



# Multi-Level Contrastive Learning for Cross-Lingual Alignment

Beiduo Chen<sup>1</sup>, Wu Guo<sup>1</sup>, Bin Gu<sup>1</sup>, Quan Liu<sup>2</sup>, Yongchao Wang<sup>2</sup>

<sup>1</sup>National Engineering Research Center for Speech and Language Information Processing, University of Science and Technology of China

<sup>2</sup>State Key Laboratory of Cognitive Intelligence, iFLYTEK Research

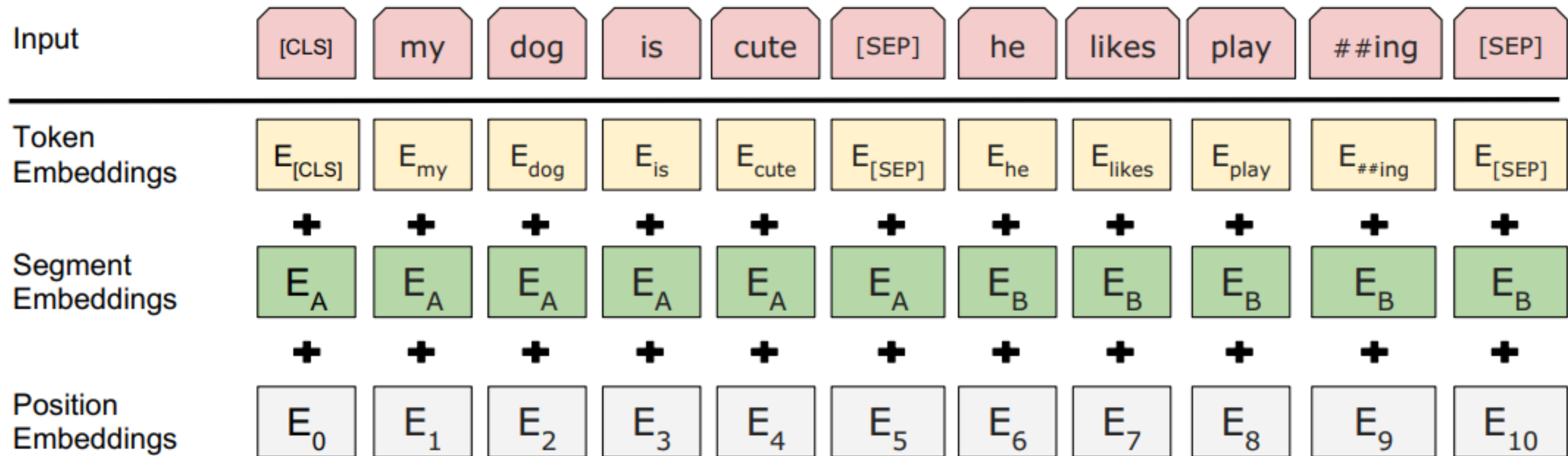


# Outline

- **Introduction**
- Multi-level contrastive learning (ML-CTL)
- Cross-zero noise contrastive estimation loss (CZ-NCE)
- Experiments and analyses
- Conclusion

