

Brian Yan, Chunlei Zhang, Meng Yu, Shi-Xiong Zhang, Siddharth Dalmia, Dan Berrebbi, Chao Weng, Shinji Watanabe, Dong Yu

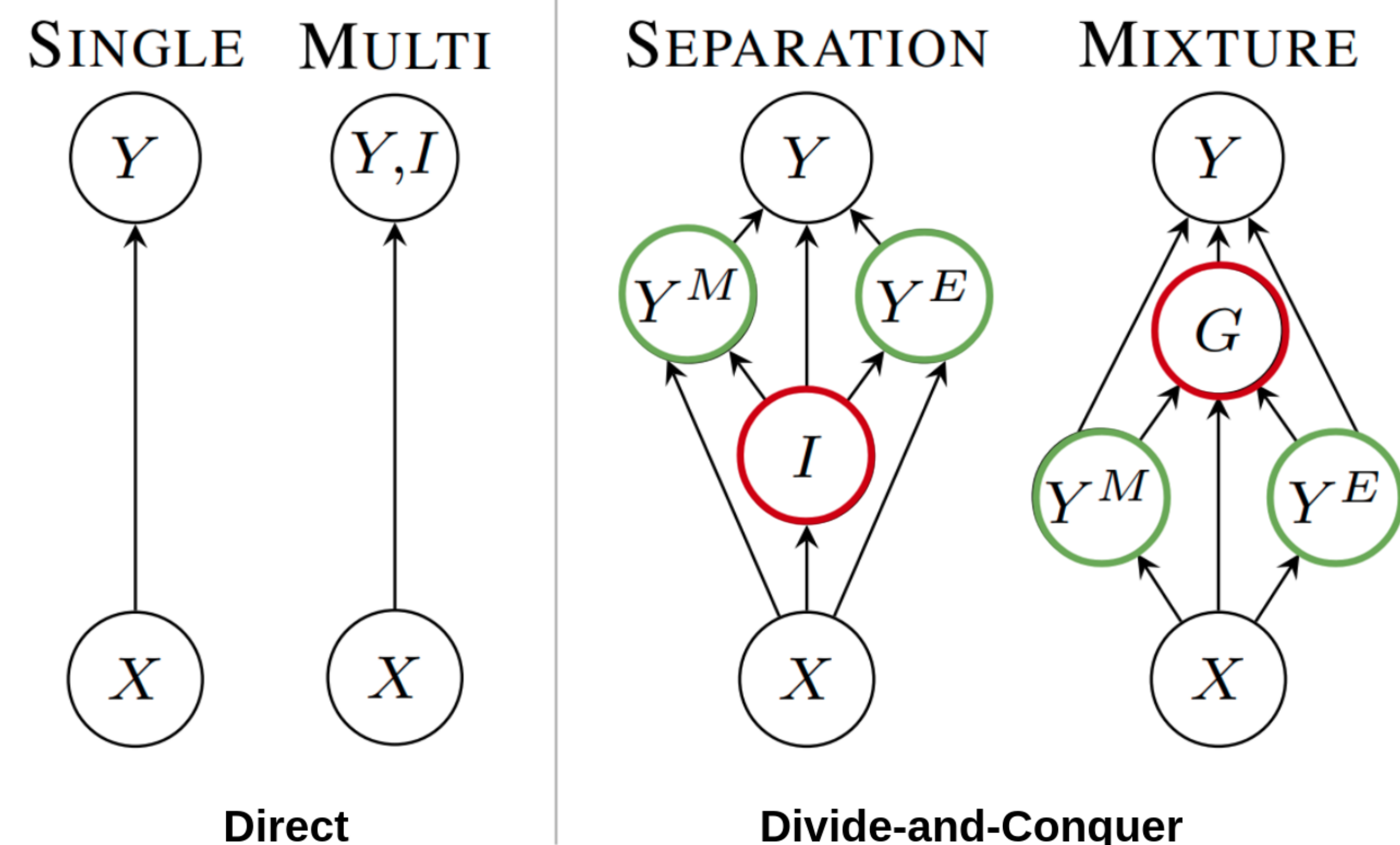
At a Glance

- We first propose a formulation of the bilingual ASR problem as a **conditionally factorized joint model** of monolingual and CS ASR where the final output is obtained given only monolingual label-to-frame synchronized information
- We then apply an end-to-end neural network, which we call the **Conditional RNN-T**, to model our conditional joint formulation

Code-Switching is a Subset of Bilingualism

In this work: we evaluate our models on not only Mandarin-English CS Mixed Error Rate (MER) but also Monolingual English Character Error Rate (CER), and Monolingual English Word Error Rate (WER).

