

# Neural Machine Translation with Source Dependency Representation

Kehai Chen<sup>1</sup>, Rui Wang<sup>2</sup>, Masao Utiyama<sup>2</sup>, Lemao Liu<sup>3</sup>, Akihiro Tamura<sup>4</sup>,  
Eiichiro Sumita<sup>2</sup> and Tiejun Zhao<sup>1</sup>

<sup>1</sup>*Harbin Institute of Technology, China*

<sup>2</sup>*National Institute of Information and Communications Technology, Japan*

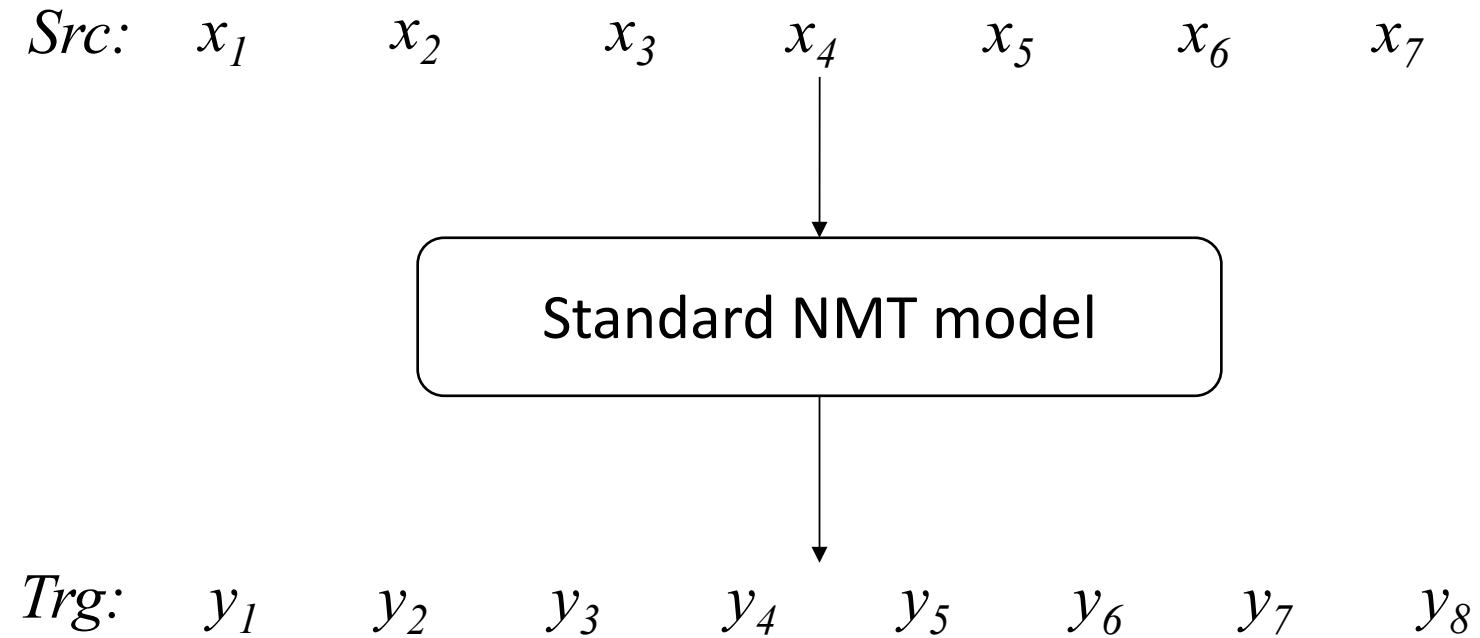
<sup>3</sup>*Tencent AI Lab, China*

<sup>4</sup>*Ehime University, Japan*

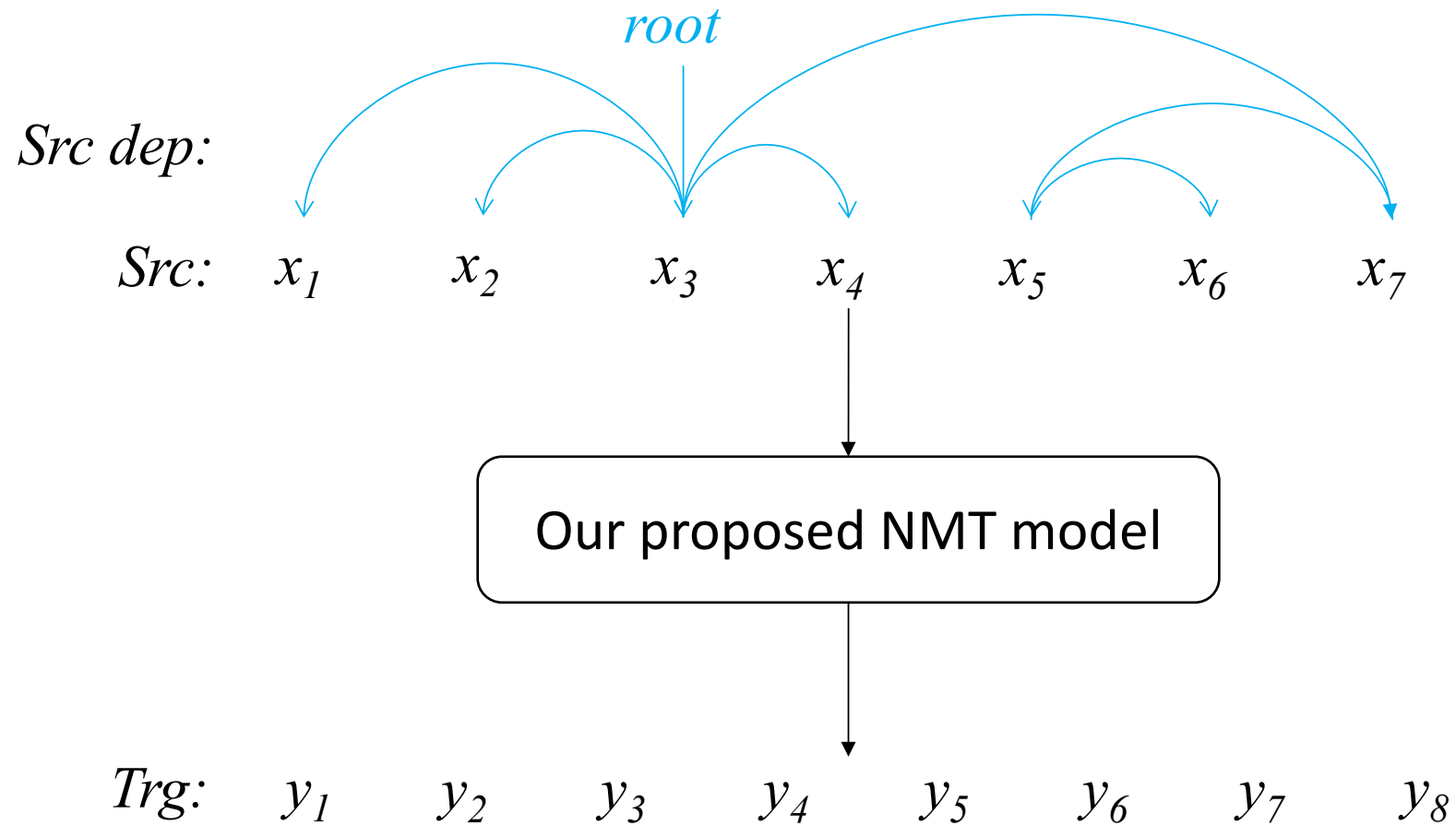


# Overview

- Traditional NMT Model



# Overview



- Our proposed NMT model

Inspired by the syntax knowledge in SMT, we want to explicitly integrate source dependency information into NMT

