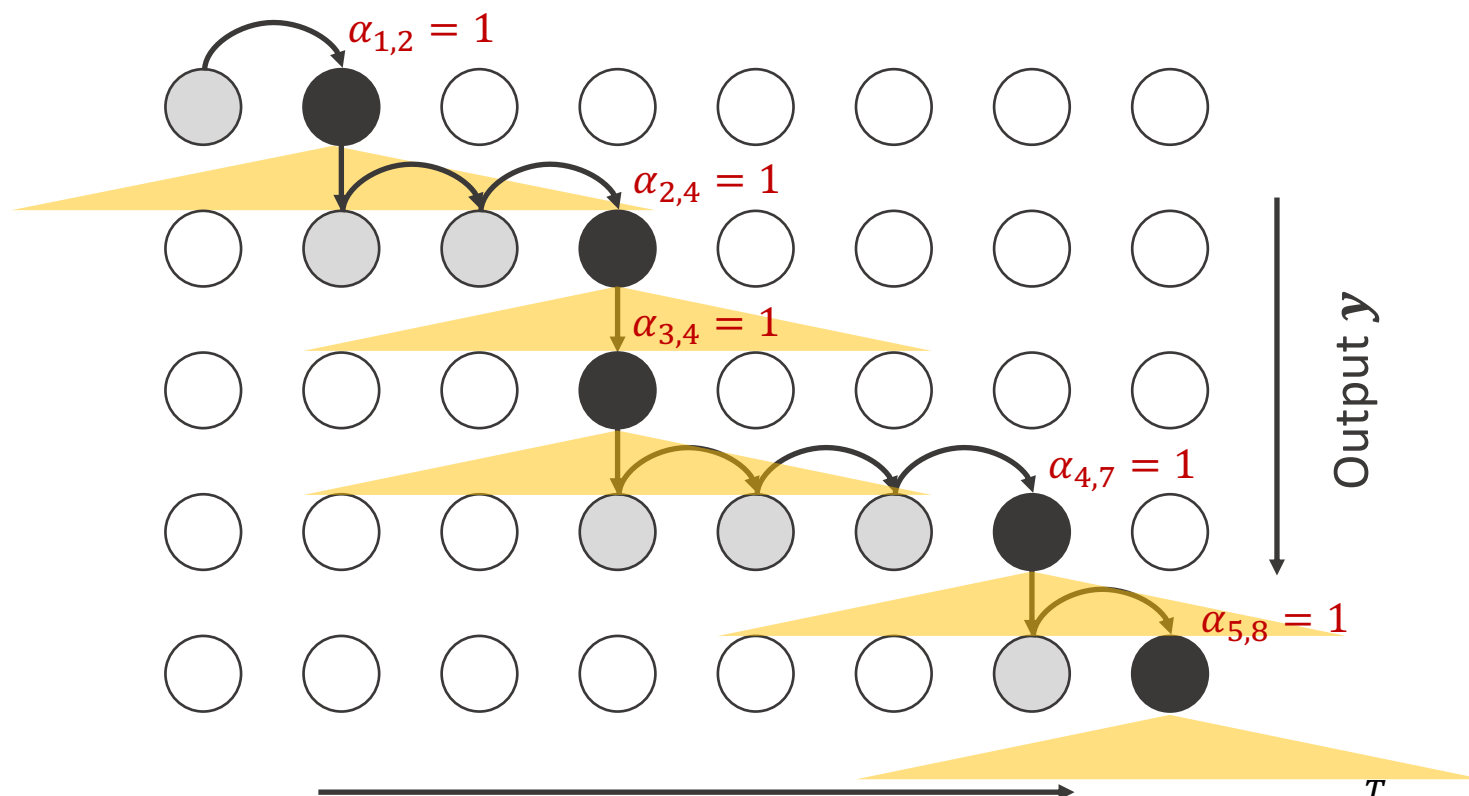
# Experimental condition

Data	Train: Cortana voice assistant (3.4k hours) Validation: Sampled disjoint 4k utterances from the training set Test: 5.6k utterances
Feature	80-dim log-mel fbank (3 frame stacked, 30ms per frame)
Output unit	Mixed units (34k vocabulary)
Architecture	Offline: 512-dim (per direction) 6-layer BiGRU encoder Streaming: 1024-dim 6-layer GRU encoder Decoder: 512-dim 2-layer GRU
Optimization	Adam
Decoding	Beam width: 8, no LM

