

CTC Alignments Improve Autoregressive Translation

Brian Yan, Siddharth Dalmia, Yosuke Higuchi, Graham Neubig, Florian Metze, Alan W Black, Shinji Watanabe

CTC vs. Attention

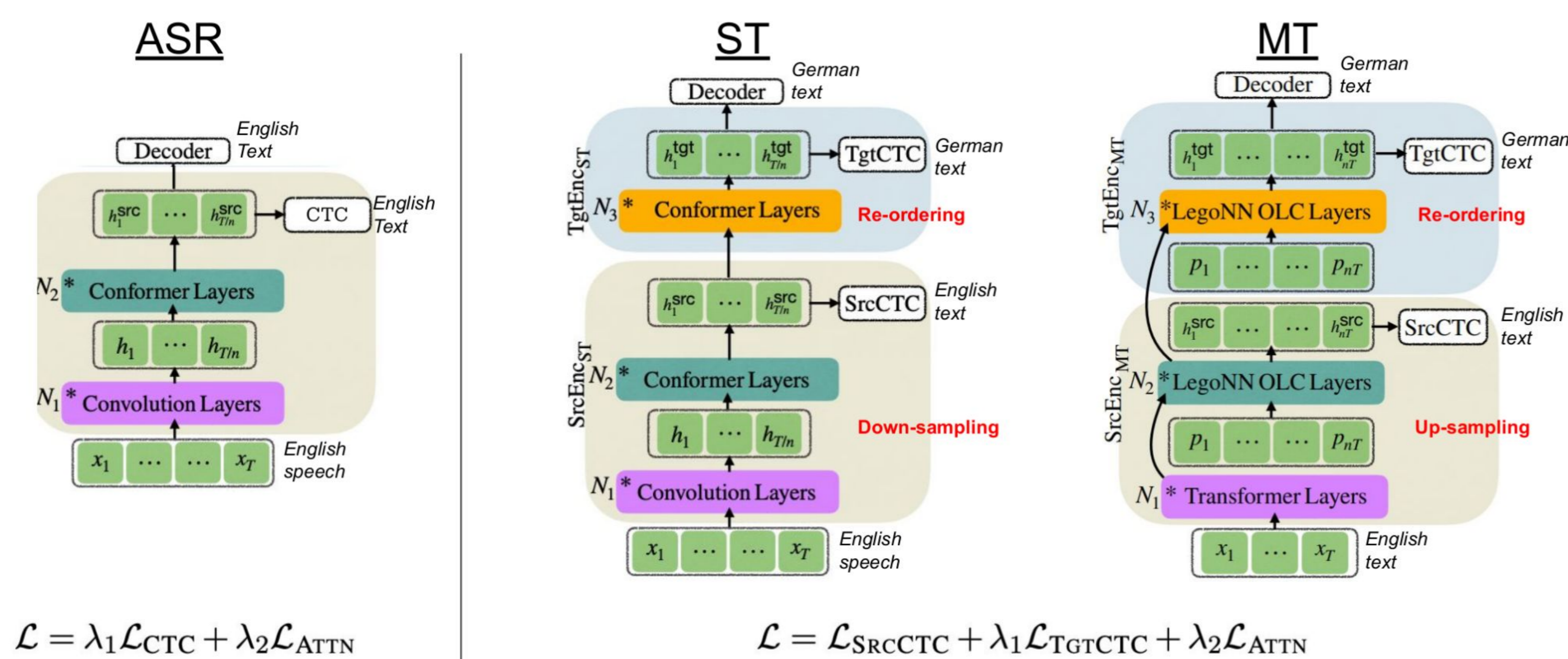
CTC	ATTENTION
$P_{CTC}(Y X) \triangleq \sum_{Z \in \mathcal{Z}} \prod_{t=1}^T P(z_t X, z_{1:t-1})$	$P_{Attn}(Y X) \triangleq \prod_{l=1}^L P(y_l y_{1:l-1}, X)$
Hard Alignment Criterion only allows monotonic alignments of inputs to outputs	Soft Alignment Flexible attention-based input-to-output mappings may overfit to irregular patterns
Conditional Independence Assumes that there are no dependencies between each output unit given the input	Conditional Dependence Locally normalized models with output dependency exhibit label/exposure biases
Input-Synchronous Emission Each input representation emits exactly one blank or non-blank output token	Autoregressive Generation Need to detect end-points and compare hypotheses of different length in beam search