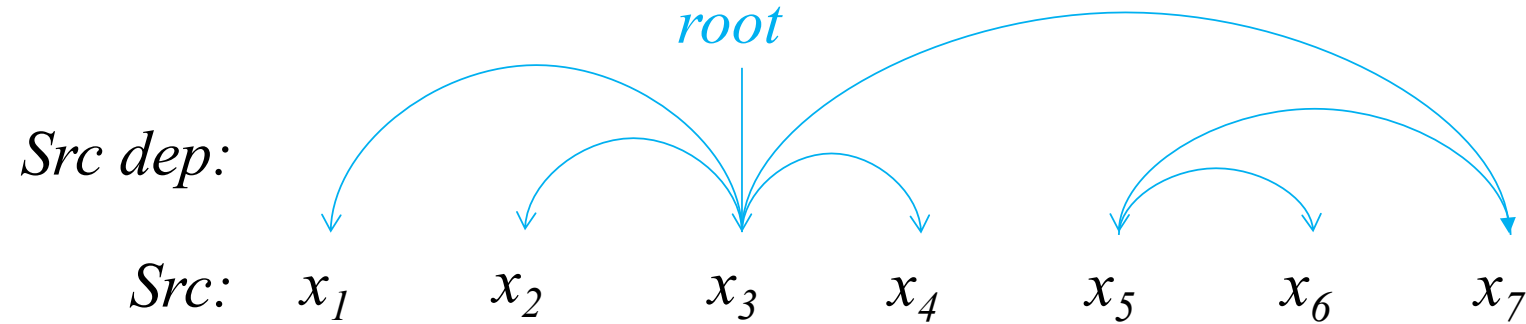
# Related Work

- NMT with source syntax information
  - Tree2seq (Eriguchi et al., 2016; Li et al., 2017; +other)  
Tree-based neural network is used to encode source phrase structures
  - Extending source inputs with syntax labels (Sennrich et al., 2016; Chen et al., 2017; +other)  
Dependency labels are concatenated to source word

# Related Work

- NMT with source syntax information
  - Tree2seq (Eriguchi et al., 2016; Li et al., 2017; +other)  
Tree-based neural network is used to encode source phrase structures
  - Extending source inputs with syntax labels (Sennrich et al., 2016; Chen et al., 2017; +other)  
Dependency labels are concatenated to source word
- Our work
  - A compromise between the two kinds of works
  - A novel double context approach to utilizing source dependency constraints

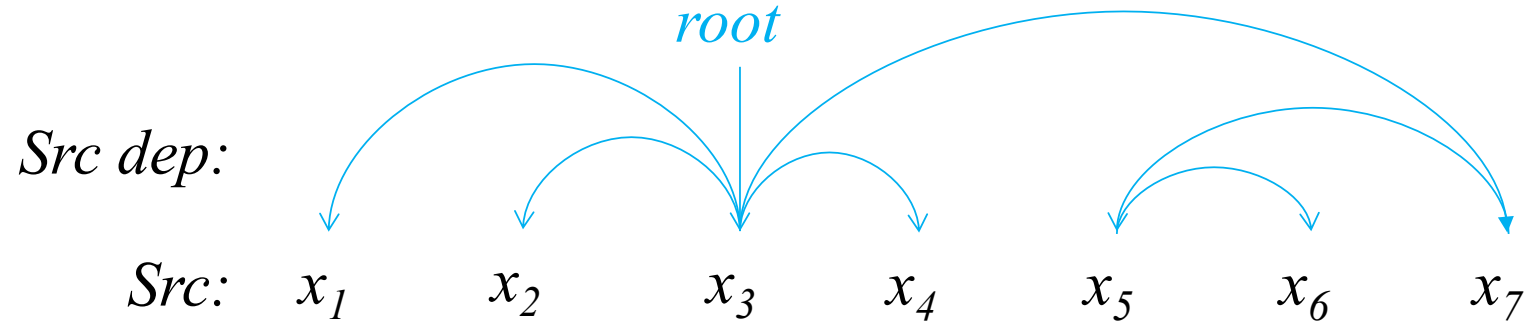
# Source Dependency Representation (SDR)



- Extracting a dependency unit for each source word to capture source long-distance dependency constraints:

$$U_j = \langle PA_{x_j}, SI_{x_j}, CH_{x_j} \rangle$$

# Source Dependency Representation (SDR)



- Extracting a dependency unit for each source word to capture source long-distance dependency constraints:

$$U_j = \langle PA_{x_j}, SI_{x_j}, CH_{x_j} \rangle$$

Where  $PA_{x_j}$ ,  $SI_{x_j}$ , and  $CH_{x_j}$  denote the parent, siblings and children words of source word  $x_j$  in a dependency structure.