Probabilistic Graphical Models of Prior Works



Direct Graphical Models

- Pro: Simple to only model a direct dependency
- Con: Dependency becomes too complex for unrelated languages

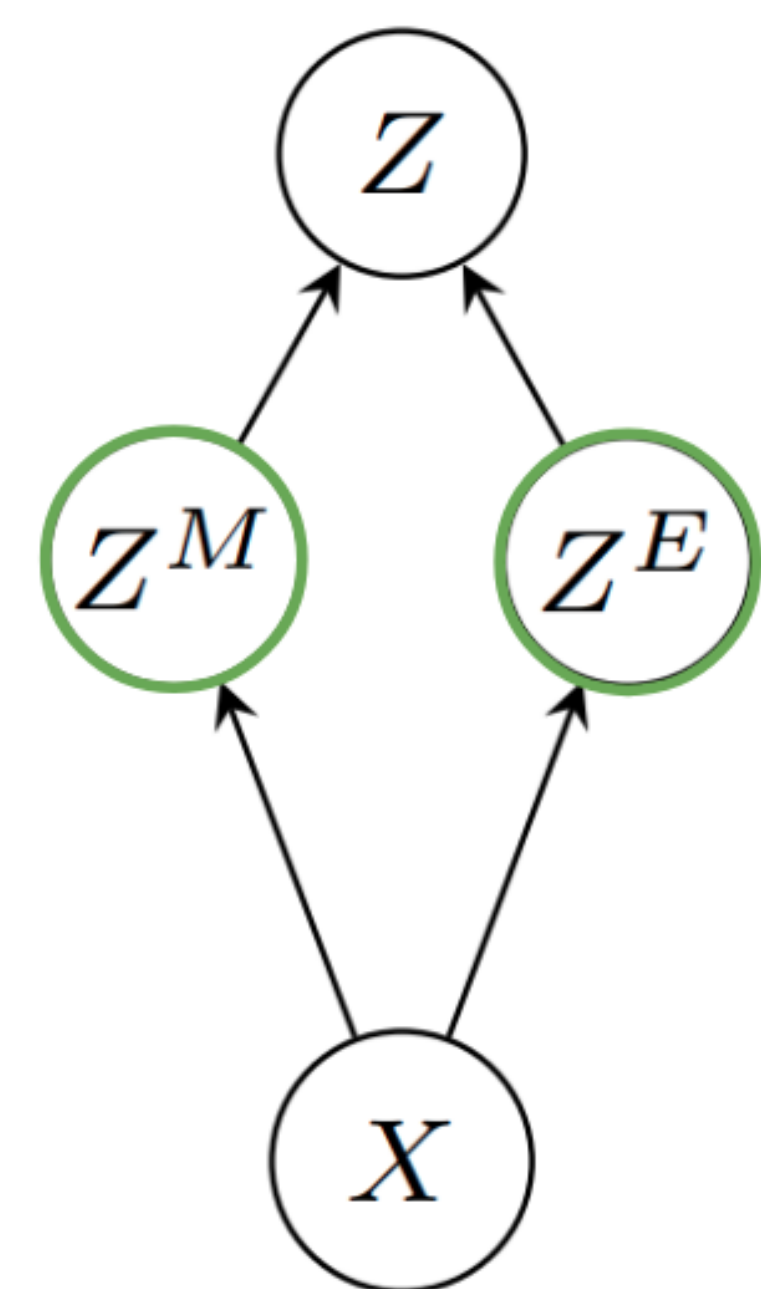
Divide-and-Conquer Models

- Pro: Division of monolingual subtasks → simpler dependencies, more compatible with monolingual data
- Con: Dependence on quality of “divider” module → risk of error propagation, increased complexity

Motivation: Simplest PGM that still has monolingual subtasks?