- Word-level alignments:  $A = (a_1, \dots, a_T)$  ( $\{a_j\}_{j=1, \dots, T}$ : one-hot vector)
  - Divide duration based on the ratio of the character length of each subword
- Start DeCoT and MinLT from the baseline MoChA (warm start training)

# Enhance monotonic attention with 1D convolution

e.g.,  $k = 5$  (lookahead: 2, 60ms)



## ◆ Motivation

- Leveraging the surrounding frames for robust binary decision

$$e_{i,j}^{\text{mono}} = g \frac{v^T}{\|v\|} \text{ReLU}(W_h h'_j + W_s s_i + b) + r$$


$$h'_j = W_c * h_j \quad (W_c \in \mathbb{R}^{d \times d \times k})$$

$k$ : kernel size,  $d$ : unit size



# Results: Baseline

Model		WER [%]
Offline	BiGRU global attention	7.01
	UniGRU global attention	8.44
	BiGRU MoChA (chunk: 4)	8.09
Streaming	UniGRU MoChA (chunk:4)	10.37
	+ 1D-convolution (baseline)	<b>9.93</b>



4.24% WERR

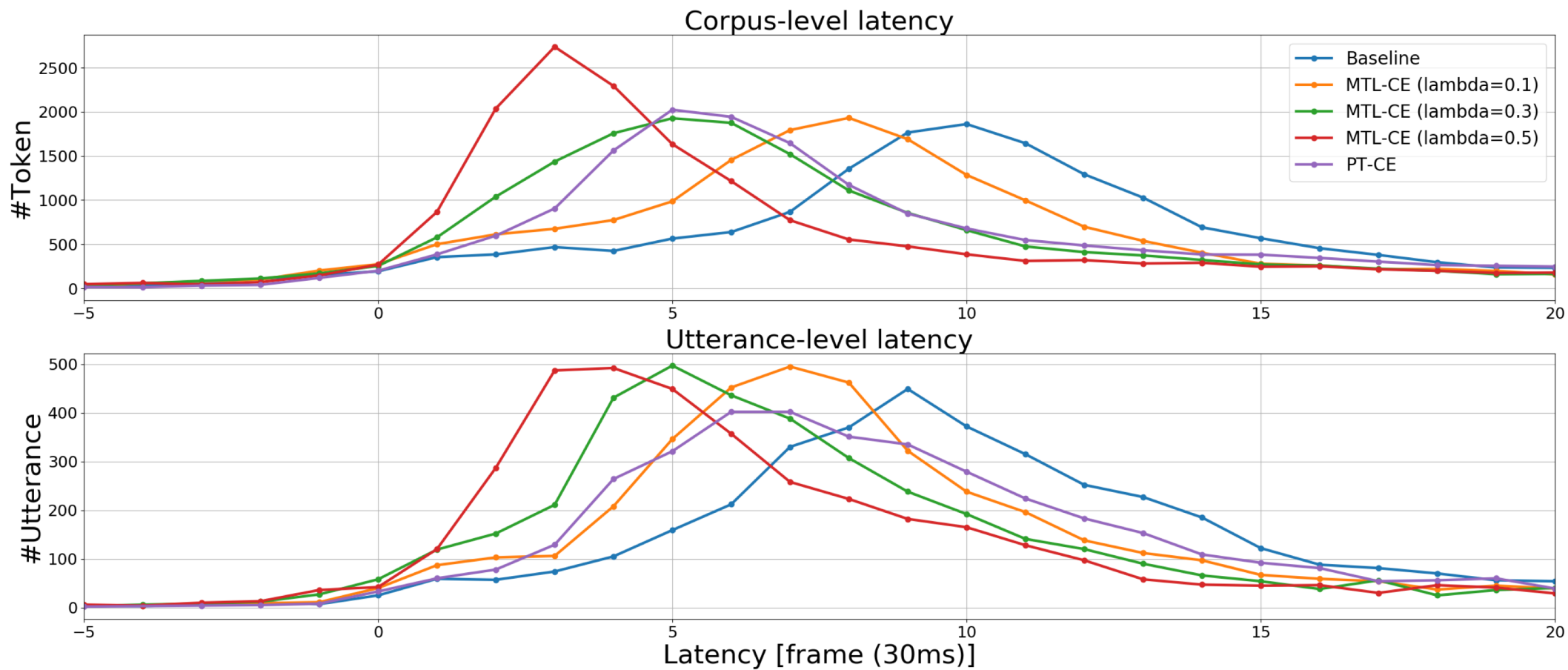
- Huge gaps between (1) bidirectional  $\leftrightarrow$  unidirectional  
(2) offline  $\leftrightarrow$  streaming S2S
- 1D-convolution layer improved the streaming MoChA by **4.24%** relatively

