

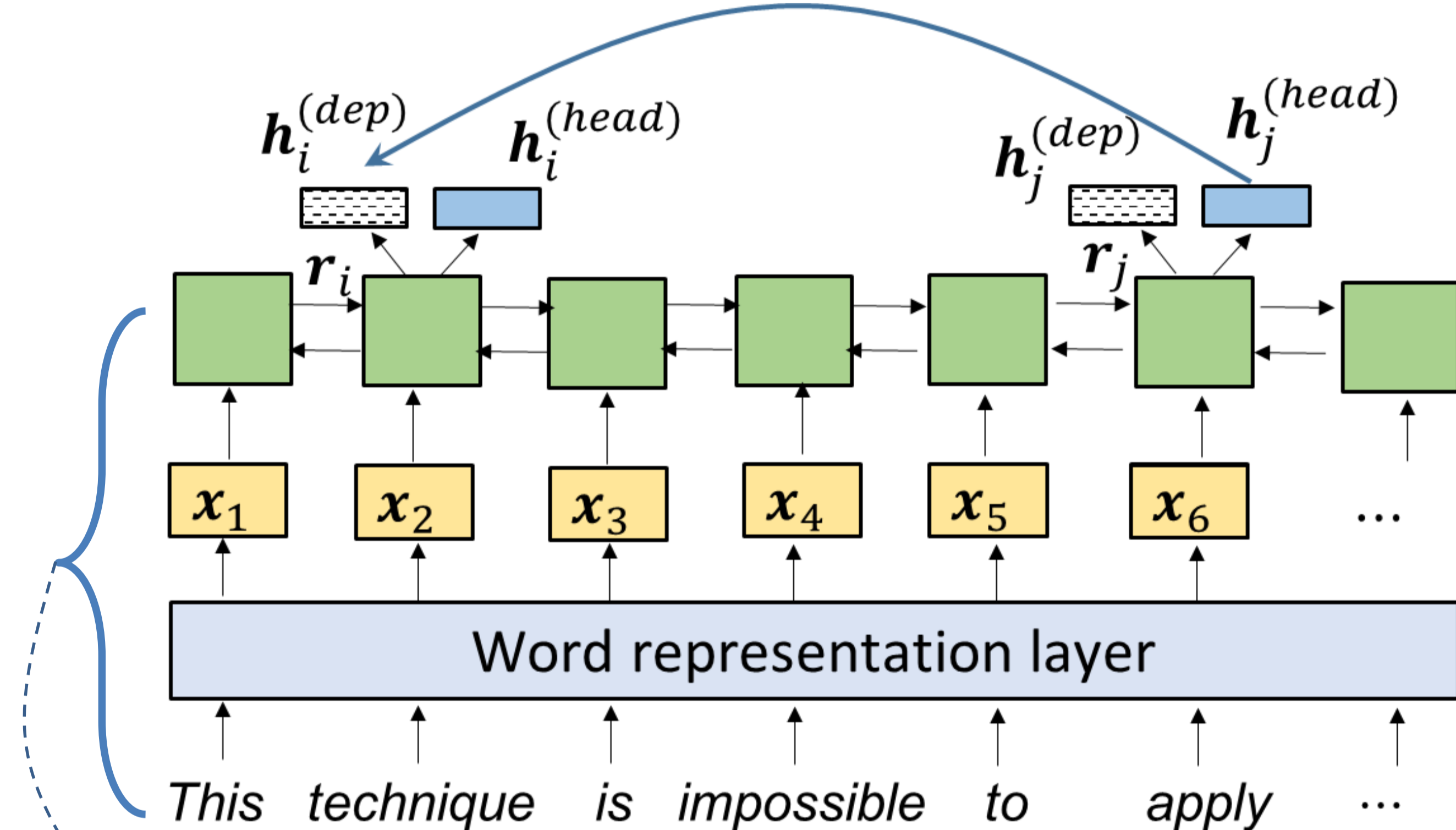
Introduction

- Multi-task learning for DM/PSD/UCCA: Propose a unified neural model for DM/PSD/UCCA frameworks based on **biaffine attention** used in (Dozat and Manning, 2017, 2018; Zhang et al., 2019)

