Research and Practice of Simultaneous Speech Translation in Huawei Noah's Ark Lab

Prof. Dr. Qun Liu

Chief Scientist of Speech and Language Computing



Huawei Noah's Ark Lab

10 July 2020, AutoSimTrans Workshop

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

2 A General SST Framework

3 Adapting ASR & NMT to SST

Implementation and Optimization

5 Video Demonstrations



・ロト・西ト・ヨト・ヨー シック

Content

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

A General SST Framework

3 Adapting ASR & NMT to SST

Implementation and Optimization

5 Video Demonstrations

HUAWEI

・ロト・「日・・日・・日・ のへの

Content

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

• Along with the success of deep learning in MT, ASR and TTS, Simultaneous Speech Translation (SST) became commercially feasible:



Microsoft Skype Translator







Tencent Conference Interpretor



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Potential SST Scenarios in Huawei

- Huawei has a very long product line.
- There are many potential scenarios where SST can be applied:



Mobile phones





Conference systems

Customer service systems



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References

SST Research in Huawei

- We started our research on SST in 2019.
- We propose a General SST Framework
- We adopt techniques to adapt our existing ASR and NMT systems to a SST system

▲□▶▲□▶▲□▶▲□▶ □ の000

- We implement and optimize our SST systems in the following scenarios:
 - · Conference speech translator on the cloud
 - · Travel assistant on mobile phones



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

2 A General SST Framework

3 Adapting ASR & NMT to SST

Implementation and Optimization

5 Video Demonstrations

HUAWEI

▲□▶▲□▶▲□▶▲□▶ □ ● ●

Content

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References

A General Framework for Simultaneous NMT [1]

- Motivation:
 - Adapt an existing NMT system into a simultaneous translation system, without retraining the NMT model.

▲□▶▲□▶▲□▶▲□▶ □ の000

- Components:
 - An external ASR system
 - A pretrained NMT System
 - An input buffer
 - An output buffer
- Preprint: https://arxiv.org/abs/1911.03154



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

A General Framework for Simultaneous NMT [1]

- Procedure as two nested loops:
 - Outer loop: update input buffer with ASR streaming output
 - Inner loop: update output buffer with NMT model.
 - The input buffer is updated only when the inner loop stops to update the output buffer.



▲□▶▲□▶▲□▶▲□▶ □ の000

• Preprint: https://arxiv.org/abs/1911.03154



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References

Key Problems

- How to continue translation
 - Streaming source prefix \rightarrow rebuild encoder hidden states
 - Force decoding with the NMT model
- How to terminate
 - Terminate when predicting the EOS token
 - Terminate when the system is not confident and prefer to wait for more input tokens

▲□▶▲□▶▲□▶▲□▶ □ の000

- A Controller is used to determinate when to terminate
- · We define two controllers for this purpose



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References

Length and EOS (LE) Controller

- The inner loop stops to generate output tokens when:
 - A EOS token is generated;
 - The output buffer is full;
 - The delay between the output sequence and the input sequence is smaller than a predefined threshold (e.g. 3 words).

▲□▶▲□▶▲□▶▲□▶ □ の000

inner step	input buffer	output buffer	
0	晓莹 你		
1	晓莹 你	xiaoying	→ Stop
2	晓莹 你	xiaoying you	



Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

HUAWEI

Trainable (TN) Controller

- The inner loop stops to generate output tokens when:
 - The output buffer is full;
 - The TN Controller output a CONTINUE signal.



