

CTC Alignments Improve Autoregressive Translation

Brian Yan*¹ **Siddharth Dalmia***¹ **Yosuke Higuchi**²

Graham Neubig¹ **Florian Metze**¹ **Alan W Black**¹ **Shinji Watanabe**^{1,3}

¹Language Technologies Institute, Carnegie Mellon University, USA

²Department of Communications and Computer Engineering, Waseda University, Japan

³Human Language Technology Center of Excellence, Johns Hopkins University, USA

{byan, sdalmia}@cs.cmu.edu



Carnegie Mellon University
Language Technologies Institute



CMU-LTI WAV Lab

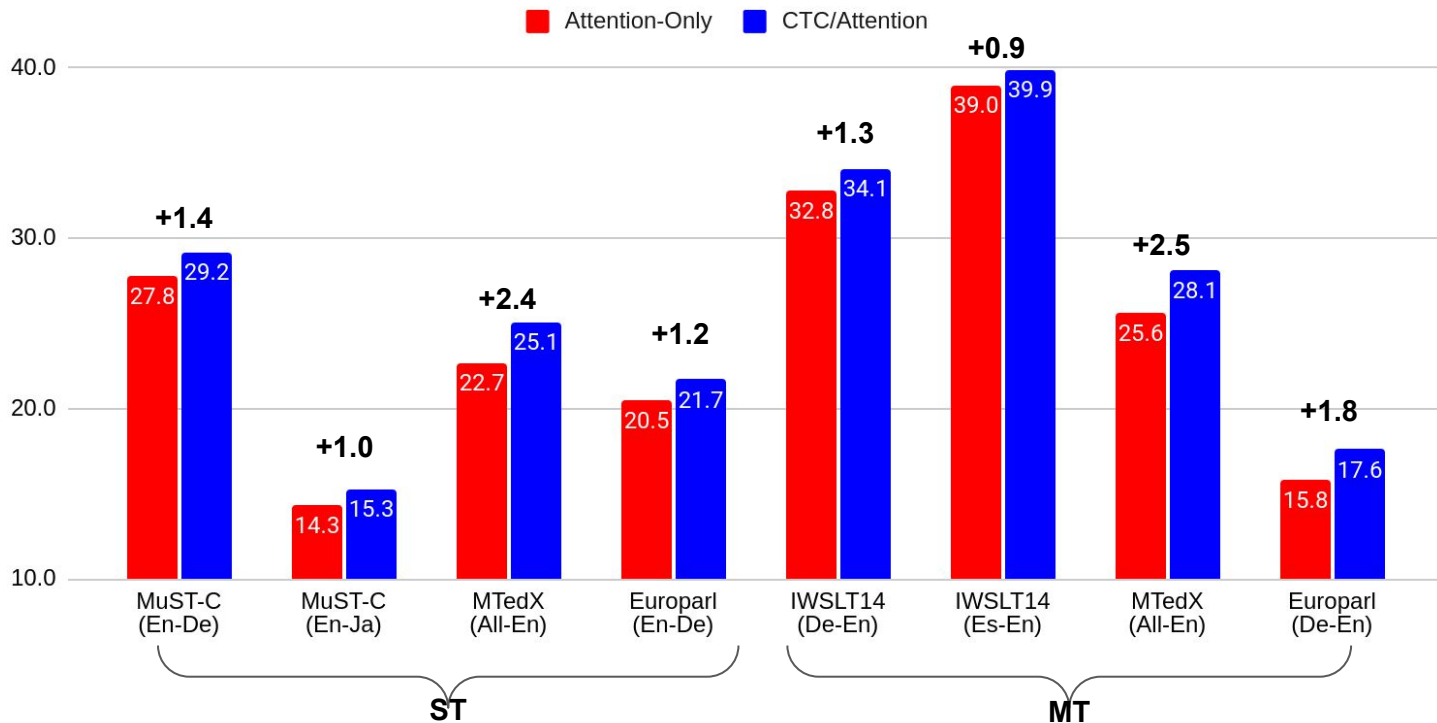
Does CTC make sense for translation?*

- Part 1: CTC vs. Attentional Encoder-Decoder
- Part 2: Joint CTC/Attention

* Considering translation quality only

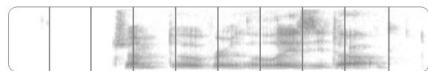
Results (Preview)

