

# JOINT MODELING OF CODE-SWITCHED AND MONOLINGUAL ASR VIA CONDITIONAL FACTORIZATION

*Brian Yan<sup>1</sup>, Chunlei Zhang<sup>2</sup>, Meng Yu<sup>2</sup>, Shi-Xiong Zhang<sup>2</sup>, Siddharth Dalmia<sup>1</sup>, Dan Berrebbi<sup>1</sup>,  
Chao Weng<sup>3</sup>, Shinji Watanabe<sup>1</sup>, Dong Yu<sup>2</sup>*

<sup>1</sup>Carnegie Mellon University, USA, <sup>2</sup>Tencent AI Lab, USA, <sup>3</sup>Tencent AI Lab, China



Carnegie Mellon University  
Language Technologies Institute



Tencent  
AI Lab



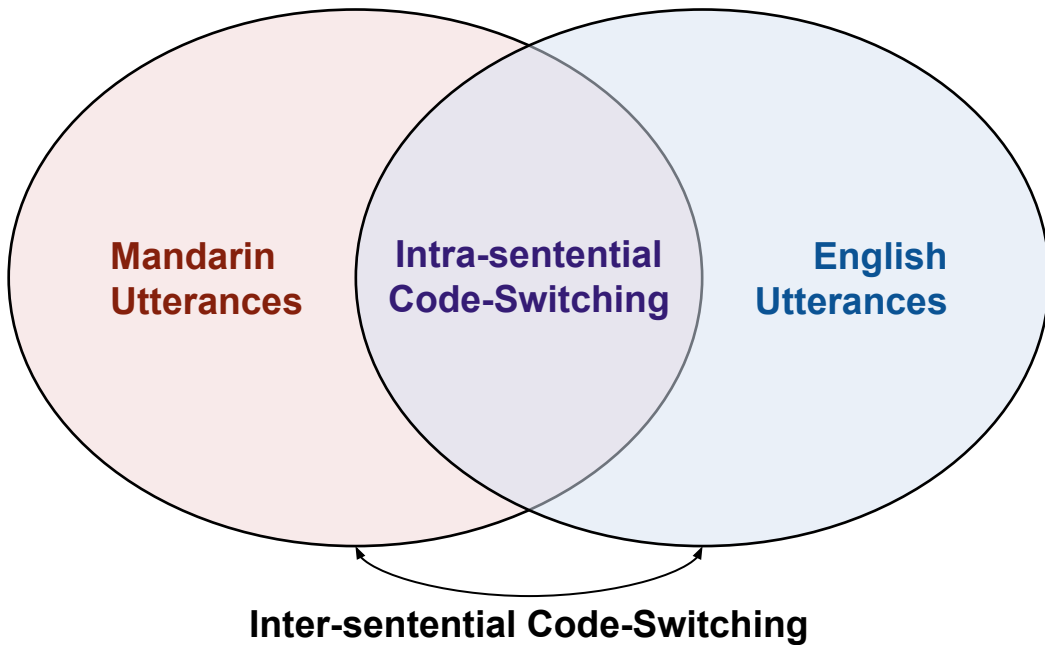
Session: SPE-14: Multi-lingual ASR

Session Time: Sunday, 8 May, 23:00 - 23:45 (Singapore Time, UTC +8)

# Code-switching (CS) $\subset$ Bilingualism

Intra-sentential CS is a **subset** of bilingual conversation, which is often 1 language at a time

Our objective is to model the **entire bilingual task**:



# Bilingual Speech Recognition

Let ...

$X = \{\mathbf{x}_t \in \mathbb{R}^D | t = 1, \dots, T\}$  denote speech features and

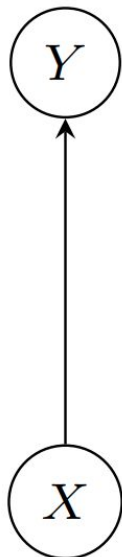
$Y = \{y_t \in (\mathcal{V}^M \cup \mathcal{V}^E) | n = 1, \dots, L\}$  denote bilingual transcriptions.

Note that  $Y$  may be purely monolingual or code-switched.

We wish to predict  $Y$  given  $X$ .

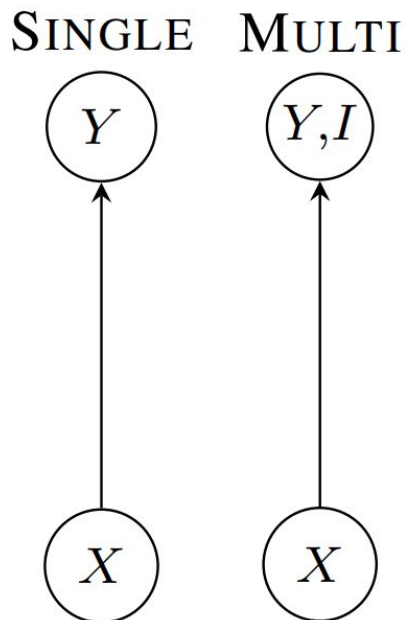
# Direct Formulations of Bilingual ASR

SINGLE



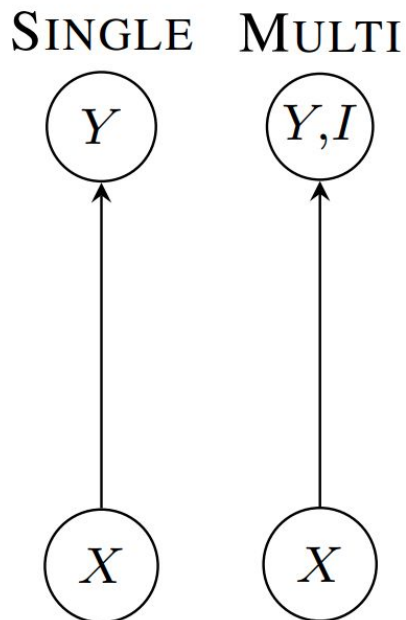
- Model Y as a **single** conditionally dependent variable
  - Hybrid: Phone merging (Sivasankaran 2018)
  - E2E: LID token method (Zhang 2020)