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

TN Controller Training: Policy Gradient

• Fix the NMT model and train the TN controller with policy gradient:

 $\begin{array}{ll} \textbf{Objective:} \qquad \mathcal{J}_{rl} = \mathbb{E}_{\pi_{\phi}}(\sum_{t=1}^{T}r_{t}) \\ \textbf{Gradient:} \qquad \nabla_{\phi}\mathcal{J}_{rl} = \mathbb{E}_{\pi_{\phi}}\left[\sum_{t=1}^{T}R_{t}\nabla_{\phi}\log\pi(\sigma_{t}|\cdot;\phi)\right] \\ \textbf{where} \qquad R_{t} = \sum_{k=t}^{T}r_{k} \end{array}$

▲□▶▲□▶▲□▶▲□▶ □ の000



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References

TN Controller Training: Reward

• The reward function trades off quality and delay:

Reward $r_t = r_t^Q$	$r^{D} + \alpha \cdot r^{D}_{t}$
Quality $r_t^Q = \begin{cases} \Delta \text{BLEU}^0(\mathbf{y}, \hat{\mathbf{y}}, t) & t < T \\ \text{BLEU}(\mathbf{y}, \hat{\mathbf{y}}) & t = T \end{cases}$ where $\Delta \text{BLEU}^0(\mathbf{y}, \hat{\mathbf{y}}, t)$ $= \text{BLEU}^0(\mathbf{y}^t, \hat{\mathbf{y}}) - \text{BLEU}^0(\mathbf{y}^{t-1}, \hat{\mathbf{y}})$	$\begin{split} \mathbf{Delay} \\ 0 < d^{AL}\left(\mathbf{x}, \mathbf{y}\right) &= \frac{1}{\tau_e} \sum_{\tau=1}^{\tau_e} l(\tau) - \frac{\tau - 1}{\lambda}, \\ r_t^D &= \left\{ \begin{array}{ll} 0 & t < T \\ -\lfloor d^{AL}(\mathbf{x}, \mathbf{y}) - d^* \rfloor_+ & t = T \end{array} \right\} \end{split}$



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへぐ

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video

Demonstrations

References



Experiments



<□▶ < □▶ < □▶ < □▶ < □▶ = □ の Q ()

- baselines:
 - test-time-waitk: Ma et al. [2018]
 - SL: Zheng et al. [2019]
 - RWAgent: Gu et al. [2017]

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection an Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Background

A General SST Framework

3 Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

Content

・ロト・日本・日本・日本・日本・日本

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection as Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

. . .

Content

3 Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation in NMT Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

・ロト・国ト・ヨト・ヨー うへぐ

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection an Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video

Demonstrations

References

Terminology Translation in NMT

- Terminology (including user-defined terms, entities, acronyms, etc.) translation is crutial in many scenarios.
- It is not a trivial task in NMT because there is no explicit phrase table used in NMT process.
- Solutions:
 - · Without word alignment: Constrained Decoding
 - Time consuming or inaccurate [6; 9].
 - With word alignment: terminology replacement in automatic post-editing
 - Word alignment induced from NMT is not accurate enough [5; 15; 8; 3].

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection ar Correction

Data Augmentation for Robust Speech Translatio

Implementation and Optimization

Video Demonstrations

References

Accurate Alignment for Terminology Translation [2]

- Our observation:
 - The hidden states at decoding step i + 1 are better representing target word y_i than the hidden states at step *i* for inducing word alignment



= nar

• Preprint: https://arxiv.org/abs/2004.14837

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection ar Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Accurate Alignment for Terminology Translation [2]

- Our work:
 - We introduce SHIFT-ATT, a pure interpretation method to induce alignments from attention weights of vanilla Transformer.
 - To select the best layer to induce alignments, we propose a surrogate layer selection criterion, to ensure the induced word alignments agree best in both translation directions, without manually labelled word alignments.

