

Summary

- We evaluated if training data from other language pairs are helpful for the end-to-end speech translation (ST) task
- Directly translate source speech to target languages with a single sequence-to-sequence (S2S) model
 - Many-to-many (M2M)
 - One-to-many (O2M)
- Outperformed the bilingual end-to-end/pipeline speech translation models
- Performed transfer learning to a very low-resource ST task: Mboshi->French (4.4h)
 - Shared representations obtained from multilingual E2E-ST were more effective than those from the bilingual one

Background

◆ End-to-end speech-to-text translation (E2E-ST)

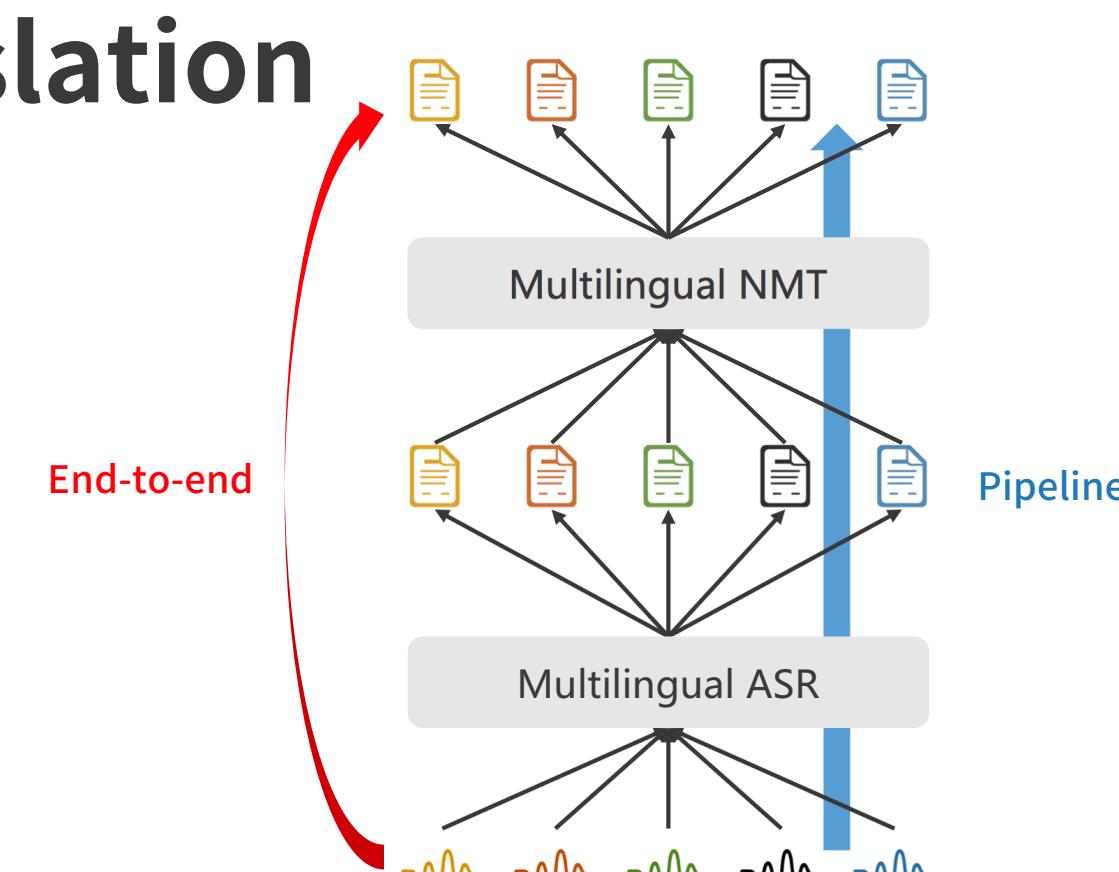
- Simplified architecture (no ASR decoder, LM, and MT encoder) -> reduced number of parameters
- Avoid error propagations from ASR

◆ Multilingual speech-to-text translation

- Applications
 - One-to-many (O2M)
 - e.g., lecture, news reading etc.
 - Many-to-many (M2M)
 - e.g., dialogue, meeting etc.
- Conventional pipeline ST system
 - Language identification->ASR->text normalization->MT

Problems of multilingual pipeline ST systems

- Data sparseness issue for low-resource directions
- Mis-identification of source languages in ASR
- Increased number of parameters
- Need text normalization per source language



◆ Main results

Fisher-CallHome (Es->En)

Model		BLEU	
		Fisher-test	CH-evltest
MT	Bi-NMT	59.6	28.9
	(M2Ma) Multi-NMT	49.5	22.8
	(M2Mb) Multi-NMT	56.7	27.7
	(M2Mc) Multi-NMT	56.2	27.7
E2E-ST	Bi-ST	41.5	14.2
	+ ASR pre-training	45.2	15.4
	(M2Ma) Multi-ST	41.3	15.2
	(M2Mb) Multi-ST	44.2	15.8
	(M2Mc) Multi-ST + ASR pre-training	45.2	16.2
Pipeline-ST	Mono ASR -> Bi NMT	38.6	16.5
	(M2Ma) Multi-ASR -> Bi-NMT	39.2	17.2
	(M2Mb) Multi-ASR -> Bi-NMT	38.9	17.0
	(M2Mc) Multi-ASR -> Bi-NMT	38.5	16.9