Weaker for translation
(Gu+ 2021, Huang+ 2022)

But attention is not perfect
(Murray+ 2018, Hannun+ 2019, Watanabe+ 2018)

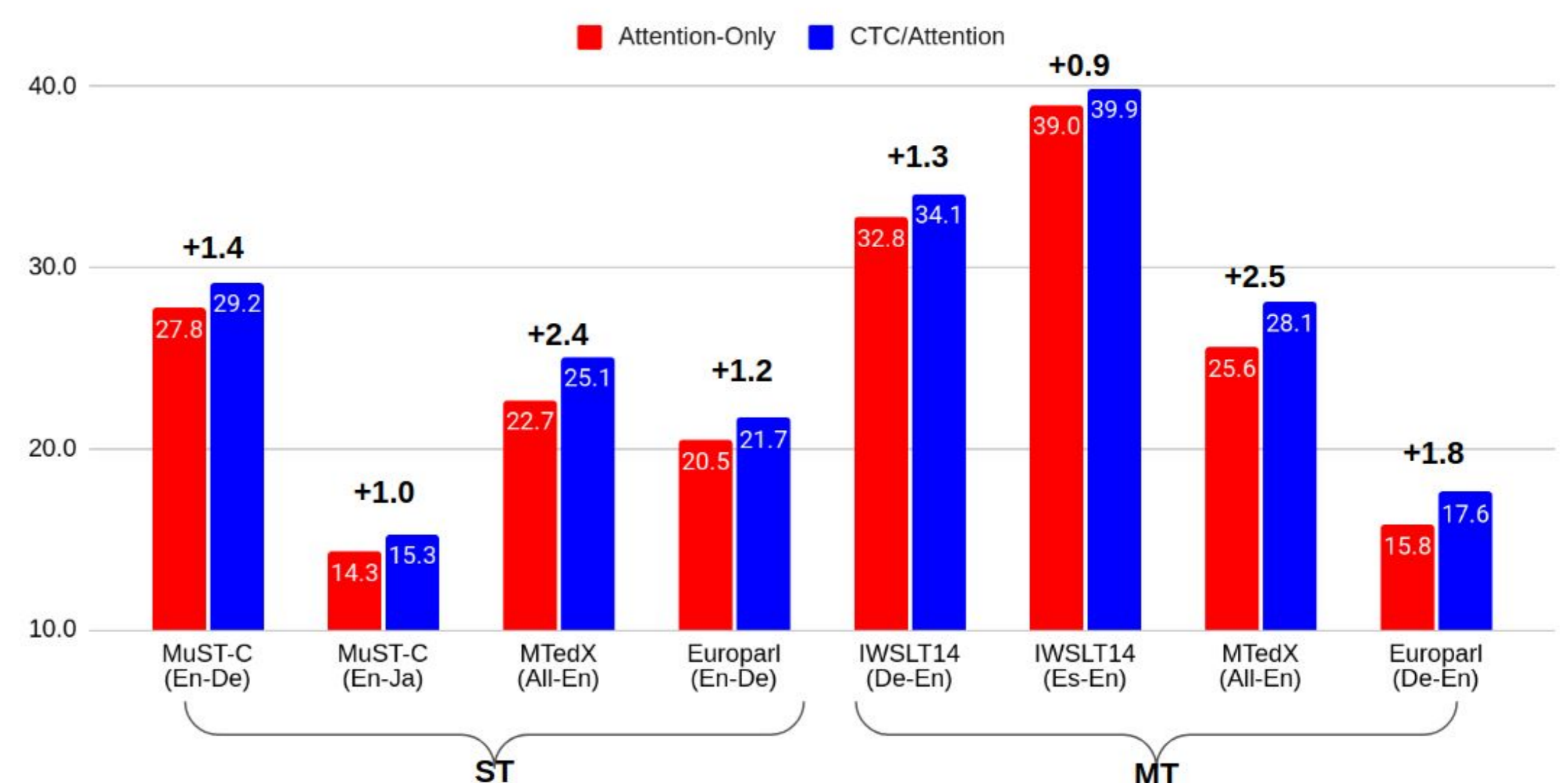
Joint Training with Hierarchical CTC

- 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**



Joint CTC/Attention: Results

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



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 .