▲□▶▲□▶▲□▶▲□▶ □ の000

- We further propose SHIFT-AET, which extracts alignments from an additional alignment module.
- Preprint: https://arxiv.org/abs/2004.14837

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection a Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Experiments and Results

Mathod		Fulle	de-en		fr-en		ro-en				
Metriod	miler.	Func	$de \rightarrow en$	$en{\rightarrow}de$	bidir	fr→en	$en{\rightarrow} fr$	bidir	$ro \rightarrow en$	$en{\rightarrow}ro$	bidir
Statistical Methods											
FAST-ALIGN (Dyer et al., 2013)	-	Y	28.5	30.4	25.7	16.3	17.1	12.1	33.6	36.8	31.8
GIZA++ (Brown et al., 1993)	-	Y	18.8	19.6	17.8	7.1	7.2	6.1	27.4	28.7	26.0
			Neura	Methods							
NAIVE-ATT (Garg et al., 2019)	Y	Ν	33.3	36.5	28.1	27.5	23.6	16.0	33.6	35.1	30.9
SMOOTHGRAD (Li et al., 2016)	Y	Ν	36.4	45.8	30.3	25.5	27.0	15.6	41.3	39.9	33.7
SD-SMOOTHGRAD (Ding et al., 2019)	Y	Ν	36.4	43.0	29.0	25.9	29.7	15.3	41.2	41.4	32.7
PD (Li et al., 2019)	Y	Ν	38.1	44.8	34.4	32.4	31.1	23.1	40.2	40.8	35.6
ADDSGD (Zenkel et al., 2019)	N	Ν	26.6	30.4	21.2	20.5	23.8	10.0	32.3	34.8	27.6
MTL-FULLC (Garg et al., 2019)	N	Y	-	-	20.2	-	-	7.7	-	-	26.0
Our Neural Methods											
Shift-Att	Y	N	20.9	25.7	17.9	17.1	16.1	6.6	27.4	26.0	23.9
Shift-AET	N	Ν	16.1	19.3	15.5	10.3	10.5	5.0	22.4	23.7	21.2

Table 1: AER on the test set with different alignment methods. *bidir* are symmetrized alignment results. The column Inter. represents whether the method is an interpretation method that can extract alignments from a pretrained vanilla Transformer model. The column Fullc denotes whether full target sentence is used to extract alignments at test time. The lower AER, the better. We mark best bidir interpretation results of vanilla Transformer with underlines, and best bidir results among all with boldface.

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection a Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Content

3 Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation in NMT Punctuation Prediction

Disfluency Detection and Correction Data Augmentation for Robust Speech Translati

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 善臣 - のへで

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection an Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Punctuation Prediction

- Models
 - Baseline: BiLSTM Yi et al. [13]
 - BERT with Fine-tune
 - TinyBERT [7] (Distilled from BERT)
- Data: IWSLT2011 TED Corpus
- Labels:

4 labels	3 labels
Comma(,)	
Period(.)	Period(.)
Question(?)	Question(?)
None(O)	None(O)

uh, gripping the onion like a you know, like a tennis ball, holding it together in place.



uh gripping the onion like a you know like a tennis ball holding it together in place

▲□▶▲□▶▲□▶▲□▶ □ の000

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video

Demonstrations

References

Punctuation Prediction

	Model	Comma			Period		Question			Overall			
		P(%)	R(%)	$F_{1}(\%)$	P(%)	R(%)	$F_{l}(\%)$	P(%)	R(%)	$F_{I}(\%)$	P(%)	R(%)	$F_{I}(\%)$
	DNN	58.1	35.8	44.3	62.1	64.8	63.4	60.5	48.9	54.1	60.2	49.8	53.9
Dof	T-BRNN-pre [17]	65.5	47.1	54.8	73.3	72.5	72.9	70.7	63.0	66.7	70.0	59.7	64.4
Ref.	BLSTM-CRF	58.9	59.1	59.0	68.9	72.1	70.5	71.8	60.6	65.7	66.5	63.9	65.1
	Teacher-Ensemble	66.2	59.9	62.9	75.1	73.7	74.4	72.3	63.8	67.8	71.2	65.8	68.4
	DNN	47.5	32.3	38.5	58.3	60.5	59.4	57.1	46.8	51.4	54.3	46.5	49.8
ASD	T-BRNN-pre [17]	59.6	42.9	49.9	70.7	72.0	71.4	60.7	48.6	54.0	66.0	57.3	61.4
ASK	BLSTM-CRF	55.7	56.8	56.2	68.7	71.5	70.1	63.8	53.4	58.1	62.7	60.6	61.5
	Teacher-Ensemble	60.6	58.3	59.4	71.7	72.9	72.3	66.2	55.8	60.6	66.2	62.3	64.1