LibriSpeech (En->Fr)

Model		BLEU
		test
MT	Bi-NMT	18.3
	(O2M) Multi-NMT	16.2
	(M2Ma) Multi-NMT	12.2
	(M2Mc) Multi-NMT	14.8
E2E-ST	Bi-ST	15.7
	+ ASR pre-training	16.3
	(O2M) Multi-ST	17.2
	(M2Ma) Multi-ST	16.4
Pipeline-ST	(M2Mc) Multi-ST + ASR pre-training	17.3
	Mono ASR -> Bi NMT	15.8
	(O2M) Multi-ASR -> Bi-NMT	16.7
	(M2Ma) Multi-ASR -> Bi-NMT	16.4
	(M2Mc) Multi-ASR -> Bi-NMT	16.7

ST-TED (En->De)

Model		BLEU
		test
MT	Bi-NMT	23.0
	(O2M) Multi-NMT	18.9
	(M2Mb) Multi-NMT	17.5
	(M2Mc) Multi-NMT	17.2
E2E-ST	Bi-ST + ASR pre-training	16.0
	(O2M) Multi-ST	17.6
	(M2Mb) Multi-ST	16.7
	(M2Mc) Multi-ST + ASR pre-training	18.6
Pipeline-ST	Mono ASR -> Bi NMT	18.1
	(O2M) Multi-ASR -> Bi-NMT	18.5
	(M2Mb) Multi-ASR -> Bi-NMT	17.7
	(M2Mc) Multi-ASR -> Bi-NMT	18.1

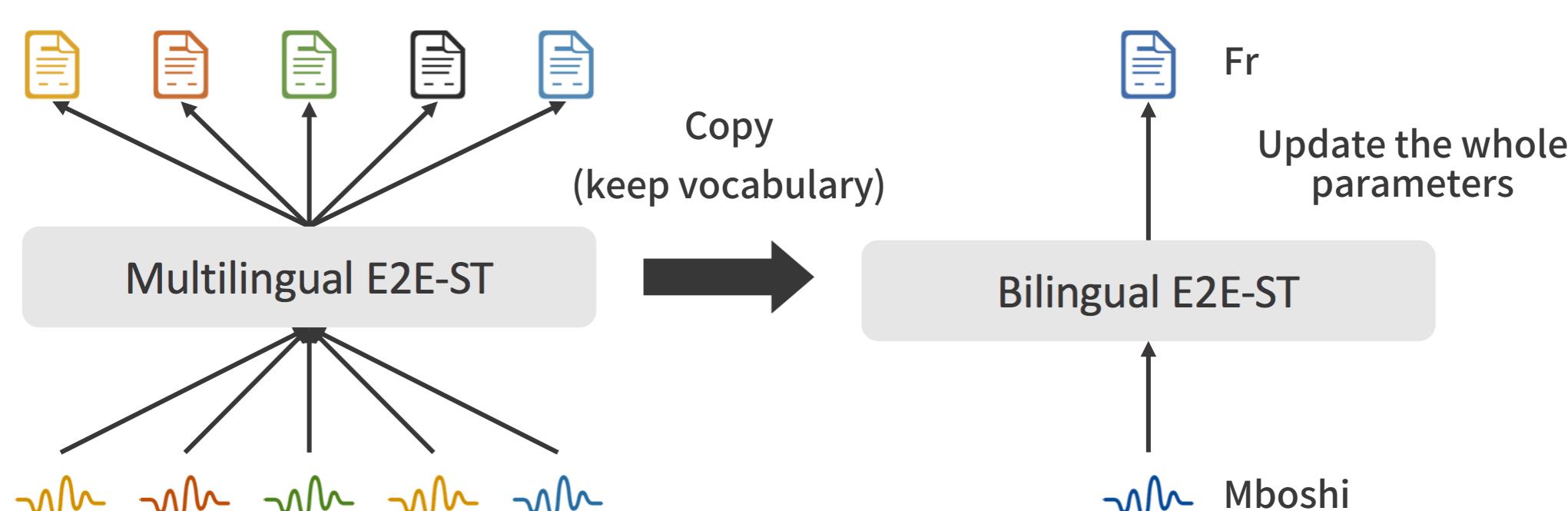
Take-home-message (many-to-many scenario)

- Additional training data from other language pairs was effective
- Multilingual E2E-ST models were not affected by the domain mismatch
- E2E-ST models got more gains from multilingual training than the pipeline systems

Take-home-message (one-to-many scenario)

- O2M training was more effective than M2M training from the perspective of data efficiency
- However, using all training data (M2Mc) got a further small gain
- O2M multilingual training benefits from not only additional English speech data but also the direct optimization

◆ Transfer learning to a very low-resource



Mboshi->Fr: 4.4h

Model	BLEU
Bi-ST (LibriSpeech)	4.55
O2M-ST	6.92
M2Ma-ST	5.50
M2Mc-ST	6.52

Take-home-message

- Multilingual E2E-ST seed was more effective than the bilingual one
- O2M seed showed the best performance among all models