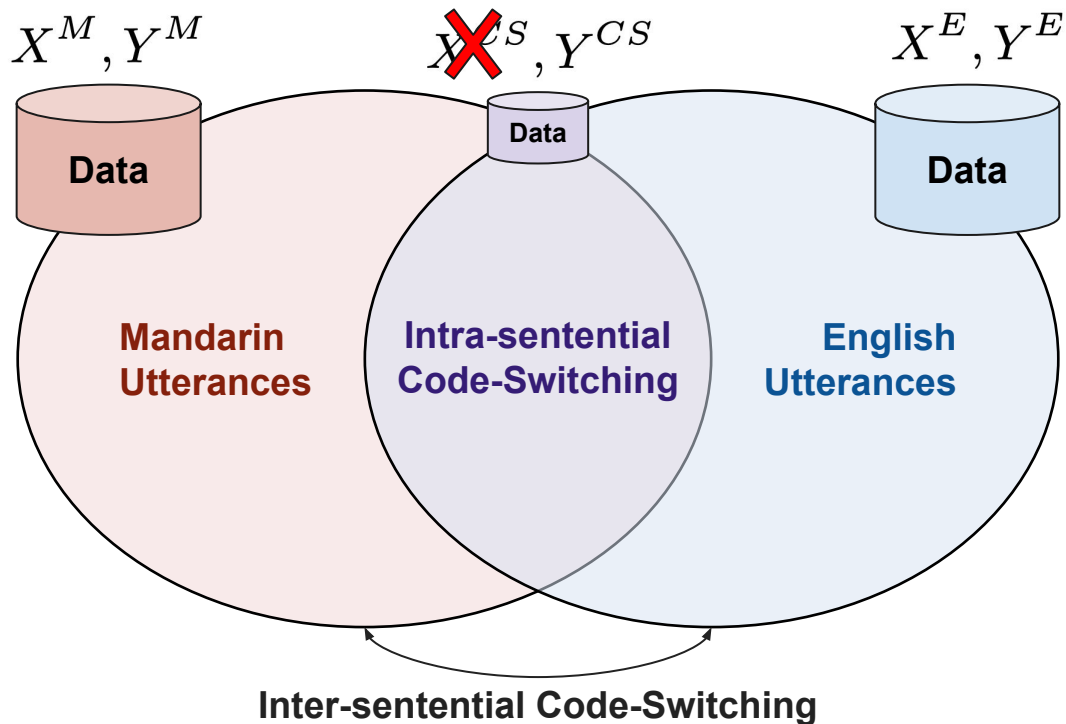# Code-Switched Modeling

**Brian Yan, Matthew Wiesner, Ondrej Klejch, Preethi Jyothi, Shinji Watanabe**
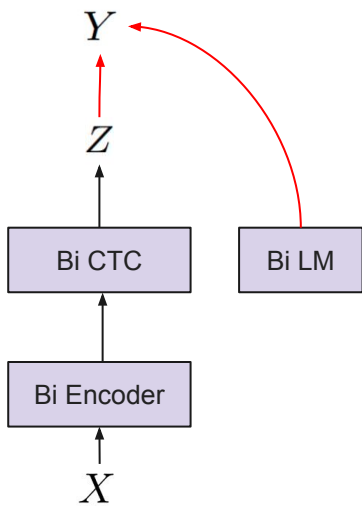
# Code-switching (CS) ⊂ Bilingualism