Model architecture

- BERT-BiLSTM sentence encoder (shared across frameworks)**: an input sentence is fed to a word representation layer using BERT, resulting in a sequence of word embedding vectors, which are then given to the BiLSTM layer to produce a sentence representation
- Biaffine attention decoder (framework-specific)**: additional feed-forward layers are applied to obtain role-dependent representations for head and dependent roles, which are then forwarded to the biaffine attention.

Decoder: Biaffine attention



Encoder: BERT-BiLSTM

Encoder: BERT-BiLSTM

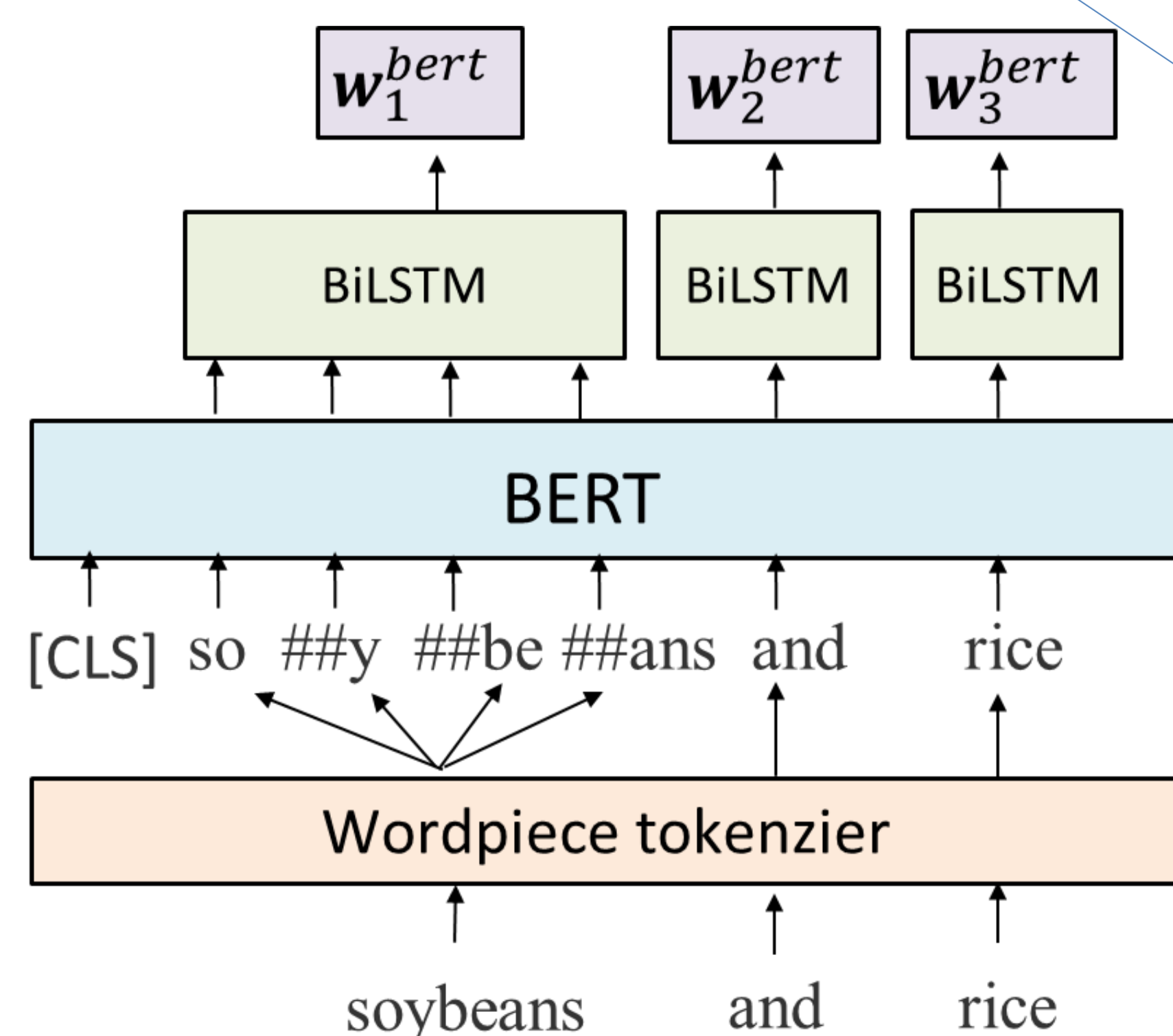
- Word representation layer using BERT**: Given a sentence, the **BERT encoder** is applied to its wordpieces and the encoded wordpiece-level representations are composed to the word-level embeddings based on BiLSTM.