Conditionally Factorized Joint Model

CONDITIONAL



X : speech features

Z : bilingual transcript

Z^M : mandarin sub-task

Z^E : english sub-task

Z is label-to-frame alignment:

$$z_t = \begin{cases} z_t^M, & \text{if } z_t^M \in \mathcal{V}^M \text{ and } z_t^E = \emptyset \\ z_t^E, & \text{if } z_t^E \in \mathcal{V}^E \text{ and } z_t^M = \emptyset \\ \emptyset, & \text{if } z_t^M = \emptyset \text{ and } z_t^E = \emptyset \end{cases}$$

Jointly model CS and Monolingual parts, w/ **conditional factorization**:

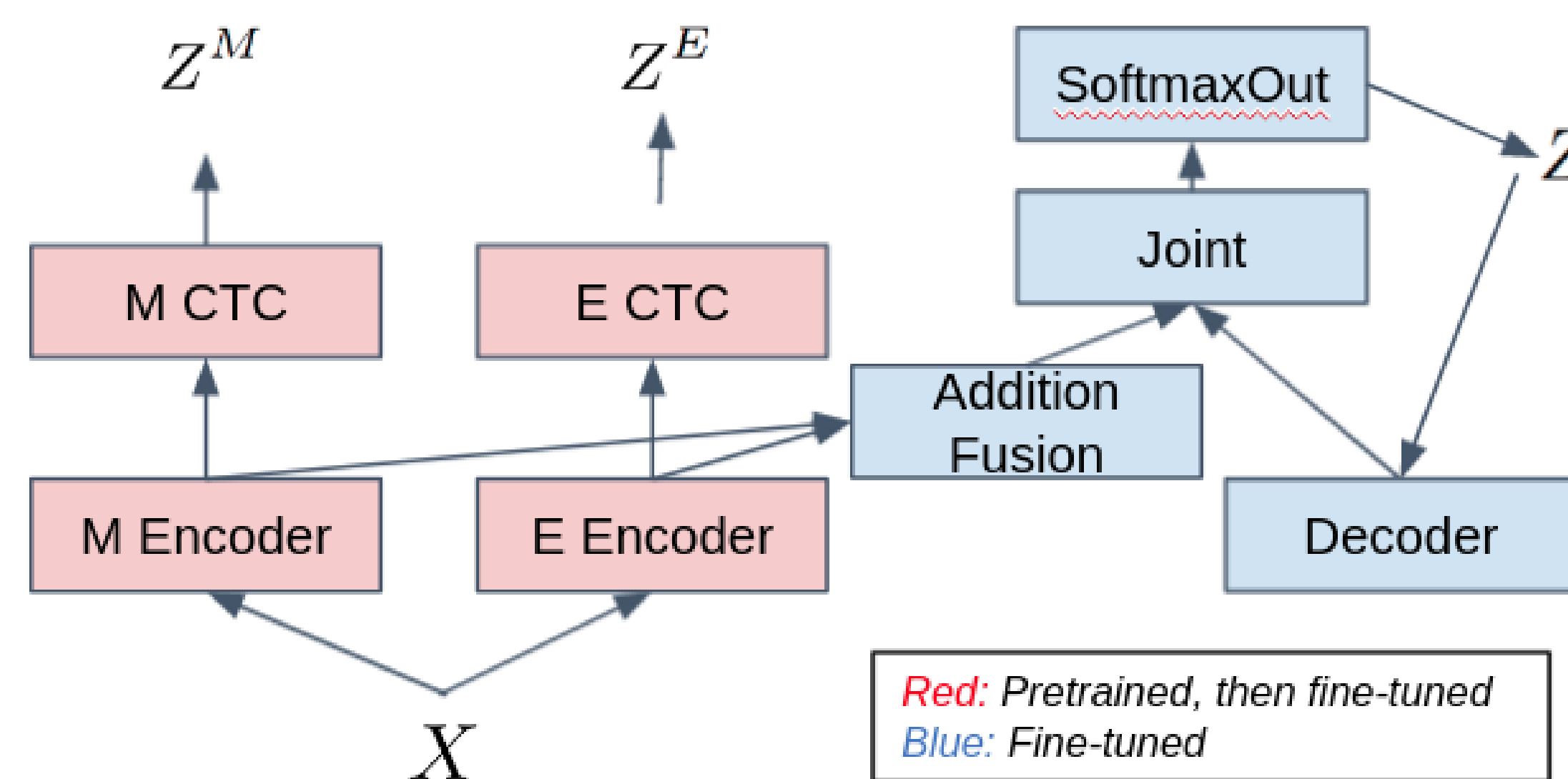
$$\begin{aligned} p(Z|X) &= p(Z, Z^M, Z^E|X) \\ &= p(Z|Z^M, Z^E, X)p(Z^M, Z^E|X) \\ &\approx p(Z|Z^M, Z^E, X)p(Z^M|X)p(Z^E|X) \end{aligned}$$