# Results: Alignments on the encoder side

Model	WER [%]	Corpus-level [frame (30ms)]			
		Ave.	Med.	90th	99th
Baseline MoChA	9.93	11.65	10.00	21.39	<b>44.29</b>
MTL-CE ( $\lambda_{\text{CE}} = 0.1$ )	10.21 <span>5.6%</span>	9.84	8.00 <span>40%</span>	<b>19.42</b>	46.54
MTL-CE ( $\lambda_{\text{CE}} = 0.3$ )	10.48	8.78	6.00	19.69	47.96
MTL-CE ( $\lambda_{\text{CE}} = 0.5$ )	11.11	<b>8.36</b>	<b>5.00</b>	21.21	49.86
PT-CE	12.74	10.49	7.00	22.90	48.65

- MTL-CE reduced latency in proportion to  $\lambda_{\text{CE}}$  while degrading WER slightly
- PT-CE also reduced latency but degraded WER too much
- Contrastive results to previous works using CTC + framewise CE objective
  - MoChA is a label-synchronous model

# Visualization of latency distribution (encoder)

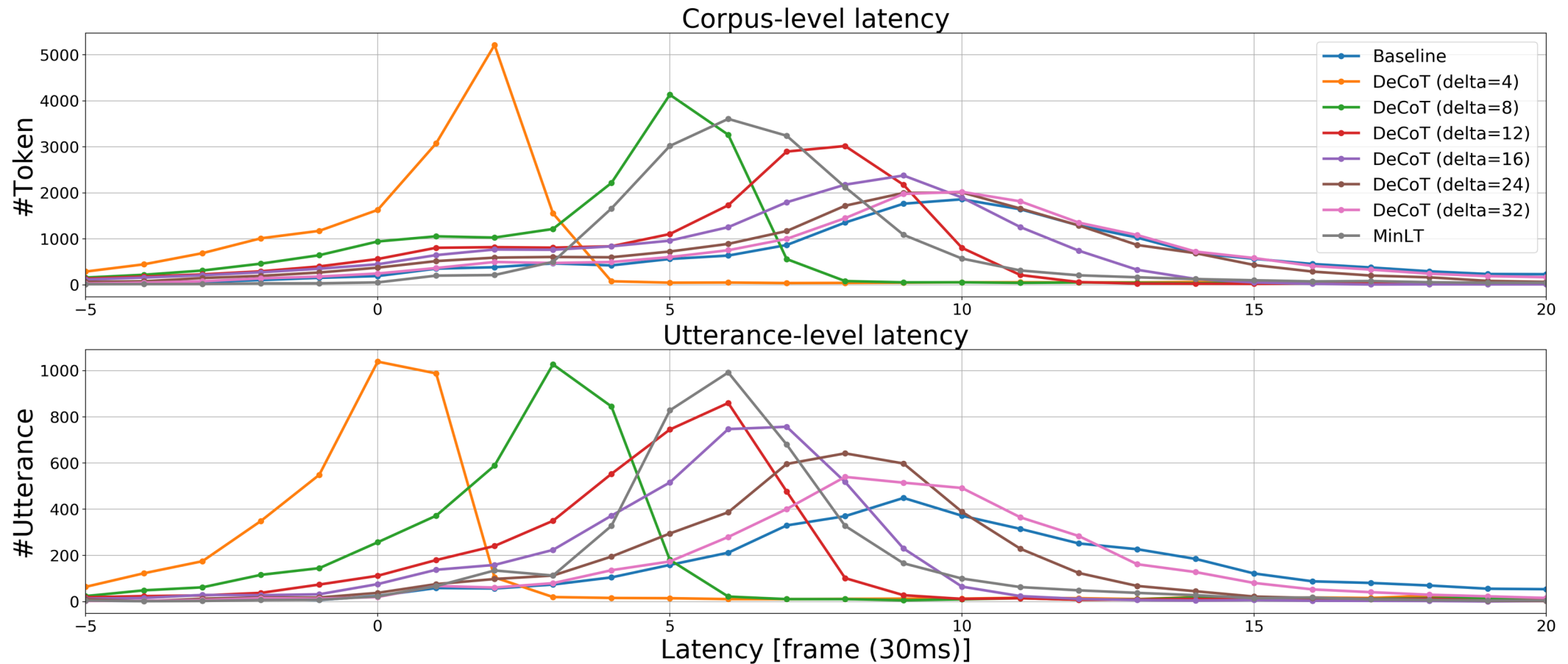


# Results: Alignments on the decoder side

Model	WER [%]	Corpus-level [frame (30ms)]			
		Ave.	Med.	90th	99th
UniGRU global attention (offline)	8.44	N/A	N/A	N/A	N/A
Baseline MoChA	9.93	11.65	10.00	21.39	44.29
DeCoT ( $\delta = 4$ )	20.25	3.66	1.00	9.56	62.27
DeCoT ( $\delta = 8$ )	14.35	4.60	5.00	7.00	47.04
DeCoT ( $\delta = 12$ )	11.40	6.02	7.00	9.92	35.58