BERT word-level embedding

$$\mathbf{x}_i = [\mathbf{w}_i^{bert}; \mathbf{e}_i^{glove}; \mathbf{e}_i^{POS}]$$

BiLSTM sentence encoding

$$\mathbf{r}_i = BiLSTM_i(\mathbf{x}_1 \cdots \mathbf{x}_n)$$



Biaffine attention decoder

- Biaffine attention is performed on the role-dependent representations to predict the existence of an edge and its labels.

Biaffine attention

$$BiAff_m(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{U}^{[1:m]} \mathbf{y} + \mathbf{V} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} + \mathbf{b}$$

$$s_{i,j}^{(edge)} = BiAff_1^{(edge)}(\mathbf{h}_i^{(dep)}, \mathbf{h}_j^{(head)})$$

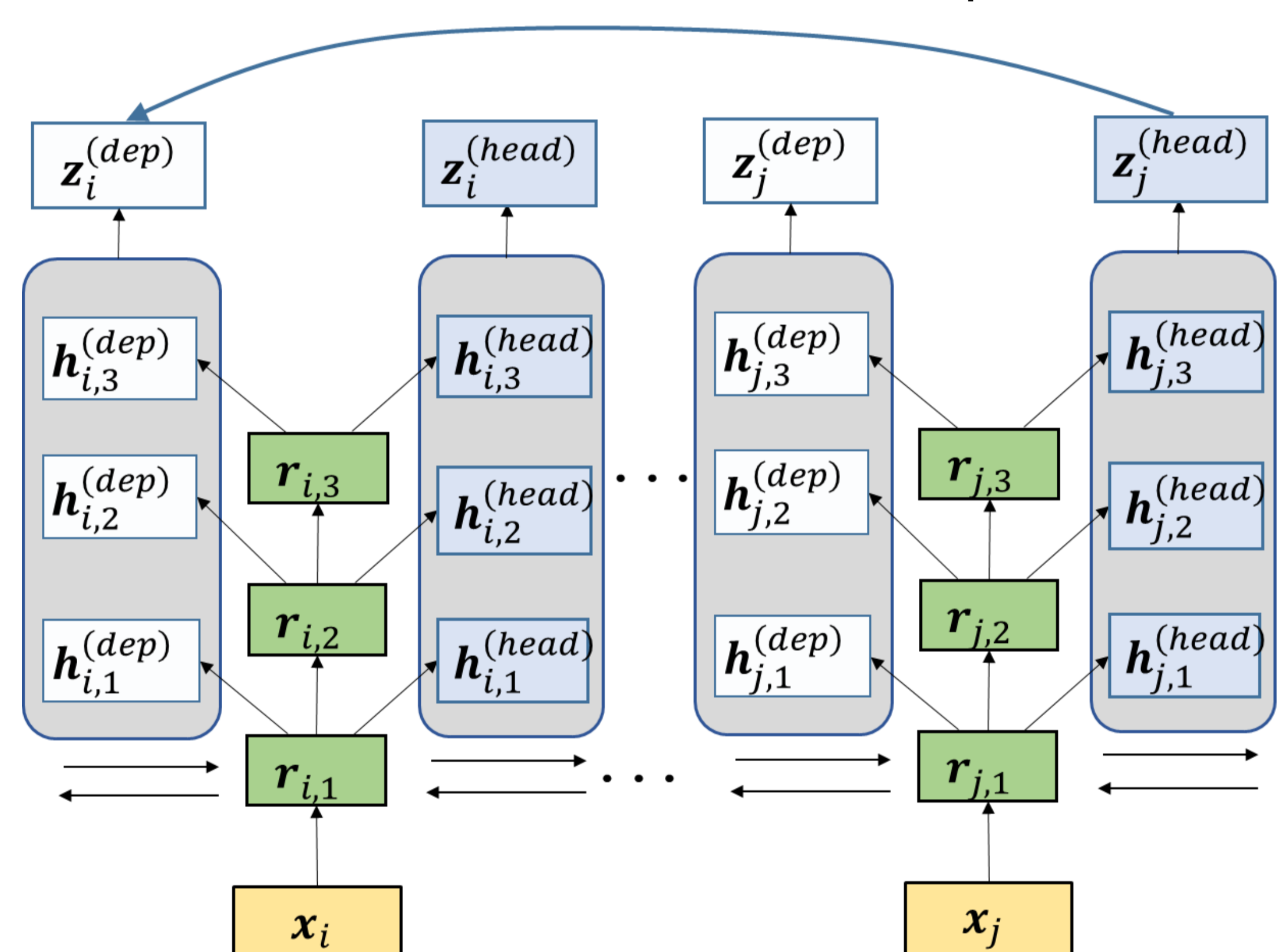
$$s_{i,j}^{(label)} = BiAff_k^{(label)}(\mathbf{h}_i^{(l-dep)}, \mathbf{h}_j^{(l-head)})$$