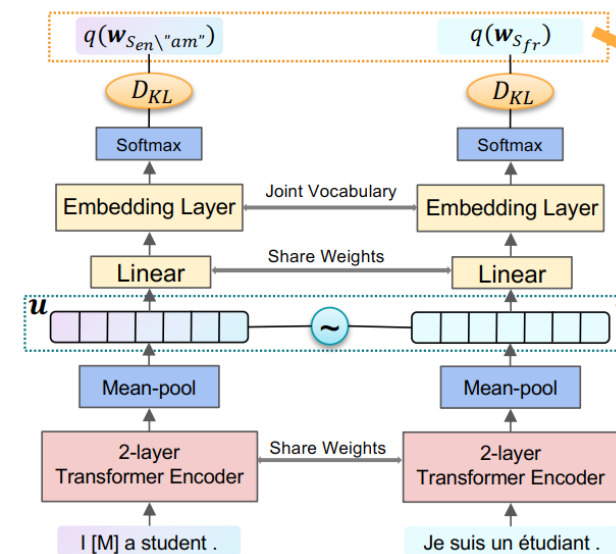
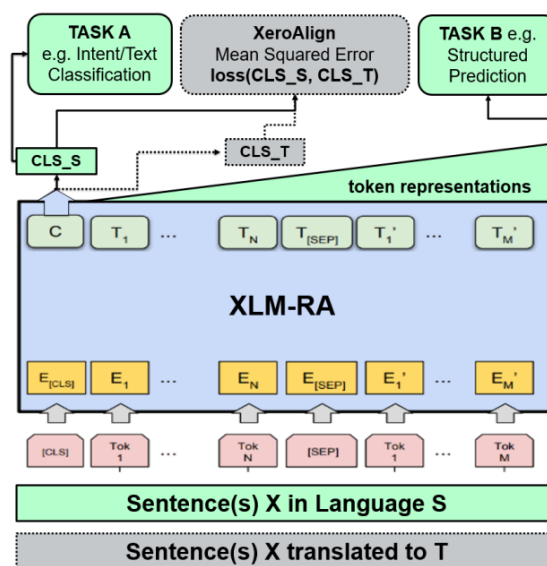
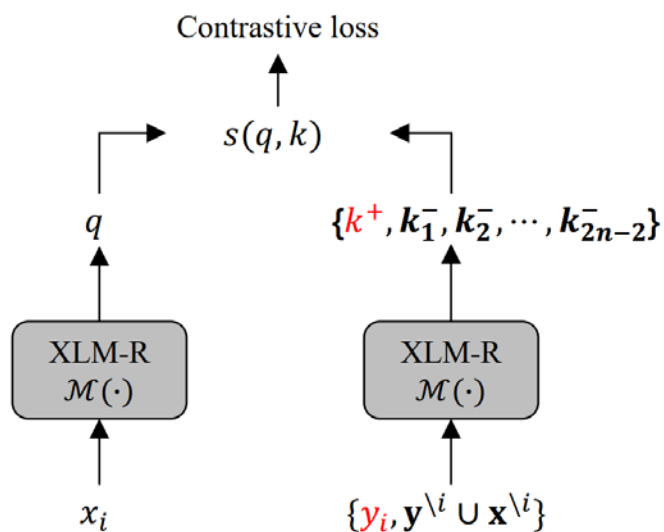
# Cross-lingual Language Model

- Transformer based: mBERT, XLM
- MLM: mask language modeling



# Contrastive Learning

- Issue: no explicit cross-lingual alignment
- Exist solution: contrastive learning (CTL)





# Contrastive Learning

- Issue: no explicit cross-lingual alignment
- Exist solution: contrastive learning (CTL)
  - New issue: only sentence-level CTL

sentence-level → **sentence**-level + **word**-level

**Multi-Level** Contrastive Learning



# Computational Resources

- Issue: high demand for computational resources
- InfoNCE:

$$L_{infoNCE} = -\log\left(\frac{e^{s^+}}{e^{s^+} + \sum e^{s_i^-}}\right)$$

Especially serious in contrastive learning

→ **Cross-zero** NCE loss

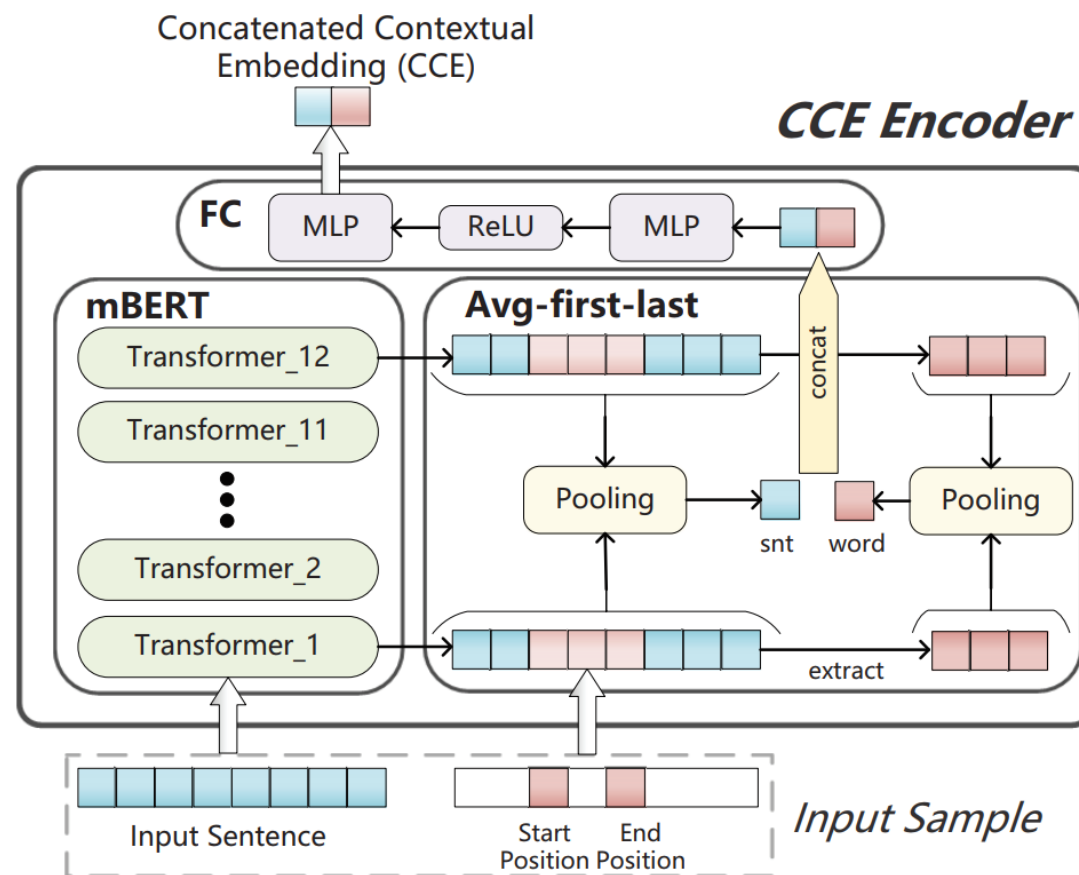


# Outline

- Introduction
- **Multi-level contrastive learning (ML-CTL)**
- Cross-zero noise contrastive estimation loss (CZ-NCE)
- Experiments and analyses
- Conclusion

# ML-CTL

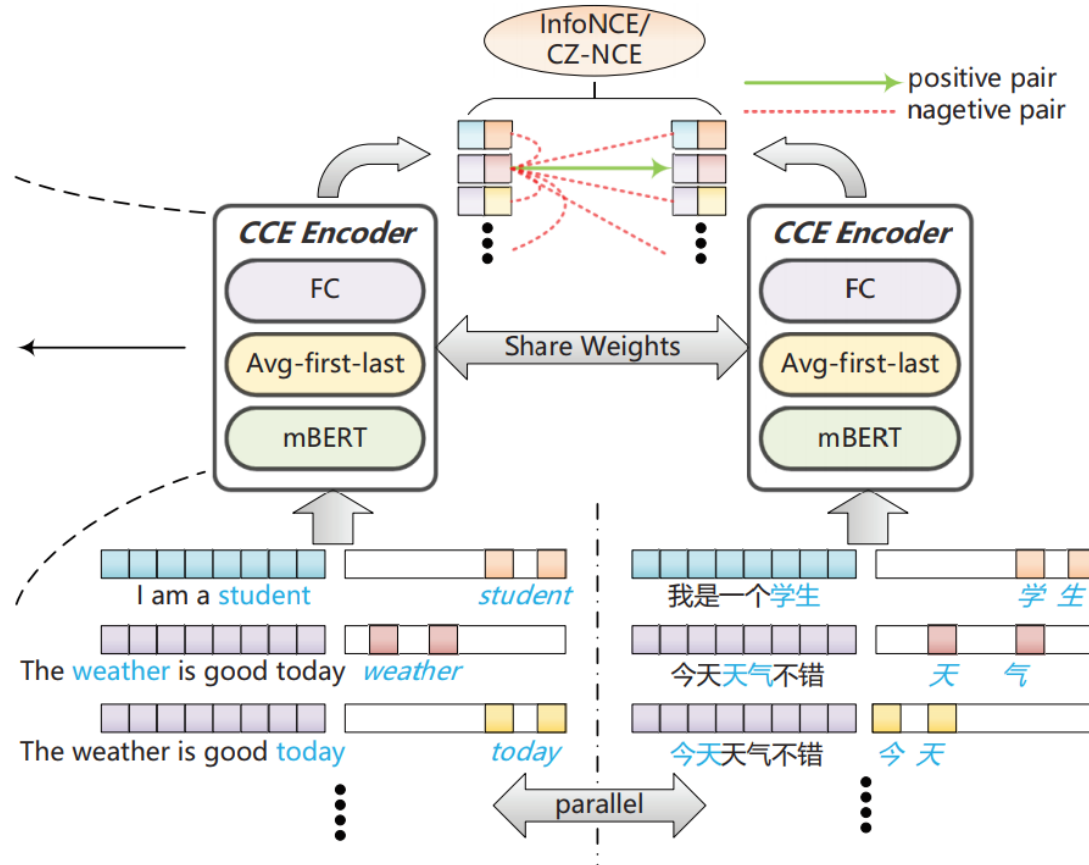
- Concatenated Contextual Embedding Encoder