- Joint CTC/attention outperforms pure-attention by an average of **+1.6 BLEU**



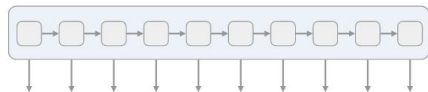
* "Attention" refers to autoregressive encoder-decoder models with cross-attention mechanisms, optimized via cross-entropy

Connectionist Temporal Classification (CTC)



Input (speech signal)

$$X$$



h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o	o
ε	ε	ε	ε	ε	ε	ε	ε	ε	ε

Frame-level posteriors
(aka alignment posteriors)

$$P(z_t | X, \underline{z}_{1:t-1})$$

h	e	ε	l	l	ε	l	l	o	o
h	h	e	l	l	ε	ε	l	ε	o
ε	e	ε	l	l	ε	ε	l	o	o

Alignment sequence likelihoods

$$\prod_{t=1}^T P(z_t | X, \underline{z}_{1:t-1})$$

h	e	l	l	o
e	l	l	o	
h	e	l	o	

Label sequence likelihoods

$$\sum_{Z \in \mathcal{Z}} \prod_{t=1}^T P(z_t | X, \underline{z}_{1:t-1})$$

Properties of CTC

CTC

$$P_{\text{CTC}}(Y|X) \triangleq \sum_{Z \in \mathcal{Z}} \prod_{t=1}^T P(z_t | X, z_{1:t-1})$$

Hard Alignment

Criterion only allows monotonic alignments of inputs to outputs

Conditional Independence

Assumes that there are no dependencies between each output unit given the input

Input-Synchronous Emission

Each input representation emits exactly one blank or non-blank output token

CTC vs. Attention

CTC

$$P_{\text{CTC}}(Y|X) \triangleq \sum_{Z \in \mathcal{Z}} \prod_{t=1}^T P(z_t | X, z_{1:t-1})$$

Hard Alignment

Criterion only allows monotonic alignments of inputs to outputs

Conditional Independence

Assumes that there are no dependencies between each output unit given the input

Input-Synchronous Emission

Each input representation emits exactly one blank or non-blank output token

Weaker for translation

(Gu+ 2021, Huang+ 2022)

ATTENTION

$$P_{\text{Attn}}(Y|X) \triangleq \prod_{l=1}^L P(y_l | y_{1:l-1}, X)$$

Soft Alignment

Conditional Dependence

Autoregressive Generation

CTC vs. Attention

CTC

$$P_{\text{CTC}}(Y|X) \triangleq \sum_{Z \in \mathcal{Z}} \prod_{t=1}^T P(z_t | X, z_{1:t-1})$$

Hard Alignment

Criterion only allows monotonic alignments of inputs to outputs

Conditional Independence

Assumes that there are no dependencies between each output unit given the input

Input-Synchronous Emission

Each input representation emits exactly one blank or non-blank output token

ATTENTION

$$P_{\text{Attn}}(Y|X) \triangleq \prod_{l=1}^L P(y_l | y_{1:l-1}, X)$$

Soft Alignment

Flexible attention-based input-to-output mappings may overfit to irregular patterns

Conditional Dependence

Locally normalized models with output dependency exhibit label/exposure biases

Autoregressive Generation

Need to detect end-points and compare hypotheses of different length in beam search

Attention is not perfect

(Murray+ 2018, Hannun+ 2019, Watanabe+ 2018)

CTC and Attention are complementary

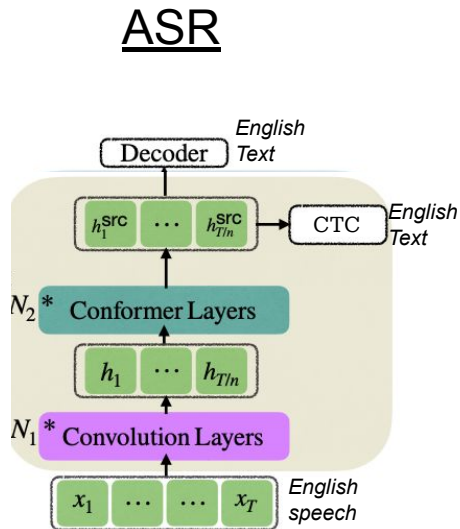
- Joint CTC/attention is excellent for ASR (Kim+ 2017, Watanabe+ 2018)
- Joint CTC/attention should also benefit MT/ST due to positive interactions:
 - Hard alignment + soft alignment
 - Conditional independence + conditional dependence
 - Input synchronous emission + autoregressive generation

CTC/Attention for MT/ST (1)