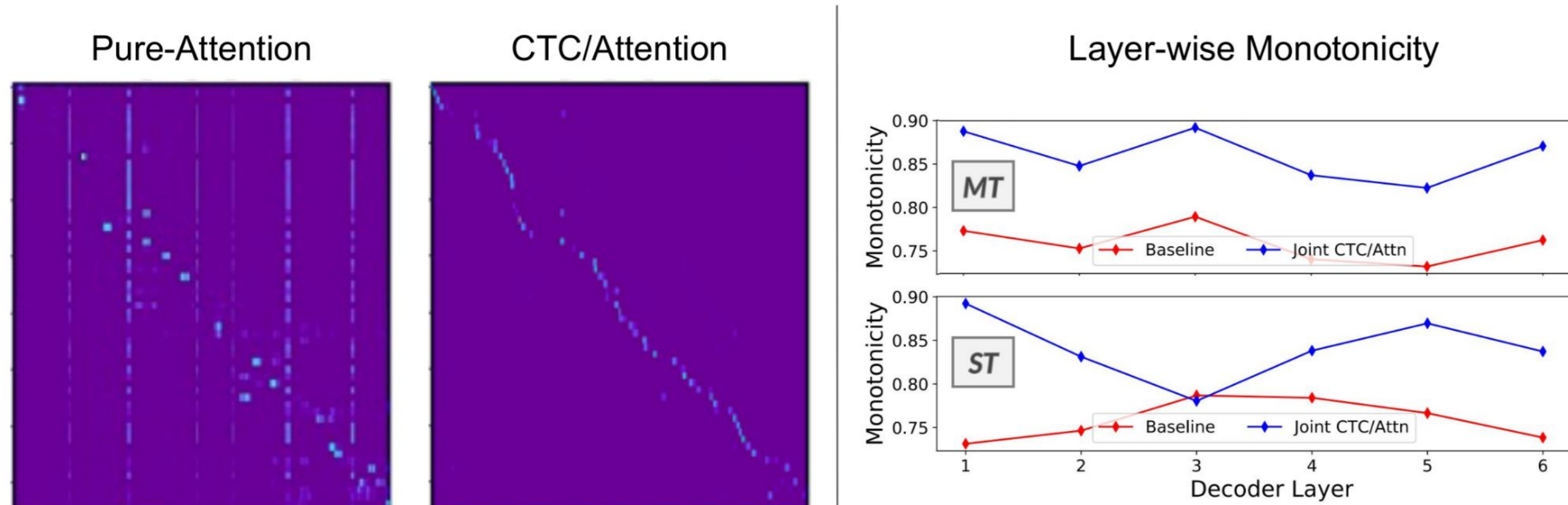
```

1: procedure OUTPUTSTEP(prtHs, X, l, p, maxL)
2:   newPrts = {}; endHs = {}
3:   for  $y_{1:l-1} \in \text{prtHs}$  do
4:     [attnCnds] = top-k( $P_{\text{attn}}(y_l|X, y_{1:l-1})$ , k = p)
5:     for  $c \in \text{attnCnds}$  do
6:        $y_{1:l} = y_{1:l-1} \oplus c$ 
7:        $\alpha_{\text{ctc}} = \text{CTCScore}(y_{1:l}, X_{1:l})$ 
8:        $\alpha_{\text{attn}} = \text{AttnScore}(y_{1:l}, X_{1:l})$ 
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 maxL) then
12:        endHs[ $y_{1:l}$ ] =  $P_{\text{beam}}(\cdot)$ 
13:      else
14:        newPrts[ $y_{1:l}$ ] =  $P_{\text{beam}}(\cdot)$ 
15:      end if
16:    end for
17:  end for
18:  return newPrts, endHs
19: end procedure
    
```

Hypothesis Expansion: Choose top- p candidates per hypothesis (e.g. $p=1.5^b$) from attention posteriors
 Joint Scoring: Interpolate label sequence likelihoods from attention and CTC
 End Detection: End loop based on conditions from attention (CTC's role in end-detection is implicit)