(Yi, Jiangyan, et al. Interspeech 2017, baseline)

	comma	period	question	other		Overa	11
	acc	acc	acc	acc	pre	rec	f1(macro)
4 labels	0.716	0.84	0.664	0.983	0.755	0.745	0.75
3 labels	_	0.92	0.71	0.99	0.893	0.887	0.894
3 labels vs. 4 labels					18%	19%	19%
3 labels vs. baseline-ref					135%	↑42%	131%

	precision	recall	f1(macro)	train speed	infer speed
teach model (bert base)	0.893	0.887	0.894		
student model (tiny-bert)	0.857	0.837	0.867	X12	X5
			-3%		

(3 labels VS 4 labels VS baseline, 1 indicate relative improvements) (Knowledge Distillation for Inference acceleration)

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Content

3 Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation in NMT Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video

Demonstrations

References

Disfluency Detection and Correction

• English Switchboard Corpus:

I want a flight [$\underbrace{to Boston}_{RM} + \underbrace{\{um\}}_{IM}$ $\underbrace{to Denver}_{RP}$]

Figure 1: A sentence from the English Switchboard corpus with disfluencies annotated. RM=Reparandum, IM=Interregnum, RP=Repair. The preceding RM is corrected by the following RP.

(from Wang et al. [11])

▲□▶▲□▶▲□▶▲□▶ □ の000

- IWSLT2020 TED Corpus (Chinese):
 - ASR: 黑胡桃木, 黑胡桃木是我们店内最好的木材。
 - REF:黑胡桃木是我们店内最好的木材。
 - ASR: 啊它的性质很好不容易开裂,
 - REF: 它的性质很好不容易开裂,

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Disfluency Detection and Correction

- BERT with Multitask Fine-tuning
 - follow Wang et al. [11]
- Task1: Sequence labeling
 - RM and IM \rightarrow D (DELETE)
 - RP and Others \rightarrow O (OTHER)
- Task2: Sentence pair classification
 - Smooth/Disfluent \rightarrow DF_0
 - Disfluent/Smooth \rightarrow DF_1
- Training Data: English Switchboard
- Data augmentation:
 - Randomly insert IM
 - Swap RM & RP





▲□▶▲□▶▲□▶▲□▶ □ の000

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection and Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video

Demonstrations

References

Disfluency Detection and Correction

• Experiment (English Switchboard Corpus):

	Р	R	F1
UBT [12]	90.3	80.5	85.1
Semi-CRF [4]	90	81.2	85.4
Bi-LSTM [14]	91.8	80.6	85.9
Transition-based [10]	91.1	84.1	87.5
MTL Self-supervised [11]	93.4	87.3	90.2
Our MTL-DA	91.3	87.7	89.4

• Experiment (IWSLT2020 TED Corpus, Chinese):

TER (Translation Error Rate)	Train	Dev	Test
Baseline	35.2	36.3	33.1
+ Expert Rules (pre-processing)	31.8(-3.4)	31.9(-4.4)	30.1(-3.0)
+ BERT Fine-tuning	28.3(-3.5)	28.3(-3.6)	56.7(-3.4)
+ Dictionary	26.4(-1.9)	25.9(-2.4)	24.9(-1.8)

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection a Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Content

3 Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation in NMT Punctuation Prediction Disfluency Detection and Correction Data Augmentation for Robust Speech Translation

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ● ●

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video

Demonstrations

References

Motivation

- ASR outputs always contain errors which severely harms the quality of the downstream MT system.
- Ideally, fine-tuning the MT system with the pairs of noisy ASR outputs and their translations will benefit.
- However, there is very little parallel data with ASR-output in the source side.
- Our solution:
 - Using ASR training data (SPEECH, TEXT) and an existing ASR system to train a GPT-2 system to generate ASR-like texts from normal texts.
 - Using the obtained system to transfer the source text of a bilingual courpus to ASR-like text.