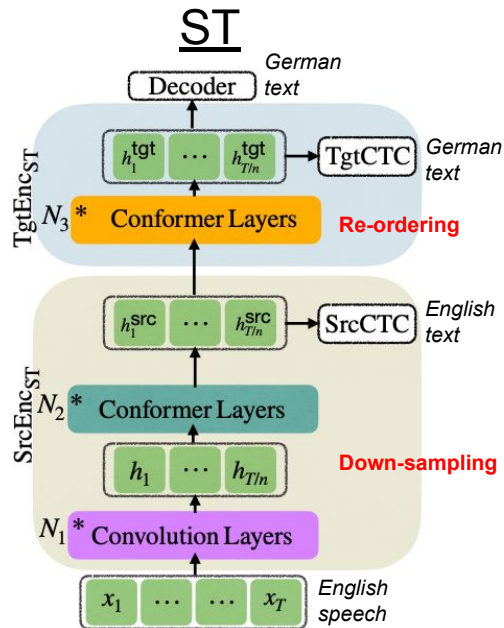
- **Hard alignment** + **soft alignment**
 - Conjecture: hard alignment objective produces stable encoder representations allowing the decoder to **more easily learn soft alignment patterns during training**
 - **Can CTC encoders perform input-to-output mappings for translation?**
 - Outputs may be longer than inputs (Libovicky+ 2018, Dalmia+ 2022)
 - Input-to-output re-ordering (Chuang+ 2021)
- Let's incorporate these requirements into our encoder architecture

Joint Training with Hierarchical CTC

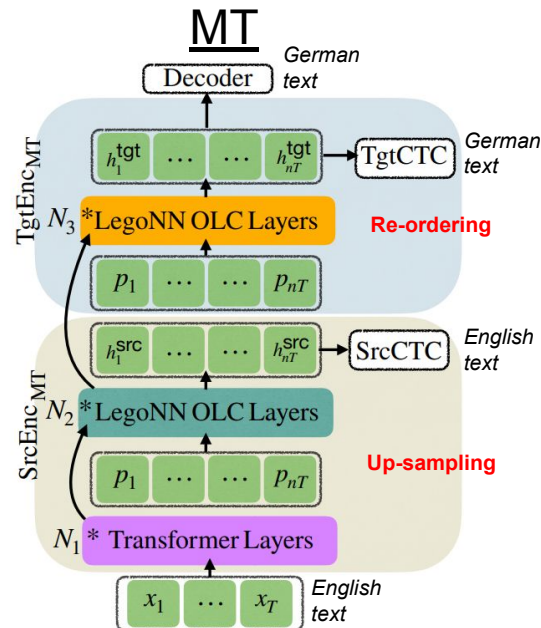
- For ASR, CTC and attention simply share a monolithic encoder
- For MT/ST, we decompose the encoder into 2 stages: **1) length-adjustment 2) re-ordering**



$$\mathcal{L} = \lambda_1 \mathcal{L}_{CTC} + \lambda_2 \mathcal{L}_{ATTN}$$



$$\mathcal{L} = \mathcal{L}_{SRCCTC} + \lambda_1 \mathcal{L}_{TGTCTC} + \lambda_2 \mathcal{L}_{ATTN}$$