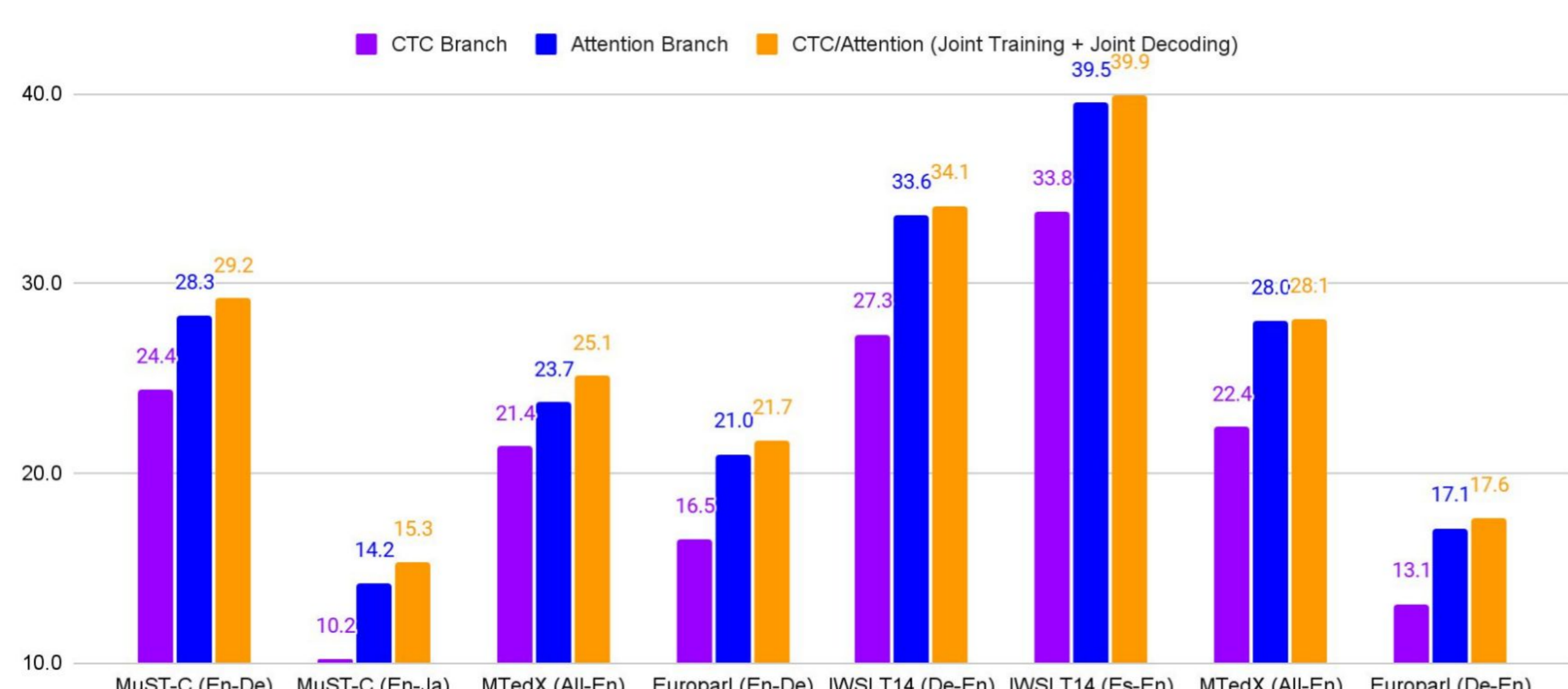
Reduced Soft Alignment Burden

- Hard alignment + soft alignment**
 - Conjecture*: hard alignment objective produces stable encoder representations allowing the decoder to **more easily learn soft alignment patterns during training**
 - Concern*: can CTC encoders perform input-to-output mappings for translation?
 - Finding*: joint training results in more regular, diagonal source-attention patterns