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detecti Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Robust MT Training



◆□ > ◆□ > ◆臣 > ◆臣 > ○ 国 ○ ○ ○ ○

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection a Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Noisy Data Filtering

- The noisy data generated by the GPT-2 system may be TOO noisy.
- We use an edit rate threshold on pronunciations to filter the GPT-2 outputs.
 - Sentence → pronunciation using cmudict¹.
 - Edit Rate: Edit distance normalized by the original length.

clean: The priest **tied the knot**. ['DH', 'AH0', 'P', 'R', 'IY1', 'S', 'T', 'T', 'AY1', 'D', 'DH', 'AH0', 'N', 'AA1', 'T']

noisy1: The priest **told the knot**. ['DH', 'AHO', 'P', 'R', 'IY1', 'S', 'T', 'T', 'OW1', 'L', 'D', 'DH', 'AHO', 'N', 'AA1', 'T'] D=2 ER = 0.4

noisy2: The priest to you, you ['DH', 'AHO', 'P', 'R', 'IY1', 'S', 'T', 'T', 'UW1', 'Y', 'UW1', 'Y', 'UW1'] D=7 ER = 1.4

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Accurate Alignment for Terminology Translation

Punctuation Prediction

Disfluency Detection Correction

Data Augmentation for Robust Speech Translation

Implementation and Optimization

Video Demonstrations

References

Experiments

- \bullet Data: EN ${\rightarrow}$ ZH, Huawei STW meeting dataset and MSLT v1.1 2
- Models:
 - transformer-small, emb size:256, intermediate size:1024, attention heads:4, training data: 2M pairs
 - transformer-big, emb size:1024, intermediate size:4096, attention heads:16, training data: 1.4B pairs

		ST	W	MSLT			
		dev	test	dev	test		
amall	clean	36.29	36.84	29.61	31.21		
small	robust	36.93(+0.64)	37.69(+0.85)	30.69(+1.08)	32.56(+1.35)		
hia	clean	44.18	44.03	33.75	34.29		
big	robust	45.85(+1.67)	45.46(+1.43)	34.65(+0.9)	35.23(+0.96)		

2https://github.com/MicrosoftTranslator/MSLT-Corpus 😽 🚛 🖉 🔊 🧟

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

5 Video Demonstrations



Content



System Architecture



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへで



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References



Adapting ASR from Cloud to Mobile Devices



Encoder

介

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Conv-Transformer Transducer

• Streamble: unidirectional Transformer for audio encoding







Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Conv-Transformer Transducer

• Low Frame Rate: interleaved convolutions for gradually downsampling





◆□▶ ◆□▶ ◆三▶ ◆三▶ ○三 のへで

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Conv-Transformer Transducer

• Reduced Computation Cost: self-attention with fixed context window







Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Conv-Transformer Transducer

 Reduced Computation Cost: hidden state reuse with relative positional encoding







A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References





Adapting NMT from Cloud to Mobile Devices

- 1. Use RNN decoder instead of Transformer decoder.
- 2. Decrease layer number to compress model size.
- 3. Share parameter for different language.
- 4. Use smaller vocab size.
- Model quantization for model storage to lower hard disk usage.
- Model quantization during training to limit parameter value range which enable more stable quantized inference.
- 7. Model quantization during inference to enhance using experience.
- 8. Use cloud model to generate more data for training on-device model.
- 9. Powered by Bolt framework.



On-device Model



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Model Simplication on Model Devices

- Use RNN layers instead of masked multi-head self-attention layers;
- Decrease number of decoder model stacks, from N layers on cloud, to M layers on devides;

▲□▶▲□▶▲□▶▲□▶ □ のQ@

- Share model parameters for different languages;
- Use smaller vocab size, from x to y.