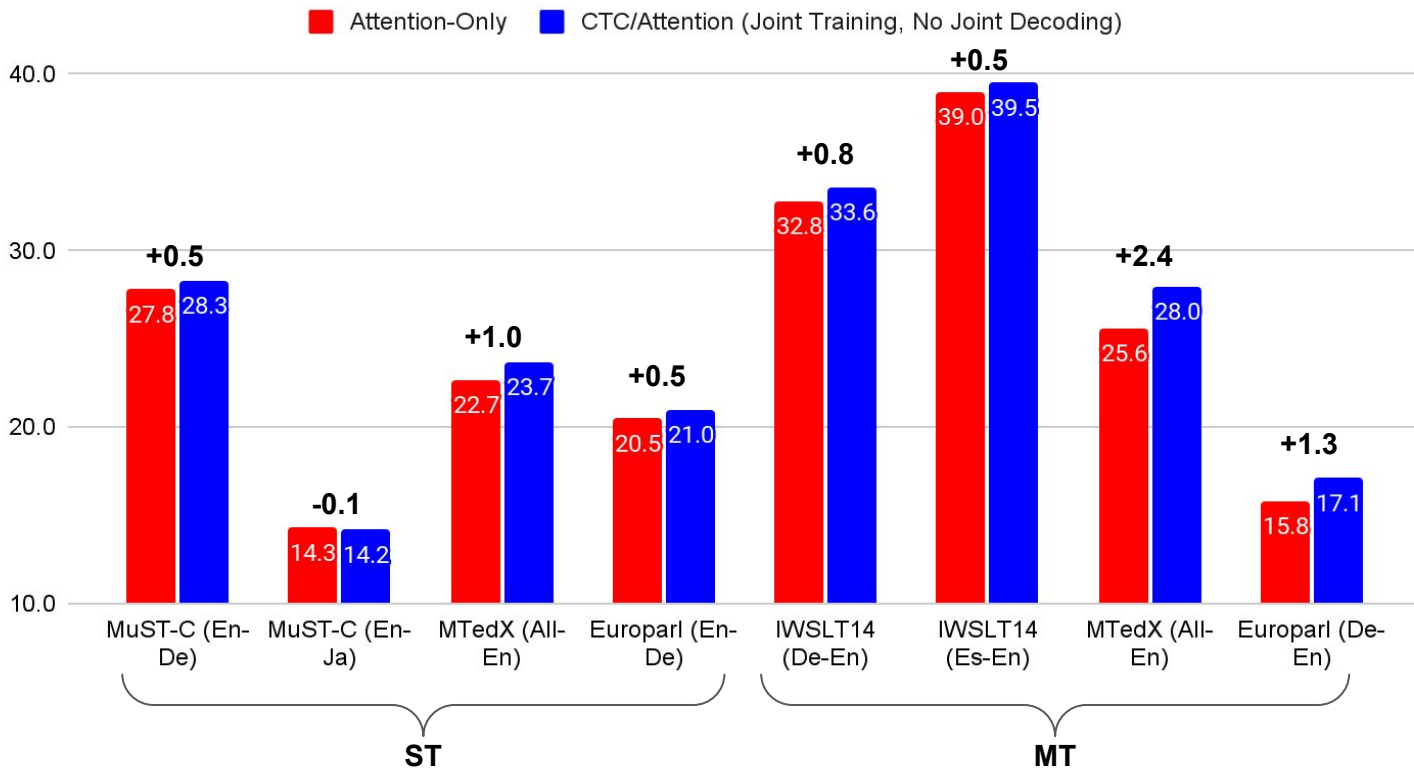
Joint Training with Hierarchical CTC

- Ablation shows that separating length-adjustment and re-ordering is beneficial

		MT (DE-EN)	ST (EN-DE)
SRCCTC	TGTCTC	IWSLT14	MuST-C-v2
X	X	32.1	27.7
✓	X	34.1	27.8
X	✓	33.3	28.1
✓	✓	34.8	28.3

Attention w/ CTC Joint Training vs. Pure-Attention

- Joint training yields an average of **+0.9** BLEU improvement

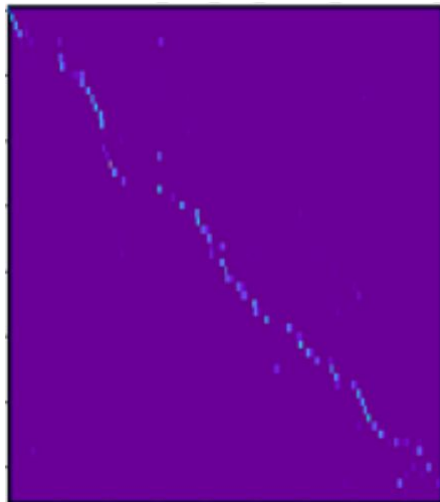
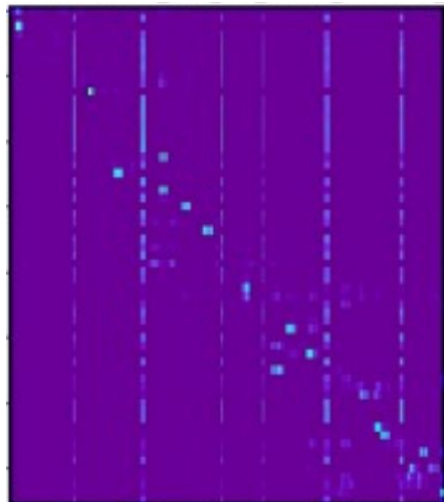