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



$X^M, Y^M$   $X^{CS}, Y^{CS}$   $X^E, Y^E$

Data   Data   Data

**Mandarin Utterances**   **Intra-sentential Code-Switching**   **English Utterances**

**Inter-sentential Code-Switching**

# Joint Modeling of Monolingual and CS ASR

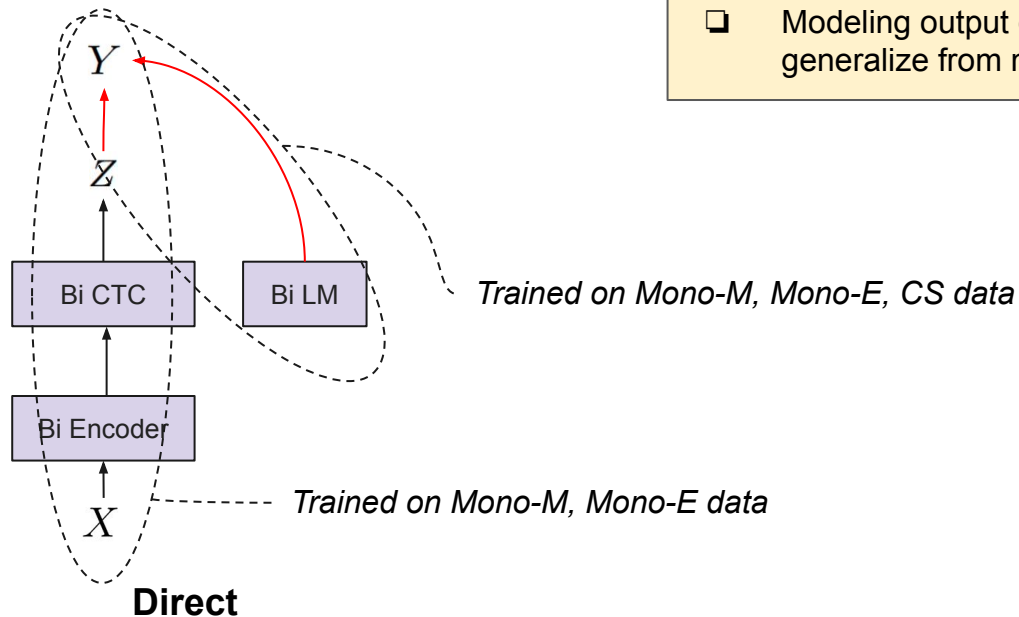$$p(Y|X) \approx \underbrace{p(Y)}_{\triangleq p_{\text{Bi\_LM}}(Y)} \underbrace{\sum_{\mathcal{Z}} p(Z|X)}_{\triangleq p_{\text{Bi\_CTC}}(Y|X)}$$

**Bilingual Modules**:
handle speech/text which may be Mandarin-only, English-only, or code-switched

$Y$

$Z$

| Bi CTC | | Bi LM |
| --- | --- | --- |

| Bi Encoder |
| --- |

$X$

**Direct**

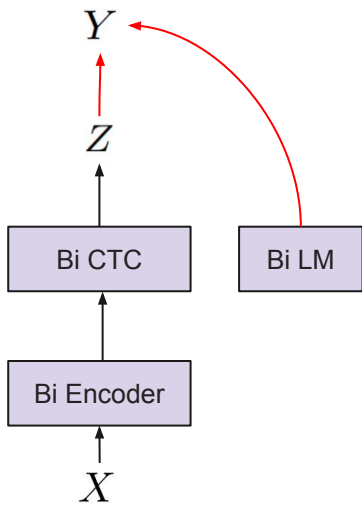# Joint Modeling of Monolingual and CS ASR

$$p(Y|X) \approx \underbrace{p(Y)}_{\triangleq p_{\text{Bi\_LM}}(Y)} \underbrace{\sum_{\mathcal{Z}} p(Z|X)}_{\triangleq p_{\text{Bi\_CTC}}(Y|X)}$$

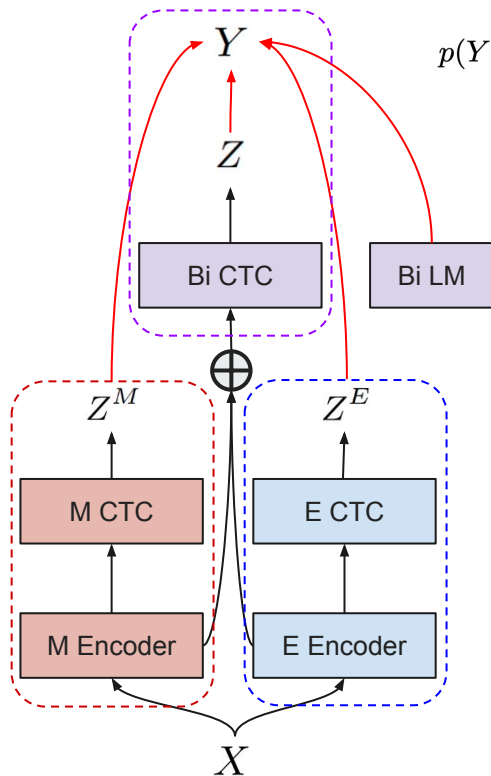❏ Modeling output dependency via an external LM → generalize from monolingual ASR training to CS testing



*Trained on Mono-M, Mono-E, CS data*

*Trained on Mono-M, Mono-E data*

**Direct**

# Joint Modeling of Monolingual and CS ASR

$$p(Y|X) \approx \underbrace{p(Y)}_{\triangleq p_{\text{Bi\_LM}}(Y)} \underbrace{\sum_{\mathcal{Z}} p(Z|X)}_{\triangleq p_{\text{Bi\_CTC}}(Y|X)}$$