Conditional independence

Independence

Conditional RNN-T

$$p(Y|X) \approx \underbrace{\sum_Z p(Z|Z^M, Z^E)}_{\triangleq p_{\text{mnt}}(Y|Z^M, Z^E)} \underbrace{\sum_{Z^M} p(Z^M|X)}_{\triangleq p_{\text{ctc}}(Y^M|X)} \underbrace{\sum_{Z^E} p(Z^E|X)}_{\triangleq p_{\text{ctc}}(Y^E|X)}$$



Red: Pretrained, then fine-tuned
Blue: Fine-tuned

$$\mathcal{L}_{\text{LS}} = \lambda \mathcal{L}_{\text{RNNT}} + (1 - \lambda)(\mathcal{L}_{\text{M.CTC}} + \mathcal{L}_{\text{E.CTC}})$$

Original Bilingual g.t. (Z)
Masked Mandarin g.t. (Z^M)
Masked English g.t. (Z^E)

什么是 Code-Switching
什么是 <en>
<zh> Code-Switching

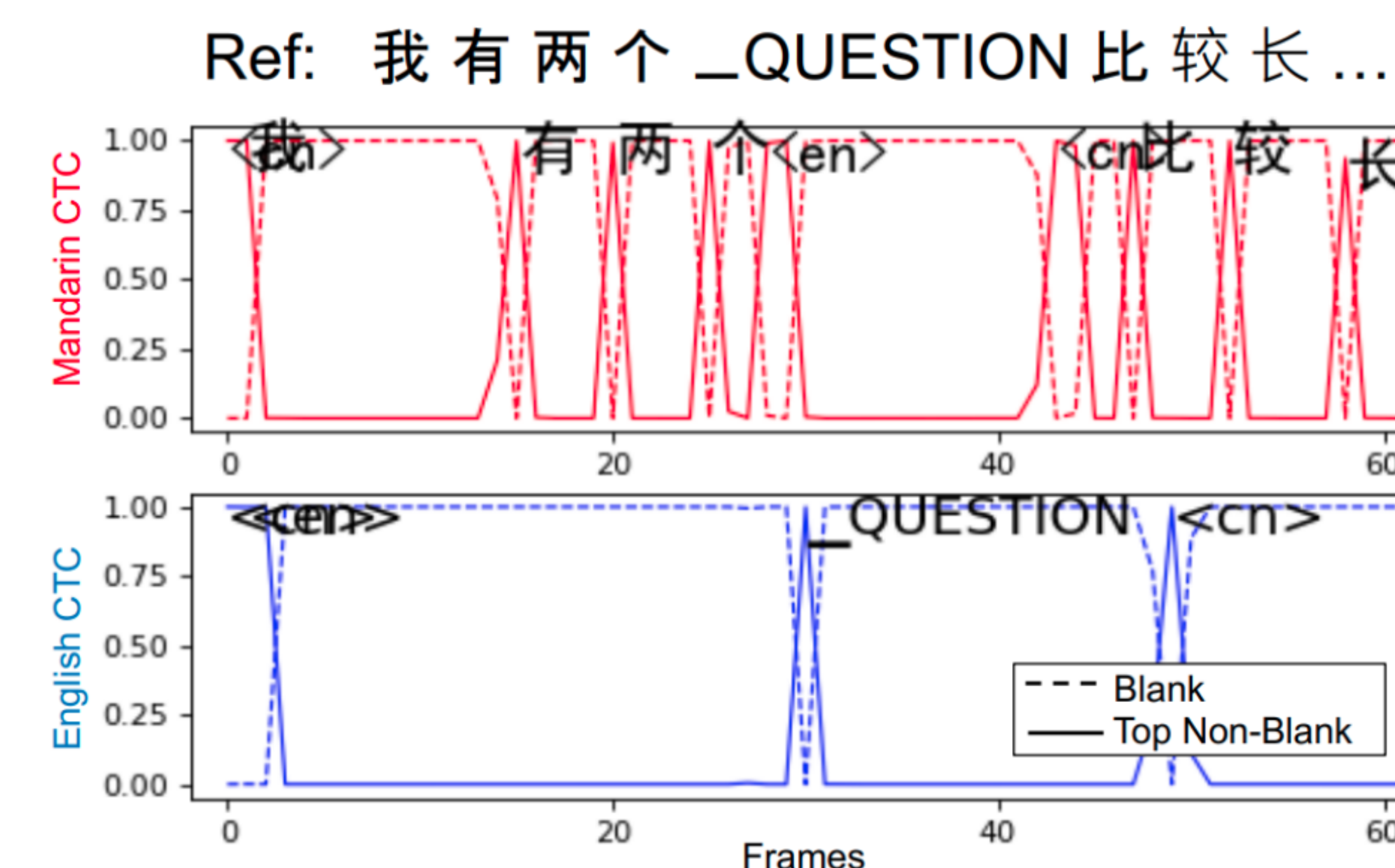
Code-Switched and Monolingual Results

