

# Toward low-latency and accurate simultaneous interpretations from speech

Hirofumi Inaguma

Ph.D. candidate, Kyoto University, Japan

12/09/2020



# Agenda

## ◆ Streaming end-to-end automatic speech recognition (ASR)

- *Monotonic chunkwise attention (MoChA)* [Chiu+ 2018]
- *How to reduce latency with alignment information?*
- Where to apply? (encoder/decoder)
  - *Minimum Latency Training Strategies for Streaming Sequence-to-Sequence ASR [ICASSP 2020]*
- Leverage CTC alignment (hybrid ASR-free)
  - *CTC-synchronous Training for Monotonic Attention Model [Interspeech2020]*

## ◆ Non-autoregressive end-to-end speech translation: A first study

- *Conditional masked language model (CMLM)* [Ghazvininejad+ 2019]
- How to estimate **target lengths** from speech directly?
  - Orthros: Non-autoregressive End-to-end Speech Translation with Dual-decoder [under review]

# Background: Hybrid ASR system

- Traditional approach (still dominant in production system)

Acoustic model (AM) Language model (LM)

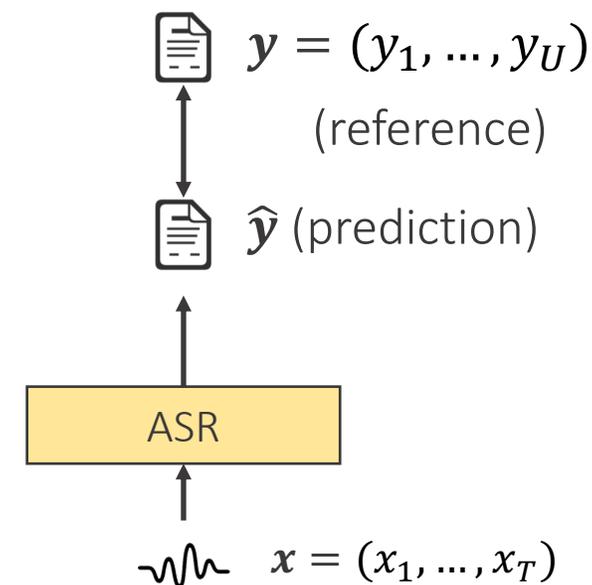
$$P(\mathbf{y}|\mathbf{x}) = \frac{P(\mathbf{x}|\mathbf{y})P(\mathbf{y})}{P(\mathbf{x})}$$

- Inference

$$\begin{aligned}\hat{\mathbf{y}} &= \arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) \\ &= \arg \max_{\mathbf{y}} \underline{P(\mathbf{x}|\mathbf{y})} P(\mathbf{y})\end{aligned}$$

$\mathbf{y}$  (word)  $\rightarrow$   $\mathbf{p}$  (pronounce)  $\rightarrow$   $\mathbf{s}$  (HMM state)

- 😊 Rare words, low-resource, module update (customization)
- 😓 Expertized knowledge



# Background: End-to-end ASR system

- Learn a direct mapping function  $\varphi(x)$  to maximize  $P(y|x)$

😊 Quick development, scalability

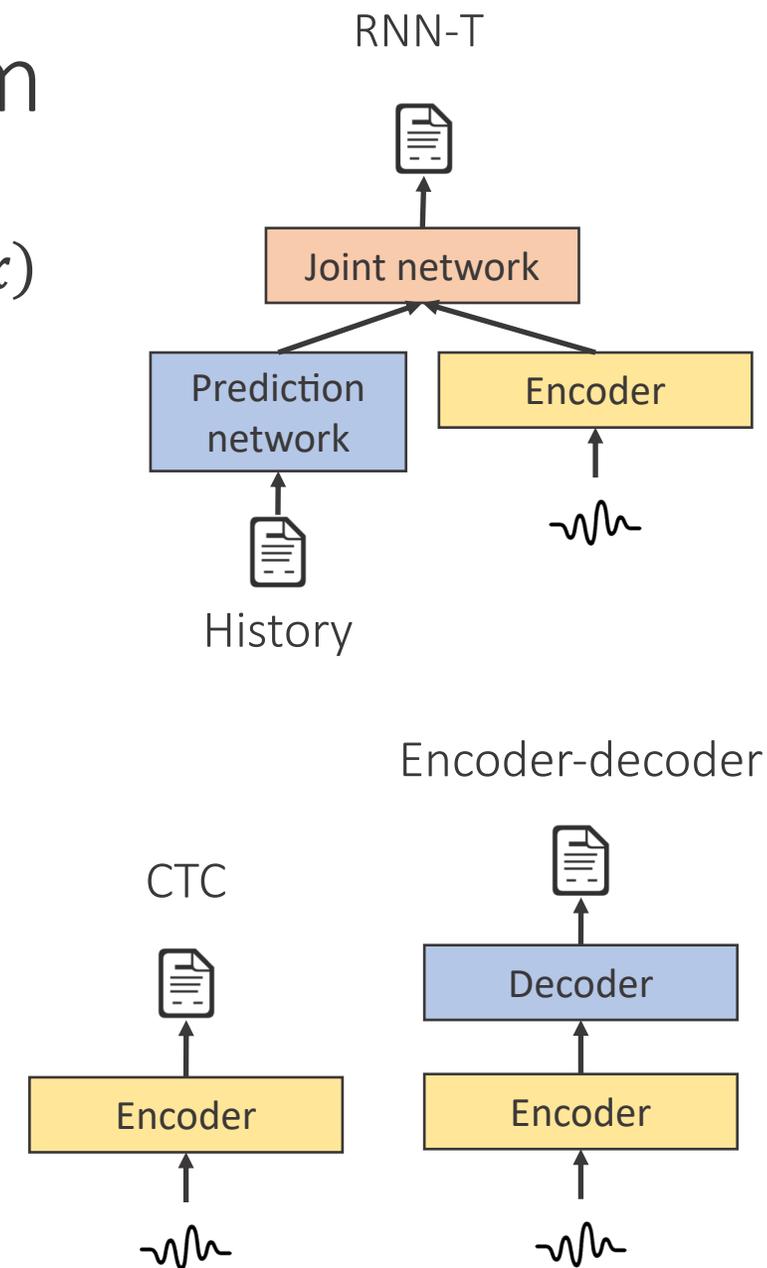
😓 Rare words, low-resource, customization

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

- *Connectionist temporal classification (CTC)* [Graves+ 2006]
- *RNN-Transducer (RNN-T)* [Graves+ 2013]
- *Recurrent neural aligner (RNA)* [Sak+ 2017]

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

- *Attention-based RNN encoder-decoder* [Bahdanau+ 2016]
- *Transformer* [Vaswani+ 2017]



# Streaming ASR

- Transcribe speech before a speaker finalizes their turn
- Applications
  - ✓ Live captioning
  - ✓ Dialogue system
  - ✓ Simultaneous translation
- RNN-T is dominant in the industry
  - 😊 Stable inference thanks to frame-wise prediction
  - 😓 Memory-consuming training (-> small mini-batch size)
    - ✓ Distributed training (a log of GPUs)
    - ✓ Efficient implementation (not publicly available in general)
    - ✓ Small vocabulary size etc. are required
  - 😓 Large search space due to frame-wise predictions (slow inference)

# Challenges in label-synchronous streaming ASR

- Why label-sync. models instead of RNN-T?
  - Small memory consumption
  - Small search space (fast inference)
- Challenges in label-sync. streaming models
  1. Need to modify the decoding scheme
    - The whole encoder outputs are required to generate the first token in general seq2seq models
  2. Poor performance for long-form speech
    - Exposure bias (not occur in frame-synchronous models such as RNN-T)

# Streaming attention-based encoder-decoder models

Learn when to generate the next token (segment audio) *on the encoder side*

Flexible audio segmentation policy

◆ **Triggered attention** [Moritz+ 2018]

Global attention over past encoder outputs truncated by CTC spikes

◆ **Adaptive computation steps (ACS)** [Li+ 2018]

Learn how many tokens to generate with encoder outputs (halting mechanism)

◆ **Continuous Integrate-and-Fire (CIF)** [Dong+ 2019]

Fine-grained version of ACS

Fixed audio segmentation policy

◆ **Neural Transducer** [Jailty+ 2015]

Attention mechanism for a fixed size of block

- Simple framework
- Good results
- Efficient training
- Linear time inference

Learn when to generate the next token (segment audio) *on the decoder side*

◆ **Monotonic chunkwise attention (MoChA)** [Raffel+ 2017], [Chiu+ 2018]

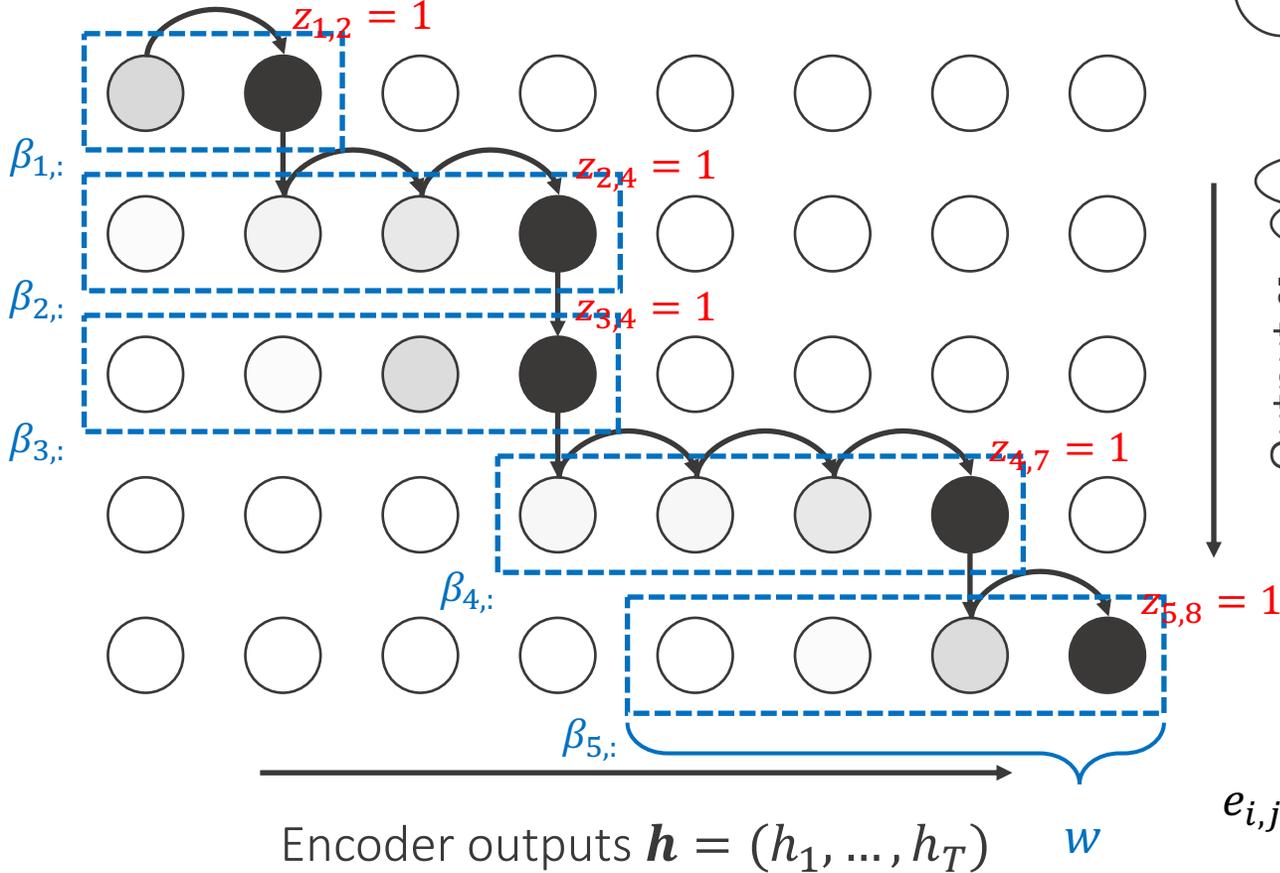
Learn to detect token boundaries via stochastic binary decision

And more...

- Windowing approaches
- Incremental decoding
- Reinforcement learning

# MoChA (test time) [Chiu+ 2018]

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



: Attend ( $z_{i,j} = 1$ )  
 : Not attend ( $z_{i,j} = 0$ )

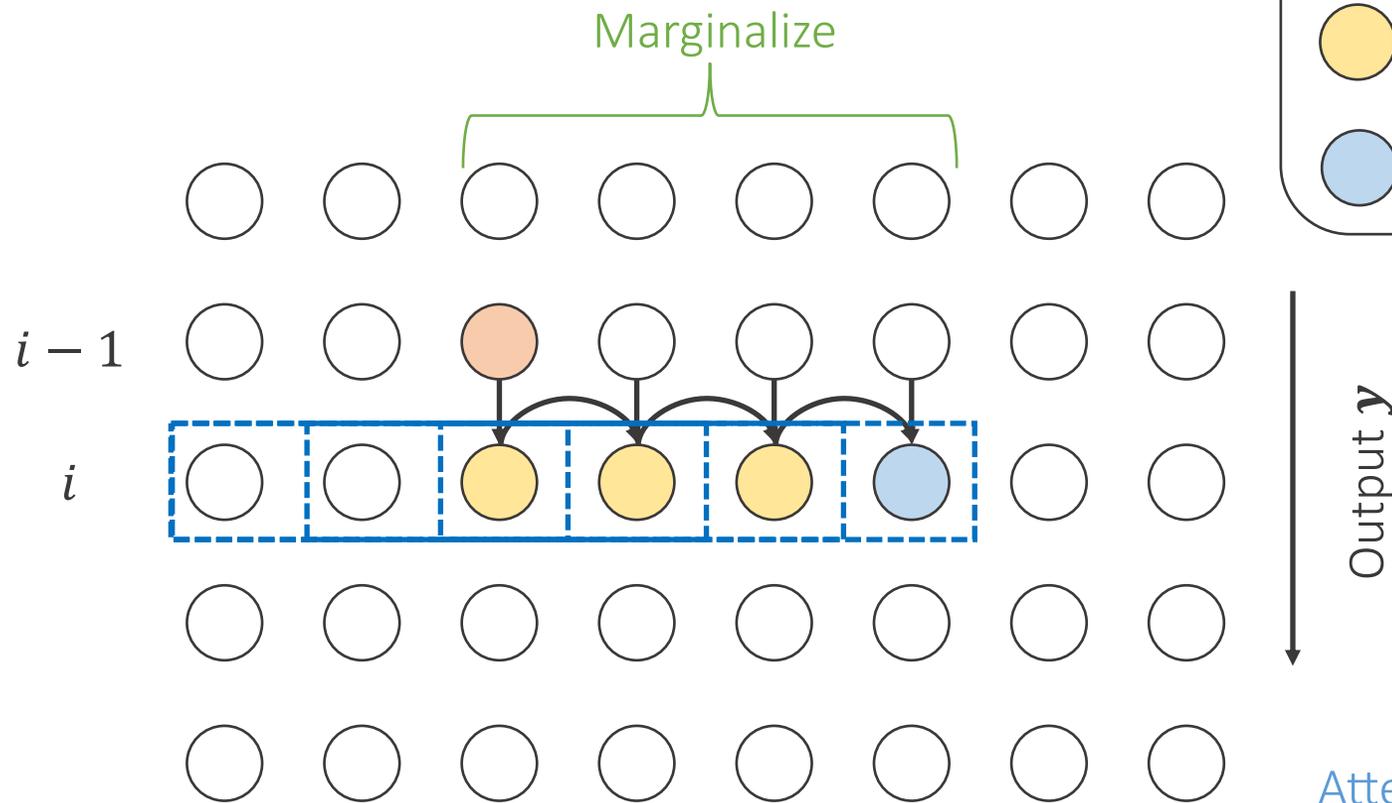
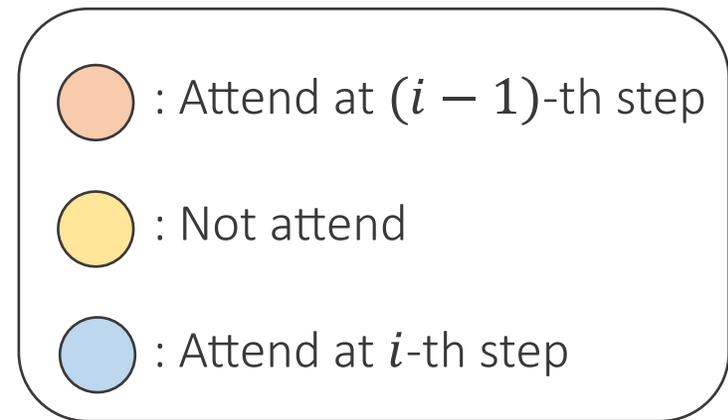
☹️ Not differentiable