# Direct Formulations of Bilingual ASR



- Model **Y** as a **single** conditionally dependent variable
  - Hybrid: Phone merging (Sivasankaran 2018)
  - E2E: LID token method (Zhang 2020)
- Model **multiple** dependents:  $Y$  and language ID,  $I$ 
  - E2E: joint LID and ASR (Zeng 2019)

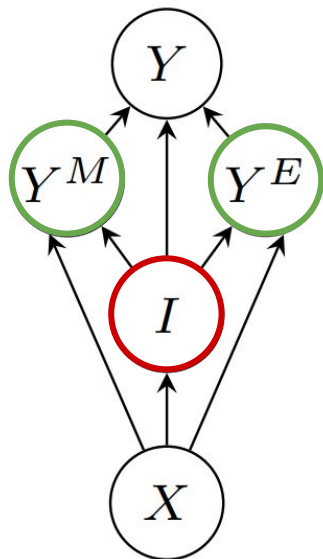
# Direct Formulations of Bilingual ASR



- Model **Y** as a **single** conditionally dependent variable
  - Hybrid: Phone merging (Sivasankaran 2018)
  - E2E: LID token method (Zhang 2020)
- Model **multiple** dependents:  $Y$  and language ID,  $I$ 
  - E2E: joint LID and ASR (Zeng 2019)
- Combining 2 unrelated languages = more complex

# Divide-and-Conquer Formulations of Bilingual ASR

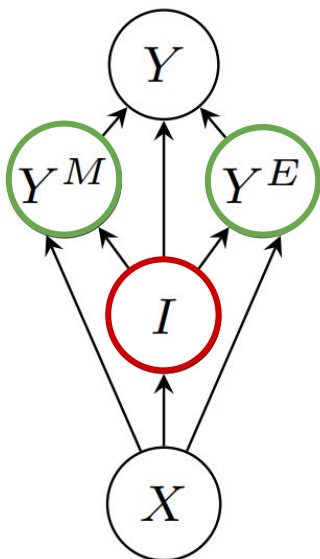
SEPARATION



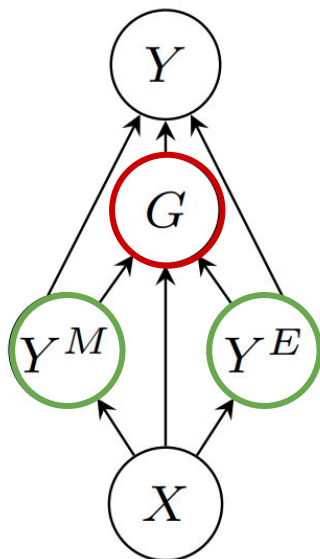
- **Separate** then Recognize
  - Hybrid: LID to monolingual ASR cascade (Chan 2004)

# Divide-and-Conquer Formulations of Bilingual ASR

SEPARATION



MIXTURE

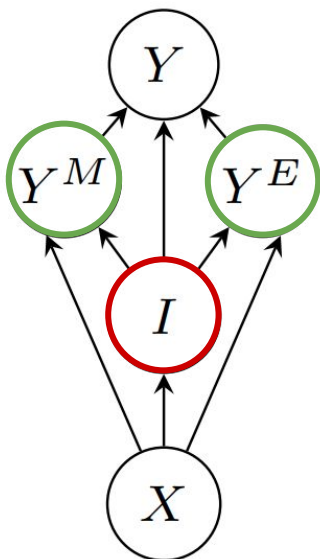


- **Separate** then Recognize
  - Hybrid: LID to monolingual ASR cascade (Chan 2004)
- **Mixture** of Monolingual Experts
  - Hybrid: Frame-level posterior weighting (Weiner 2012)
  - E2E: Mixture of experts (Lu 2020, Dalmia 2021)

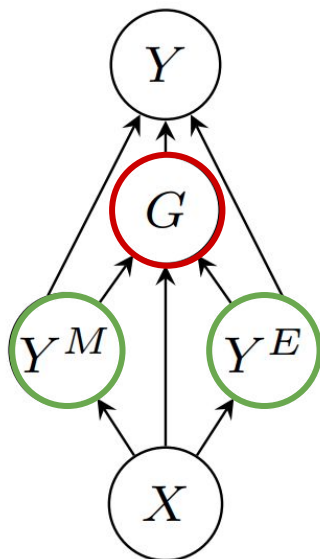


# Divide-and-Conquer Formulations of Bilingual ASR

SEPARATION



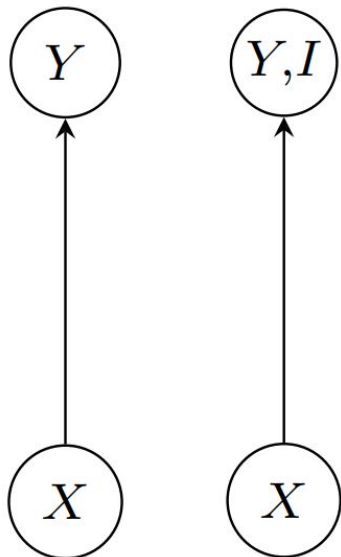
MIXTURE



- **Separate** then Recognize
  - Hybrid: LID to monolingual ASR cascade (Chan 2004)
- **Mixture** of Monolingual Experts
  - Hybrid: Frame-level posterior weighting (Weiner 2012)
  - E2E: Mixture of experts (Lu 2020, Dalmia 2021)
- **Division of monolingual tasks**
  - simpler, more compatible with monolingual data
- **Dependence on quality of “divider” module**
  - Risk of error propagation, increased complexity