# ML-CTL

- Multi-level contrastive learning framework



$$L_{info}(\mathbf{x}_{ci}) = -\log\left(\frac{e^{s(\mathbf{x}_{ci}, \mathbf{y}_{ci})}}{e^{s(\mathbf{x}_{ci}, \mathbf{y}_{ci})} + \sum_{\mathbf{k}_j^-} e^{s(\mathbf{x}_{ci}, \mathbf{k}_j^-)}}\right)$$

$$L_{info\_batch} = \frac{\sum_i^n (L_{info}(\mathbf{x}_{ci}) + L_{info}(\mathbf{y}_{ci}))}{2n}$$

$$L_{multi_1} = L_{info\_batch} + \alpha * L_{MLM}$$



# Outline

- Introduction
- Multi-level contrastive learning (ML-CTL)
- **Cross-zero noise contrastive estimation loss (CZ-NCE)**
- Experiments and analyses
- Conclusion



# CZ-NCE

$$L_{infoNCE} = -\log\left(\frac{e^{s^+}}{e^{s^+} + \sum e^{s_i^-}}\right) \Rightarrow L_{CZ-NCE} = -\log\left(\frac{e^{s^+}}{\sum e^{s_i^-}}\right)$$

# CZ-NCE

- Demonstration

$$L_{CZ-NCE} = -\log\left(\frac{e^{s^+}}{\sum e^{s_i^-}}\right) = \log\left(\sum e^{s_i^- - s^+}\right) = \log(\varphi)$$

$$\nabla_{\theta} L_{CZ-NCE} = \nabla_{\theta} \log(\varphi) = \frac{\nabla_{\theta} \varphi}{\varphi} = \nabla_{\theta} \left(\frac{\varphi}{\text{sg}(\varphi)}\right)$$

$$L_{multi2} = L_{CZ-NCE\_batch} + \alpha * L_{MLM}$$



# Outline

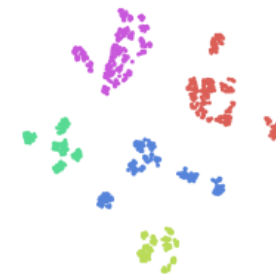
- Introduction
- Multi-level contrastive learning (ML-CTL)
- Cross-zero noise contrastive estimation loss (CZ-NCE)
- **Experiments and analyses**
- Conclusion

# Experiments and analyses

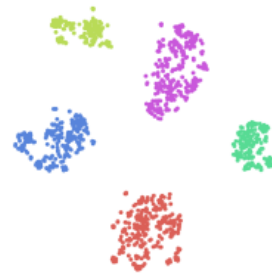
Task	XNLI	PAWS-X	POS	NER	BUCC	TATOEBA
Model\Metric	Acc. (%)	Acc. (%)	F1 (%)	F1 (%)	F1 (%)	Acc. (%)
<i>Main results compared to strong baselines</i>						
mBERT (base)	65.4	81.9	70.3	62.2	56.7	38.7
XLM	<b>69.1</b>	80.9	70.1	61.2	56.8	32.6
MMTE	67.4	81.3	72.3	58.3	59.8	37.9
ML-CTL-CZ	67.8	<b>85.3</b>	<b>72.3</b>	<b>62.9</b>	<b>78.4</b>	<b>43.4</b>
<i>Results of ablation study</i>						
mBERT (base)	65.4	81.9	70.3	62.2	56.7	38.7
<i>info-snt</i>	66.255	84.092	71.544	62.157	76.426	41.148
<i>CZ-snt</i>	66.862	84.485	71.733	62.337	77.403	41.751
ML-CTL-CZ	<b>67.750</b>	<b>85.321</b>	<b>72.289</b>	<b>62.865</b>	<b>78.440</b>	<b>43.389</b>



(a) mBERT



(b) info-snt



(c) CZ-snt



(d) ML-CTL-CZ



# Outline

- Introduction
- Multi-level contrastive learning (ML-CTL)
- Cross-zero noise contrastive estimation loss (CZ-NCE)
- Experiments and analyses
- **Conclusion**



# Conclusion

- ML-CTL is proposed to improve the cross-lingual alignment ability of pre-trained language models by applying contrastive learning on concatenated contextual embeddings which contain information of both sentences and words.
- CZ-NCE is proposed to alleviate the impact of the floating-point error with a small training batch size.