$$p(Y|X) \approx \underbrace{p(Y)}_{\triangleq p_{\text{Bi\_LM}}(Y)} \underbrace{\sum_{\mathcal{Z}} p(Z|Z^M, Z^E)}_{\triangleq p_{\text{Bi\_CTC}}(Y|Z^M, Z^E)} \underbrace{\sum_{\mathcal{Z}^M} p(Z^M|X)}_{\triangleq p_{\text{M\_CTC}}(Y^M|X)} \underbrace{\sum_{\mathcal{Z}^E} p(Z^E|X)}_{\triangleq p_{\text{E\_CTC}}(Y^E|X)}$$
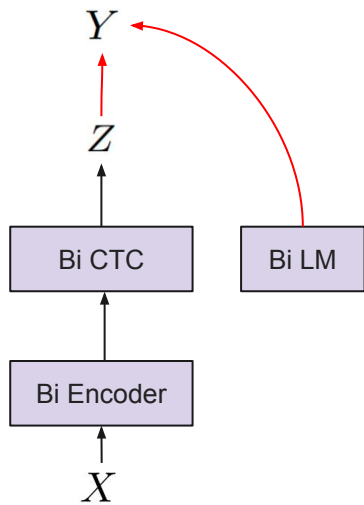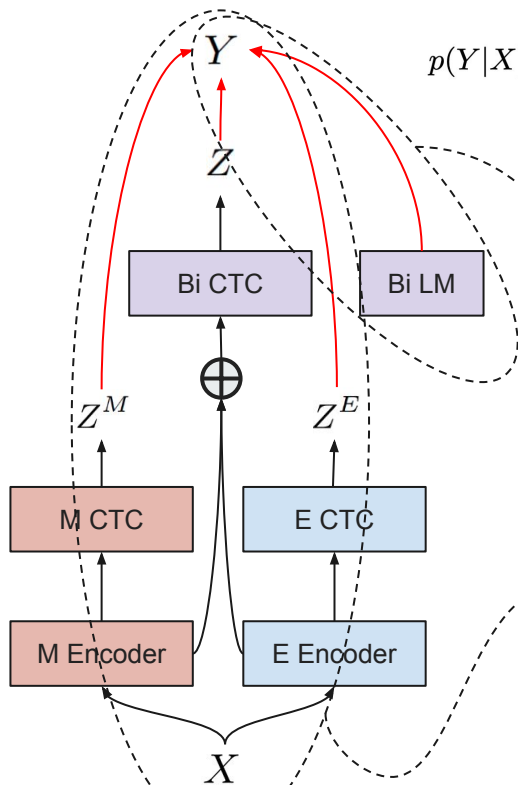


Bilingual Modules +
Monolingual English Expert +
Monolingual Mandarin Expert

**Direct**

**Conditionally Factorized**

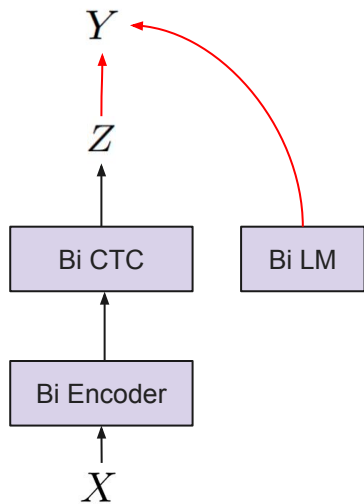2

# Joint Modeling of Monolingual and CS ASR



Trained on Mono-M, Mono-E, CS data

Trained on Mono-M, Mono-E data

**Direct**

**Conditionally Factorized**

# Joint Modeling of Monolingual and CS ASR



**Direct**

**Conditionally Factorized**

❏ Dedicated monolingual sub-components → data efficient training

❏ Re-framed the bilingual task → choosing the language per $z_i$ given monolingual information

2

# Joint Modeling of Monolingual and CS ASR



*Training Scheme*

$$Y|X^E = \text{\_account} \quad \text{\_ing} \qquad Y|X^M = \text{还} \quad \text{有}$$

$$Y^M|X^E = \quad \text{[null]} \quad \text{[null]} \qquad Y^M|X^M = \text{还} \quad \text{有}$$

$$Y^E|X^E = \text{\_account} \quad \text{\_ing} \qquad Y^E|X^M = \text{[null]} \quad \text{[null]}$$

$$\mathcal{L}_{\text{LS}} = \lambda \mathcal{L}_{\text{Bi\_CTC}} + (1 - \lambda)(\mathcal{L}_{\text{M\_CTC}} + \mathcal{L}_{\text{E\_CTC}})$$

# Joint Modeling of Monolingual and CS ASR



*Training Scheme*

$$Y|X^{CS} = \quad \text{\_account} \quad \text{\_ing} \quad \quad 还 \quad 有$$

$$Y^M|X^{CS} = \quad \text{[null]} \quad \text{[null]} \quad \quad 还 \quad 有$$

$$Y^E|X^{CS} = \quad \text{\_account} \quad \text{\_ing} \quad \quad \text{[null]} \quad \text{[null]}$$