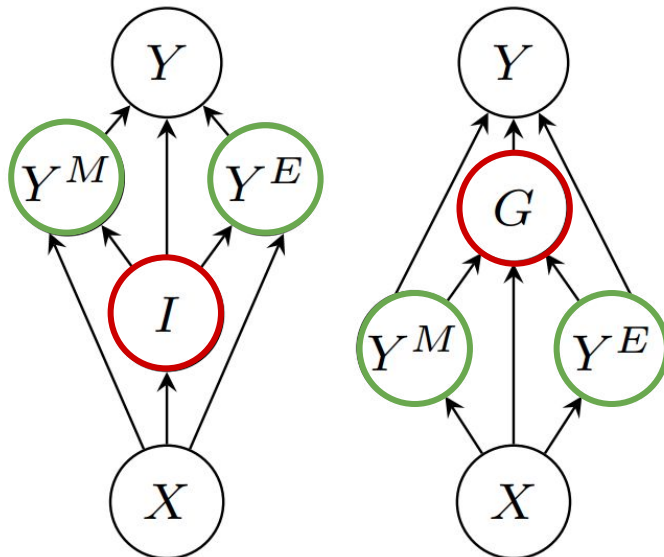
# Our Motivation

SINGLE MULTI



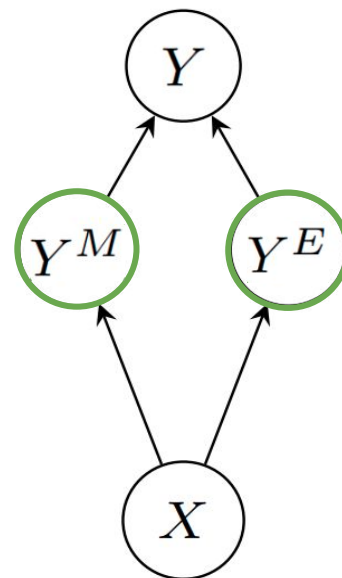
Direct

SEPARATION MIXTURE

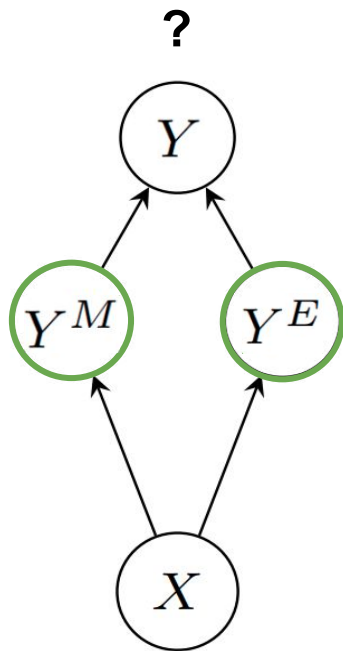


Divide-and-Conquer

?

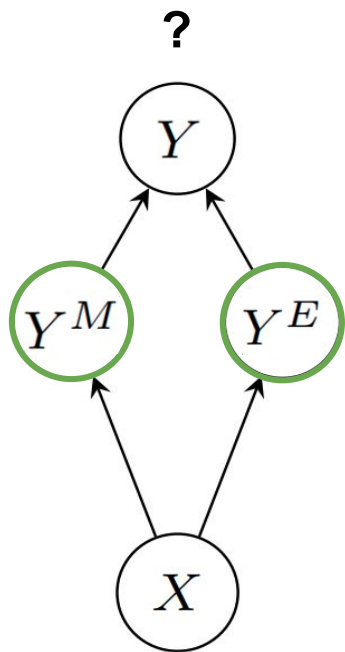


# Desiderata



1. Can we build CS + bilingual ASR with **monolingual sub-components**...
2. ...where the final output is **conditioned only on those 2 sub-components** and nothing else?
3. And does such a conditional approach more **efficiently leverage monolingual and CS training data**?

# Label-to-Frame Synchronization



If ...

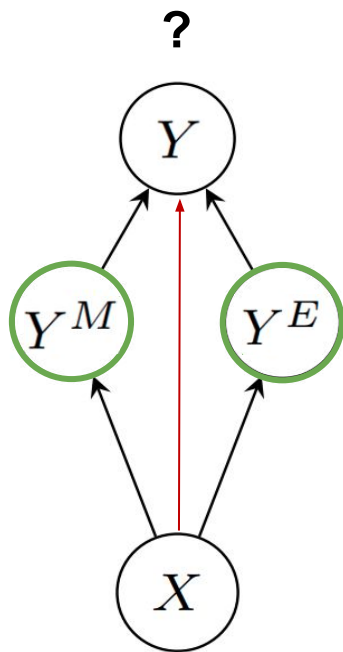
$Y^M$  = 什么是,

$Y^E$  = code-switching

Can  $Y$  be determined?

$$P(y, y^M, y^E | x) = P(y | y^M, y^E) P(y^M | x) P(y^E | x)$$

# The Need for Label-to-Frame Synchronization



If ...

$Y^M$  = 什么是,       $Y^E$  = code-switching

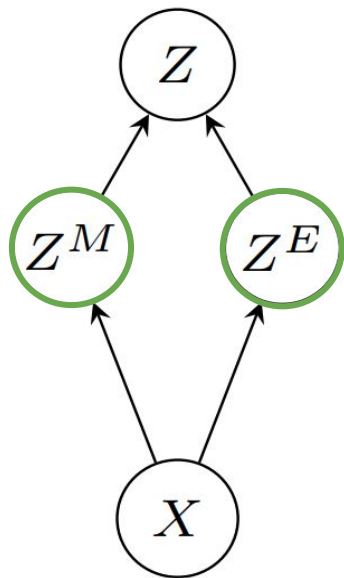
Can  $Y$  be determined?

- No! We're missing ordering information
- This formulation is not conditionally independent

$$P(y, y^M, y^E | x) = \boxed{P(y | y^M, y^E)} P(y^M | x) P(y^E | x)$$

# Conditionally Factorized Bilingual ASR

CONDITIONAL



Let ...