DeCoT ( $\delta = 16$ )	<b>9.13</b>	<b>6.63</b>	<b>8.00</b>	<b>11.71</b>	<b>16.43</b>
DeCoT ( $\delta = 24$ )	<b>8.87</b>	<b>8.37</b>	<b>9.00</b>	<b>14.45</b>	<b>21.07</b>
DeCoT ( $\delta = 32$ )	9.17	9.79	10.00	16.54	27.01
MinLT	<b>9.70</b>	<b>7.06</b>	<b>6.00</b>	<b>10.63</b>	<b>26.76</b>

- DeCoT: large WER improvement and moderate latency reduction (tail part)
- MinLT: small WER improvement and large latency reduction (median)

# Visualization of latency distribution (decoder)



# Ablation study: Decoder side

- **Combination of DeCoT and MinLT reduced the latency, but degraded WER too much**

Model	WER [%]	Corpus-level [frame (30ms)]			
		Ave.	Med.	90th	99th
DeCoT ( $\delta = 16$ )	9.13	6.63	8.00	11.71	16.43
+ MinLT	12.75	4.05	4.00	7.96	15.92
MinLT	9.70	7.06	6.00	10.63	26.76

# Ablation study: Decoder side

- **Quantity loss was essential for DeCoT but not necessary for the baseline and MinLT**

Model	WER [%]	Corpus-level [frame (30ms)]			
		Ave.	Med.	90th	99th
Baseline MoChA	9.93	11.65	10.00	21.39	44.29
+ Quantity loss	10.30	11.24	10.00	20.39	36.01
DeCoT ( $\delta = 16$ )	9.13	6.63	8.00	11.71	16.43
- Quantity loss	14.28	3.93	3.00	7.20	27.39
MinLT	9.70	7.06	6.00	10.63	26.76
+ Quantity loss	13.66	6.82	6.00	10.45	25.57

