ICT's System for AutoSimTrans 2021: Robust Char-Level Simultaneous Translation

Shaolei Zhang^{1,2}, **Yang Feng**^{1,2*}

¹Key Laboratory of Intelligent Information Processing Institute of Computing Technology, Chinese Academy of Sciences (ICT/CAS) ² University of Chinese Academy of Sciences, Beijing, China {zhangshaolei20z, fengyang}@ ict.ac.cn







Method

Experiments



Pipeline of simultaneous interpretation

□ Automatic Speech Recognition (ASR) \rightarrow simultaneous translation (ST) \rightarrow

Text-to-Speech Synthesis (TTS)

- Input of simultaneous translation:
 - Inaccurate, unsegmented.
 - Spoken language domain.

Robustness and Domain adaptability

Streaming Transcript	Translation
大	
大家	
大家好	Hello everyone!
双迎大	
次迎入豕米到込 波而子安玉和近田	
从迎入豕术到这里	nere.



For robustness

- ASR result (streaming transcription): incremental, unsegmented.
- Subword-level segmentation result of the streaming transcription is unstable.
 - Existing method: remove the last to prevent it from being incomplete.





For domain adaptability

- General domain the spoken language domain are quite different:
 - Word order
 - Punctuation
 - Modal particles
 - o ...

Our system

- Robust:
 - Propose the Char-Level Wait-k Policy
- Domain adaptation:
 - Apply data augmentation on spoken language domain.
 - Combine two training methods to enhance the predictive ability.





Method

Experiments



- Char-Level Wait-k Policy
 - Source: character sequence after char-level tokenization.
 - Target: subword sequence after subwordlevel segmentation and BPE.
 - Read / Write policy: waiting for k source characters first, and then reading and writing alternately.

Standard wait-k policy:





Input Sentence		欢迎来到UNIT系统的第12期高级课程。	
Output Sentence		welcome to the 12th advanced course on UNIT system.	
	subword-level MT	欢迎/来到/UN@@/IT/系统/的/第/12@@/期/高级/课程/。	
S.	character-level MT	欢/迎/来/到/U/N/I/T/系/统/的/第/1/2/期/高/级/课/程/。	
·	char-level tokenization	欢/迎/来/到/UNIT/系/统/的/第/12/期/高/级/课/程/。	
T.	subword-level MT	welcome/to/the/12@@/th/advanced/course/on/UNIT/system/.	

Why char-level simultaneous translation?

More robust

• avoid unstable prefixes caused by subword segmentation.

More fine-grained latency

- if one character is enough to express the meaning of a entire word, the ST system does not have to wait for the complete word.
- **Translation quality will not be affected too much**
 - only performs char-level tokenization on the source, and the target retains subword-level tokenization.



- **Domain Adaptation**
 - Depunctuation
 - Source: delete the ending punctuation.
 - Target: unchanged.

Original	各位开发者、各位朋友
	们,大家下午好!
Depunctuation	「各位开发者、各位朋友」
	们,大家下午好

	Data Augmentation
--	-------------------

- For spoken language domain corpus.
- Source: we perform 5 data augmentation operations.
- Target: unchanged.

Original	1957年我到北京上大学	
Add	1057年 我到北京上大学	
Comma	1937年,我到北永上八手	
Ādd		
Tone character		
Сору	1957年我到北北京上大学	
Character		
Replace	1957年我到北经上大学	
Homophone		
Delete	1957年我到北京上大学	
Character		



Training Methods

- Pre-training : general domain MT corpus
 - Multi-path training (Elbayad et al., 2020)
 - Future-guided training (Zhang et al., 2020b)
- **Fine-tuning** : spoken language domain corpus
 - Original training: fix k and use the original prefix-to-prefix framework for training, and train different models for different k.

Maha Elbayad, Laurent Besacier, and Jakob Verbeek. 2020. Efficient wait-k models for simultaneous machine translation. Shaolei Zhang, Yang Feng, and Liangyou Li. 2020b. Future-guided incremental transformer for simultaneous translation.





Method

Experiments



Experiments

Datasets

- **CWMT19** Chinese \rightarrow English: for pre-training.
- **Transcription**: for fine-tuning.
- Dev. Set: for evaluation.

Datasets	Domain	#Sentence Pairs
CWMT19	General	9,023,708
Transcription	Spoken	37,901
Dev. Set	Spoken	956

System setting

- **Offline:** full-sentence MT based on Transformer.
- **Standard Wait-k:** standard subword-level waitk policy.
- Standard Wait-k + rm Last Token: In the inference time, the last token after the word segmentation is remove to prevent it from being incomplete.
- Char-Level Wait-k: our proposed method.



Experiments

Main Result

- Char-Level Wait-k improves about 6 BLEU at low latency (AL=1.10).
- More stable and robust.





Experiments

Ablation Study

- Data processing: 'Depunctuation' and 'Data Augmentation'
- Training methods: 'Future-guided' and 'Multi-path'







Method

Experiments



Conclusion

- ▶ The proposed char-level wait-k policy is more robust.
- Data processing and two training methods improve the spoken language domain adaptability.

▷ For some language pairs with a large length ratio between the source (char) and the target (bpe), we can read multiple characters at each step to deal with the long char-level source. We put this into our future work.



Thanks!

Contact me with: zhangshaolei20z@ict.ac.cn