$$\mathcal{L}_{\text{LS}} = \lambda\mathcal{L}_{\text{Bi\_CTC}} + (1-\lambda)(\mathcal{L}_{\text{M\_CTC}} + \mathcal{L}_{\text{E\_CTC}})$$
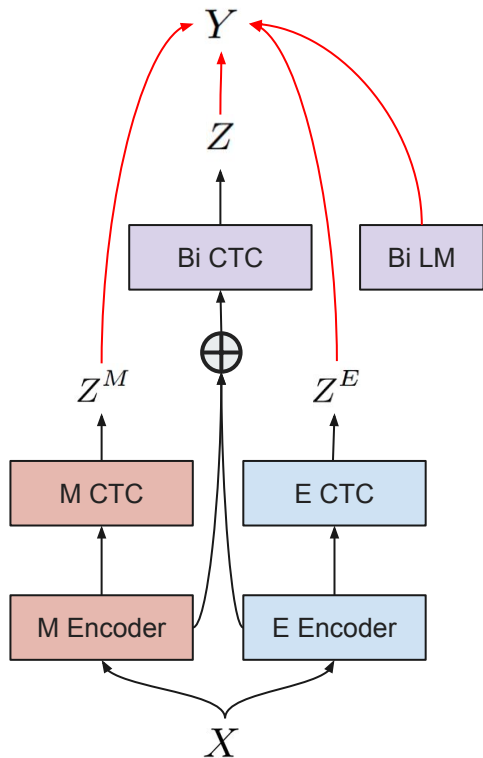
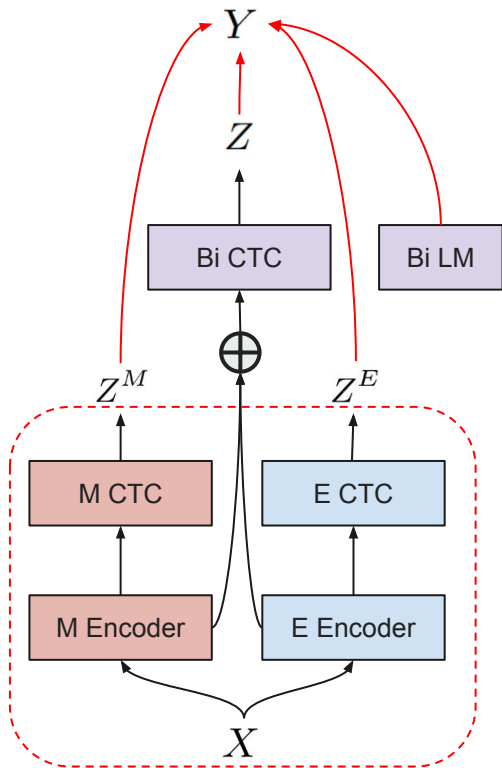# Joint Modeling of Monolingual and CS ASR



### *Training Scheme*

$$Y | X^{CS} = \text{\_account} \quad \text{\_ing} \quad 还 \quad 有$$

$$Y^M | X^{CS} = \text{[null]} \quad \text{[null]} \quad 还 \quad 有$$

$$Y^E | X^{CS} = \text{\_account} \quad \text{\_ing} \quad \text{[null]} \quad \text{[null]}$$

$$\mathcal{L}_{\text{LS}} = \lambda \mathcal{L}_{\text{Bi\_CTC}} + (1 - \lambda)(\mathcal{L}_{\text{M\_CTC}} + \mathcal{L}_{\text{E\_CTC}})$$

### *Inference Procedure*

1. Monolingual CTC modules transcribe their respective parts

2. Bilingual CTC module transcribes whole, conditioned on monolingual info.

3. Mono/bilingual CTC modules + bilingual LM jointly decode the final output sequence (e.g. via time sync beam search)

2

# Joint Modeling of Monolingual and CS ASR



_Training Scheme_

$$Y|X^{CS} = \text{\_account} \quad \text{\_ing} \quad 还 \quad 有$$

$$Y^M|X^{CS} = [\text{null}] \quad [\text{null}] \quad 还 \quad 有$$

$$Y^E|X^{CS} = \text{\_account} \quad \text{\_ing} \quad [\text{null}] \quad [\text{null}]$$

$$\mathcal{L}_{LS} = \lambda\mathcal{L}_{Bi\_CTC} + (1-\lambda)(\mathcal{L}_{M\_CTC} + \mathcal{L}_{E\_CTC})$$
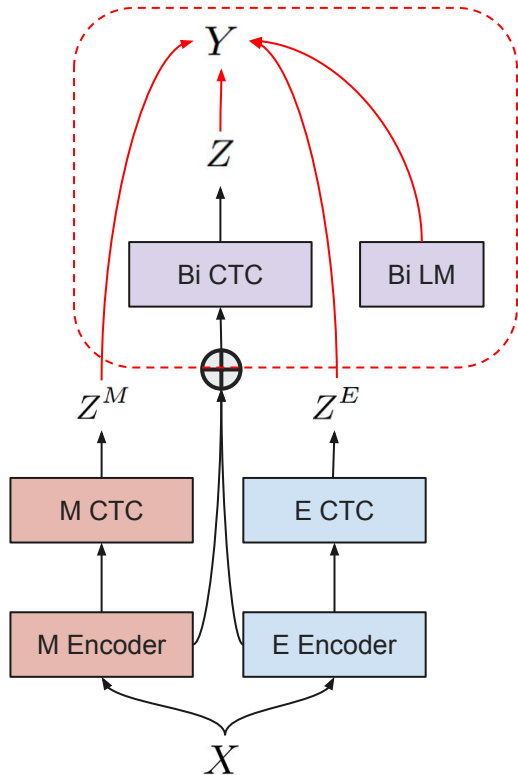
_Inference Procedure_            *Making a language segmentation decision*