Reduced Soft Alignment Burden

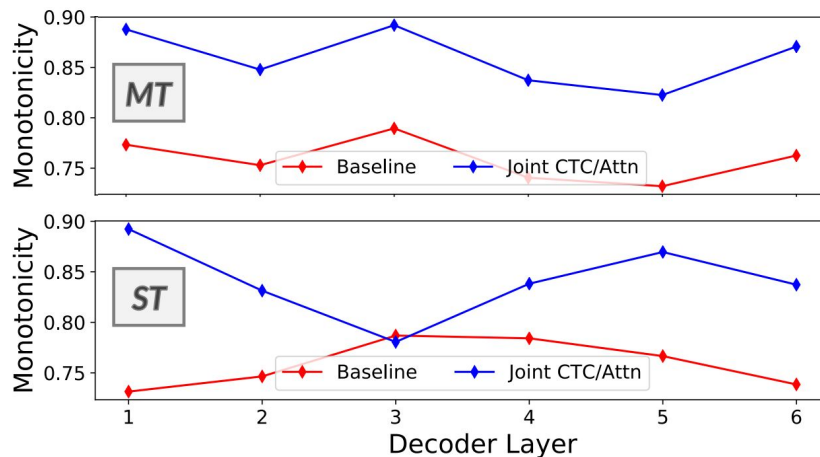
- Joint training results in more regular, diagonal source-attention patterns

Pure-Attention

CTC/Attention

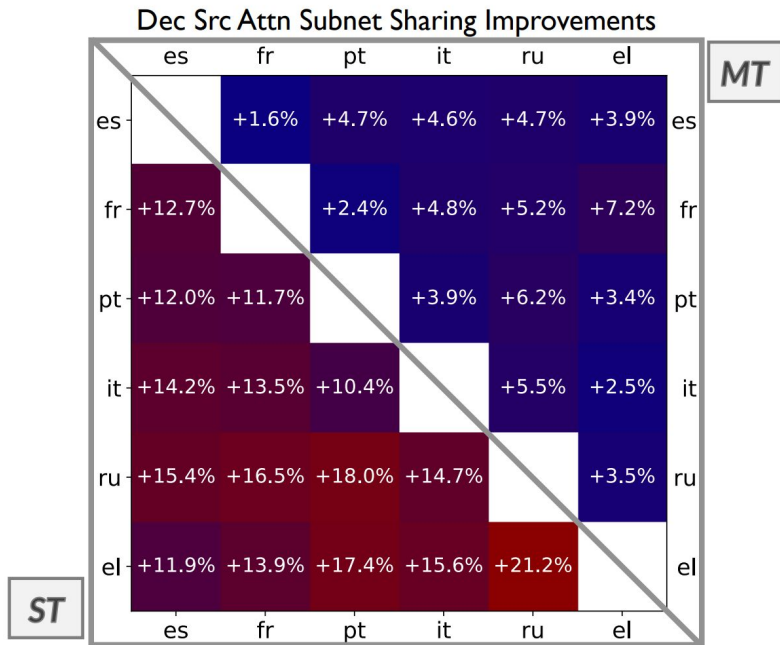


Layer-wise Monotonicity



Increased Multilingual Parameter Sharing

- For X→En models, decoder source-attention parameter sharing between languages was higher in CTC/attention vs. pure-attention



Language specific subnets extracted via Lottery Ticket Sparse Fine-Tuning (Ansell+ 2022)

CTC/Attention for MT/ST (2)

- **Conditional independence** + **conditional dependence**
 - Conjecture: use of conditionally independent likelihoods in joint scoring **eases the exposure/label biases** from conditionally dependent likelihoods **during decoding**
- **Does CTC translation quality lag too far behind attention to be useful?**
 - In our study, pure-CTC models are up to 28% worse than pure-attention models

→ Let's examine joint decoding of CTC/attention

Joint Decoding with Output-Sync Beam Search

- Attention** plays a primary role while **CTC** plays a secondary role (Watanabe+ 2018)