Take  $x_2$  as an example:  $PA_{x_2} = \langle x_3 \rangle$ , then,  $U_2 = \langle x_3, x_1, x_4, x_7, \varepsilon \rangle$

$SI_{x_2} = \langle x_1, x_4, x_7 \rangle$ ,

$CH_{x_2} = \langle \varepsilon \rangle$ ,

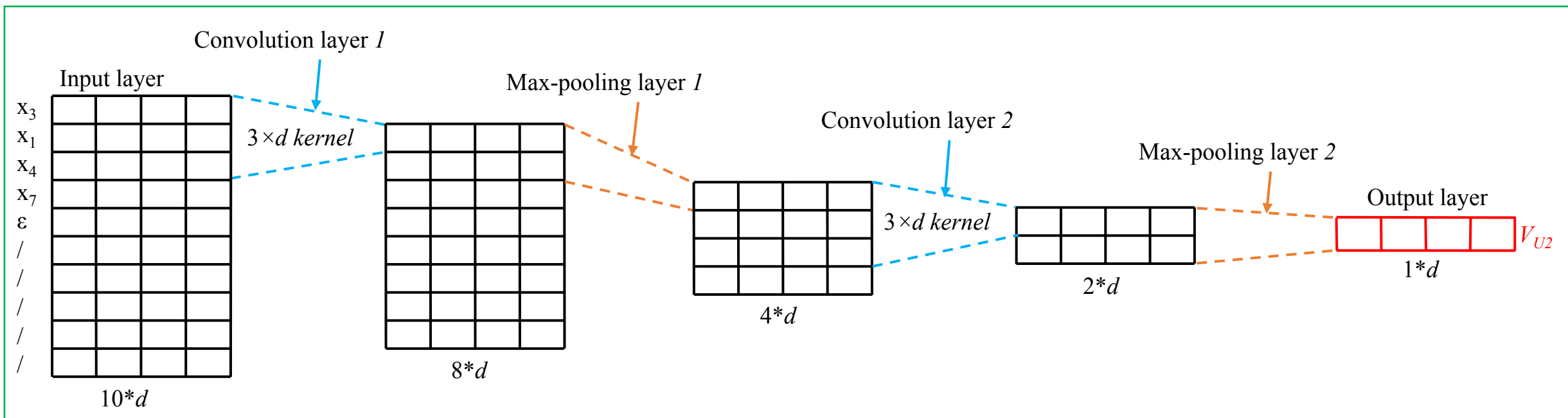
# Source Dependency Representation (SDR)

- Learn semantic representation of each dependency unit

Take  $x_2$  as an example:  $PA_{x_2} = \langle x_3 \rangle$ , then,  $U_2 = \langle x_3, x_1, x_4, x_7, \varepsilon \rangle$

$SI_{x_2} = \langle x_1, x_4, x_7 \rangle$ ,

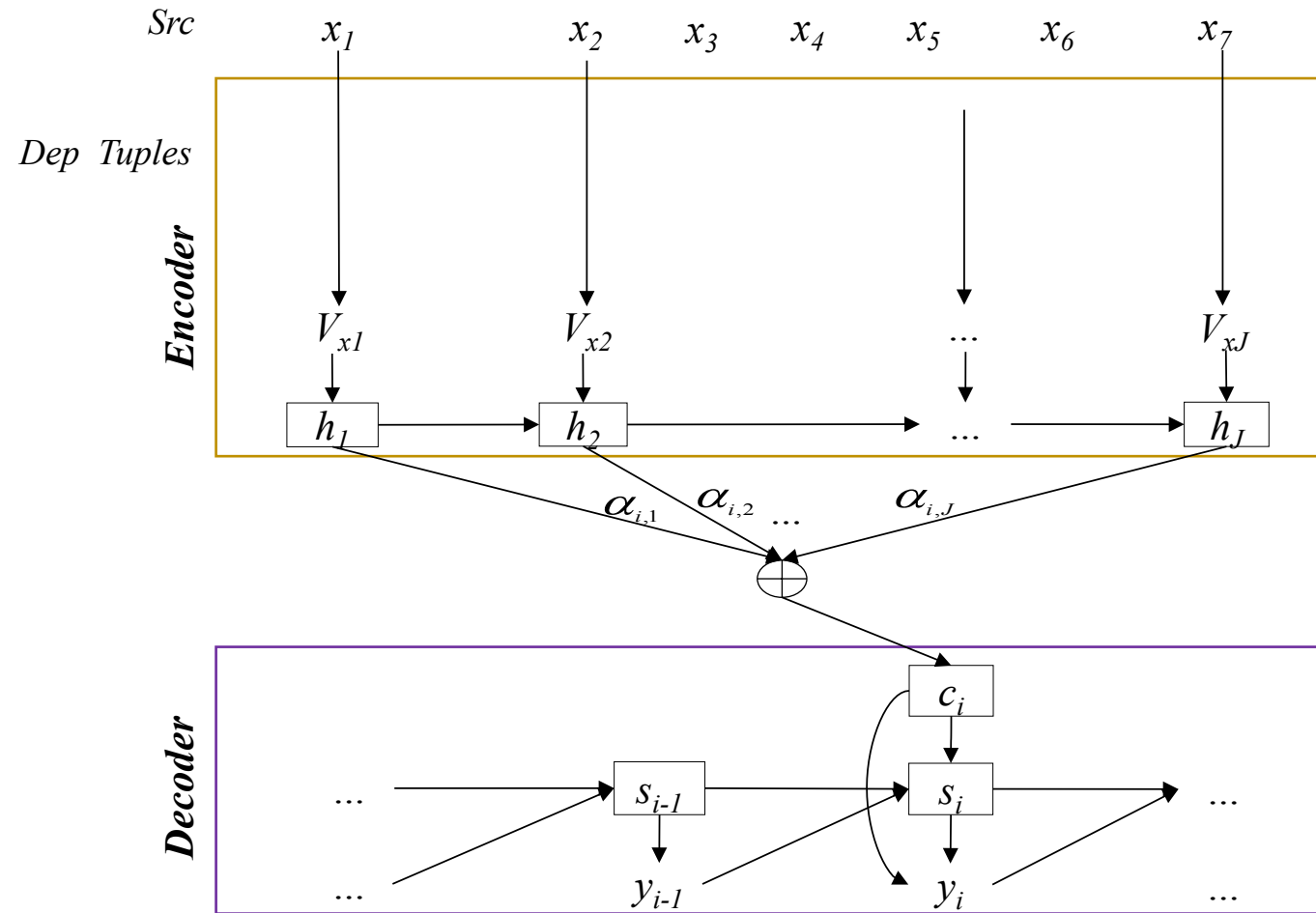
$CH_{x_2} = \langle \varepsilon \rangle$ ,