1.  Monolingual CTC modules transcribe their respective parts

2.  Bilingual CTC module transcribes whole, conditioned on monolingual info.

3.  Mono/bilingual CTC modules + bilingual LM jointly decode the final output sequence (e.g. via time sync beam search)

2

# Joint Modeling of Monolingual and CS ASR



*Language segmentation decision*

*Can we make the language segmentation decision later?*
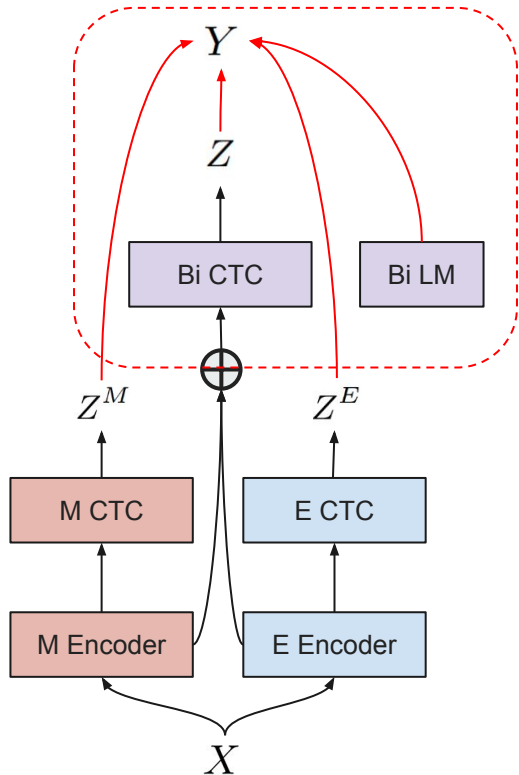
*Inference Procedure*

*Making a language segmentation decision*

1. Monolingual CTC modules transcribe their respective parts

2. Bilingual CTC module transcribes whole, conditioned on monolingual info.

3. Mono/bilingual CTC modules + bilingual LM jointly decode the final output sequence (e.g. via time sync beam search)
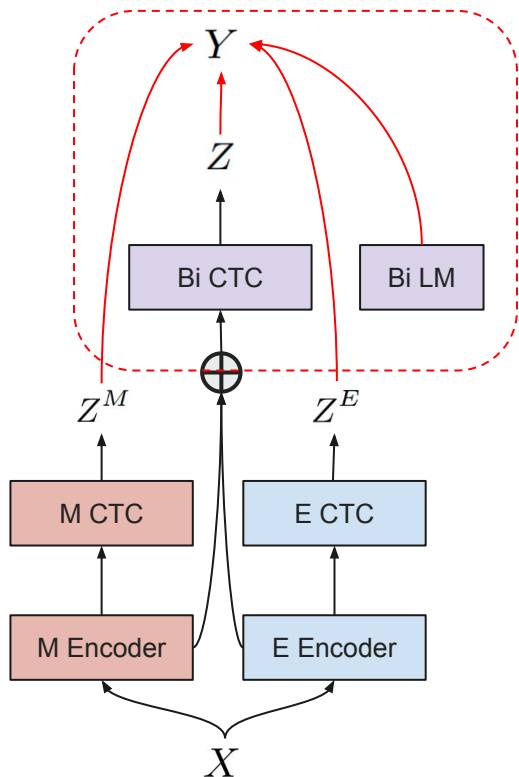
2

# Joint Modeling of Monolingual and CS ASR



❏ Encourage monolingual modules to transcribe the opposite language → leave language segmentation decision to bilingual modules (CTC, LM)