Algorithm 1 *Output-Synchronous Step Function*:
 attentional decoder proposes candidates to expand hypotheses which are all of l -length at step l .

```

1: procedure OUTPUTSTEP(prtHs,  $X$ ,  $l$ ,  $p$ , max $L$ )
2:   newPrtHs = {}; endHs = {}
3:   for  $y_{1:l-1} \in$  prtHs do
4:     attnCnds = top-k( $P_{\text{Attn}}(y_l|X, y_{1:l-1})$ , k =  $p$ )
5:     for  $c \in$  attnCnds do
6:        $y_{1:l} = y_{1:l-1} \oplus c$  ← Hypothesis Expansion
7:        $\alpha_{\text{CTC}} = \text{CTCScore}(y_{1:l}, X_{1:T})$ 
8:        $\alpha_{\text{Attn}} = \text{AttnScore}(y_{1:l}, X_{1:T})$  ← Joint Scoring
9:        $\beta = \text{LengthPen}(y_{1:l})$ 
10:       $P_{\text{Beam}}(y_{1:l}|X) = \alpha_{\text{CTC}} + \alpha_{\text{Attn}} + \beta$ 
11:      if ( $c$  is <eos>) or ( $l$  is max $L$ ) then ← End Detection
12:        endHs[ $y_{1:l}$ ] =  $P_{\text{Beam}}(\cdot)$ 
13:      else
14:        newPrtHs[ $y_{1:l}$ ] =  $P_{\text{Beam}}(\cdot)$ 
15:      end if
16:    end for
17:  end for
18:  return newPrtHs, endHs
19: end procedure
  
```

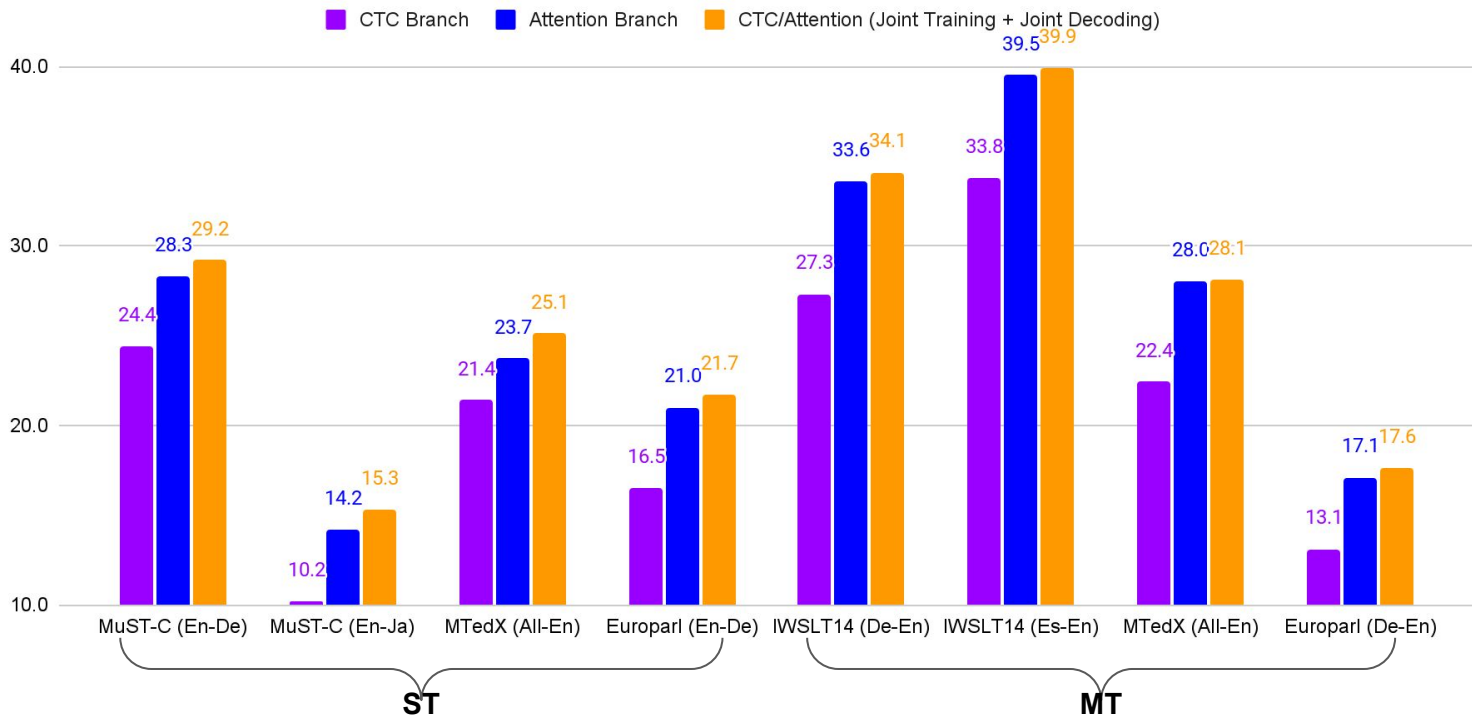
Choose top- p candidates per hypothesis (e.g. $p=1.5*b$) **from attention posteriors**

Interpolate label sequence likelihoods **from attention and CTC**

End loop based on conditions **from attention** (CTC's role in end-detection is implicit)