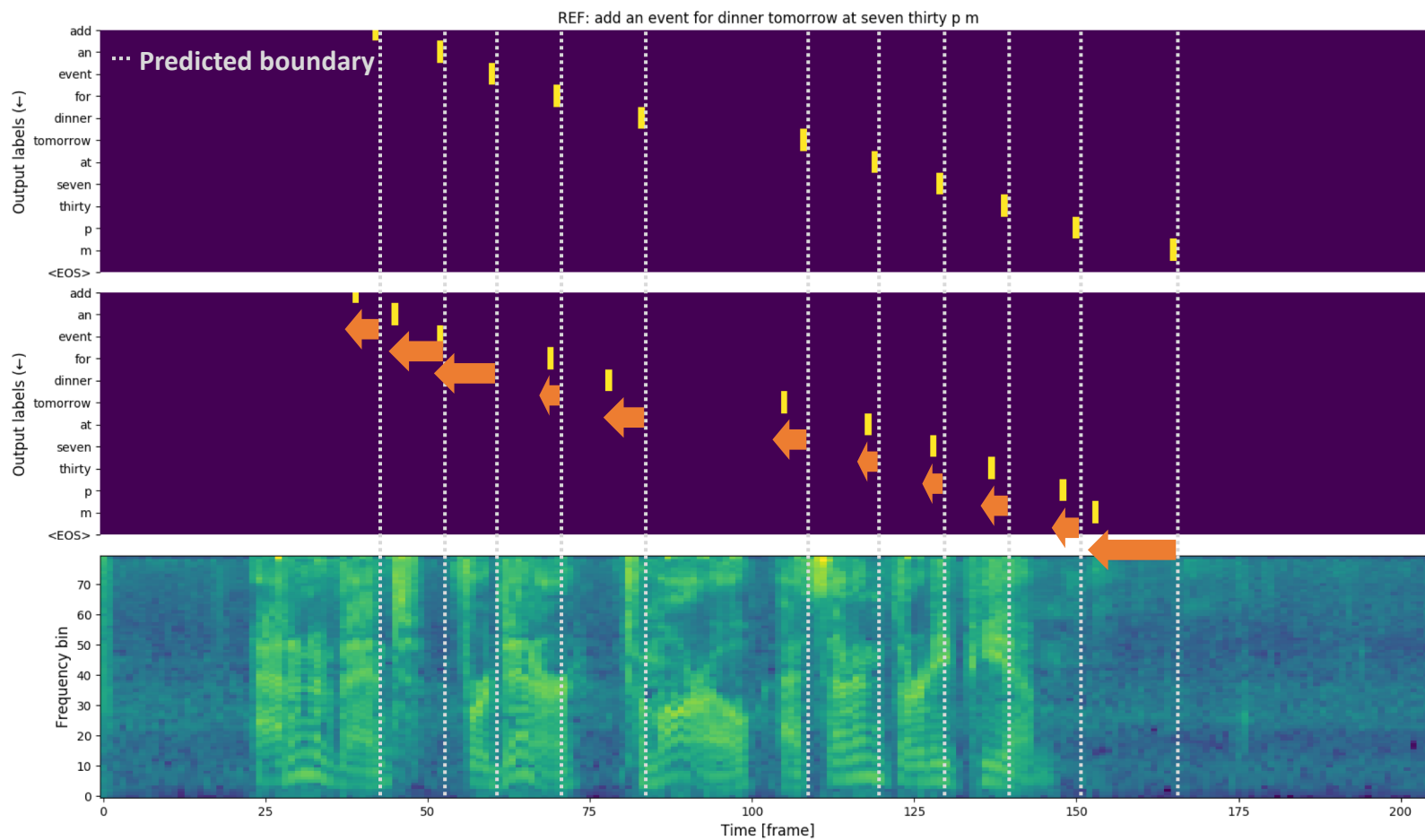
# Ablation study: Decoder side

- **Warm start training from the baseline was effective for DeCoT and MinLT**

Model	WER [%]	Corpus-level [frame (30ms)]			
		Ave.	Med.	90th	99th
Baseline MoChA	9.93	11.65	10.00	21.39	44.29
+ Warm start training	9.21	12.27	11.00	22.23	43.16
DeCoT ( $\delta = 16$ )	9.13	6.63	8.00	11.71	16.43
- Warm start training	10.72	6.28	7.00	11.12	36.03
MinLT	9.70	7.06	6.00	10.63	26.76