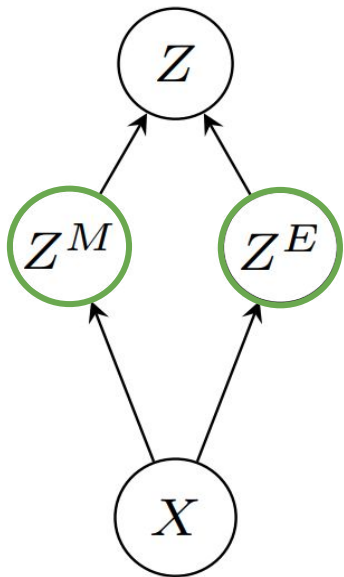
$X = \{\mathbf{x}_t \in \mathbb{R}^D | t = 1, \dots, T\}$  denote speech features and

$Z = \{z_t \in \mathcal{V}^M \cup \mathcal{V}^E \cup \{\emptyset\} | t = 1 \dots T\}$  denote label-to-frame alignments.

*Denotes null emission as in CTC*

# Conditionally Factorized Bilingual ASR

CONDITIONAL



Let ...

$X = \{\mathbf{x}_t \in \mathbb{R}^D \mid t = 1, \dots, T\}$  denote speech features and

$Z = \{z_t \in \mathcal{V}^M \cup \mathcal{V}^E \cup \{\emptyset\} \mid t = 1 \dots T\}$  denote label-to-frame alignments.

*Denotes null emission as in CTC*

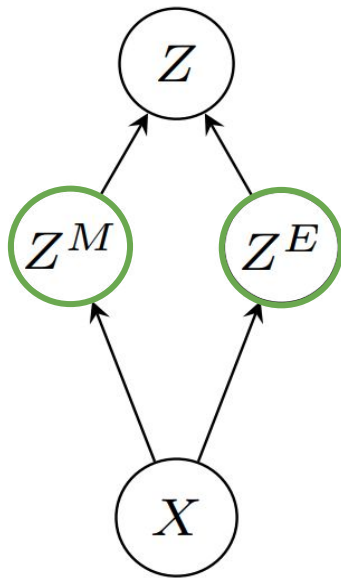
Bilingual label-to-frame alignments can be specified in terms of its constituent monolingual label-to-frame alignments.

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

*By definition, at least one side is null for a given  $t$*

# Conditionally Factorized Bilingual ASR

CONDITIONAL



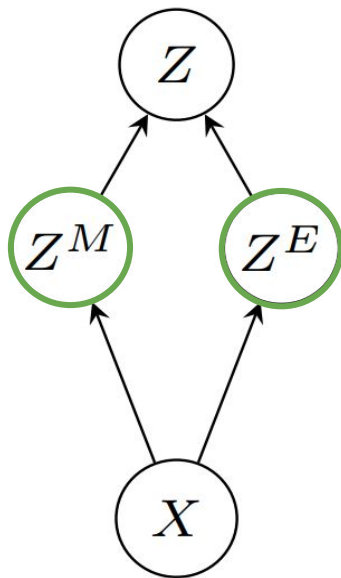
Formulate likelihood in terms of label-to-frame alignments:

$$p(Y|X) = \sum_{Z \in \mathcal{Z}(Y)} p(Z|X)$$



# Conditionally Factorized Bilingual ASR

CONDITIONAL



Formulate likelihood in terms of label-to-frame alignments:

$$p(Y|X) = \sum_{Z \in \mathcal{Z}(Y)} p(Z|X)$$

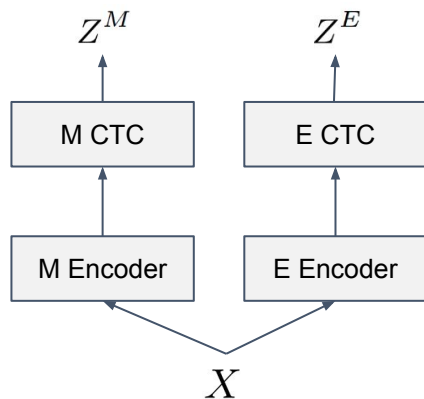
Jointly model CS and Monolingual parts, w/ **conditional factorization**:

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

*Independence*

*Conditional independence*

# Conditional RNN Transducer

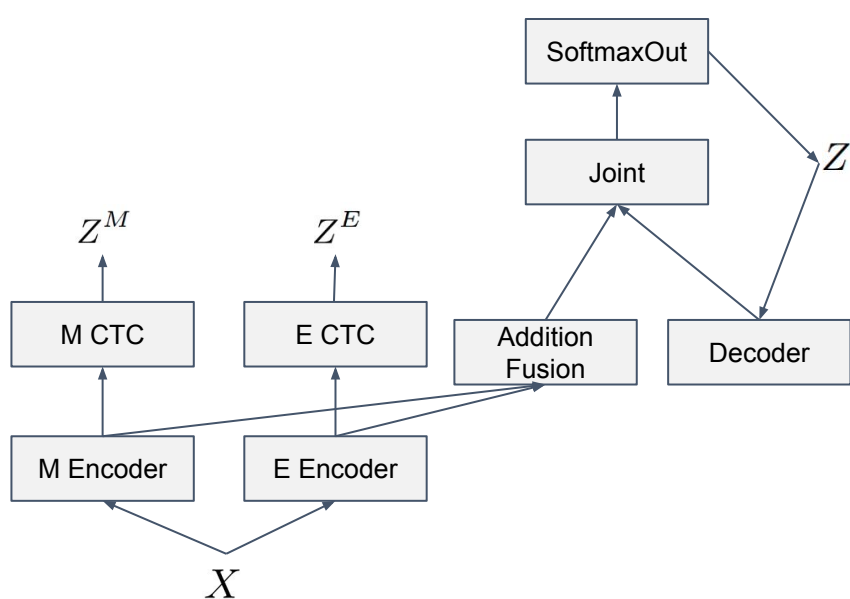


Monolingual CTC Modules:

$$p(Z^M | X) \approx \prod_{t=1}^T p(z_t^M | X, \cancel{z_{1:t-1}^M})$$

$$p(Z^E | X) \approx \prod_{t=1}^T p(z_t^E | X, \cancel{z_{1:t-1}^E})$$

# Conditional RNN Transducer



Monolingual CTC Modules:

$$p(Z^M | X) \approx \prod_{t=1}^T p(z_t^M | X, \cancel{z_{1:t-1}^M})$$

$$p(Z^E | X) \approx \prod_{t=1}^T p(z_t^E | X, \cancel{z_{1:t-1}^E})$$

Bilingual RNN-T Module:

$$p(Z | Z^M, Z^E) = \prod_{i=1}^{T+L} p(z_i | Z^M, Z^E, z_{1:i-1})$$

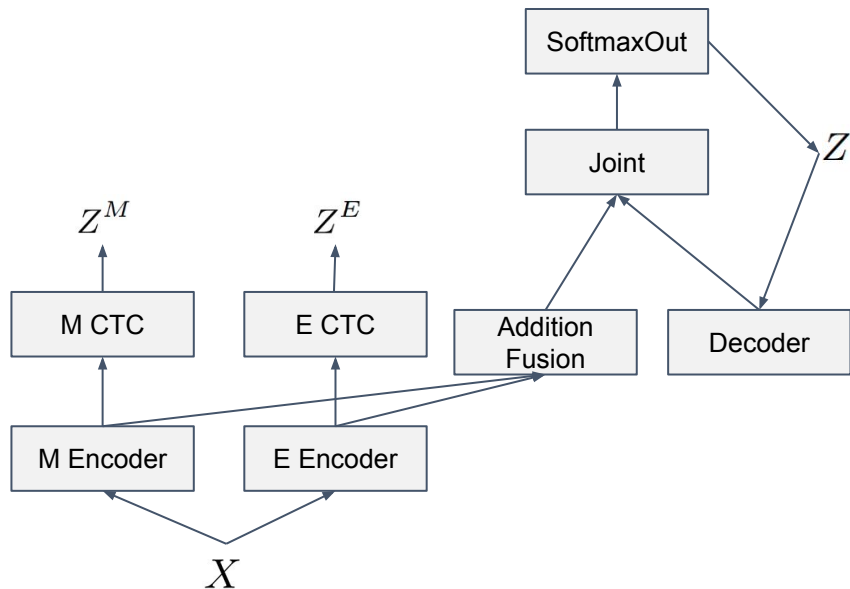
$$\mathbf{h}_t^{\text{ENC}} = \mathbf{h}_t^M + \mathbf{h}_t^E$$

$$\mathbf{h}_l^{\text{DEC}} = \text{DECODER}(z_{1:l-1})$$

$$\mathbf{h}_{t,l}^{\text{JNT}} = \text{JOINT}(\mathbf{h}_t^{\text{ENC}}, \mathbf{h}_l^{\text{DEC}})$$

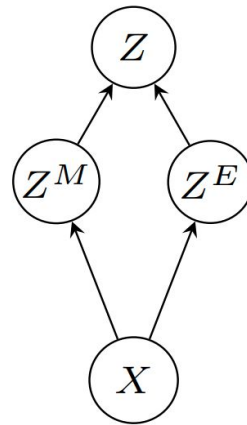
$$p(z_i | \mathbf{h}^M, \mathbf{h}^E, z_{1:i-1}) = \text{SOFTMAXOUT}(\mathbf{h}_{t,l}^{\text{JNT}})$$

# Conditional RNN Transducer

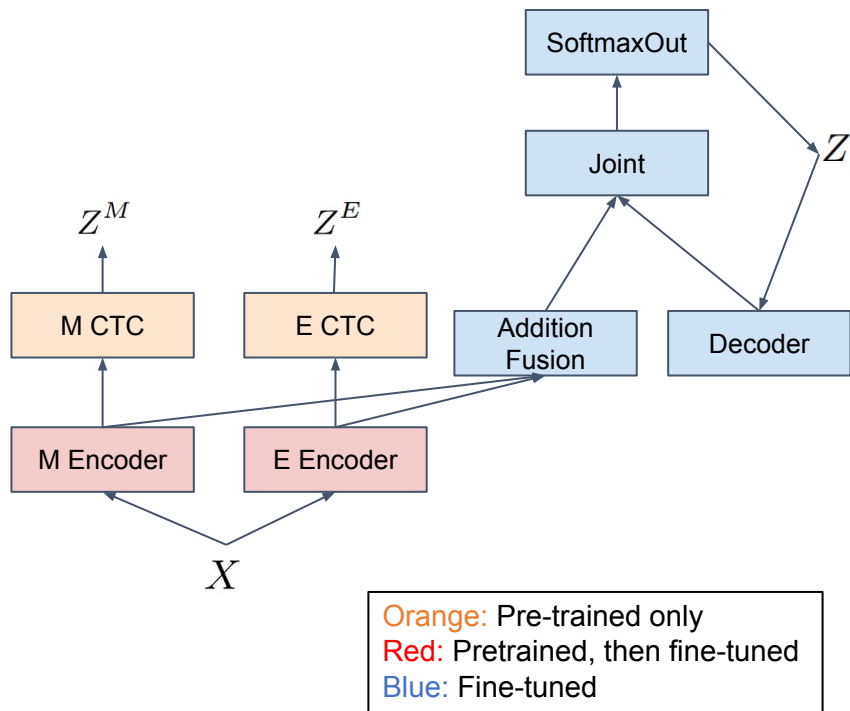


Full Network:

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



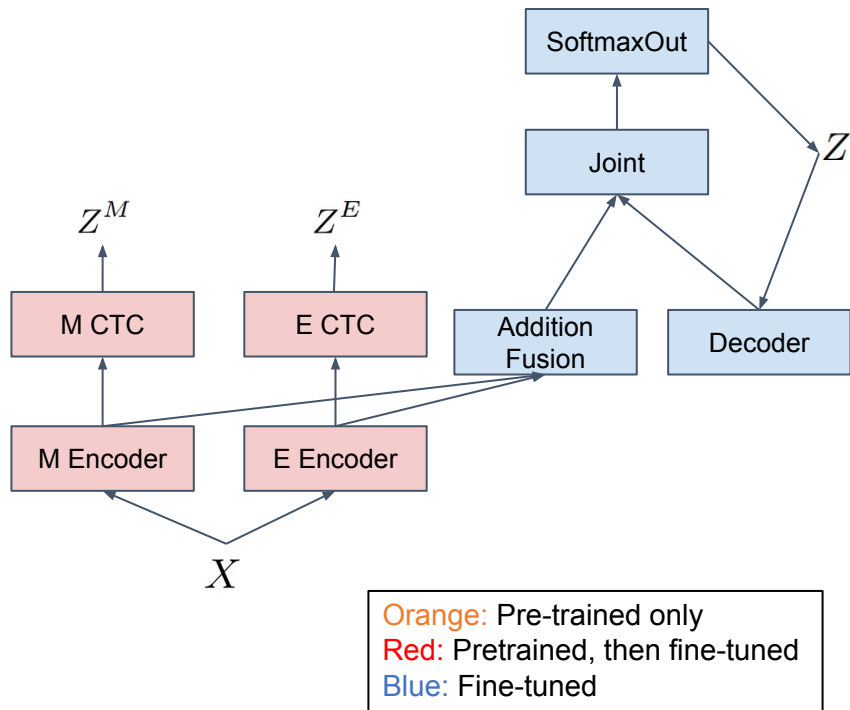
# Conditional RNN Transducer



Implicit conditioning via pre-training:

- Pre-train:  
 $\mathcal{L}_{M.CTC}$ ,  $\mathcal{L}_{E.CTC}$
- Fine-tune:  
 $\lambda \mathcal{L}_{RNNT}$

# Conditional RNN Transducer



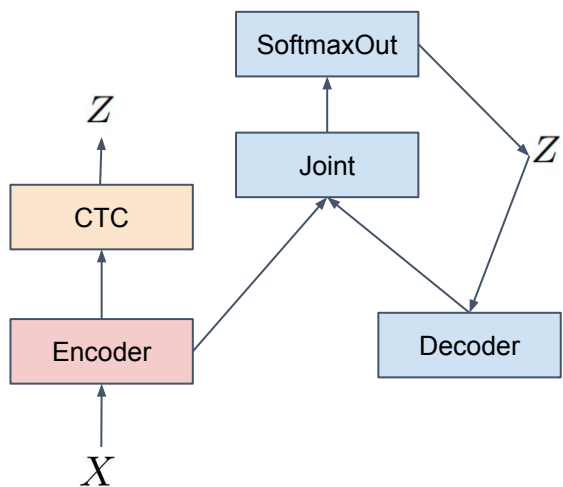
Explicit conditioning via masking:

- Pre-train:  
 $\mathcal{L}_{M.CTC}$  ,  $\mathcal{L}_{E.CTC}$
- Fine-tune:  
 $\mathcal{L}_{LS} = \lambda \mathcal{L}_{RNNT} + (1 - \lambda)(\mathcal{L}_{M.CTC} + \mathcal{L}_{E.CTC})$
- Monolingual ground truths are obtained via language-specific masking of the bilingual ground truth → referred to as **Language-Separation**

Original Bilingual g.t.  
Masked Mandarin g.t.  
Masked English g.t.

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

# Single RNN-T Baseline

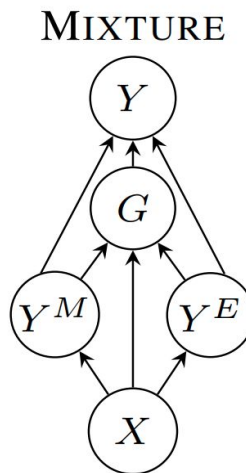
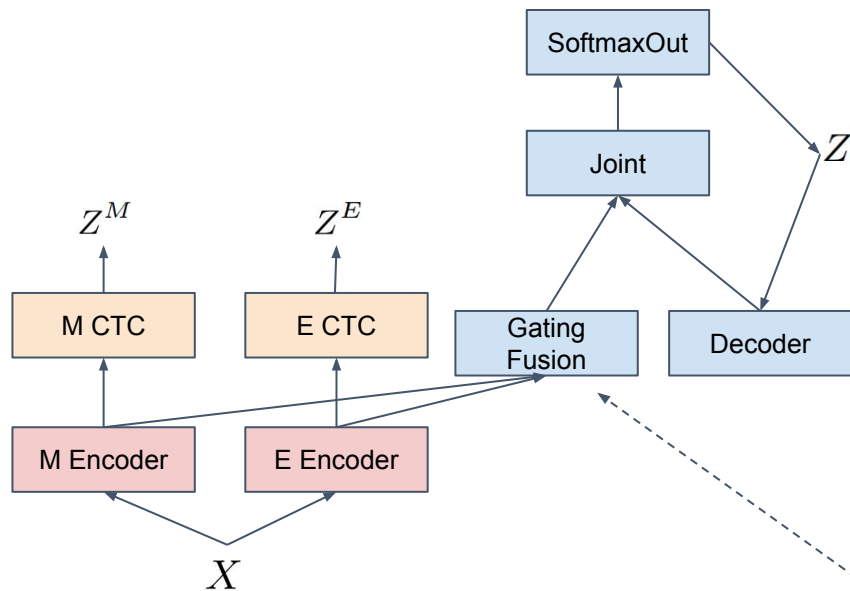


SINGLE



Orange: Pre-trained only  
Red: Pretrained, then fine-tuned  
Blue: Fine-tuned

# Gating RNN-T Baseline



Orange: Pre-trained only  
Red: Pretrained, then fine-tuned  
Blue: Fine-tuned

$$\alpha_{\text{enc}} = \text{sigmoid}(w_{\alpha}(\tanh(w_{\text{enc}}^{\text{man}}(h^{\text{man}}) + w_{\text{enc}}^{\text{en}}(h^{\text{en}}))))$$
$$h = \alpha_{\text{enc}} * h^{\text{man}} + (1 - \alpha_{\text{enc}}) * h^{\text{en}}$$



