

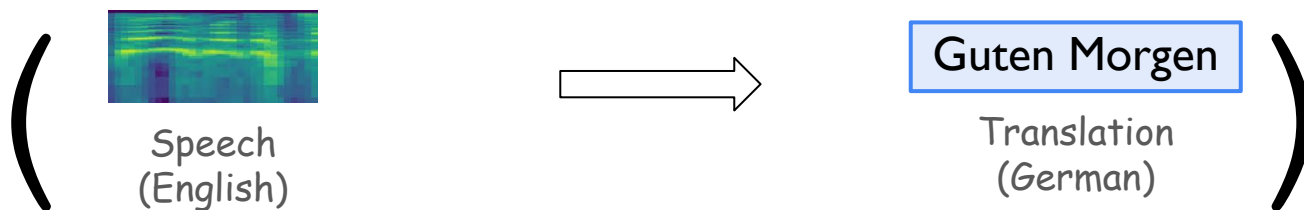
# Revisiting End-to-End Speech-to-Text Translation From Scratch

Biao Zhang<sup>1</sup>, Barry Haddow<sup>1</sup>, Rico Sennrich<sup>2,1</sup>

<sup>1</sup>University of Edinburgh <sup>2</sup>University of Zurich



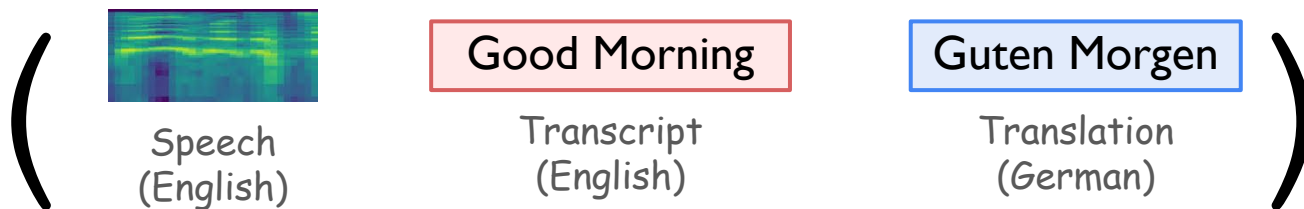
# End-to-End Speech-to-Text Translation Is Challenging



E2E ST aims at translating speech directly to a foreign text **without any intermediate outputs**, e.g., transcript

- ✗ Implicitly modeling ASR and MT via a single model is difficult
- ✗ Performance of the direct ST baseline lags far behind the cascade

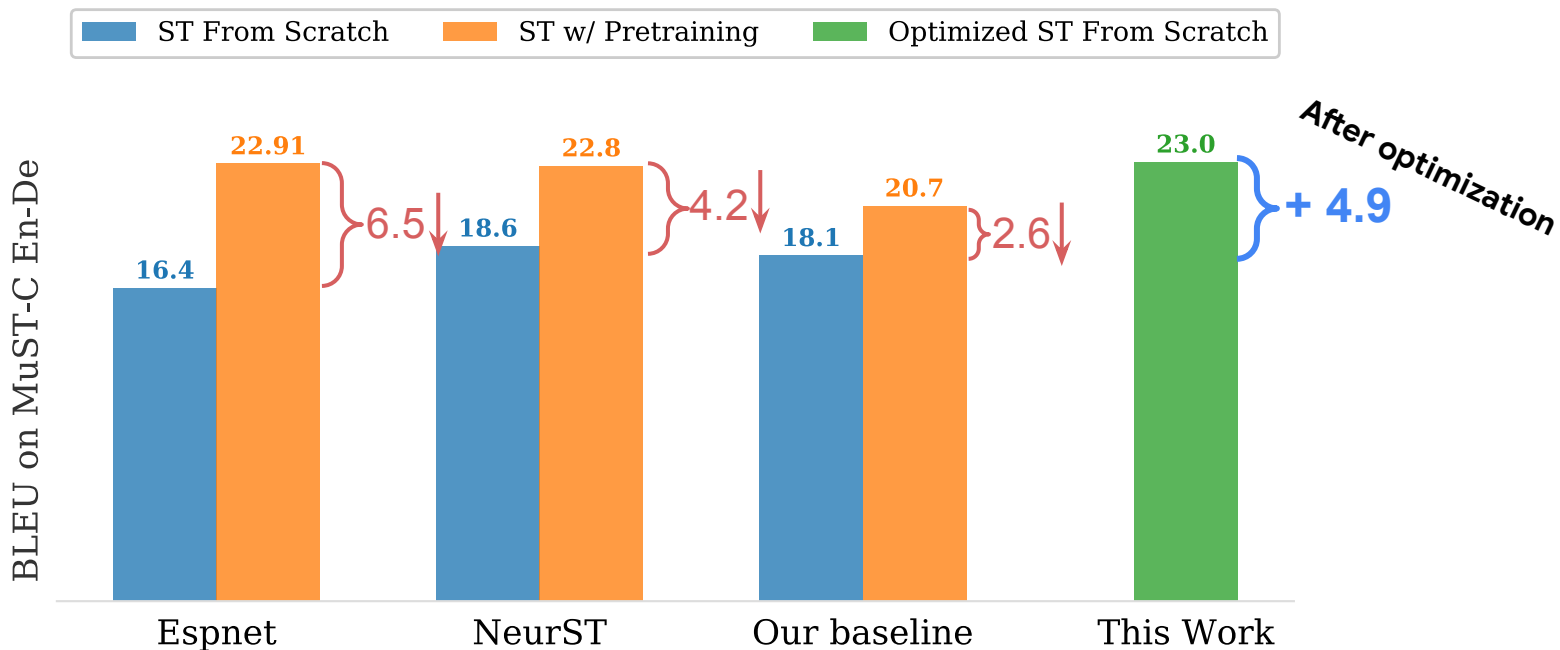
# Rescue: ASR/MT Pretraining and ST Finetuning