- Warm start training	13.63	11.83	10.00	21.41	45.06



# Alignment visualization

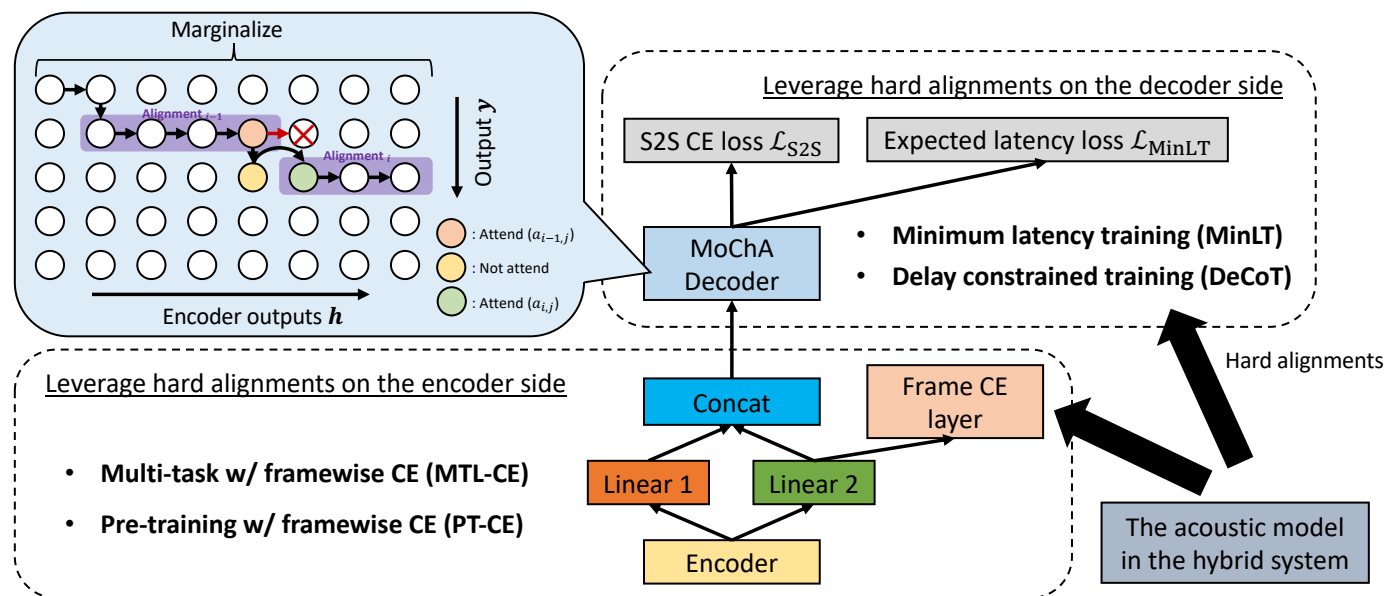


Baseline MoChA

DeCoT ( $\delta = 16$ )

# Conclusion

- Explored to leverage frame-level hard alignments extracted from the hybrid system to reduce user perceived latency
- Alignments were effective for latency reduction on both sides, and also improved ASR performance when applying on the decoder side



Question?

E-mail: inaguma [at] sap.ist.i.kyoto-u.ac.jp