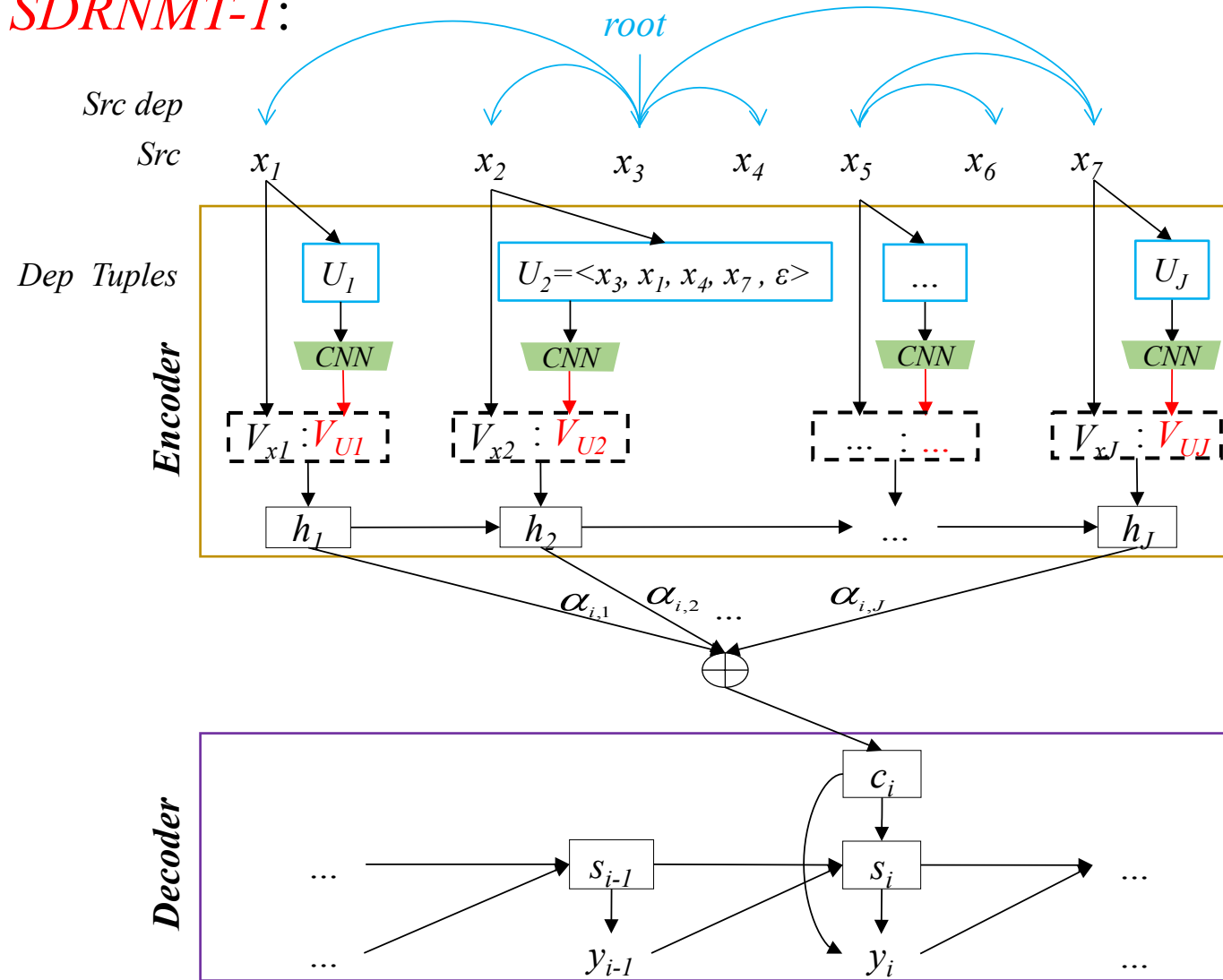
# Neural Machine Translation with SDR

*SDRNMT-1:*



# Neural Machine Translation with SDR

*SDRNMT-1:*