**Step 1:** pretrain ST encoder/decoder with ASR/MT using **transcripts**

**Step 2:** finetune the model on direct speech-translation pairs

# E2E ST Without Transcript Performs Poorly, Really?



ST from Scratch: train ST on **speech-translation pairs alone** and from scratch

# Why Revisit E2E ST *From Scratch*?

- ✓ Improve our understanding of pretraining in E2E ST
  - when and where ASR/MT pretraining really matters
- ✓ Transcript is not always available
  - > 3000 languages in our world have no written form
- ✓ Simplify the training pipeline, and develop useful inductive biases
  - using two-stage training complicates the modeling

# Improving E2E ST Towards Training From Scratch

## CTC Regularization

- Use translation as CTC labels
- No transcripts are used

## Parameterized Distance Penalty (PDP)

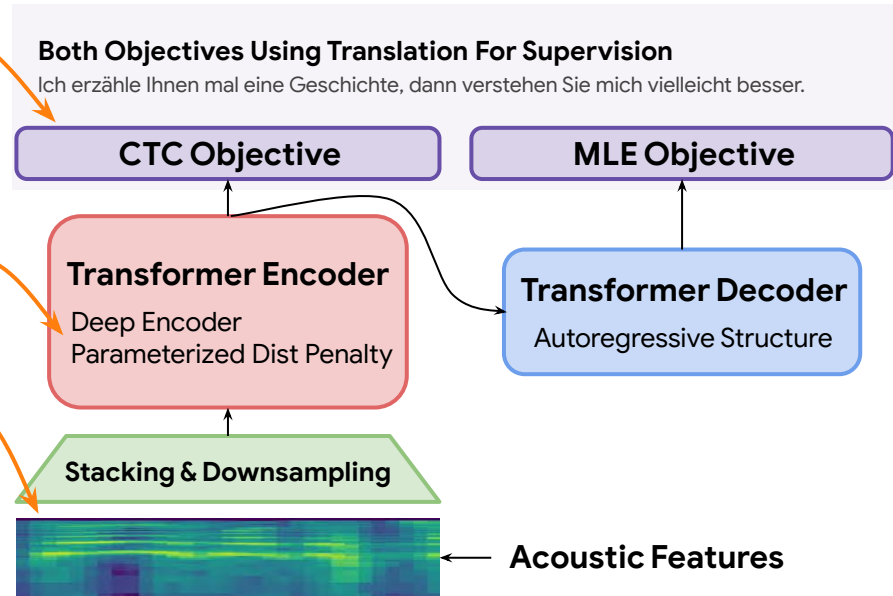
- Add freedom in local attention modeling

## Neural Acoustic Feature Modeling

- Use raw waveform to retain local details

## Hyperparameter Tuning

- Beam search; Model depth/width



Using speech-translation pairs alone **with no transcripts**

# Improved Results: Different Techs Are Complementary

System	BLEU		Avg
NeurST (pretrain-finetune)	22.8		24.9
Baseline	18.1		-
+ hyperparameter tuning	21.1	+3.0	-
+ PDP (R=512)	21.8	+0.7	-
+ CTC regularization	22.7	+0.9	-
+ neural acoustic model	<b>23.0</b>	+0.3	<b>25.2</b>

Test performance on MuST-C En-De and average results on the other language pairs

Note all our models are trained with speech-translation pairs alone

# To Summarize

- The quality gap between ST with and without transcripts is overestimated
- We figure out a set of practices for improving ST from scratch
  - deep post-LN encoder, wider feed-forward layer, ST-based CTC regularization and parameterized distance penalty, neural acoustic feature modeling
- Pretraining still matters
  - low-resource regime and large-scale external data available

Paper: <https://arxiv.org/abs/2206.04571>

Code: [https://github.com/bzhangGo/st\\_from\\_scratch](https://github.com/bzhangGo/st_from_scratch)