$h_j$ : encoder state  
 $s_i$ : decoder state

$e_{i,j} = \text{MonotonicEnergy}(h_j, s_i)$   
 $p_{i,j} = \sigma(e_{i,j})$  (selection probability)  
 $z_{i,j} \sim \text{Bernoulli}(p_{i,j})$

1. **Hard monotonic attention** [Raffel+ 2017]: whether to attend or not
2. **Chunkwise attention**: soft attention over a small window of size  $w$

# MoChA (training time) [Chiu+ 2018]



Previous attention

Attend                      Not attend

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

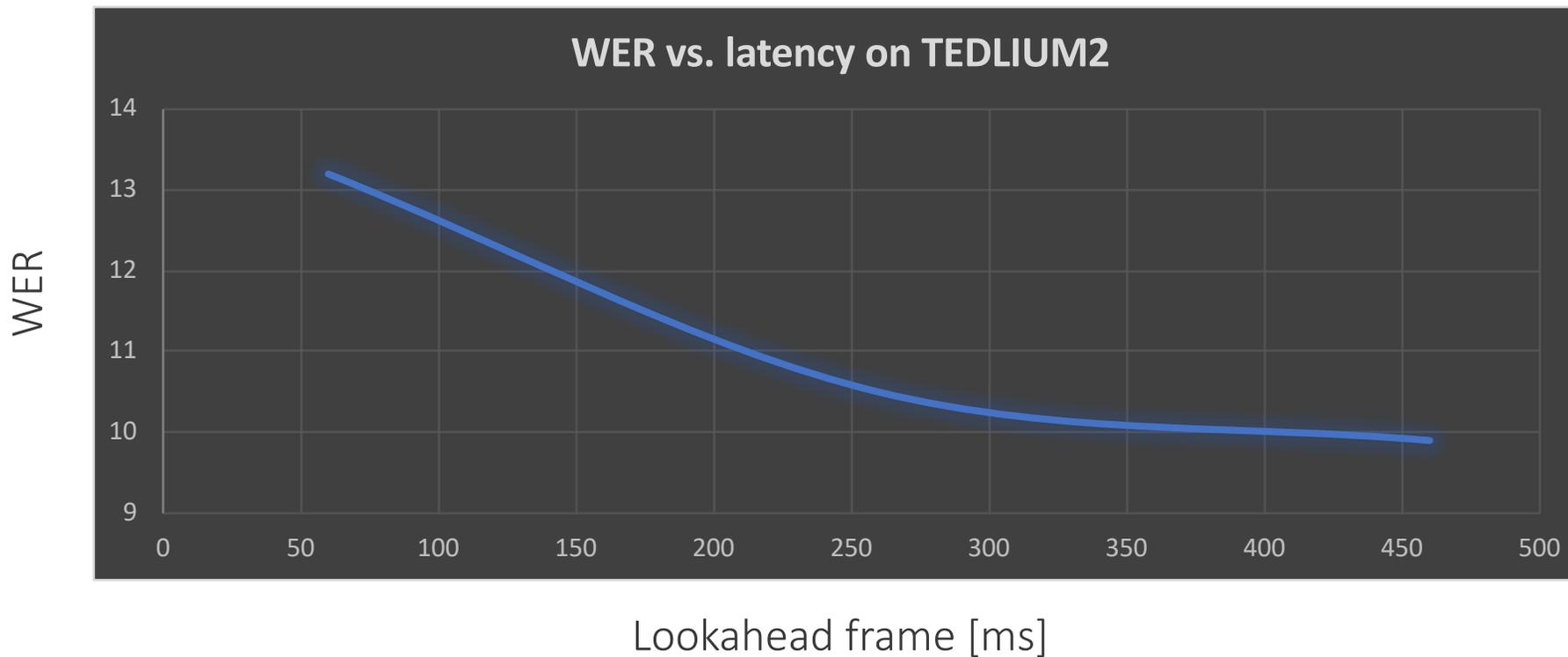
$$= (1 - p_{i,j-1}) \frac{\alpha_{i,j-1}}{p_{i,j-1}} + \alpha_{i-1,j}$$

Calculate expected alignments  $\alpha$

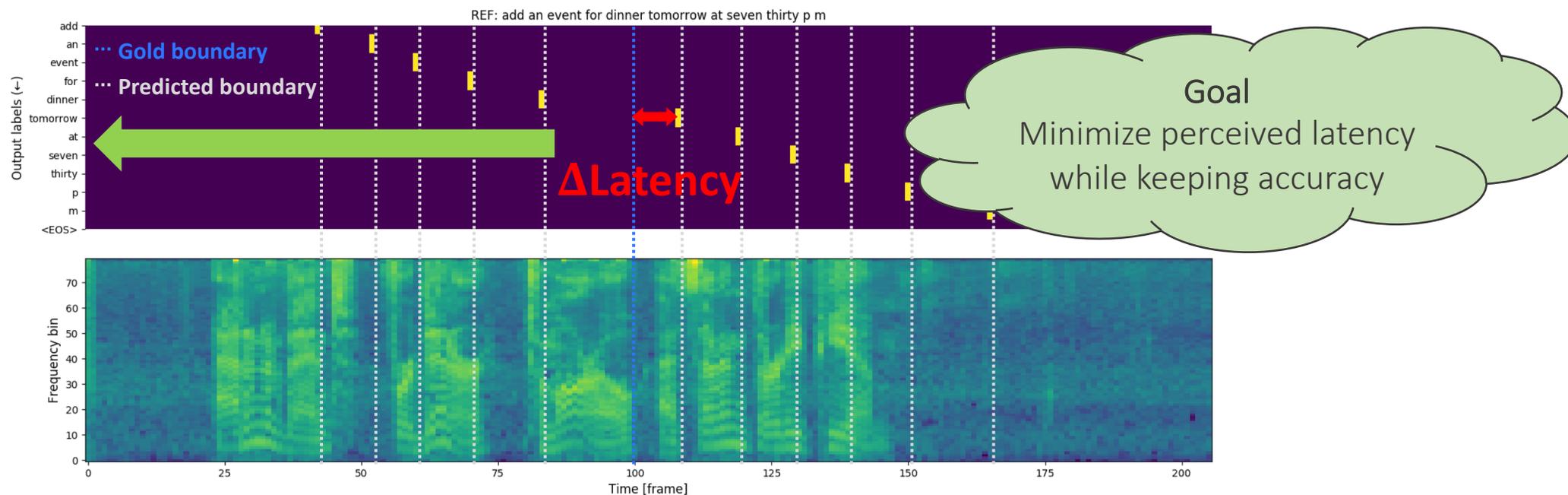
Can be implemented efficiently in parallel with  $j$

# Lookahead latency and accuracy trade-off in streaming ASR

- Future information (lookahead) is very important to improve accuracy
- Large lookahead leads to large **algorithmic** latency
  - Can be controlled on demand



# Delayed token generation problem

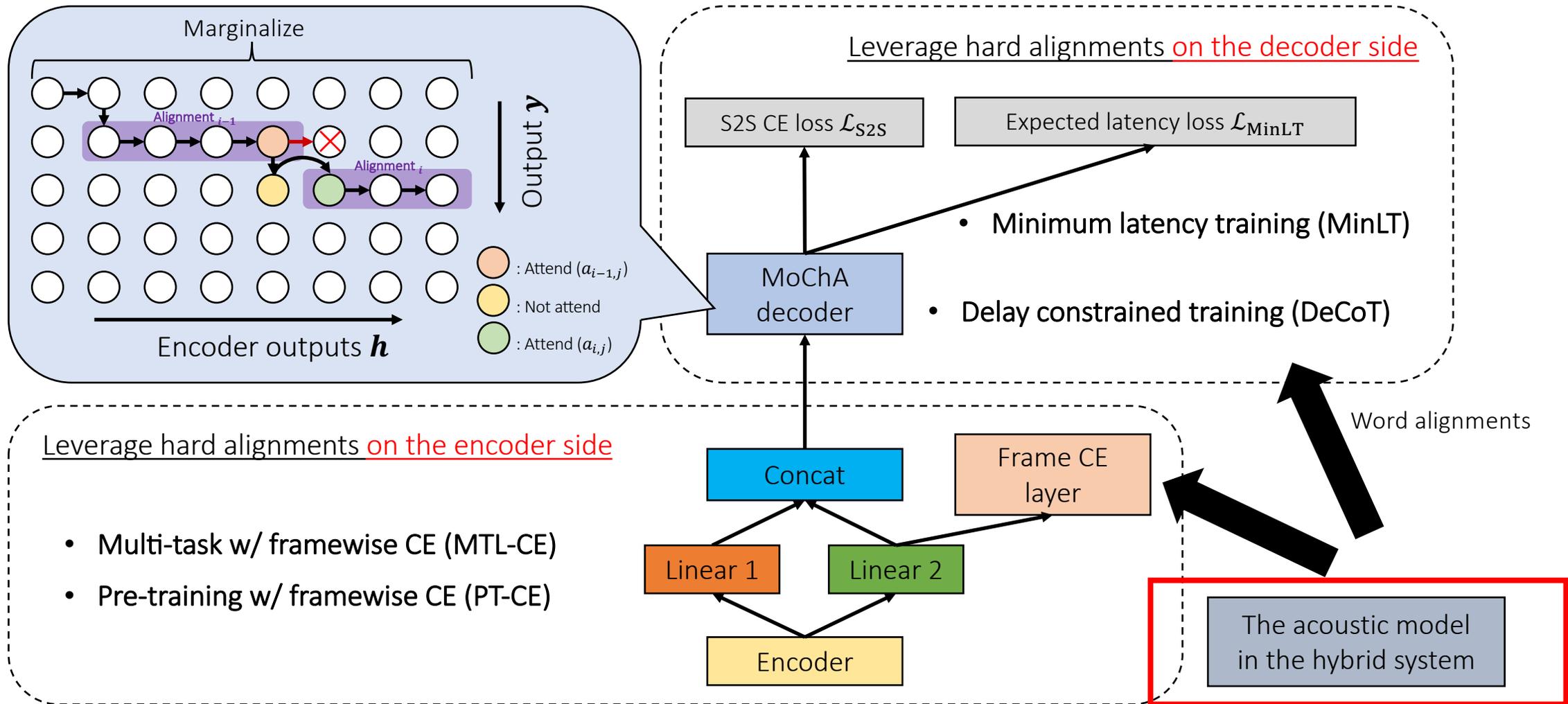


- Decision boundaries (yellow dots) are delayed from the corresponding acoustic boundary
  1. Unidirectional encoder (lacking the future information)
  2. Sequence-level criterion (utilizing as many future frames as possible to maximize the log-likelihood)
- Increase user perceived latency
  - Similar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

# Proposed methods

- Leverage external frame-level alignments extracted from the hybrid ASR system
- Investigate where to apply alignment information to streaming encoder-decoder model
  - Encoder side
    1. Multi-task learning with frame-wise CE
    2. Pre-training with frame-wise CE
  - Decoder side
    3. Delay constrained training (DeCoT)
    4. Minimum latency training (MinLT)

# Overview



Leveraging word alignments extracted from the hybrid system

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

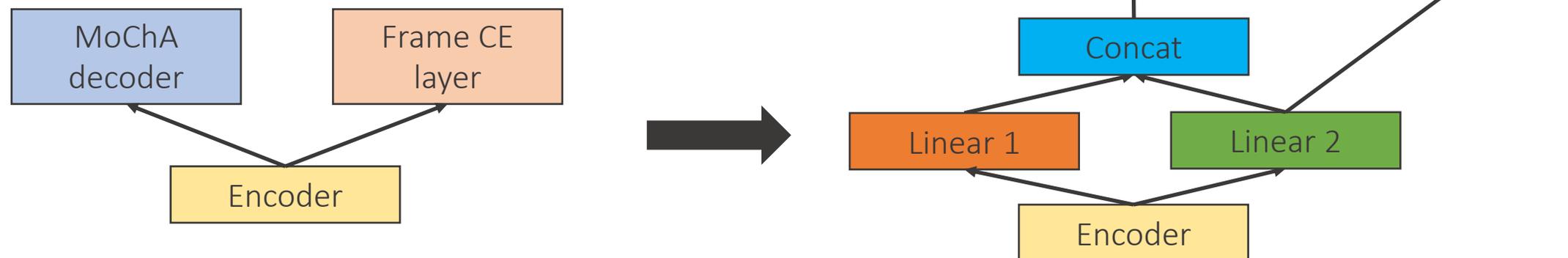
## ◆ Objective function

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

- Motivation: align encoder outputs to the true acoustic location

## ◆ Insert linear bottleneck layers

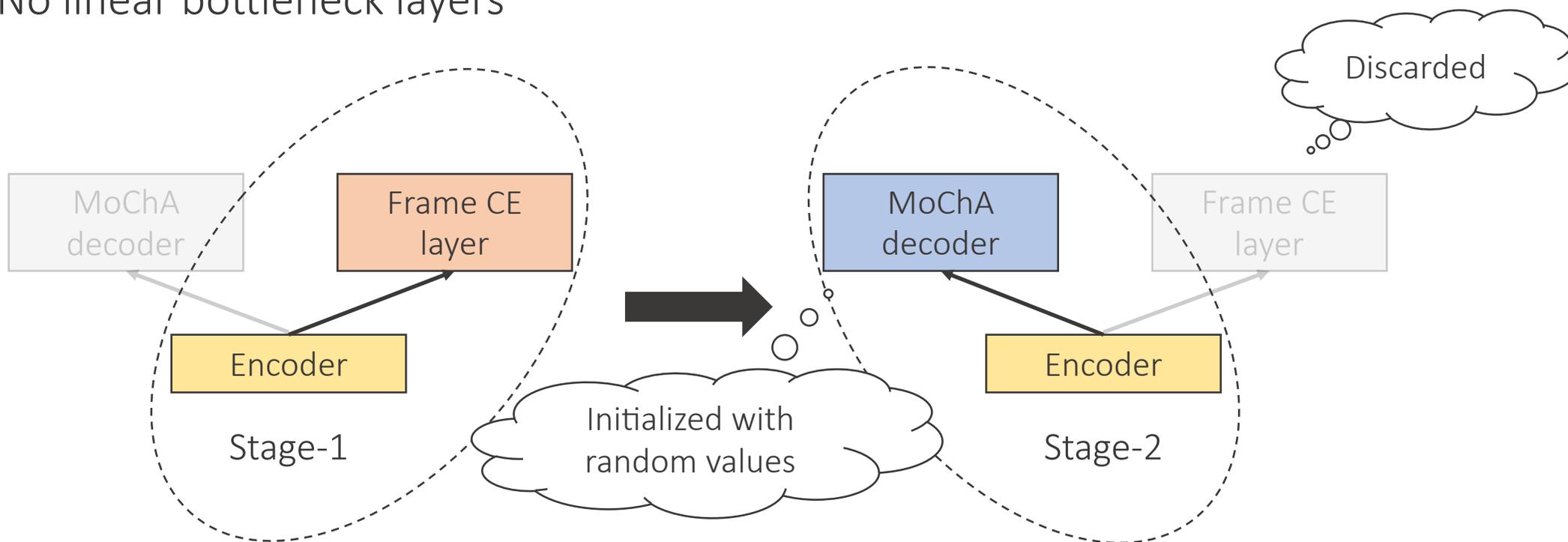
- Inspired by the CTC acoustic model [Yu+ 2018]