# Joint Modeling of Monolingual and CS ASR

*Language segmentation decision*
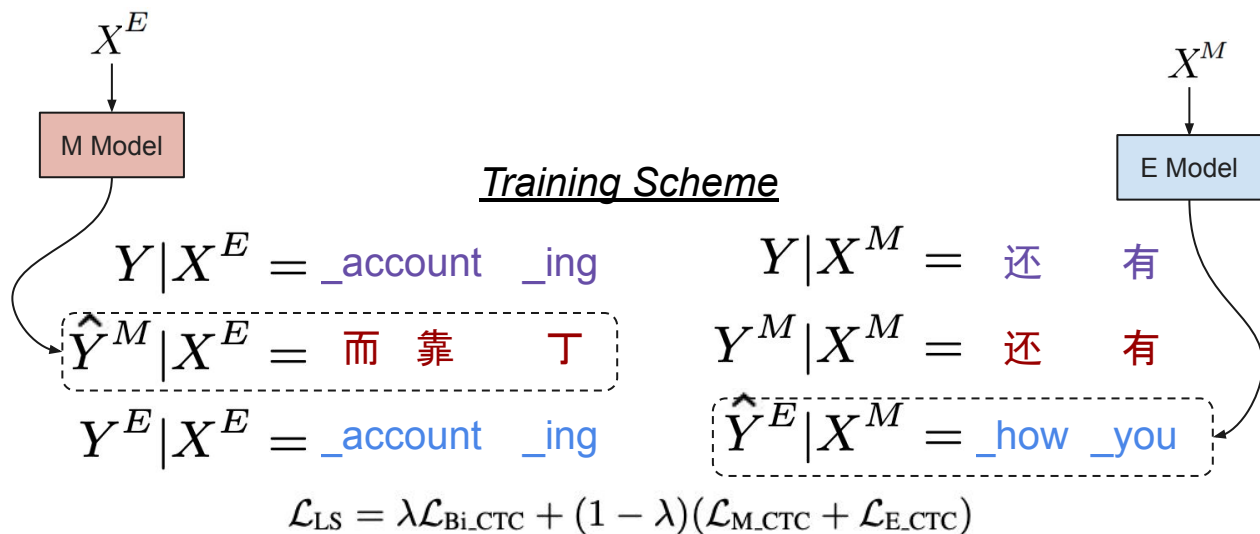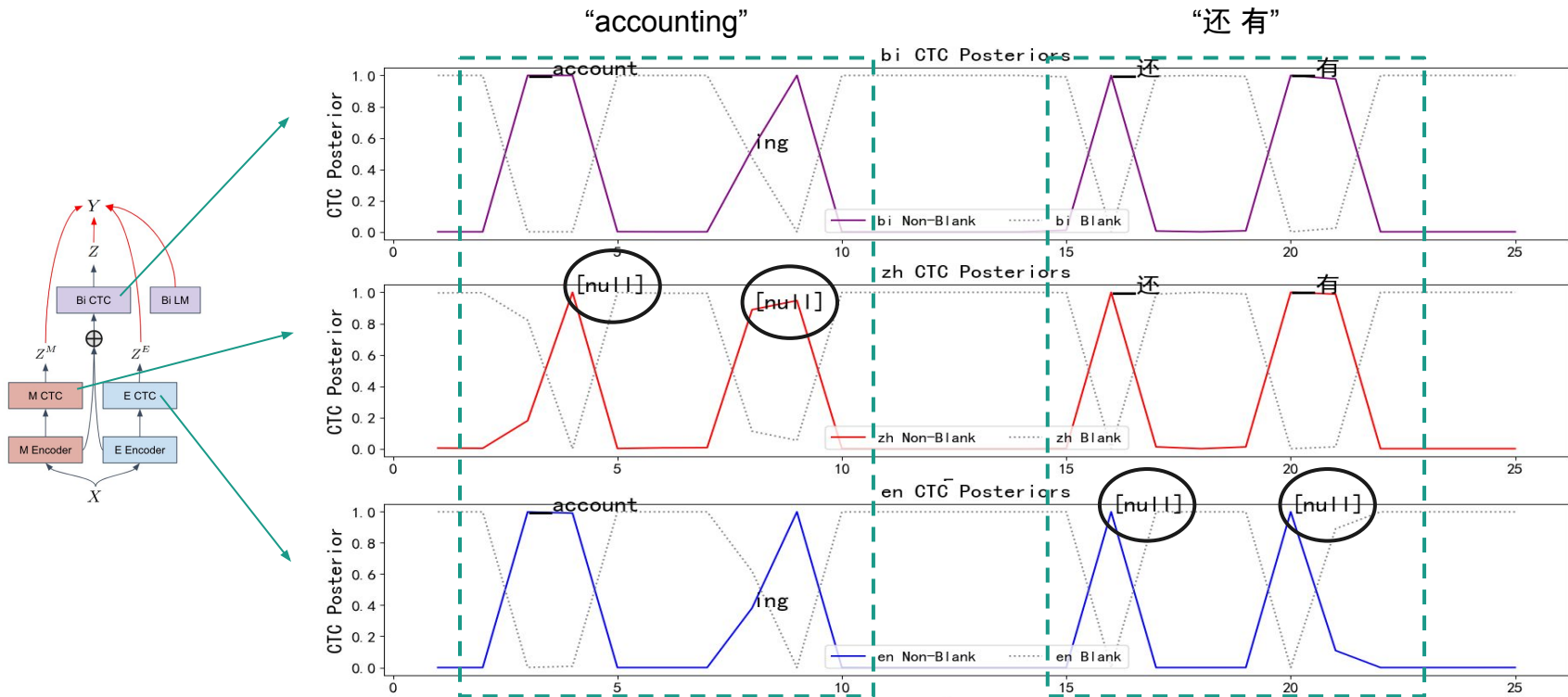