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Background

A General SST Framework

3 Adapting ASR & NMT to SST

Implementation and Optimization

5 Video Demonstrations



Content

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References



Video Demonstrations





▲□▶▲圖▶▲圖▶▲圖▶ 圖 のQ@

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References



References I

- [1] Yun Chen, Liangyou Li, Xin Jiang, Xiao Chen, and Qun Liu. How to do simultaneous translation better with consecutive neural machine translation?, 2019.
- [2] Yun Chen, Yang Liu, Guanhua Chen, Xin Jiang, and Qun Liu. Accurate word alignment induction from neural machine translation, 2020.
- [3] Shuoyang Ding, Hainan Xu, and Philipp Koehn. Saliency-driven word alignment interpretation for neural machine translation. In Proceedings of WMT 2019), pages 1–12, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5201. URL https://www.aclweb.org/anthology/W19-5201.
- [4] James Ferguson, Greg Durrett, and Dan Klein. Disfluency detection with a semi-markov model and prosodic features. In Proceedings of the 2015 Conference of NAACL-HLT 2015, pages 257–262, 2015.
- [5] Sarthak Garg, Stephan Peitz, Udhyakumar Nallasamy, and Matthias Paulik. Jointly learning to align and translate with transformer models. In *Proceedings of EMNLP-IJCNLP 2019*, pages 4453–4462, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1453. URL https://www.aclweb.org/anthology/D19-1453.
- [6] Chris Hokamp and Qun Liu. Lexically constrained decoding for sequence generation using grid beam search. arXiv preprint arXiv:1704.07138, 2017.
- [7] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. In proceedings of ICLR 2020, 2020.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

[8] Xintong Li, Guanlin Li, Lemao Liu, Max Meng, and Shuming Shi. On the word alignment from neural machine translation. In Proceedings of ACL 2019, pages 1293–1303, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1124. URL https://www.aclweb.org/anthology/P19-1124.

Qun Liu

Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

References II

- [9] Matt Post and David Vilar. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. arXiv preprint arXiv:1804.06609, 2018.
- [10] Shaolei Wang, Wanxiang Che, Yue Zhang, Meishan Zhang, and Ting Liu. Transition-based disfluency detection using lstms. In Proceedings of EMNLP 2017, pages 2785–2794, 2017.
- [11] Shaolei Wang, Wanxiang Che, Qi Liu, Pengda Qin, Ting Liu, and William Yang Wang. Multi-task self-supervised learning for disfluency detection. In Proceedings of AAAI 2020, 2020.
 - [12] Shuangzhi Wu, Dongdong Zhang, Ming Zhou, and Tiejun Zhao. Efficient disfluency detection with transition-based parsing. In Proceedings of ACL-IJCNLP 2015, pages 495–503, 2015.
 - [13] Jiangyan Yi, Jianhua Tao, Zhengqi Wen, Ya Li, et al. Distilling knowledge from an ensemble of models for punctuation prediction. In Proceedings of Interspeech 2017, 2017.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

- [14] Vicky Zayats, Mari Ostendorf, and Hannaneh Hajishirzi. Disfluency detection using a bidirectional lstm. *arXiv preprint arXiv:1604.03209*, 2016.
- [15] Thomas Zenkel, Joern Wuebker, and John DeNero. End-to-end neural word alignment outperforms giza++, 2020.



Background

A General SST Framework

Adapting ASR & NMT to SST

Implementation and Optimization

Video Demonstrations

References

Thank for listening!



Team members:

CHEN Xiao, CHEN Yun, CUI Tong, DENG Yao, FU Xu, ZHANG Guchun, GUO Weisheng, HU Wenchao, HUANG Wenyong, JIANG Xin, LI Liangyou, LIN Fuguo, LIN Xia, LIU Jianfeng, LIU Kai, LIU Qun, LIU Xin, PENG Wei, WANG Minghan, XIAO Jinghui, XIA Hairong, YANG Hao, ZHOU Huan.