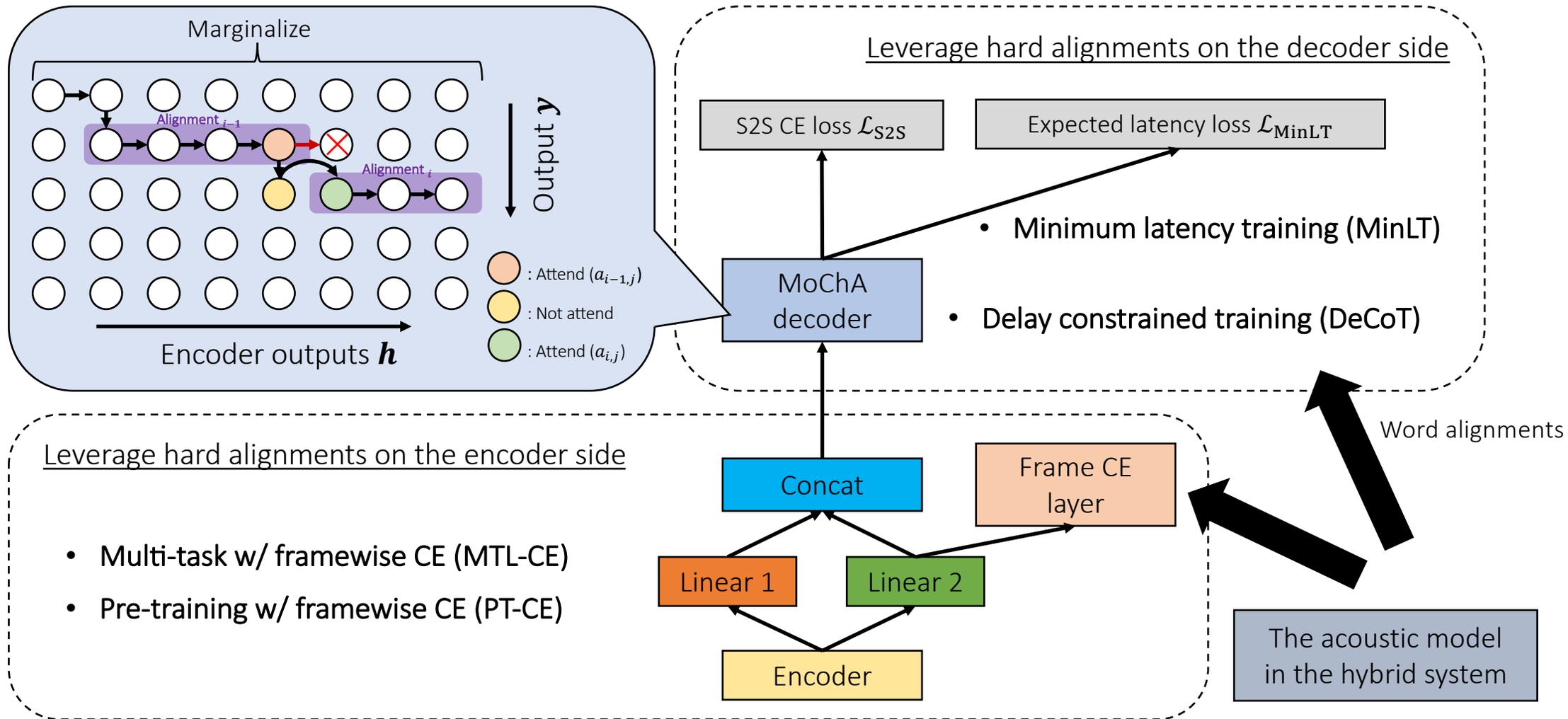
## 2. Pre-training with framewise CE (PT-CE)

### ◆ 2-staged training

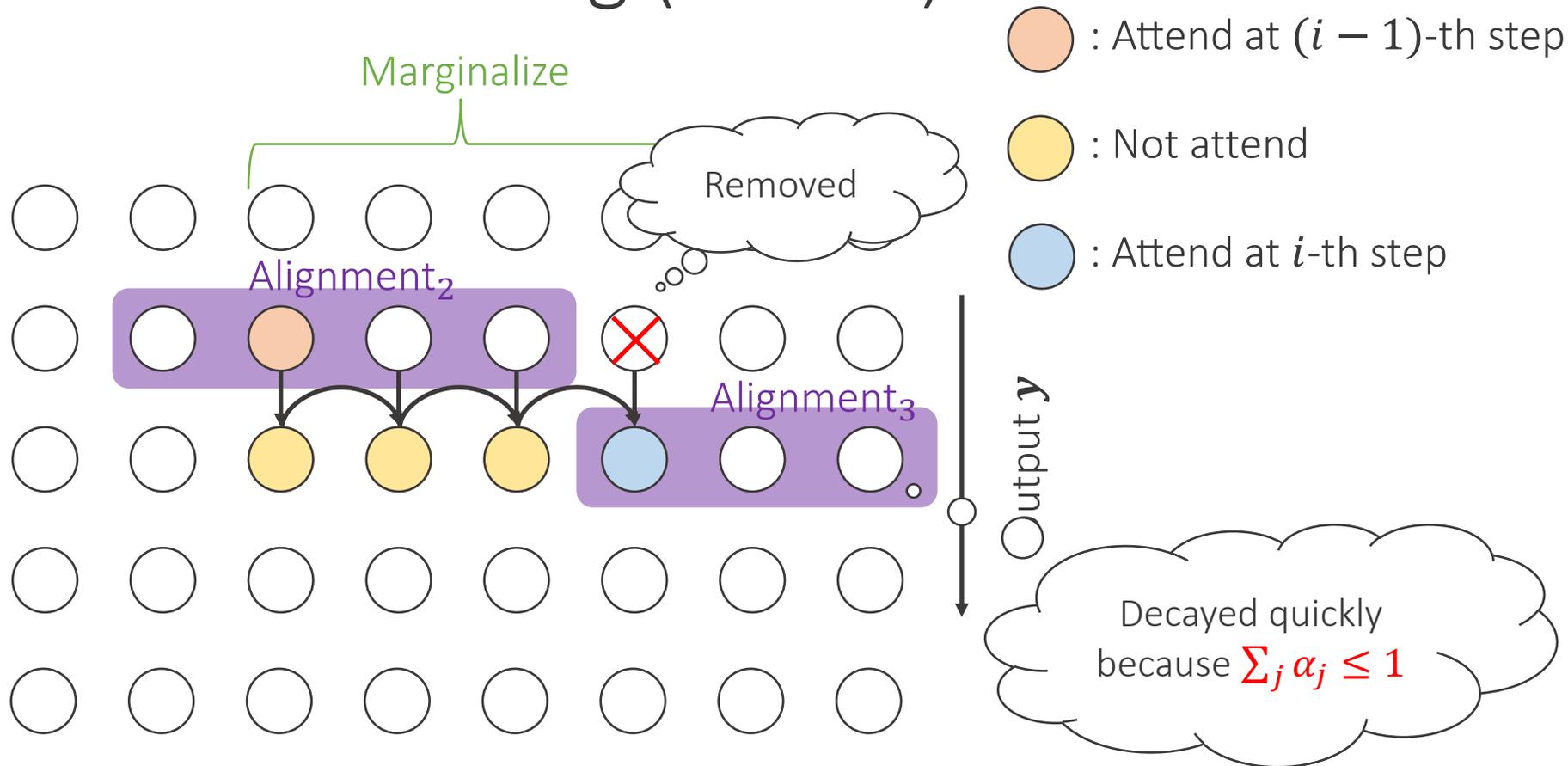
- Motivation
  - Start training from well-aligned encoder representations
  - Do not have to tune the framewise CE weight  $\lambda_{\text{CE}}$
- No linear bottleneck layers



# Overview



# 3. Delay constrained training (DeCoT)



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

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

$\mathbf{b}_i$ : gold boundary

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

### 3. Delay constrained training (DeCoT)

#### Quantity regularization

- Add a regularization term to keep  $\sum_j \alpha_{i,j} = 1$
- Originally proposed in CIF [Dong+ 2019] with a different motivation

$U$ : the number of tokens in the reference

$$\mathcal{L}_{\text{QUA}} = \left| U - \sum_{i=1}^U \sum_{j=1}^T \alpha_{i,j} \right| \quad (\text{quantity loss})$$

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

## 4. Minimum latency training (MinLT)

### ◆ Objective function

- Directly minimize the expected latency  $\mathcal{L}_{\text{MinLT}}$

Expected boundary

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

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

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

### ◆ Related work

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

# Experimental condition

Data	Train: Microsoft Cortana voice assistant (3.4k hours) Validation: Sampled disjoint 4k utterances form the training set Test: 5.6k utterances
Feature	80-dim log-mel fbank (3 frames stacked, 30ms per frame)
Output unit	Mixed units (34k)
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-> subword-level alignments
  - Divide duration per word by the ratio of the character length of each subword
- Warm start training
  - Start DeCoT and MinLT from the baseline MoChA to stabilize training

# Evaluation metric: Token emission latency

- Averaged time difference between a predicted boundary  $\widehat{b}_i^k$  and the gold boundary  $b_i^k$

*Corpus-level* latency (averaged per token)

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

- Report 50-th ( $TEL@50$ ) and 90-th percentile ( $TEL@90$ )
- Perform teacher-forcing when calculating latency to match the sequence lengths

# Results: Alignments on the encoder side

Model	WER [%] (↓)	Corpus-level latency [ms] (↓)	
		TEL@50	TEL@90
Baseline MoChA	9.93	300	642
+ MTL-CE ( $\lambda_{CE} = 0.1$ )	10.21 <small>5.6%</small>	240 <small>40%</small>	583
+ MTL-CE ( $\lambda_{CE} = 0.3$ )	10.48	180	591
+ MTL-CE ( $\lambda_{CE} = 0.5$ )	11.11	150	637
+ PT-CE	12.74	210	687

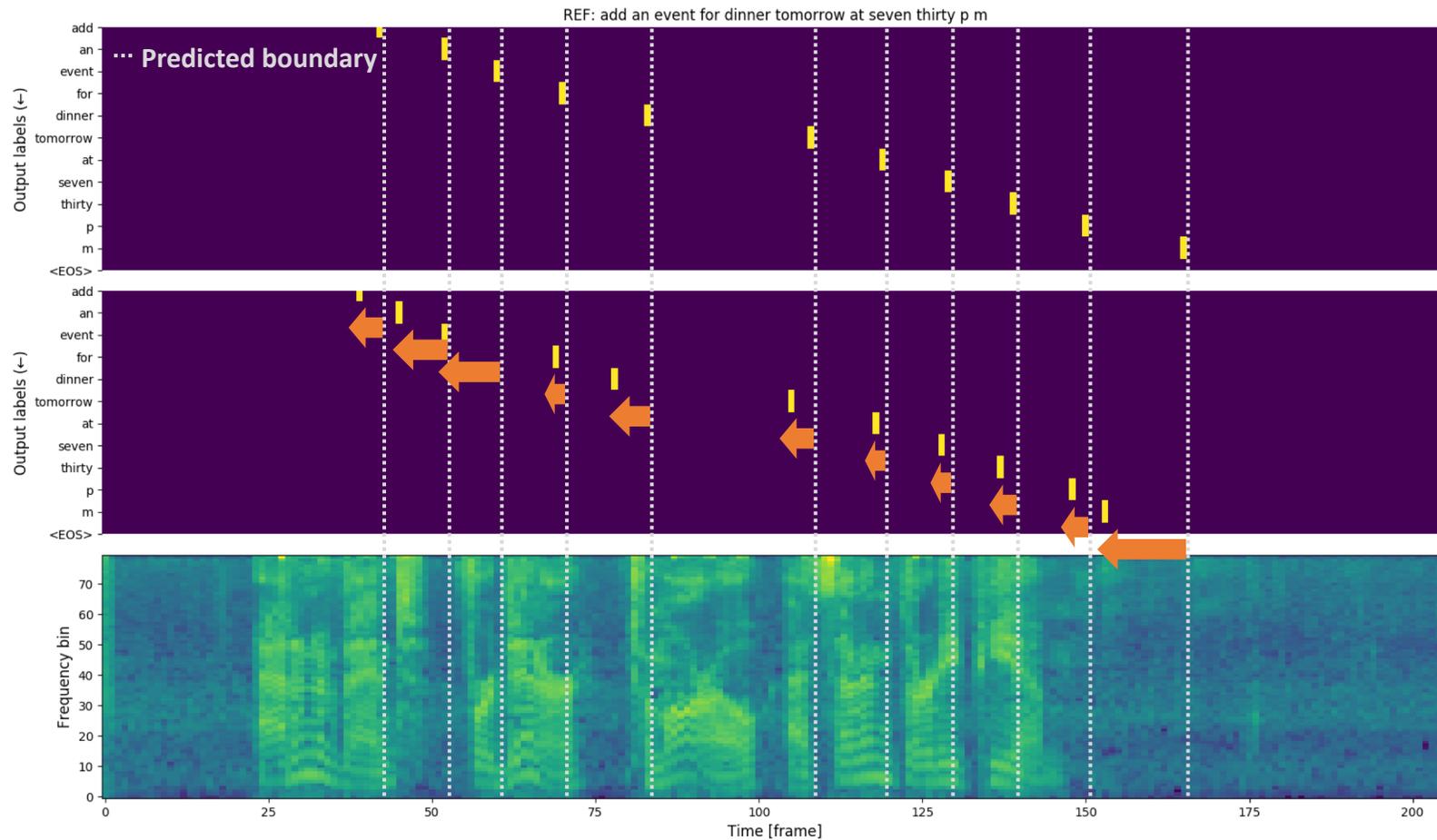
- MTL-CE reduced latency in proportion to  $\lambda_{CE}$  while degrading WER slightly
- PT-CE also reduced latency but degraded WER too much
- Contrastive results to previous works using CTC + framewise CE objective
  - MoChA is a label-synchronous model
  - Frame-wise CE on the encoder is not compatible with label-wise CE on the decoder

# Results: Alignments on the decoder side

Model	WER [%] (↓)	Corpus-level latency [ms] (↓)	
		TEL@50	TEL@90
Global attention (offline)	8.44	N/A	N/A
Baseline MoChA	9.93	300	642
+ DeCoT ( $\delta = 4$ , 120ms)	20.25	30	287
+ DeCoT ( $\delta = 8$ , 240ms)	14.35	150	210
+ DeCoT ( $\delta = 12$ , 360ms)	11.40	210	298
+ DeCoT ( $\delta = 16$ , 480ms)	9.13	240	352
+ DeCoT ( $\delta = 24$ , 720ms)	8.87	270	434
+ DeCoT ( $\delta = 32$ , 960ms)	9.17	300	497
+ MinLT	9.70	180	319
+ DeCoT ( $\delta = 16$ )	12.75	120	239

- DeCoT: **large WER reduction** and **moderate latency reduction (tail part)**
- MinLT: **small WER reduction** and **large latency reduction (entire)**
- Combination of DeCoT and MinLT reduced latency further, but degraded WER too much

# Alignment visualization



Baseline MoChA

DeCoT ( $\delta = 16$ )

# Summary: alignment information from hybrid ASR

- Alignment information is beneficial when applying it **on the decoder side**
  - 🥵 This is NOT purely end-to-end
- Can we remove the dependency to hybrid ASR system for alignment extraction?
  - **CTC alignment**

# Optimization problem

Recap

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

$$p_{i,j} = \sigma(e_{i,j})$$

## 1. $\sum_j \alpha_{i,j} = 1$ is not satisfied during training

- $\alpha_{i,j}$  is NOT globally normalized over the whole encoder outputs  $\{h_j\}_{j=1,\dots,T}$ 
  - $\alpha_{i,j}$  is not a valid probability distribution
  - $\alpha_{i,j}$  attenuates quickly during marginalization
  - Selection probability  $p_{i,j}$  is not learnt well
- Enlarge the mismatch between training and test time