❏ Encourage monolingual modules to transcribe the opposite language → leave language segmentation decision to bilingual modules (CTC, LM)

*Training Scheme*

$$Y|X^E = \text{\_account \_ing} \qquad Y|X^M = 还 \quad 有$$

$$\hat{Y}^M|X^E = 而 \quad 靠 \quad 丁 \qquad Y^M|X^M = 还 \quad 有$$

$$Y^E|X^E = \text{\_account \_ing} \qquad \hat{Y}^E|X^M = \text{\_how \_you}$$

$$\mathcal{L}_{\text{LS}} = \lambda\mathcal{L}_{\text{Bi\_CTC}} + (1-\lambda)(\mathcal{L}_{\text{M\_CTC}} + \mathcal{L}_{\text{E\_CTC}})$$
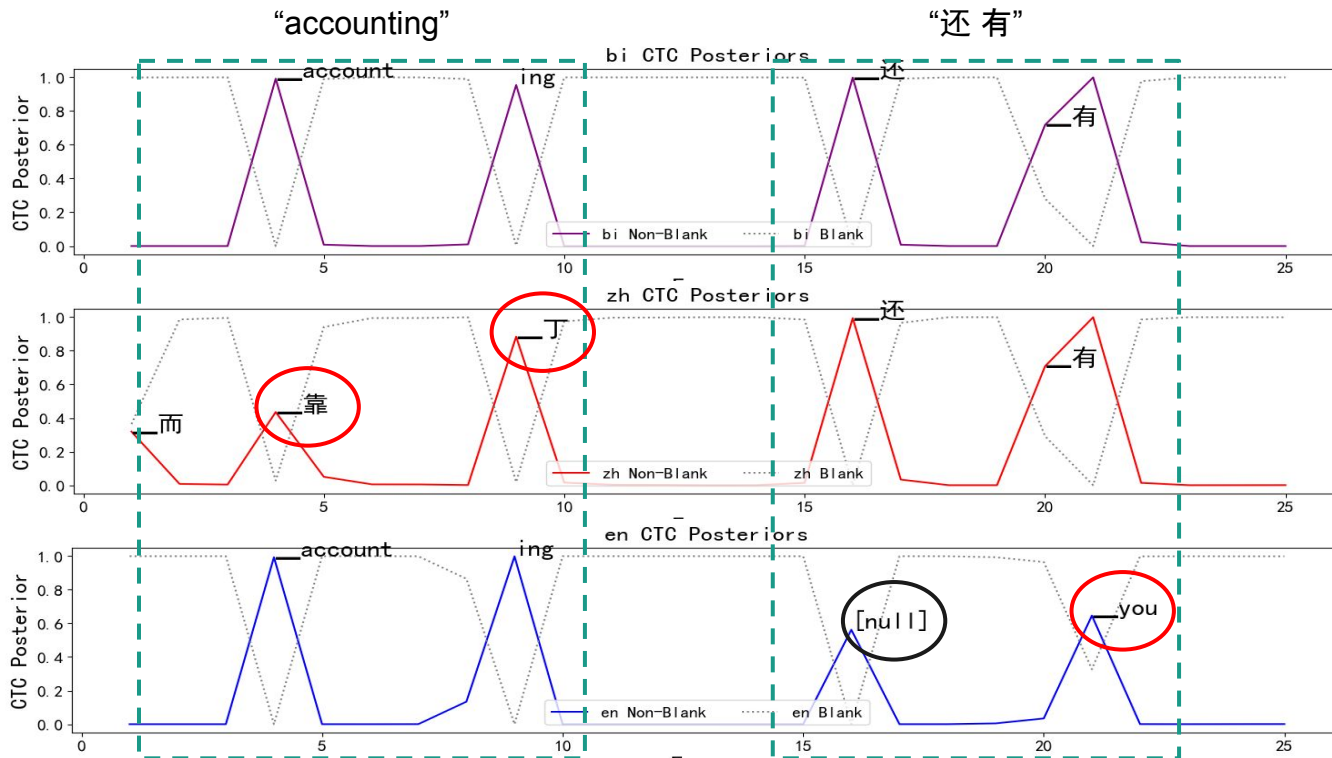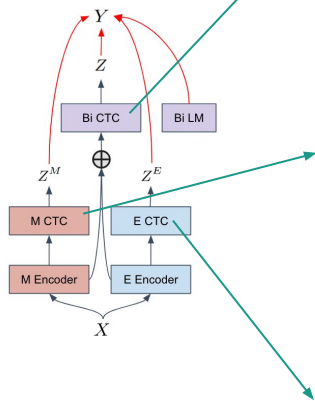
# Qualitative Example: Conditional CTC Posteriors