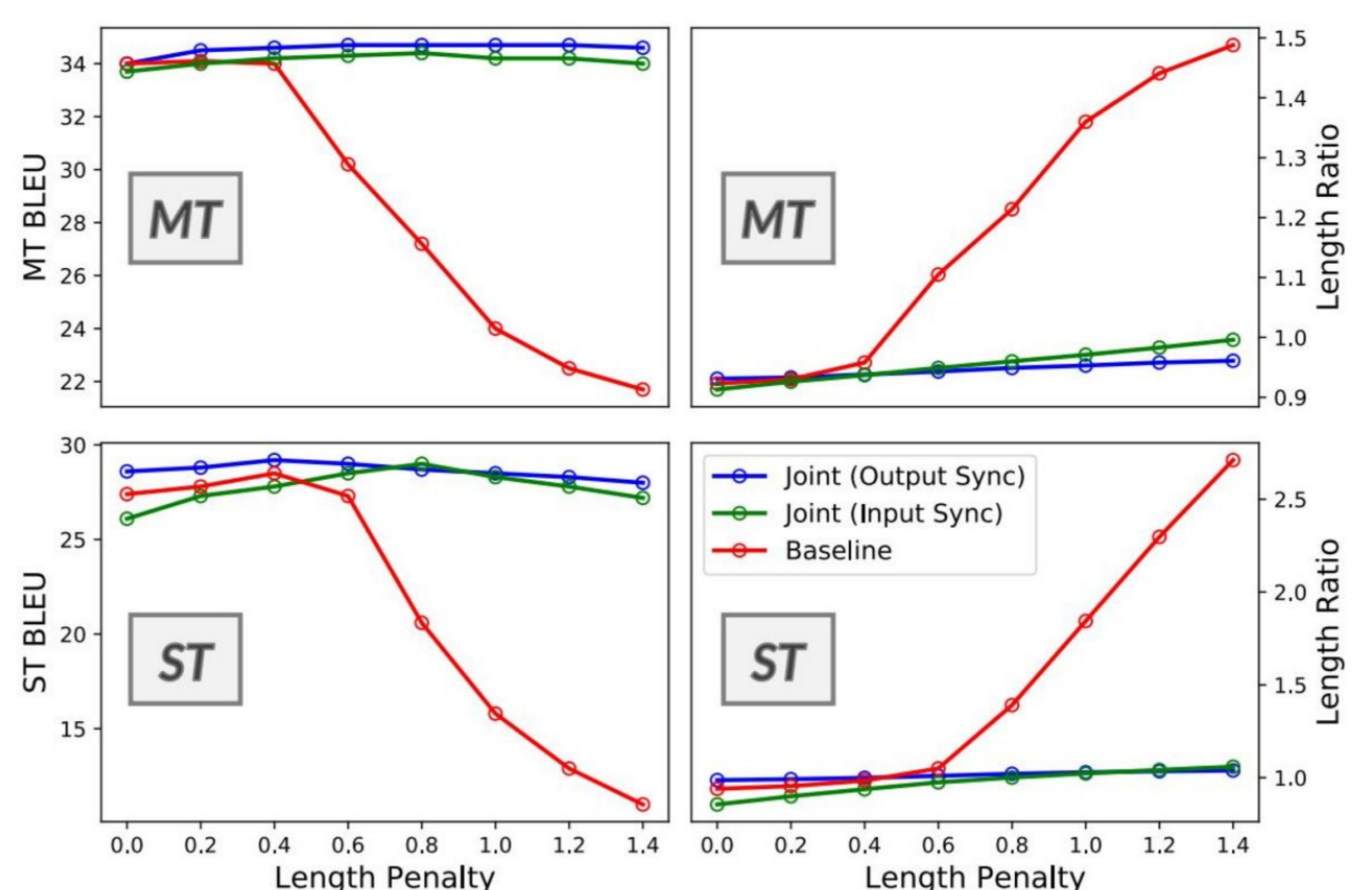
Positive Ensembling Effect

- Conditional independence + conditional dependence**
 - Conjecture*: use of conditionally independent likelihoods in joint scoring **eases the exposure/label biases** from conditionally dependent likelihoods **during decoding**
 - Concern*: does CTC translation quality lag too far behind attention to be useful?
 - Finding*: even weak CTC models provide a positive ensembling effect
 - 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**



Robust End-Detection

- Input synchronous emission + autoregressive generation**
 - Conjecture*: input-synchronous emission determines output length based on input length **counteracting the autoregressive end-detection problem during decoding**
 - Concern*: is the alignment information from CTC translation models reasonable?
 - Finding*: Joint CTC/attention decoding has low length penalty elasticity
 - Sanity check: input-sync beam search, where CTC plays the
 - Pure-attention models are highly sensitive to length penalty \rightarrow easily overtuned



TLDR

- What did we do?
 - Applied joint CTC/attention to MT/ST with only minor changes from ASR
- Why did we do it?
 - CTC may be weaker for translation, but attention is not perfect either
 - CTC and attention have complementary properties
- How can I try it too?
 - Code and recipes are open-sourced in ESPnet