## 2. Alignment errors are propagated to later token generation

- $\alpha_{i,j}$  depends on past alignments
- Backward algorithm cannot be used for  $\alpha_{i,j}$ 
  - $\alpha_{i,j}$  is not a valid probability distribution
  - Autoregressive decoder
- Model needs to learn (1) a proper scale of  $\alpha_{i,j}$  and (2) accurate decision boundaries ( $j$  s. t.  $\alpha_{i,j} = 1$ ) at the same time

Problematic for long and noisy speech utterances

# Related work: Joint CTC-attention [Kim+ 2017]

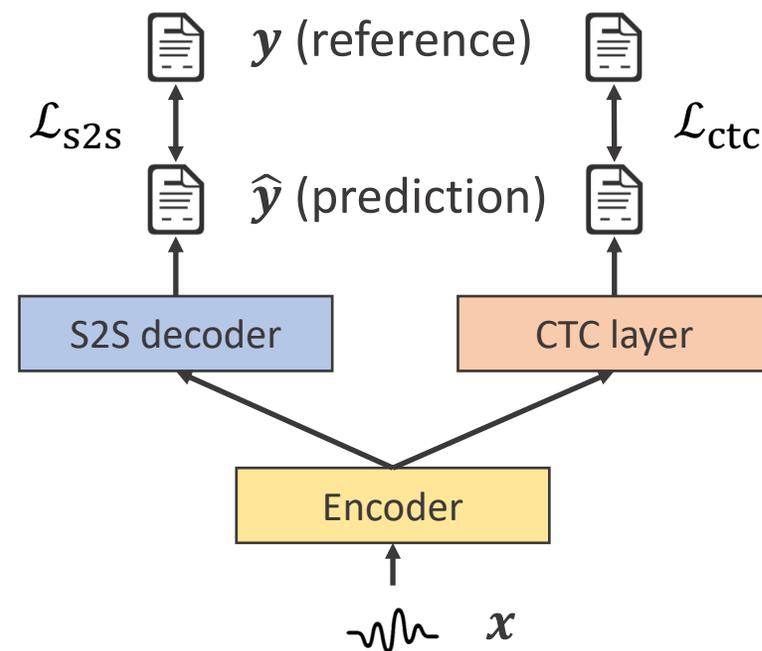
- Auxiliary CTC loss encourages the monotonicity between input and output alignments

## Objective function of encoder-decoder model

$$\mathcal{L}_{s2s} = -\log P(y|x) = -\sum_{i=1}^U \log P(y_i | y_{<i}, x)$$

## Multitask learning with CTC objective

$$\mathcal{L}_{\text{total}} = (1 - \lambda_{\text{ctc}})\mathcal{L}_{s2s} + \lambda_{\text{ctc}} \underbrace{\mathcal{L}_{\text{ctc}}}_{\text{CTC loss}} \quad (0 \leq \lambda_{\text{ctc}} \leq 1)$$

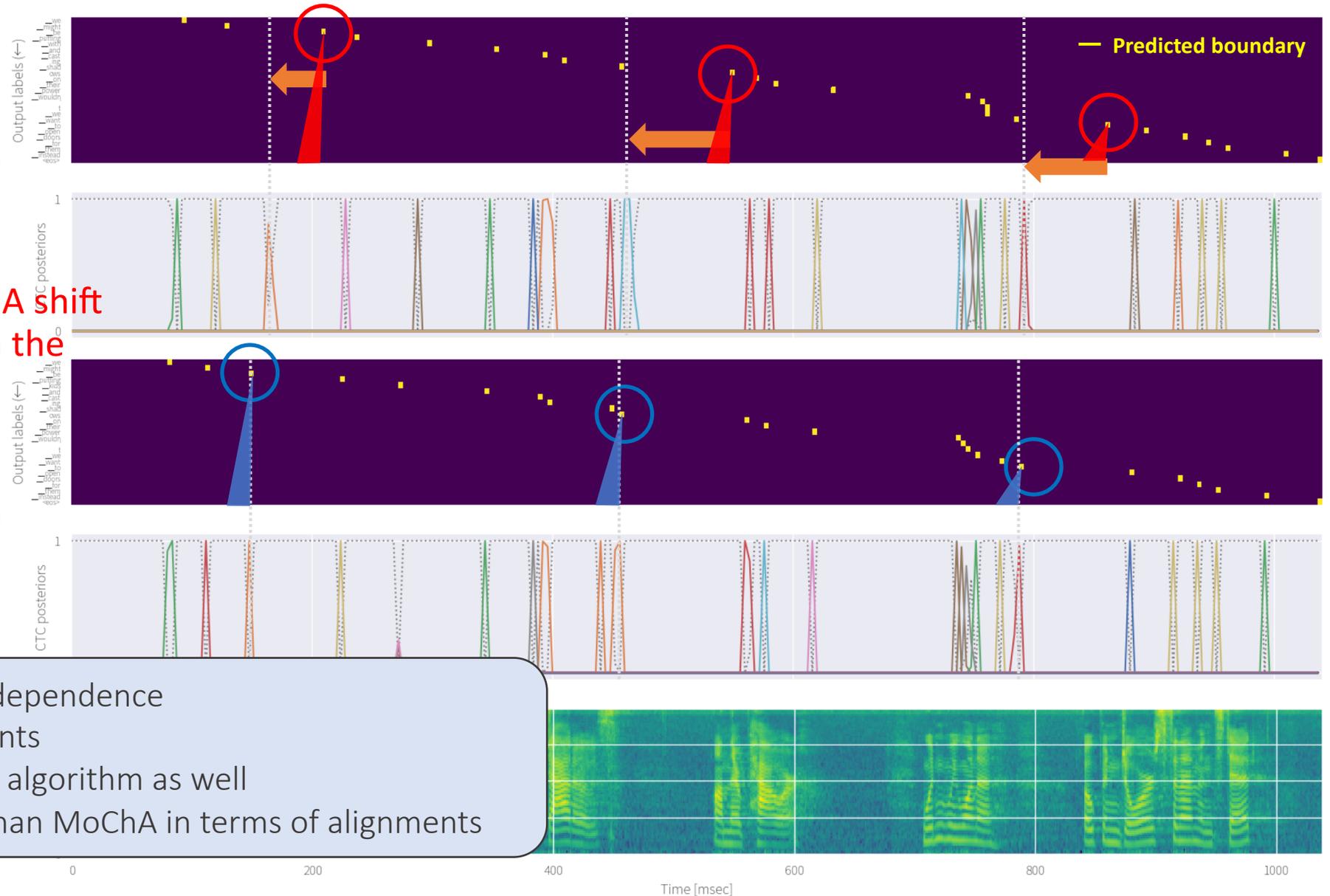


# Comparison of boundary positions: CTC vs. MoChA

Baseline

Decision boundaries of MoChA shift to the right side (future) from the corresponding CTC spikes

Proposed



- CTC assumes conditional independence
  - Robust to past alignments
- CTC leverages the backward algorithm as well
  - CTC is more accurate than MoChA in terms of alignments

# Proposed method: CTC-synchronous training (CTC-ST)

- Leverage CTC's posterior spikes as reference boundaries for MoChA
- MoChA is trained to mimic the CTC model to generate the similar decision boundaries

## Objective function

$$\mathcal{L}_{\text{sync}} = \frac{1}{U} \sum_{i=1}^U \left| \mathbf{b}_i^{\text{ctc}} - \sum_{j=1}^T j \alpha_{i,j} \right|$$

CTC boundary    Expected MoChA boundary

$$\mathcal{L}_{\text{qua}} = \left| U - \sum_{i=1}^U \sum_{j=1}^T \alpha_{i,j} \right|$$

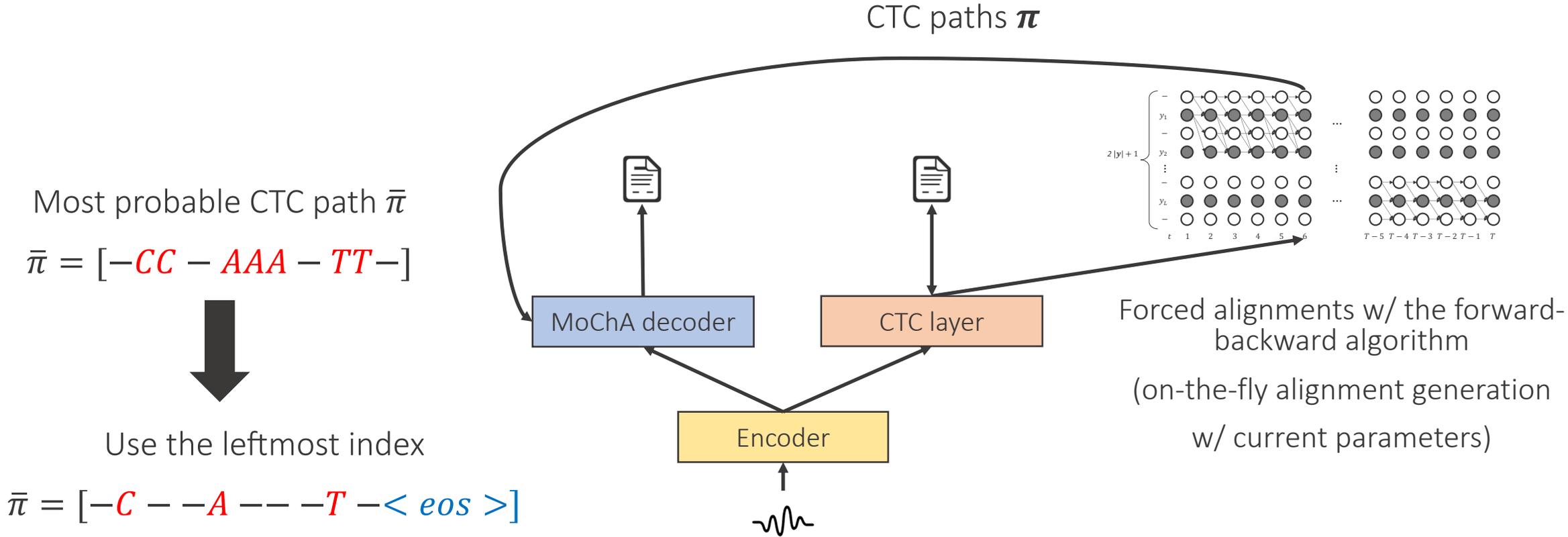
Important regularization  
for baseline model

$$\mathcal{L}_{\text{total}} = (1 - \lambda_{\text{ctc}}) \mathcal{L}_{\text{mocha}} + \lambda_{\text{ctc}} \mathcal{L}_{\text{ctc}} + \lambda_{\text{qua}} \mathcal{L}_{\text{qua}} + \lambda_{\text{sync}} \mathcal{L}_{\text{sync}} \quad (\lambda_{\text{sync}} \geq 0)$$

- Unless otherwise noted,  $\lambda_{\text{qua}}$  is set to 0 when using CTC-ST

# Extraction of CTC alignments

- Encoder network is shared between both branches
- Both branches are jointly optimized
- CTC alignments are extracted via forced alignment over the transcription



# Curriculum learning strategy

- Applying CTC-ST from scratch is inefficient because  $\sum_{j=1}^T \alpha_{ij} \ll 1$  in the early training stage
  - Difficult to estimate the expected boundaries  $\sum_{j=1}^T j \alpha_{i,j}$  accurately
  - Propose curriculum learning strategy composed of two stages

## Stage-1 (expected to learn a proper scale of $\alpha_{ij}$ )

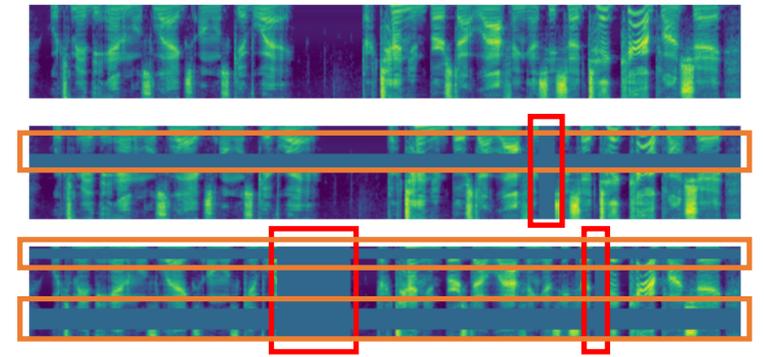
- Train **BLSTM encoder + MoChA** with quantity regularization until convergence

## Stage-2 (expected to learn boundary location)

- Initialize with model parameters in stage-1
- Train **latency-controlled BLSTM (LC-BLSTM) encoder + MoChA** with CTC-ST

NOTE: When using the unidirectional LSTM encoder, the same encoder is used in both stages

# Combination with SpecAugment



## SpecAugment [Park+ 2019]

- On-the-fly data augmentation method over input log-mel filterbank features
- Zero out successive frames in **time** and **frequency** bins

## Problem of SpecAugment for MoChA

- Recurrency of  $\alpha_{i,j}$  can be easily collapsed after the masked region
- The naïve MoChA did not obtain any gains with SpecAugment
- CTC can estimate boundaries accurately even right after the masked region thanks to the conditional independence assumption per frame
- CTC-ST is expected to improve the effectiveness of SpecAugment for MoChA

Recap

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

# Experimental condition

Corpus	TEDLUM2 (210h, lecture), Librispeech (960h, read)
Feature	80-dim log-mel fbank
Output unit	BPE10k units
Architecture	Offline: 4-layer CNN -> 512-dim (per direction) 5-layer BLSTM encoder Streaming: 4-layer CNN -> 512-dim 5-layer LC-BLSTM encoder or 4-layer CNN -> 1024-dim 5-layer unidirectional LSTM encoder
	Decoder: 1024-dim 1-layer LSTM $w$ : 4 (window size for chunkwise attention in MoChA)
Optimization	Adam
Loss weight	$\lambda_{\text{ctc}} = 0.3, \lambda_{\text{qua}} = 1.0, \lambda_{\text{sync}} = 1.0$
Decoding	Beam width: 10, shallow fusion with external 4-layers of LSTM-LM

# Main results: TEDLIUM2 (210h)

Latency-controlled BLSTM  
 LC-BLSTM- $N_l + N_r$   
 hop size    lookahead frame  
 (ms)            (ms)

	Model	WER [%]	
Offline	BLSTM - Global attention (T1)	9.5	
	BLSTM - MoChA	12.6	22.2% (↑)
	+ Quantity regularization (T2)	9.8	
	+ CTC-ST	10.2	
Streaming	UniLSTM - MoChA	15.0	12.0% (↑)
	+ CTC-ST	13.2	
	LC-BLSTM-40+20 - MoChA	12.2	13.9% (↑)
	+ CTCT-ST	10.5	
	LC-BLSTM-40+40 - MoChA (T5)	11.3	12.3% (↑)
+ CTC-ST (T6)	9.9		
	+ Quantity regularization	10.1	

- Combination of CTC-ST and quantity regularization was not effective
  - CTC-ST has a similar effect to improve the scale of  $\alpha_{ij}$
- Curriculum learning was effective

# Results of curriculum learning

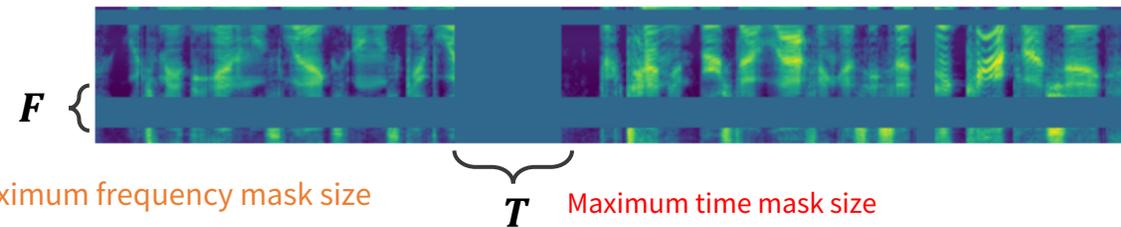
Quantity regularization      CTC-ST

$$\mathcal{L}_{\text{total}} = (1 - \lambda_{\text{ctc}})\mathcal{L}_{\text{mocha}} + \lambda_{\text{ctc}}\mathcal{L}_{\text{ctc}} + \lambda_{\text{qua}}\mathcal{L}_{\text{qua}} + \lambda_{\text{sync}}\mathcal{L}_{\text{sync}}$$