# Experimental Setup

## Data:

- **200h Mandarin-English CS data** from ASRU'19 challenge
- **500h monolingual Mandarin data** from ASRU'19 challenge
  - 200h subset used for fine-tuning
- **700h monolingual accented English data** from King-ASR-190
  - 200h subset used for fine-tuning

## Evaluation:

- **CS set** as measured by Mixed Error Rate (MER)
- **Monolingual Mandarin set** as measured by Character Error Rate (CER)
- **Monolingual English set** as measured by Word Error Rate (WER)

# Main Results

Model Type	Model Name	Pre-trained Encoder(s)	Fine-tuning Data	CODE-SWITCHED			MONO-MAN	MONO-ENG
				MER	CER	WER	CER	WER
Direct	Vanilla RNN-T [21, 24]	✓	CS	12.3	9.9	34.3	17.9	81.4
Mixture	Gating RNN-T [22, 24]	✓	CS	11.5	9.1	33.0	17.7	<b>78.3</b>
Conditional	Our Proposed Model	✓	CS	11.5	9.1	33.2	15.5	82.9
Conditional	+ Language-Separation (LS)	✓	CS	<b>11.1</b>	<b>8.7</b>	<b>32.7</b>	<b>15.3</b>	82.7
Direct	Single RNN-T [21, 24]	✓	CS + M	11.3	9.3	30.8	6.5	17.8
Mixture	Gating RNN-T [22, 24]	✓	CS + M	11.2	8.8	34.7	5.7	34.6
Conditional	Our Proposed Model	✓	CS + M	10.3	8.2	29.5	5.4	16.5
Conditional	+ Language-Separation (LS)	✓	CS + M	<b>10.2</b>	<b>8.1</b>	<b>29.2</b>	<b>5.3</b>	<b>16.3</b>

- All models perform significantly better on monolingual sets when using monolingual fine-tuning

# Main Results

Model Type	Model Name	Pre-trained Encoder(s)	Fine-tuning Data	CODE-SWITCHED			MONO-MAN	MONO-ENG
				MER	CER	WER	CER	WER
Direct	Vanilla RNN-T [21, 24]	✓	CS	12.3	9.9	34.3	17.9	81.4
Mixture	Gating RNN-T [22, 24]	✓	CS	11.5	9.1	33.0	17.7	<b>78.3</b>
Conditional	Our Proposed Model	✓	CS	11.5	9.1	33.2	15.5	82.9
Conditional	+ Language-Separation (LS)	✓	CS	<b>11.1</b>	<b>8.7</b>	<b>32.7</b>	<b>15.3</b>	82.7
Direct	Single RNN-T [21, 24]	✓	CS + M	11.3	9.3	30.8	6.5	17.8
Mixture	Gating RNN-T [22, 24]	✓	CS + M	11.2	8.8	34.7	5.7	34.6
Conditional	Our Proposed Model	✓	CS + M	10.3	8.2	29.5	5.4	16.5
Conditional	+ Language-Separation (LS)	✓	CS + M	<b>10.2</b>	<b>8.1</b>	<b>29.2</b>	<b>5.3</b>	<b>16.3</b>

- All models perform significantly better on monolingual sets when using monolingual fine-tuning
- Gating RNN-T outperforms Single RNN-T on CS set, but not on monolingual English

# Main Results

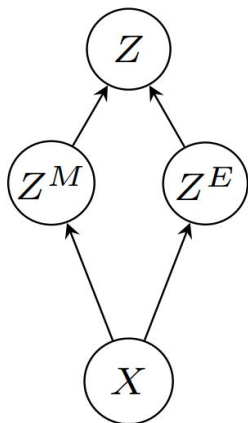
Model Type	Model Name	Pre-trained Encoder(s)	Fine-tuning Data	CODE-SWITCHED			MONO-MAN	MONO-ENG
				MER	CER	WER	CER	WER
Direct	Vanilla RNN-T [21, 24]	✓	CS	12.3	9.9	34.3	17.9	81.4
Mixture	Gating RNN-T [22, 24]	✓	CS	11.5	9.1	33.0	17.7	<b>78.3</b>
Conditional	Our Proposed Model	✓	CS	11.5	9.1	33.2	15.5	82.9
Conditional	+ Language-Separation (LS)	✓	CS	<b>11.1</b>	<b>8.7</b>	<b>32.7</b>	<b>15.3</b>	82.7
Direct	Single RNN-T [21, 24]	✓	CS + M	11.3	9.3	30.8	6.5	17.8
Mixture	Gating RNN-T [22, 24]	✓	CS + M	11.2	8.8	34.7	5.7	34.6
Conditional	Our Proposed Model	✓	CS + M	10.3	8.2	29.5	5.4	16.5
Conditional	+ Language-Separation (LS)	✓	CS + M	<b>10.2</b>	<b>8.1</b>	<b>29.2</b>	<b>5.3</b>	<b>16.3</b>

- All models perform significantly better on monolingual sets when using monolingual fine-tuning
- Gating RNN-T outperforms Single RNN-T on CS set, but not on monolingual English
- Conditional RNN-T outperforms both baselines across CS and monolingual sets

# Analysis of Language-Separation (LS) Ability

Recall:

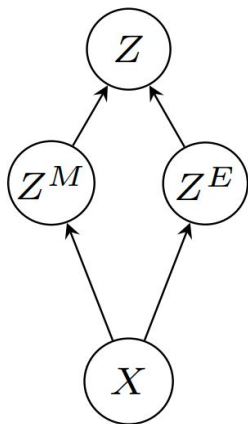
$$z_t = \begin{cases} z_t^M, & \text{if } z_t^M \in \mathcal{V}^M \text{ and } z_t^E = \emptyset \\ z_t^E, & \text{if } z_t^E \in \mathcal{V}^E \text{ and } z_t^M = \emptyset \\ \emptyset, & \text{if } z_t^M = \emptyset \text{ and } z_t^E = \emptyset \end{cases}$$