Joint Training + Decoding vs. Only Joint Training

- Joint decoding yields an average of **+0.7 BLEU** improvement over the attention branch
- For jointly trained models, attention branch outperforms CTC branch by an avg. of **+4.5 BLEU**



CTC/Attention for MT/ST (3)

- **Input synchronous emission + autoregressive generation**
 - Conjecture: input-synchronous emission determines output length based on input length **counteracting the autoregressive end-detection problem during decoding**

- **Is the alignment information from CTC translation models reasonable?**
 - Even in ASR, some alignment “drift” can occur (Kurzinger+ 2020)
 - If alignments are highly noisy, CTC’s end-detection property may not be useful

→ We address this via sanity checks

Joint Decoding with Input-Sync Beam Search

- CTC plays a primary role while attention plays a secondary role (inverse of output-sync)

Algorithm 2 *Input-Synchronous Step Function:*

CTC proposes candidates to expand hypotheses which are all produced from t input units at step t .

```
1: procedure INPUTSTEP(prtHs, X, t, p, T)
2:   newPrtHs = {}; endHs = {}
3:   CTCCnds = top-k( $P_{CTC}(z_t|X)$ , k = p)
4:   for  $y \in \text{prtHs}$  do
5:     for  $c \in \text{CTCCnds}$  do
6:       if (c is  $\emptyset$ ) or (c is  $y[-1]$ ) then
7:         Hypothesis Expansion →  $\tilde{y} = y$ 
8:       else
9:          $\tilde{y} = y \oplus c$ 
10:      end if
11:      Joint Scoring →  $\alpha_{CTC} = \text{CTCScore}(\tilde{y}, X_{1:t})$ 
12:       $\alpha_{Attn} = \text{AttnScore}(\tilde{y}, X_{1:T})$ 
13:       $\beta = \text{LengthPen}(\tilde{y})$ 
14:       $P_{\text{Beam}}(\tilde{y}|X) = \alpha_{CTC} + \alpha_{ATTN} + \beta$ 
15:      if  $t$  is  $T$  then
16:        End Detection →  $\text{endHs}[\tilde{y}] = P_{\text{Beam}}(\cdot)$ 
17:      else
18:         $\text{newPrtHs}[\tilde{y}] = P_{\text{Beam}}(\cdot)$ 
19:      end if
20:    end for
21:  end for
22:  return newPrtHs, endHs
23: end procedure
```

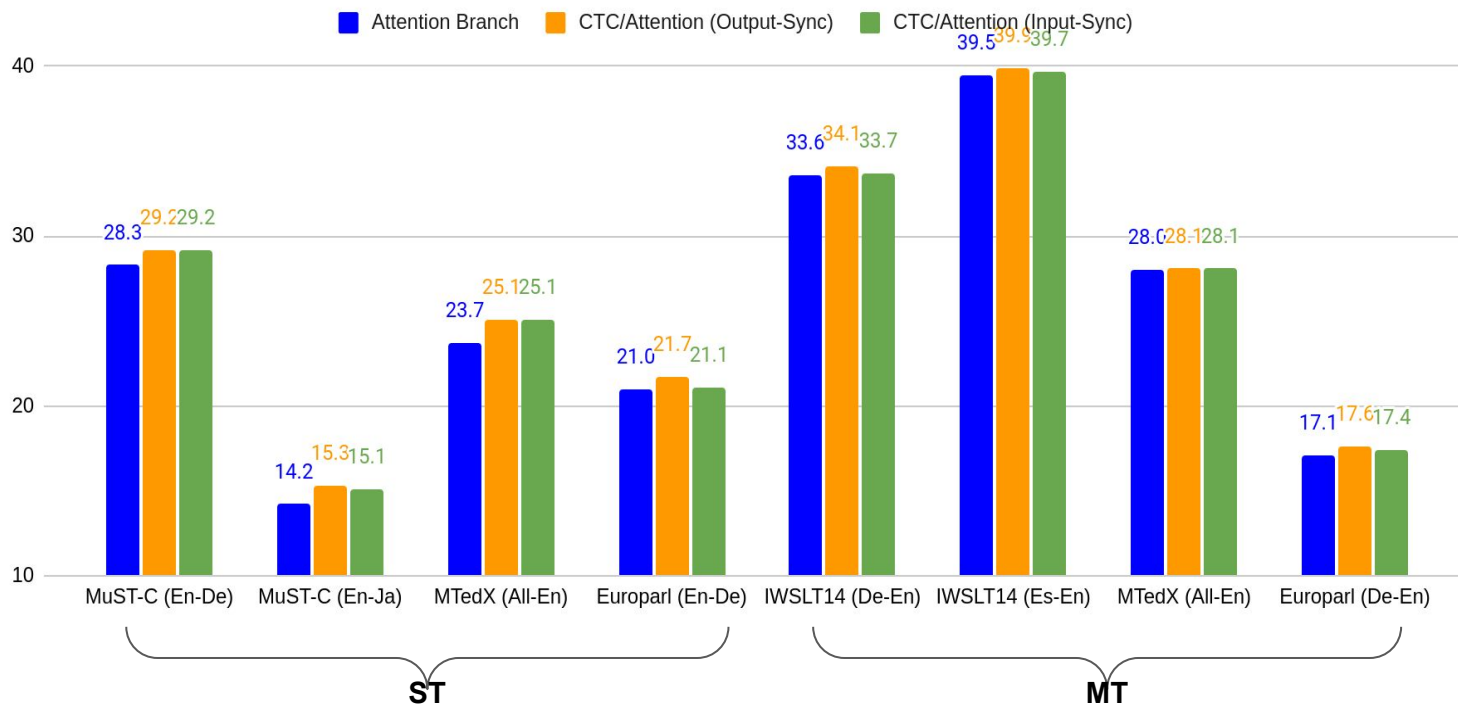
Choose top- p candidates (e.g. $p=1.5*b$) from CTC posterior