Model	Quantity regularization	CTC-ST	WER [%]
LC-BLSTM-40+40 - MoChA (from scratch)	✓	-	12.3
	-	✓	10.9
LC-BLSTM-40+40 - MoChA (from BLSTM - MoChA)	-	-	16.9
	✓	-	11.3
	-	✓	9.9
	✓	✓	10.1

- Seeding by BLSTM- MoChA was effective
- Combination of CTC-ST and quantity regularization was not effective
  - CTC-ST has a similar effect to improve the scale of  $\alpha_{ij}$
- Curriculum learning was effective
  - Quantity regularization (stage-1)-> CTC-ST (stage-2)

# Results with SpecAugment

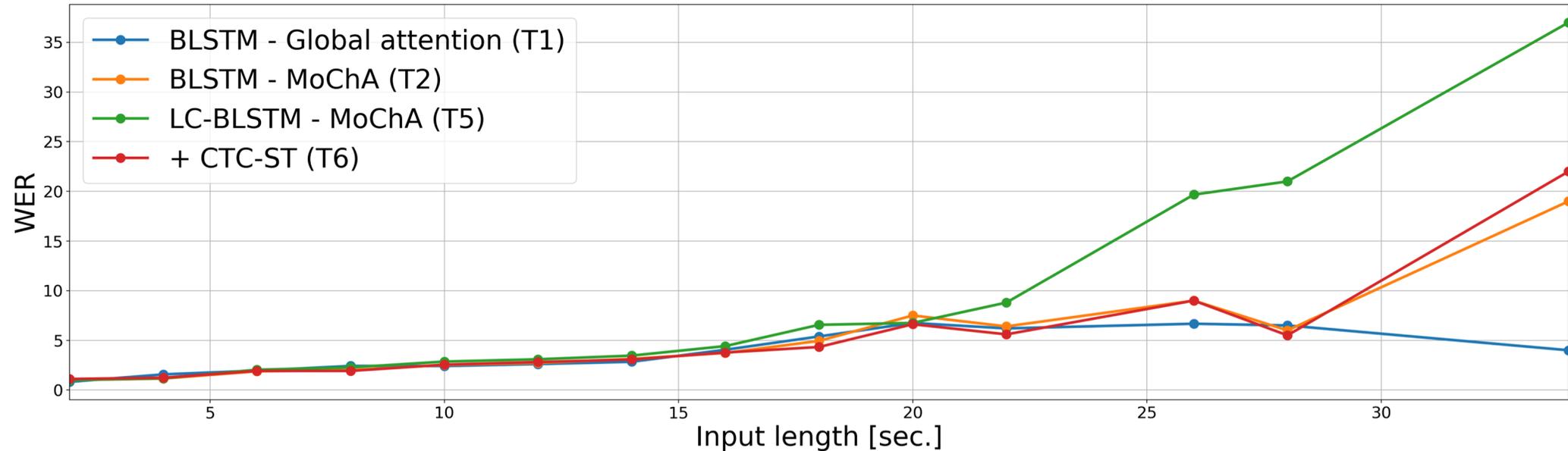


	Model	$F$	$T$	WER [%]
Offline	Transformer [Karita+ 2019]	30	40	8.1
	BLSTM - Global attention [Zeyer+ 2019]	N/A	N/A	8.8
	BLSTM - Global attention	-	-	9.5
Streaming		27	100	<b>8.1</b>
		-	-	11.3
	LC-BLSTM-40-+40 - MoChA	27	100	12.8
	(seed: BLSTM - MoChA)	27	50	11.0
		13	50	11.2
		-	-	9.9
	+ CTC-ST	27	100	9.0
	27	50	<b>8.6</b>	
	13	50	9.0	

13.1% (↑)

- MoChA did not benefit from SpecAugment w/o CTC-ST
- CTC-ST was robust to the input mask size
- Achieved the comparable performance to the offline model (8.1 vs. 8.6)

# WER vs. input sequence length



- CTC-ST improved WER for long-form utterances

# Results on Librispeech (960h)

		Model	WER [%]		
			Test-clean	Test-other	
Offline		BLSTM - global attention	3.1	9.5	
		+ SpecAugment ( $F = 27, T = 100$ )	2.8	7.6	
		BLSTM - MoChA	3.6	10.5	
		+ Quantity regularization (T2)	<b>3.3</b>	<b>10.0</b>	8.3/4.7% (↑)
Streaming	Initialization	UniLSTM - MoChA	5.3	14.5	
		+ CTC-ST	<b>4.7</b>	<b>13.6</b>	11.3/6.2% (↑)
		+ SpecAugment	<b>4.2</b>	<b>11.2</b>	
		LC-BLSTM-40+40 - MoChA	4.1	11.2	
		+ SpecAugment ( $F = 13, T = 50$ )	4.0	9.5	
		+ SpecAugment ( $F = 27, T = 50$ )	4.8	9.3	
		+ SpecAugment ( $F = 27, T = 100$ )	5.0	9.7	
		+ CTC-ST	3.9	11.2	
		+ SpecAugment ( $F = 13, T = 50$ )	3.6	9.4	
		+ SpecAugment ( $F = 27, T = 50$ )	<b>3.5</b>	<b>9.1</b>	
	+ SpecAugment ( $F = 27, T = 100$ )	3.6	9.2	10.2/18.7% (↑)	

# Comparison with previous works on Librispeech

Model	WER [%]	
	Test-clean	Test-other
LSTM - MoChA + MWER [Kim+ 2019]	5.6	15.6
LSTM - MoChA + {BPE, char}-CTC + SpecAugment [Garg+ 2019]	4.4	15.2
LSTM - MoChA + CTC-ST + SpecAugment (ours)	<b>4.2</b>	<b>11.2</b>
LC-BLSTM - sMoChA [Miao+ 2019]	6.0	16.7
LC-BLSTM - MTA [Miao+ 2020]	4.2	12.3
LC-BLSTM - MoChA + CTC-ST (ours)	3.9	11.2
+ SpecAugment	<b>3.5</b>	<b>9.1</b>

# Hybrid ASR alignment vs. CTC alignment (TEDLIUM2)

† SpecAugment is used

Alignment	Model	WER [%] (↓)	Corpus-level latency [ms] (↓)	
			TEL@50	TEL@90
-	UniLSTM MoChA	15.0	280	680
CTC	+ CTC-ST	13.2	160	360
	+ CTC-ST †	<b>11.6</b>	<b>200</b>	<b>360</b>
Hybrid ASR	+ DeCoT ( $\delta = 12, 480\text{ms}$ ) †	11.2	200	320
	+ DeCoT ( $\delta = 16, 640\text{ms}$ ) †	<b>11.0</b>	<b>280</b>	<b>440</b>
	+ DeCoT ( $\delta = 20, 800\text{ms}$ ) †	11.3	240	400
	+ DeCoT ( $\delta = 24, 960\text{ms}$ ) †	11.7	280	480
	+ MinLT †	11.7	240	360

- CTC-ST not only improves WER but also reduces token emission latency
- CTC-ST is as good as DeCoT/MinLT for latency reduction **w/o external alignment**

# Hybrid ASR alignment vs. CTC alignment (Librispeech)

† SpecAugment is used

Alignment	Model	WER [%] (↓)		Corpus-level latency [ms] (↓)	
		test-clean	test-other	TEL@50	TEL@90
-	UniLSTM MoChA	5.3	14.5	360	560
CTC	+ CTC-ST	4.7	13.6	240	400
	+ CTC-ST †	<b>4.2</b>	<b>11.2</b>	<b>280</b>	<b>400</b>
Hybrid ASR	+ DeCoT ( $\delta = 16, 640\text{ms}$ ) †	<b>4.3</b>	<b>11.5</b>	320	440
	+ MinLT †	4.7	11.8	320	480

- When training data is large, CTC alignment is very accurate and reliable

# Non-autoregressive End-to-end Speech Translation

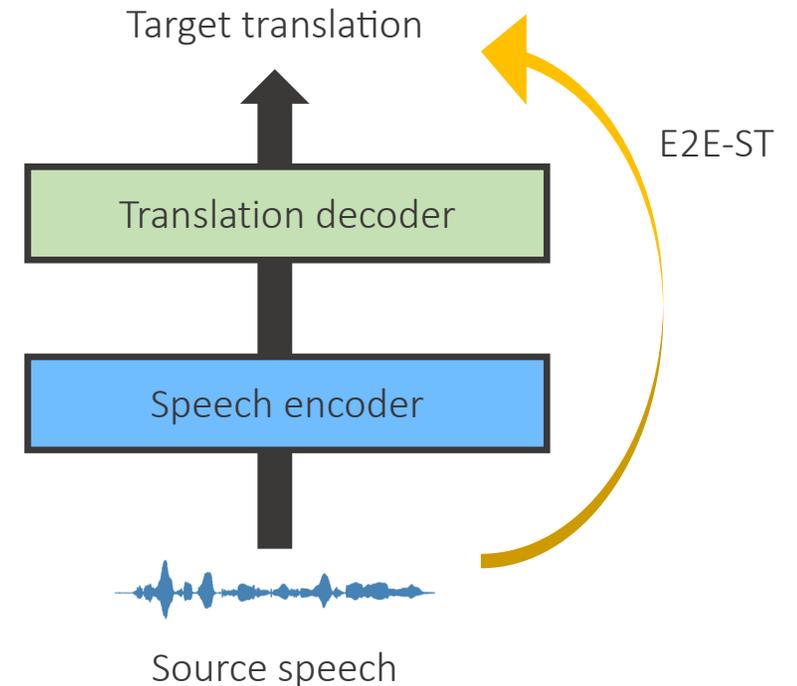
# Background: End-to-end speech translation (E2E-ST)

## Pros.

- Simplified architecture
- Avoid error propagation from ASR module
- **Low-latency inference**
- Endangered language documentation

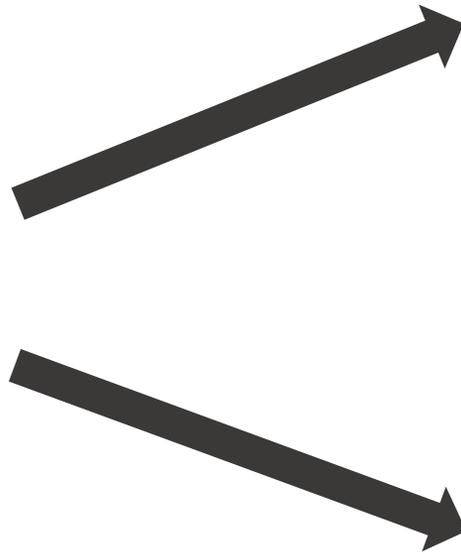
## Cons.

- Lack of supervised training data
- Most previous works focused on improving translation quality
- E2E-ST is conceptually suitable for fast decoding than cascaded systems
  - However, such evaluation has not been investigated so far



# Low-latency E2E-ST

Low-latency  
inference



Simultaneous, streaming

Decoding speed-up, offline  
(this work)

# Autoregressive (AR) sequence generation

## ◆ Notation

- $X = (x_1, \dots, x_U)$  (input speech)
- $Y = (y_1, \dots, y_N)$  (target translation)
- $Y^{\text{src}} = (y_1^{\text{src}}, \dots, y_{N_{\text{src}}}^{\text{src}})$  (source transcription)



English speech

Danke (German)

Thank you (English)

## ◆ Autoregressive decoder

- Decompose a probability distribution of  $Y$  given  $X$  into a chain of conditional probabilities from left to right

$$P(Y|X) = \prod_{i=1}^N P_{\text{ar}}(y_i | y_{<i}, X)$$

- Optimized with cross-entropy loss  $\mathcal{L}_{\text{ar}} = -\log P_{\text{ar}}(Y|X)$
- Finish decoding after generating <eos>

# Non-autoregressive (NAR) sequence generation

## ◆ Motivation

- AR left-to-right decoding still suffers from slow inference
- Incremental decoding does not enjoy the computational power of GPU/TPU
  - Toward parallel sequence generation

## ◆ Non-autoregressive decoder [\[Gu+ 2018\]](#)

- Assume conditional independence among output tokens

$$P(Y|X) = \prod_{i=1}^N P_{\text{nar}}(y_i|X)$$

- Predict target length in advance  
e.g., Fertility model, linear classifier etc.

# Modeling choice of NAR decoding

## Single forward pass model (faster but less accurate)

### Naïve model

- NAT [\[Gu+ 2018\]](#)
- NAT-REG [\[Wang+ 2019\]](#)
- bag-of-ngram loss [\[Shao+ 2020\]](#)

### Latent variable model

- FlowSeq [\[Ma+ 2019\]](#)
- Delta posterior [\[Shu+ 2020\]](#)

### Alignment model

- CTC [\[Libovický+ 2018\]](#)
- CRF [\[Sun+ 2019\]](#)

## Iterative refinement model (more accurate at the cost of speed)

### Insertion-based model

- Levenshtein Transformer [\[Gu+ 2019\]](#)
- Insertion-deletion Transformer [\[Ruis+ 2019\]](#)
- KERMIT [\[Chan+ 2019\]](#)
- InDIGO [\[Gu+ 2019\]](#)

### Energy-based model

- ENGINE [\[Tu+ 2020\]](#)

### Mask-based model

- Conditional masked language model (CMLM) [\[Ghazvininejad+ 2019\]](#)
- Semi-autoregressive training (SMART) [\[Ghazvininejad+ 2020\]](#)
- Aligned XE [\[Ghazvininejad+ 2020\]](#)
- Disentangled Context Transformer [\[Kasai+ 2020\]](#)
- Imputer [\[Saharia+ 2020\]](#)

# Modeling choice in E2E-ST

- Single-pass model requires a copy of encoder output to initialize decoder input
  - 😓 Non-silence speech frames are NOT uniformly distributed over input speech
  - 😓 Using intermediate prediction from ASR sub-module (e.g., CTC) contradicts the motivation to alleviate error propagation by E2E modeling
- Iterative refinement model can flexibly trade quality and latency during inference by changing the number of iterations
- Want to keep trainability with auxiliary tasks (ASR/MT)
  - Encoder-decoder architecture

We focus on **conditional masked language model (CMLM)** [Ghazvininejad+ 2019]

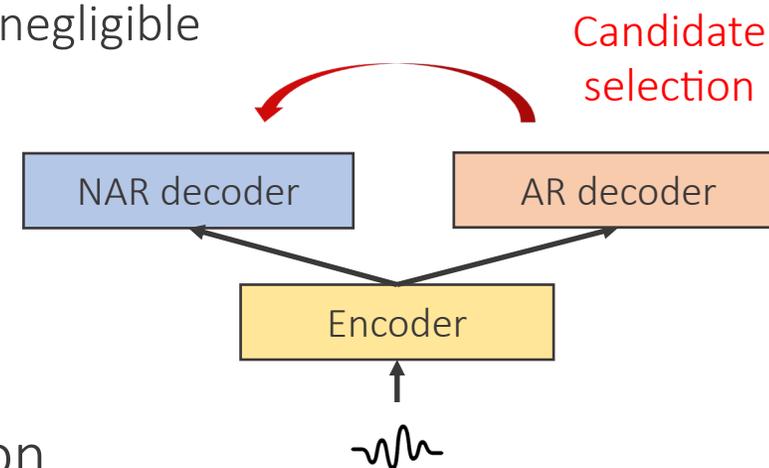
- ✓ Easy implementation
- ✓ Good translation performance