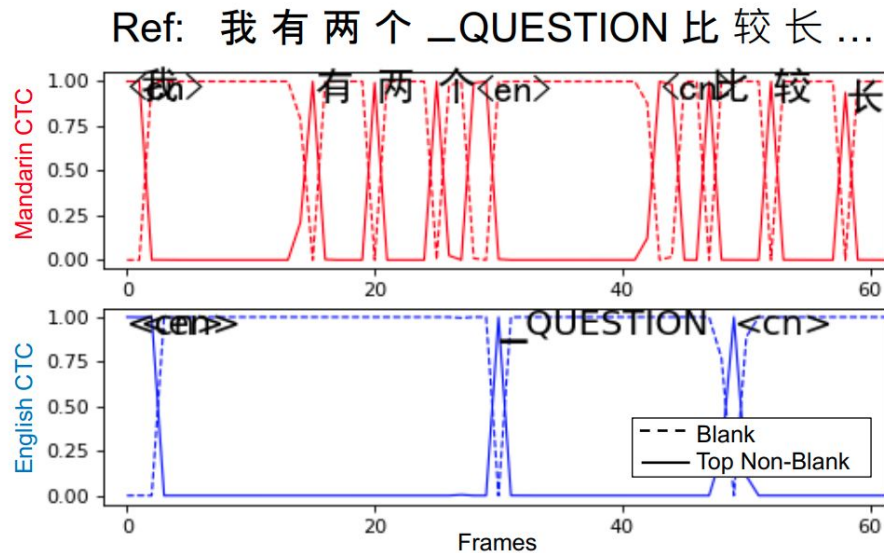
# Analysis of Language-Separation (LS) Ability

Recall:

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



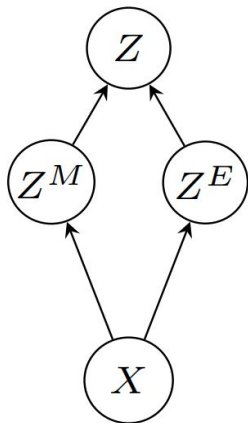
Visualizing Language Separation:



# Analysis of Language-Separation (LS) Ability

Recall:

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



Evaluating the monolingual CTC sub-nets:

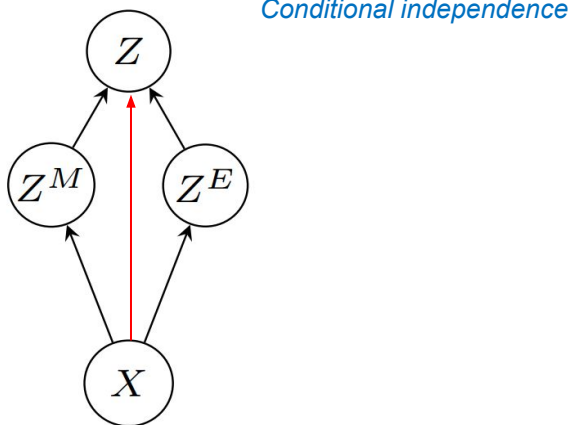
Model	MAN PORTION OF CS			ENG PORTION OF CS		
	Sub-Net	CER	INS	Sub-Net	WER	INS
Cond. RNN-T	$p(Z^M X)$	11.8	3.7	$p(Z^E X)$	42.7	7.9
Cond. RNN-T + LS	$p(Z^M X)$	<b>8.6</b>	<b>0.7</b>	$p(Z^E X)$	<b>37.1</b>	<b>4.6</b>

- We evaluate the monolingual CTC outputs against the language-specific portions of the CS reference
- Both models can perform reasonable language diarization
- Conditional RNN-T + LS has reduced insertion errors

# Conditional Independence of Bilingual Module

Recall:

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

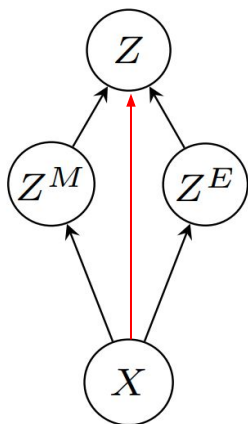




# Conditional Independence of Bilingual Module

Recall:

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



Conditional independence

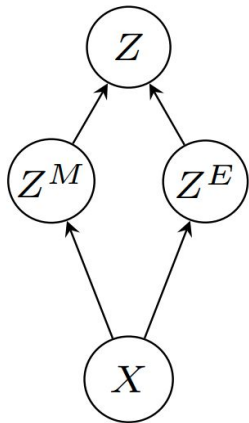
Experimental Validation:

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

- The 3-encoder variant removes the conditional independence assumption by directly encoding speech features,  $X$ , to the bilingual module
- This **dependency** adds no additional information as the monolingual alignments are enough to determine the bilingual alignment

# Recap: Did we satisfy our desiderata?

CONDITIONAL



- ✓ Can we build CS + bilingual ASR with **monolingual sub-components**...
- ✓ ...where the final output is **conditioned only on those 2 sub-components** and nothing else?
- ✓ And does such a conditional approach more **efficiently leverage monolingual and CS training data**?

# Thank You!

**Session: SPE-14: Multi-lingual ASR**

**Session Time: Sunday, 8 May, 23:00 - 23:45 (Singapore Time, UTC +8)**

*Brian Yan<sup>1</sup>, Chunlei Zhang<sup>2</sup>, Meng Yu<sup>2</sup>, Shi-Xiong Zhang<sup>2</sup>, Siddharth Dalmia<sup>1</sup>, Dan Berrebbi<sup>1</sup>,  
Chao Weng<sup>3</sup>, Shinji Watanabe<sup>1</sup>, Dong Yu<sup>2</sup>*

<sup>1</sup>Carnegie Mellon University, USA, <sup>2</sup>Tencent AI Lab, USA, <sup>3</sup>Tencent AI Lab, China

[byan@cs.cmu.edu](mailto:byan@cs.cmu.edu)



Carnegie Mellon University  
Language Technologies Institute



Tencent  
AI Lab