$$s_i^{(top)} = FFN^{(top)}(\mathbf{r}_i)$$

- The top score $S^{top}(i)$ is newly introduced in our model

Multi-level biaffine attention

- Motivation: Predicting an arc and a label may be resolved not just by single-level representation but by the combination of various levels of representations



Multi-level biaffine attention

$$\mathbf{z}_i^{(head)} = sfu^{(head)}(\mathbf{h}_{i,1}^{(head)}, \mathbf{h}_{i,2}^{(head)}, \mathbf{h}_{i,3}^{(head)})$$

$$\mathbf{z}_i^{(dep)} = sfu^{(dep)}(\mathbf{h}_{i,1}^{(dep)}, \mathbf{h}_{i,2}^{(dep)}, \mathbf{h}_{i,3}^{(dep)})$$

$$s_{i,j}^{(edge')} = BiAff_1^{(edge')}(\mathbf{z}_i^{(dep)}, \mathbf{z}_j^{(head)})$$

Property prediction based on BiLSTM

Use a simple BiLSTM with a single output layer

$$\mathbf{r}_i^{(prop)} = BiLSTM_i^{(prop)}(\mathbf{x}_1 \cdots \mathbf{x}_n) \quad s_i^{(prop)} = FFN^{(prop)}(\mathbf{r}_i^{(prop)})$$

Preliminary Experiment

method	DM			PSD			UCCA		
	Top	UF	LF	Top	UF	LF	Top	UF	LF
Biaffine	93.67	92.08	90.86	95.97	90.50	78.21	72.60	69.67	65.17
BERT+Biaffine	95.06	93.85	93.00	96.89	92.30	80.24	77.09	74.85	70.15
BERT+Multi-level Biaffine	95.09	93.86	93.02	96.76	91.95	79.76	78.12	74.42	69.81
BERT+Biaffine+MTL	N/A	93.66	92.73	N/A	92.13	79.63	N/A	75.40	70.59

- Lesson: instead of naively using the shared encoder only, other advanced multi-task learning approaches such as placing **task-specific encoding**, as detailed in (Peng et al., 2017), need to be considered