# Proposed framework: Orthros

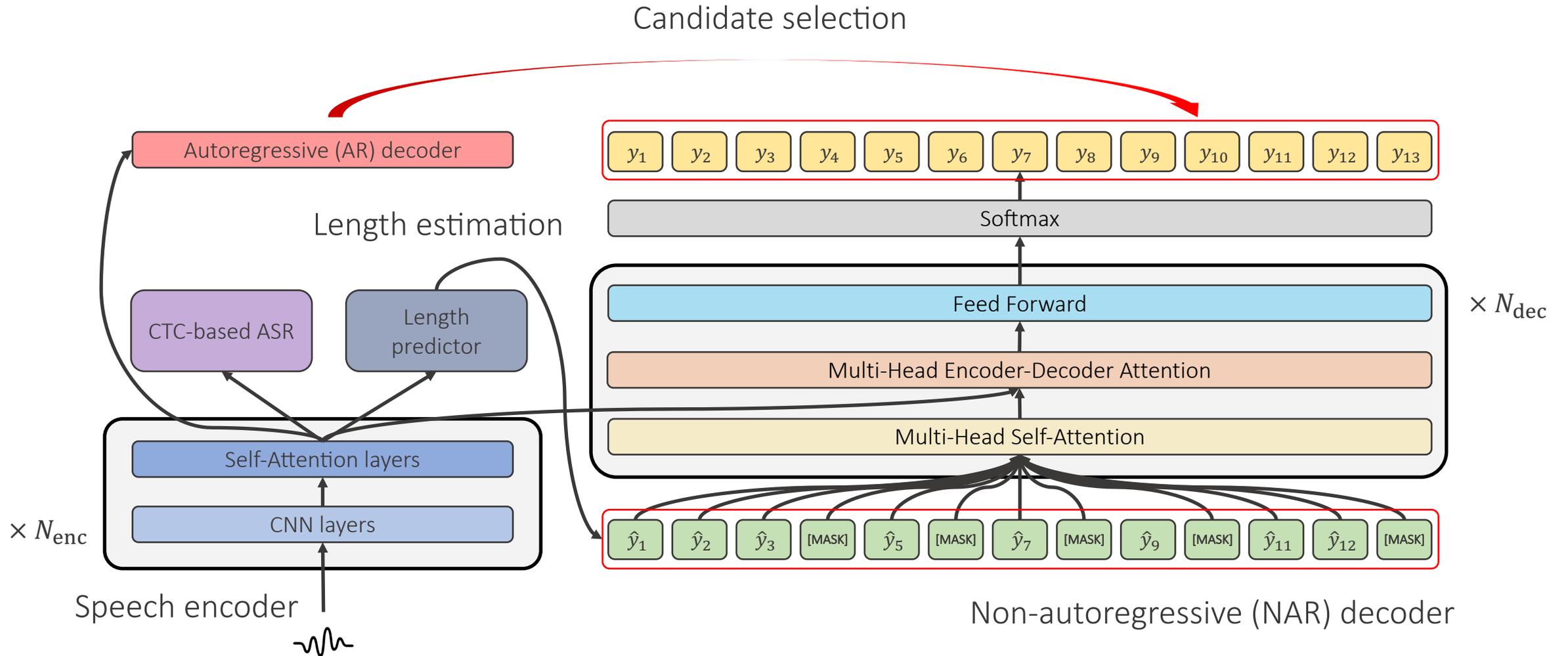
- ◆ Challenge: target length prediction from speech
  - Flexible sequence length: pause, speaking rate, language etc.
  - $|X| \gg |Y|$  even after downsampling
  - Rescoring multiple candidates from NAR model with separate AR model?
    - Extra computation for speech encoding by AR model is not negligible

- ◆ Proposed framework: *Orthros*

- AR and NAR decoders on the **shared** speech encoder
- Unified architecture, trainable in an end-to-end fashion
- Select the most probable candidate from NAR decoder **by scores from AR decoder** (AR decoder can generate scores **in parallel**)



# System overview: Orthros



# CMLM: inference

## ◆ Mask-predict algorithm [Ghazvininejad+ 2019]

- Alternate two operations (*mask*, *predict*) for a constant number of iterations  $T$
- $\hat{Y}_{\text{mask}}^{(t)} \subset Y^{(t-1)}$  (masked tokens at  $t$ -th iteration,  $1 \leq t \leq T$ )
- $\hat{Y}_{\text{obs}}^{(t)} = Y^{(t-1)} \setminus \hat{Y}_{\text{mask}}^{(t)}$  (observed tokens at  $t$ -th iteration)
- Initialize  $\hat{Y}_{\text{obs}}^{(0)}$  with [MASK]

### 1. Mask operation

Predicted target length

- Mask out  $k_t$  tokens having the lowest confidence scores ( $k_t = \lfloor \hat{N} \cdot \frac{T-t}{t} \rfloor$ )

### 2. Predict operation

- Take the most probable token at every masked position  $i$  and update  $y_i^{(t)} \leftarrow y_i^{(t-1)}$

$$y_i^{(t)} = \operatorname{argmax}_{w_i \in V} P_{\text{cmlm}}(w_i | \hat{Y}_{\text{obs}}^{(t)}, X)$$
$$p_i^{(t)} \leftarrow P_{\text{cmlm}}(y_i^{(t)} | \hat{Y}_{\text{obs}}^{(t)}, X)$$

# CMLM: inference

## ◆ Target length prediction

- Take top- $l$  sequence lengths from length distribution  $P_{lp}$

## ◆ Length parallel decoding

- Predict multiple  $l$  sequences having different lengths in parallel
  - In actual implementation, perform batch-decoding, i.e., input/output tensor size:  $[l, \hat{N}_{\max}]$
- Select the most probable sequence at the last iteration among  $l$  candidates

$$score = \frac{1}{\hat{N}} \sum_{i=1}^{\hat{N}} \log P_{i,cmlm}^{(T)}$$

# CMLM: training

## ◆ Notation

- $Y_{\text{mask}} \subset Y$  (masked tokens in ground-truth  $Y$ )
- $Y_{\text{obs}} \subset Y \setminus Y_{\text{mask}}$  (observed tokens in  $Y$ )

## ◆ Training objective

- The number of masked tokens is sampled from uniform distribution  $\mathcal{U}(1, N)$

$$\mathcal{L}_{\text{cmlm}} = - \sum_{y \in Y_{\text{mask}}} \log P_{\text{cmlm}}(y | Y_{\text{obs}}, X)$$

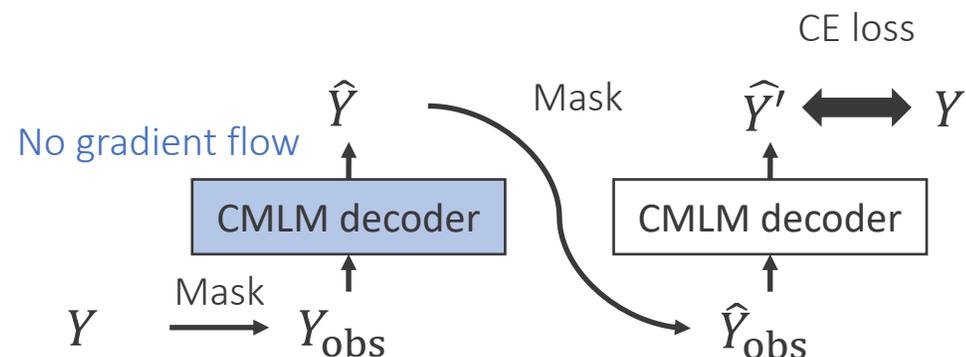
# Semi-autoregressive training (SMART) [Ghazvininejad+ 2020]

- Bridge the gap between training and test conditions by feeding output from the model to the CMLM decoder

## ◆ Procedure

1. Obtain prediction at all positions ( $\hat{Y}$ ) from the current model by feeding  $Y_{\text{obs}}$
2. Obtain new decoder input  $\hat{Y}_{\text{obs}}$  by applying random mask to  $\hat{Y}$
3. Train model to predict  $Y$  given  $\hat{Y}_{\text{obs}}$

Unlike original CMLM, cross-entropy loss is calculated at all position regardless of mask



## ◆ Training objective

$$\mathcal{L}_{\text{cmlm}} = - \sum_{y \in \hat{Y}} \log P_{\text{cmlm}}(y | \hat{Y}_{\text{obs}}, X)$$

# Orthros: training

## ◆ Training objective

$$\mathcal{L}_{\text{total}} = (1 - \lambda_{\text{cmlm}}) \underbrace{\mathcal{L}_{\text{cmlm}}(Y|X)}_{\text{NAR decoder}} + \lambda_{\text{ar}} \underbrace{\mathcal{L}_{\text{ar}}(Y|X)}_{\text{AR decoder}}$$
$$+ \lambda_{\text{lp}} \underbrace{\mathcal{L}_{\text{lp}}(N|X)}_{\text{Length prediction}} + \lambda_{\text{asr}} \underbrace{\mathcal{L}_{\text{asr}}(Y^{\text{src}}|X)}_{\text{ASR}}$$

- Length prediction:  $\mathcal{L}_{\text{lp}}(N|X) = -\log P_{\text{lp}}(N|X)$
- ASR (CTC):  $\mathcal{L}_{\text{asr}}(Y^{\text{src}}|X) = -\log P_{\text{ctc}}(Y^{\text{src}}|X)$

# Orthros: inference

1. Mask-predict for  $T$  iterations
2. Candidate selection with AR decoder
  - After the last iteration, feed outputs from the NAR decoder to the AR decoder in parallel
  - Obtain sequence-level scores from the AR decoder
  - Pick up the most probable candidate among  $l$  candidates

$$score = \frac{1}{\hat{N}} \sum_{i=1}^{\hat{N}} \log P_{i,ar}$$

# Experimental setting

## ◆ Datasets

- Must-C En-De (229k pairs, 408h), En-Fr (275k pairs, 492h)
- Fisher-CallHome Spanish (Es->En, 138k pairs, 170h)
- Libri-trans (En->Fr, 45k pairs, 100h)

## ◆ Model configuration

- Implemented with ESPnet-ST [\[Inaguma+ 2020\]](#)  **ESPnet**
- Transformer base/large ( $d_{\text{model}} = 256/512$ ,  $d_{\text{ff}} = 2048$ ,  $H = 4/8$ )
- 2-layers CNN->12-layers encoder, 6-layers decoder
- **Sequence-level knowledge distillation (Seq-KD)** [\[Kim+ 2016\]](#) from text-based AR MT model
- Vocabulary size
  - AR: BPE8k (Must-C), 1k (Fisher-CallHome, Libri-trans)
  - NAR: BPE8k

# Evaluation metric

## ◆ Translation quality

- 4-gram BLEU

## ◆ Inference speed

- GPU decoding with a NVIDIA TITAN RTX
- Decoding configuration
  - ✓ AR: beam width  $b \in \{1,4\}$
  - ✓ NAR: iteration  $T \in \{4,10\}$ , length beam width  $l = 9$
  - ✓ Batch size: 1
- Averaged over 5 runs

# Main results: Must-C En-De/En-Fr

Model		En-De			En-Fr
		BLEU	Latency [ms]	Speedup	BLEU
Autoregressive	Transformer ( $b=1$ )	21.54	175ms	1.54×	32.26
	Transformer ( $b=4$ )	23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
Non-autoregressive	CTC ( $b=1$ )	19.40	13ms	20.84×	27.38
	Orthros (CMLM, $T=4$ )	18.78	-	-	25.99
	Orthros (CMLM, $T=10$ +AR decoder)	19.62	-	-	27.77
	Orthros (CMLM $T=10$ )	20.89	-	-	28.74
	Orthros (CMLM, $T=10$ +AR decoder)	21.79	-	-	30.31
	Orthros (SMART, $T=4$ )	20.03	46	5.89×	27.22
	Orthros (SMART, $T=10$ +AR decoder)	<b>21.08</b>	<b>61</b>	<b>4.44×</b>	<b>29.30</b>
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	<b>22.27</b>	<b>117</b>	<b>2.44×</b>	<b>31.07</b>
	+ BPE8k -> BPE16k	<b>22.88</b>	117	2.31×	<b>32.20</b>
	+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24
+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>	

## Semi-autoregressive training (SMART)

- Improved BLEU significantly with no extra latency during inference

## n-De/En-Fr

		BLEU	En-De Latency [ms]	Speedup	En-Fr BLEU
Autoregressive	Transformer ( $b=1$ )	21.54	175ms	1.54×	32.26
	Transformer ( $b=4$ )	23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
Non-autoregressive	CTC ( $b=1$ )	19.40	13ms	20.84×	27.38
	Orthros (CMLM, $T=4$ )	18.78	-	-	25.99
	Orthros (CMLM, $T=10$ +AR decoder)	19.62	-	-	27.77
	Orthros (CMLM $T=10$ )	20.89	-	-	28.74
	Orthros (CMLM, $T=10$ +AR decoder)	21.79	-	-	30.31
	Orthros (SMART, $T=4$ )	20.03	46	5.89×	27.22
	Orthros (SMART, $T=10$ +AR decoder)	<b>21.08</b>	<b>61</b>	<b>4.44×</b>	<b>29.30</b>
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	<b>22.27</b>	<b>117</b>	<b>2.44×</b>	<b>31.07</b>
	+ BPE8k -> BPE16k	<b>22.88</b>	117	2.31×	<b>32.20</b>
	+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24
+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>	

### Candidates selection with AR decoder

- Improved BLEU scores is significantly
- This corresponds to performing one more iteration (about +15ms)
- CMLM does not have the ability to generate useful sentence-level scores

En-Fr

		En-De BLEU	En-De Latency [ms]	Speedup	En-Fr BLEU
Autoregressive		23.54	175ms	1.54×	32.26
		23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
Non-autoregressive	CTC ( $b=1$ )	19.40	13ms	20.84×	27.38
	Orthros (CMLM, $T=4$ )	18.78	-	-	25.99
	Orthros (CMLM, $T=10$ +AR decoder)	19.62	-	-	27.77
	Orthros (CMLM, $T=10$ )	20.89	-	-	28.74
	Orthros (CMLM, $T=10$ +AR decoder)	21.79	-	-	30.31
	Orthros (SMART, $T=4$ )	20.03	46	5.89×	27.22
	Orthros (SMART, $T=10$ +AR decoder)	<b>21.08</b>	<b>61</b>	<b>4.44×</b>	<b>29.30</b>
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	<b>22.27</b>	<b>117</b>	<b>2.44×</b>	<b>31.07</b>
	+ BPE8k -> BPE16k	<b>22.88</b>	117	2.31×	<b>32.20</b>
	+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24
	+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>

### Vocabulary size

- Large BPE vocabulary improved BLEU scores
- This was not true for AR models (shown in the later slide)

## En-De/En-Fr

		BLEU	En-De Latency [ms]	Speedup	En-Fr BLEU
Autoregressive	Transformer ( $b=1$ )	21.54	175ms	1.54×	32.26
	Transformer ( $b=4$ )	23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
Non-autoregressive	CTC ( $b=1$ )	19.40	13ms	20.84×	27.38
	Orthros (CMLM, $T=4$ )	18.78	-	-	25.99
	Orthros (CMLM, $T=10$ +AR decoder)	19.62	-	-	27.77
	Orthros (CMLM $T=10$ )	20.89	-	-	28.74
	Orthros (CMLM, $T=10$ +AR decoder)	21.79	-	-	30.31
	Orthros (SMART, $T=4$ )	20.03	46	5.89×	27.22
	Orthros (SMART, $T=10$ +AR decoder)	<b>21.08</b>	<b>61</b>	<b>4.44×</b>	<b>29.30</b>
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	<b>22.27</b>	<b>117</b>	<b>2.44×</b>	<b>31.07</b>
	+ BPE8k -> BPE16k	<b>22.88</b>	117	2.31×	<b>32.20</b>
+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24	
+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>	

# En-De/En-Fr

## Large model

- Increasing model capacity was very important for NAR models
- This was not true for AR models (shown in the later slide)