Conditional RNN-T outperforms Direct and Divide-and-Conquer RNNT-T baselines on CS and monolingual test sets.

Model Type	CODE-SWITCHED MER	MONO-MAN CER	MONO-ENG WER
Direct	11.3	6.5	17.8
Divide-and-Conquer	11.2	5.7	34.6
Conditional	10.2	5.3	16.3

Language-Separation Ability

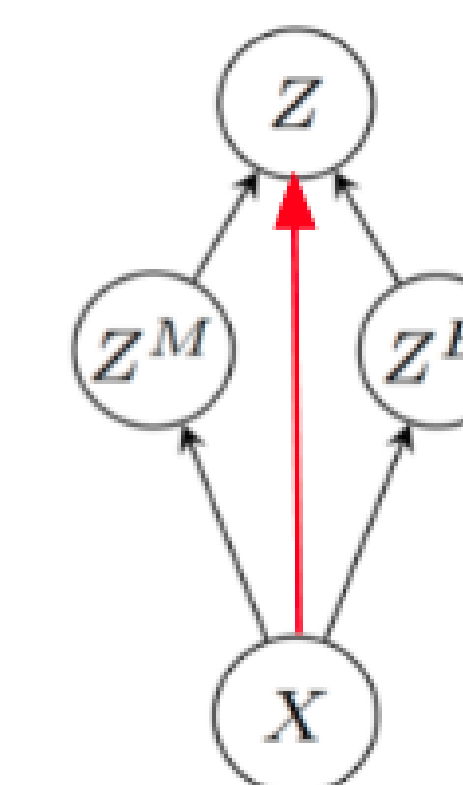
When Mandarin CTC is predicting non-blanks English CTC is predicting blanks, and vice versa.



Experimental Validation of Conditional Independence

No additional information from X is required given Z^M and Z^E .

Model	Bilingual Condition	CODE-SWITCHED MER	CER	WER
Cond. RNN-T + LS	$p(Z Z^M, Z^E)$	11.1	8.9	31.1
3-Enc. RNN-T + LS	$p(Z Z^M, Z^E, X)$	11.2	9.0	31.1



3-encoder variant passes X to the bilingual module, creating the dependency in red

$$p(Z|Z^M, Z^E, X)$$