Interpolate label sequence likelihoods from attention and CTC

End loop based on conditions from CTC

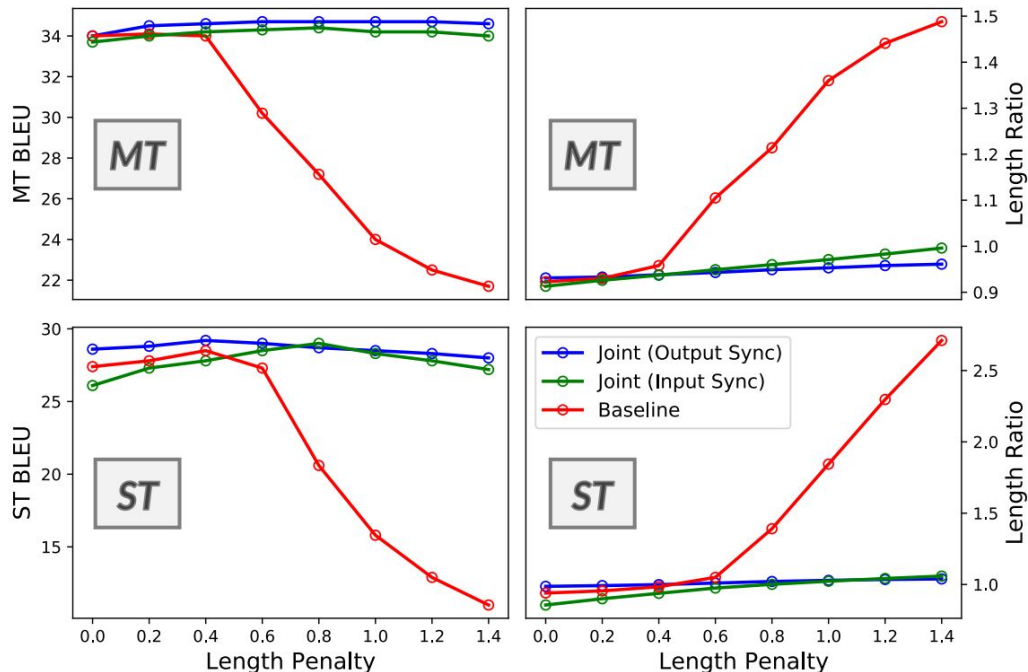
Input-Sync vs. Output-Sync Joint Decoding

- Input-sync joint CTC/attention outperforms the attention branch by **+0.5 BLEU**
- Input-sync is only **-0.2 BLEU** worse than output-sync



Robust End-Detection

- Both variants of joint CTC/attention decoding have low length penalty elasticity
- Pure-attention models are highly sensitive to length penalty → easily overtuned



Summary

- Joint CTC/attention is effectively applied to MT/ST with only minor changes from ASR
- Both joint training and joint decoding yield performance gains
- Why is joint CTC/attention better than pure-attention?
 - Simplifies the soft alignment task
 - Positive ensembling effect
 - Robust end-detection



Thank You!