		BLEU	En-De Latency [ms]	Speedup	En-Fr BLEU
Autoregressive		21.54	175ms	1.54×	32.26
		23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
Non-autoregressive	CTC ( $b=1$ )	19.40	13ms	20.84×	27.38
	Orthros (CMLM, $T=4$ )	18.78	-	-	25.99
	Orthros (CMLM, $T=10$ +AR decoder)	19.62	-	-	27.77
	Orthros (CMLM $T=10$ )	20.89	-	-	28.74
	Orthros (CMLM, $T=10$ +AR decoder)	21.79	-	-	30.31
	Orthros (SMART, $T=4$ )	20.03	46	5.89×	27.22
	Orthros (SMART, $T=10$ +AR decoder)	<b>21.08</b>	<b>61</b>	<b>4.44×</b>	<b>29.30</b>
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	<b>22.27</b>	<b>117</b>	<b>2.44×</b>	<b>31.07</b>
	+ BPE8k -> BPE16k	<b>22.88</b>	117	2.31×	<b>32.20</b>
+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24	
+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>	

# Main results: Must-C En-De/En-Fr

Model		En-De			En-Fr
		BLEU	Latency [ms]	Speedup	BLEU
Autoregressive	Transformer ( $b=1$ )	21.54	175ms	1.54×	32.26
	Transformer ( $b=4$ )	23.12	271ms	1.00×	33.84
	Transformer + Seq-KD ( $b=1$ )	23.88	-	-	33.92
	Transformer + Seq-KD ( $b=4$ )	24.43	-	-	34.57
	STG ( $T=1$ )	19.40	13ms	20.84×	27.38
		18.78	-	-	25.99
		19.62	-	-	27.77
		20.89	-	-	28.74
		21.79	-	-	30.31
		20.03	46	5.89×	27.22
		21.08	61	4.44×	29.30
	Orthros (SMART, $T=10$ )	21.25	99	2.73×	29.31
	Orthros (SMART, $T=10$ +AR decoder)	22.27	117	2.44×	31.07
	+ BPE8k -> BPE16k	22.88	117	2.31×	32.20
	+ large (SMART, $T=4$ +AR decoder, $l=7$ )	22.54	59	4.59×	31.24
	+ large (SMART, $T=10$ +AR decoder, $l=7$ )	<b>23.92</b>	113	2.39×	<b>33.05</b>

## NAR vs AR

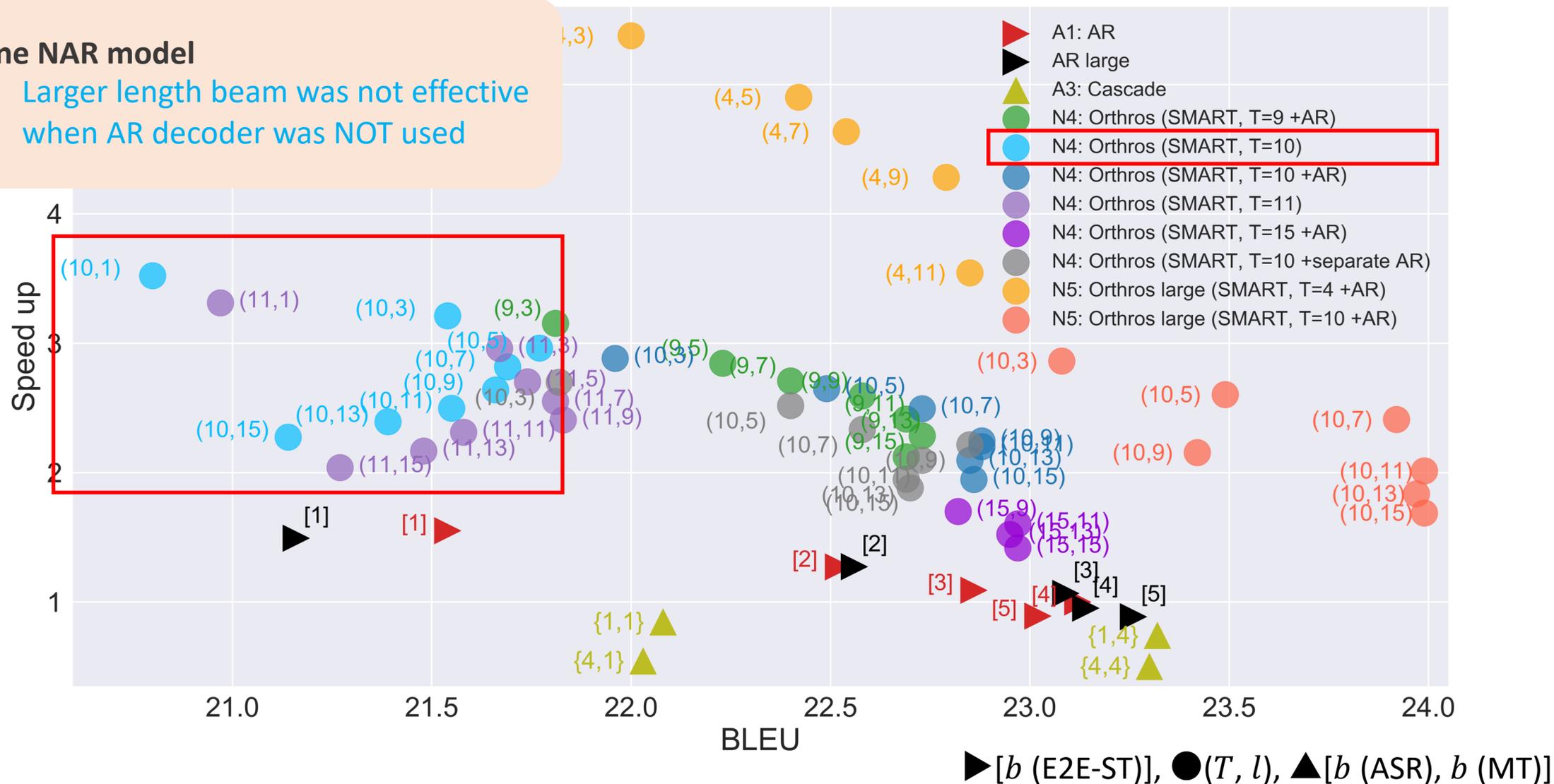
- Achieved comparable BLEU scores to baseline Transformer
- Seq-KD boosted AR model's performance further
- This differs from MT:
  - MT: large AR teacher -> small AR student
  - E2E-ST: AR MT teacher -> AR E2E-ST student



# Quality and latency trade-off on Must-C En-De

## Baseline NAR model

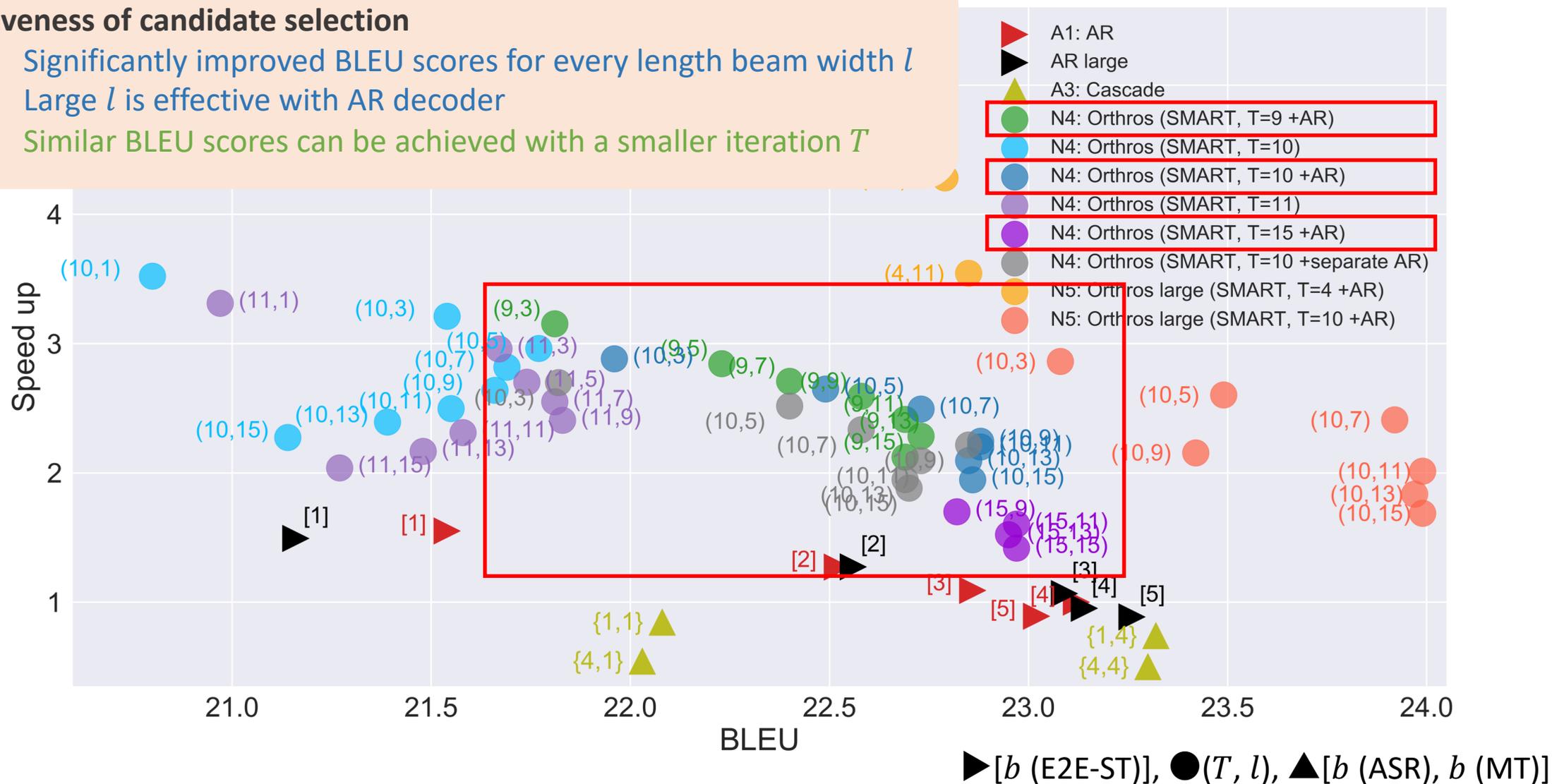
- Larger length beam was not effective when AR decoder was NOT used



# Quality and latency trade-off on Must-C En-De

## Effectiveness of candidate selection

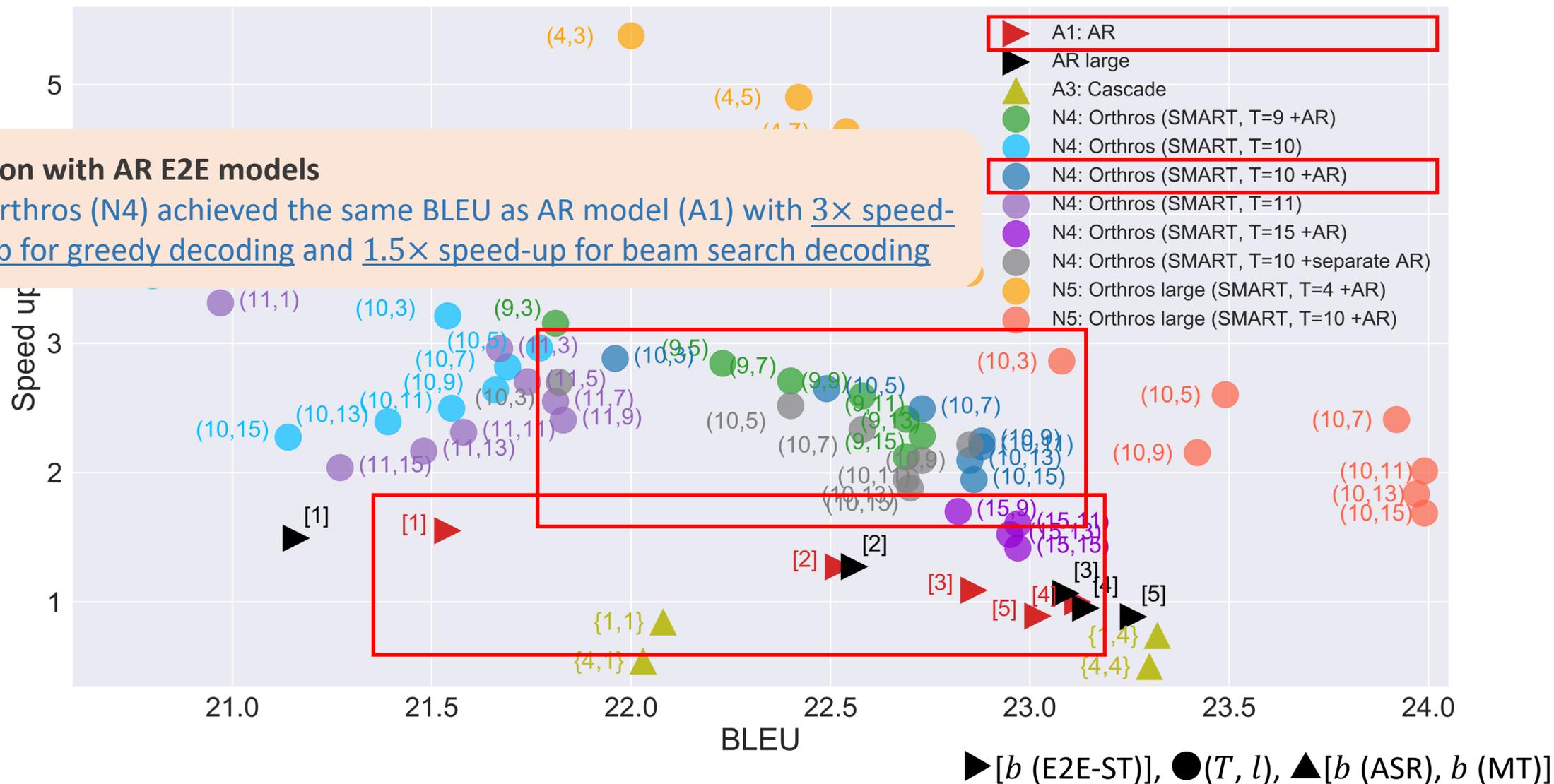
- Significantly improved BLEU scores for every length beam width  $l$
- Large  $l$  is effective with AR decoder
- Similar BLEU scores can be achieved with a smaller iteration  $T$