Given CS ASR training data, early language segmentation works well
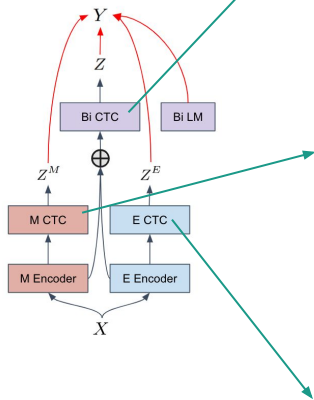
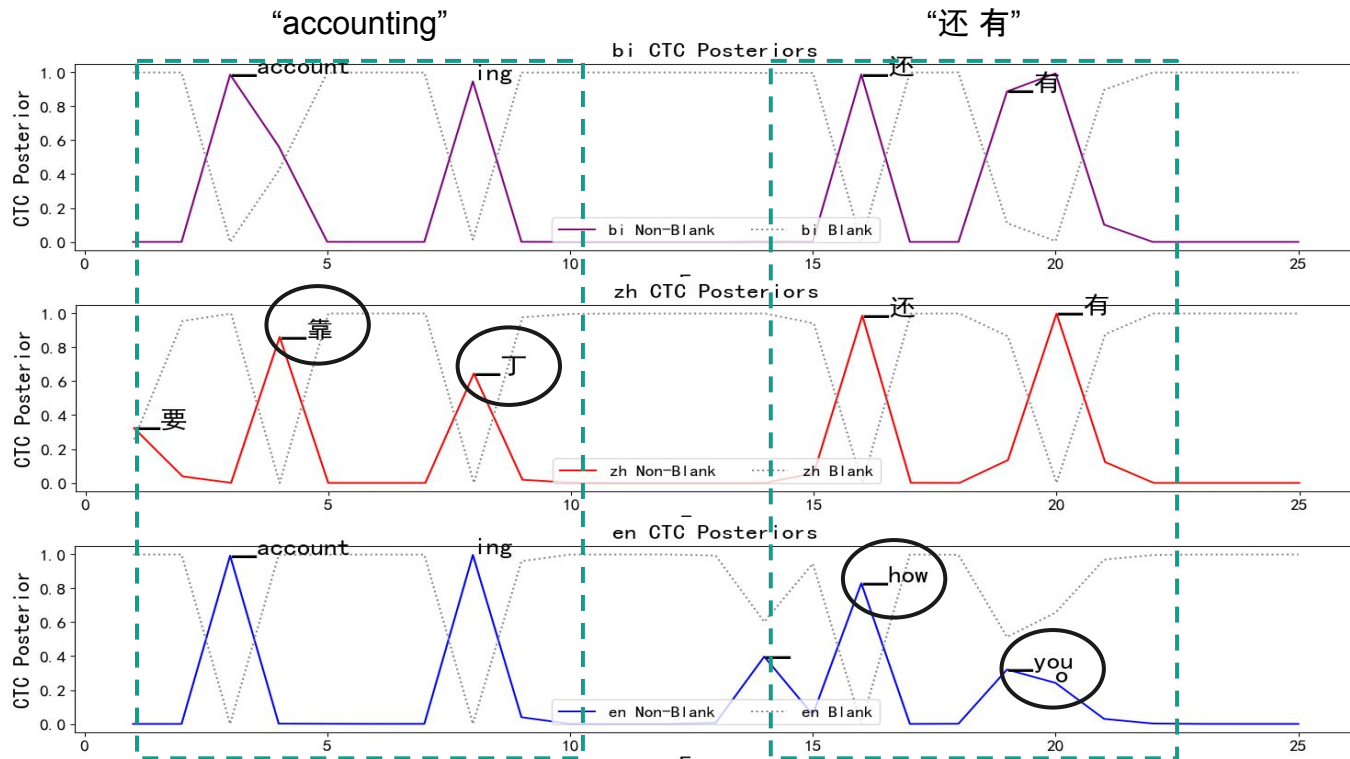# Qualitative Example: Conditional CTC Posteriors

Without CS ASR training data, early language segmentation is unreliable

# Qualitative Example: Conditional CTC Posteriors

- Monolingual modules produce smooth likelihoods for opposite lang. instead of [null]
- Language separation information is soft; LM can help decide → late decision

# Results

| Model | Language Segmentation | ASR Data | LM Data | devman MER($\downarrow$) |
|---|---|---|---|---|
| Conditional CTC | Early | CS + M | - | 17.5 |
| Conditional CTC + LM | Early | CS + M | CS + M | **16.8** |
| Conditional CTC | Early | CS | - | 32.3 |
| Conditional CTC + LM | Early | CS | CS + M | 30.1 |
| Conditional CTC | Late | CS | - | 27.9 |
| Conditional CTC + LM | Late | CS | CS + M | **25.2** |

+13.3 MER

-4.9 MER

# **Takeaways**

➢ Language segmentation of code-switched speech is hard, especially if we don't have code-switched supervision

➢ Making later decisions about language segmentation is better, allowing us to consider more information (e.g. external LM)