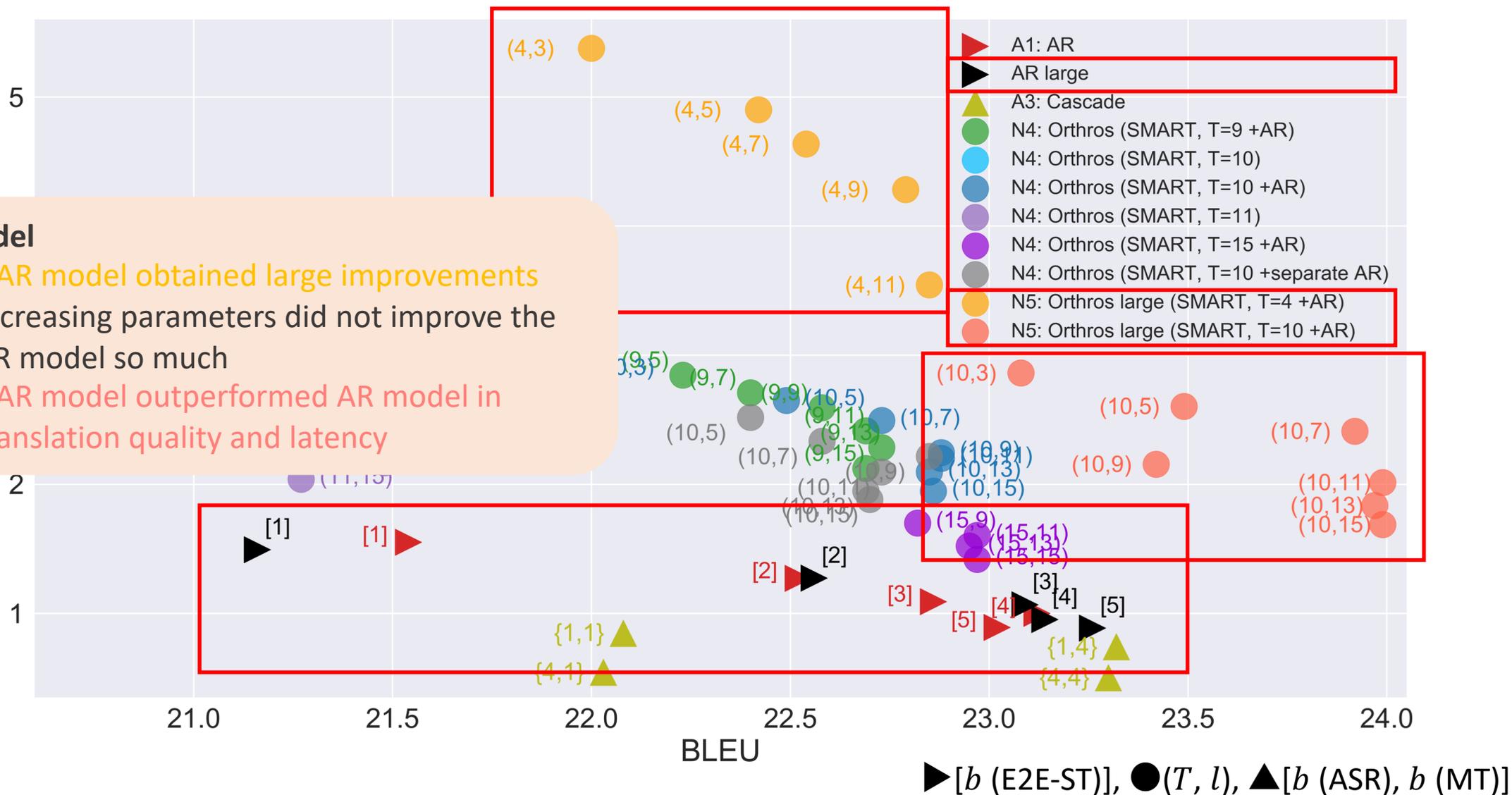
# Quality and latency trade-off on Must-C En-De

## Comparison with AR E2E models

- Orthros (N4) achieved the same BLEU as AR model (A1) with 3× speed-up for greedy decoding and 1.5× speed-up for beam search decoding



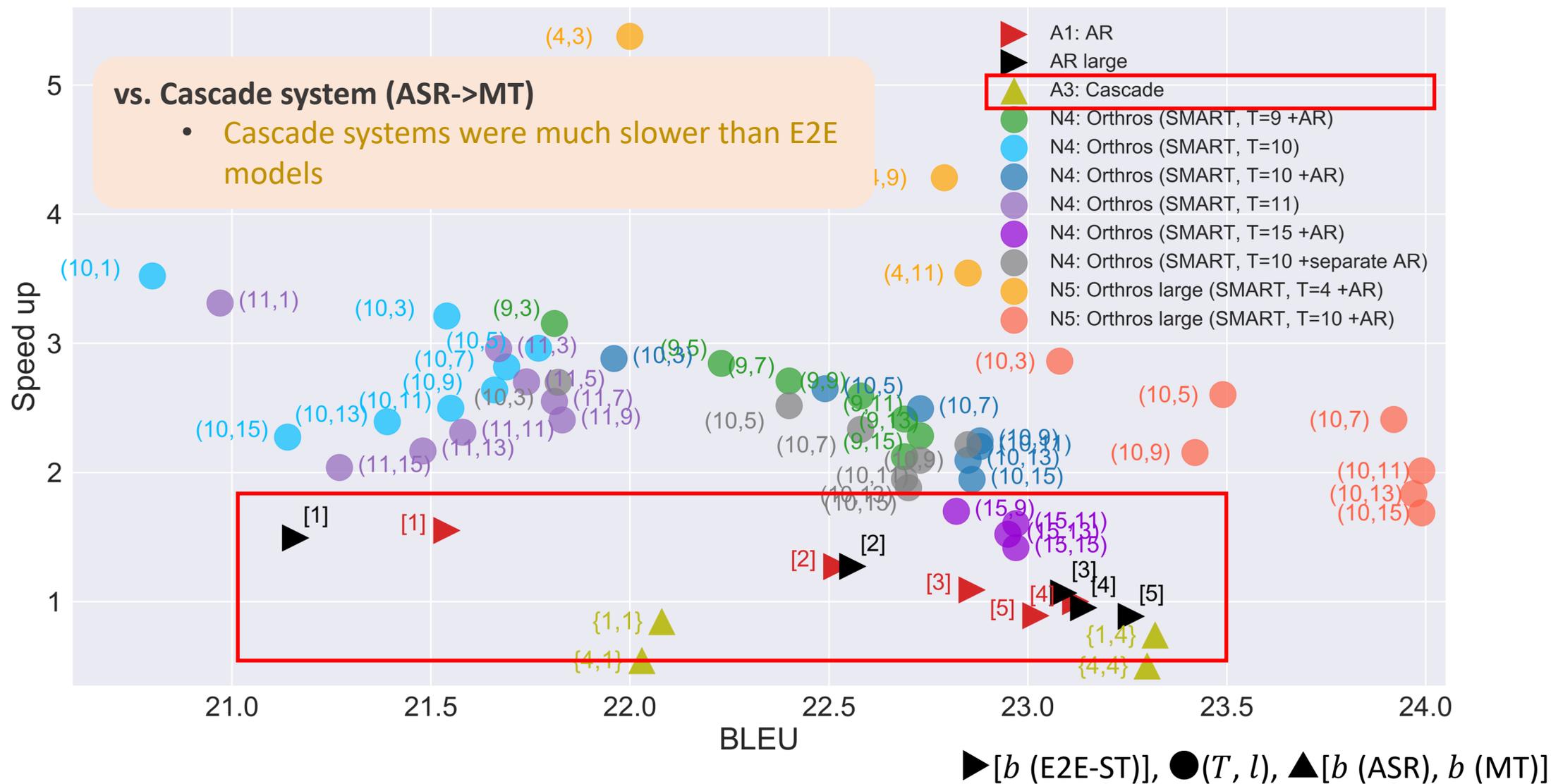
# Quality and latency trade-off on Must-C En-De



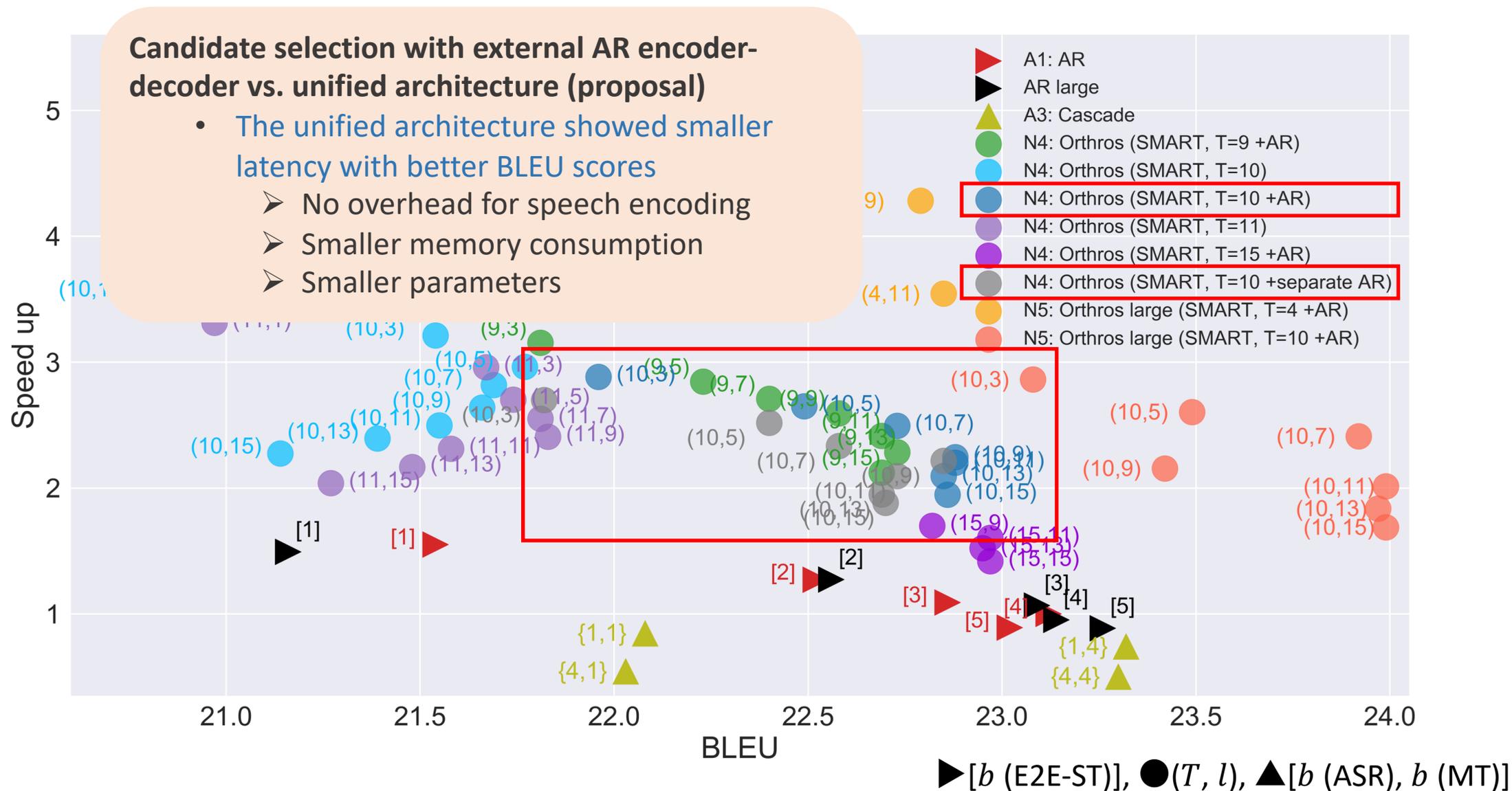
## Large model

- NAR model obtained large improvements
- Increasing parameters did not improve the AR model so much
- NAR model outperformed AR model in translation quality and latency

# Quality and latency trade-off on Must-C En-De



# Quality and latency trade-off on Must-C En-De



# Results: Fisher-CallHome Spanish/Libri-trans

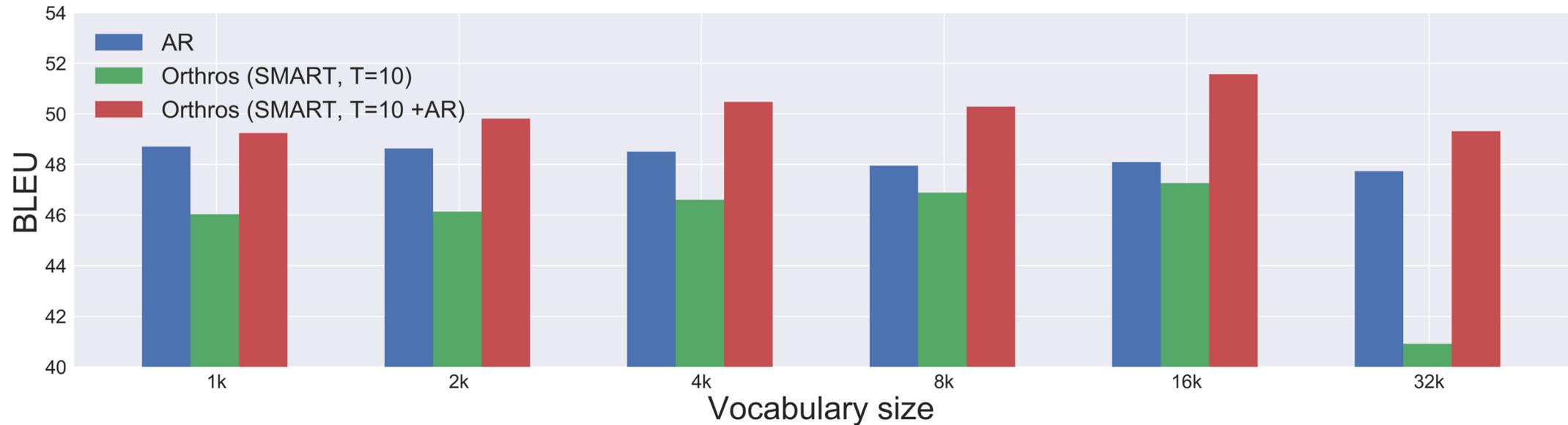
Model		BLEU		
		Fisher-CallHome Spanish		
		Fisher-test	CallHome-evltest	Libri-trans
Autoregressive	Transformer ( $b=1$ )	48.38	18.07	16.52
	Transformer ( $b=4$ )	48.49	18.90	16.84
	Transformer + Seq-KD ( $b=1$ )	50.34	19.09	15.91
	Transformer + Seq-KD ( $b=4$ )	50.32	19.81	16.44
Non-autoregressive	CTC ( $b=1$ )	45.97	15.91	12.10
	Orthros (CMLM, $T=4$ )	46.03	16.71	12.90
	Orthros (CMLM, $T=10$ +AR decoder)	47.80	18.28	13.69
	Orthros (CMLM $T=10$ )	48.56	18.60	14.68
	Orthros (CMLM, $T=10$ +AR decoder)	49.98	19.71	15.43
	Orthros (SMART, $T=4$ )	45.89	17.39	14.17
	Orthros (SMART, $T=10$ +AR decoder)	<b>48.73</b>	<b>19.25</b>	<b>14.99</b>
	Orthros (SMART, $T=10$ )	47.09	18.25	15.11
	Orthros (SMART, $T=10$ +AR decoder)	<b>50.07</b>	<b>20.10</b>	<b>16.08</b>
	+ BPE8k -> BPE16k	<b>50.18</b>	19.88	<b>16.22</b>

# Ablation study on Fisher-CallHome dev set

Model	BLEU			
	$T = 4$		$T = 10$	
	w/o AR	w/ AR	w/o AR	w/ AR
Orthros BPE8k	45.76	49.01	46.88	50.28
- Seq-KD	44.36	47.42	44.25	49.50
- AR decoder	45.53	-	46.94	-
+ length prediction w/ CTC	45.41	48.18	46.79	50.05

- Seq-KD was beneficial (multi-modality problem was alleviated)
- Joint training with AR decoder itself had no impact on BLEU scores
- Linear classifier-based length prediction was better than the CTC-based one
  - CTC-based length prediction:  $[\hat{N} - \lfloor \frac{l}{2} \rfloor, \hat{N} + \lfloor \frac{l}{2} \rfloor]$ , where  $\hat{N} = \lfloor \alpha \hat{N}_{src} \rfloor$   
( $\alpha$ : hyperparameter,  $\hat{N}_{src}$ : ASR hypothesis length obtained by CTC greedy decoding)

# Effect of vocabulary size on Fisher-CallHome dev set



- AR models have a peak around BPE1k (due to data sparseness, 170h)
- Candidate selection with AR decoder is always effective regardless of BPE size
- Orthros + candidate selection continued to improve until BPE16k
  - Most tokens in vocabulary are “complete” word
  - Complementary effect on the conditional independence assumption

# Conclusion and future work

- Perceived latency reduction for streaming encoder-decoder ASR
  - Alignment information is effective **on the decoder side**
  - CTC alignment is as good as alignment from hybrid ASR system
- Fast non-autoregressive decoding for E2E-ST, *Orthros*
  - AR decoder + NAR decoder on shared speech encoder
  - Candidate selection with AR decoder was very effective to estimate target lengths
  - Reached comparable translation quality to SOTA AR E2E-ST models with more than **2× latency reduction**