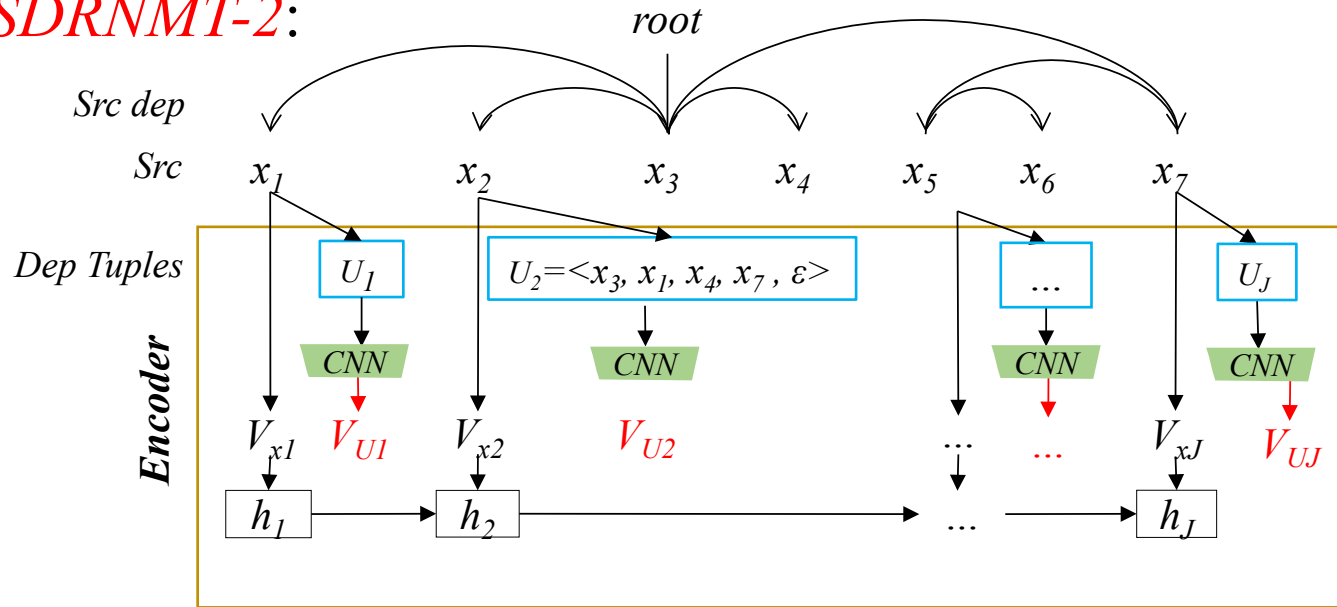


$$h_j = f_{enc}(V_{x_j} : V_{U_j}, h_{j-1})$$

Where the  $V_{x_j}$  is 360-dim and the learned  $V_{U_j}$  is 260-dim.

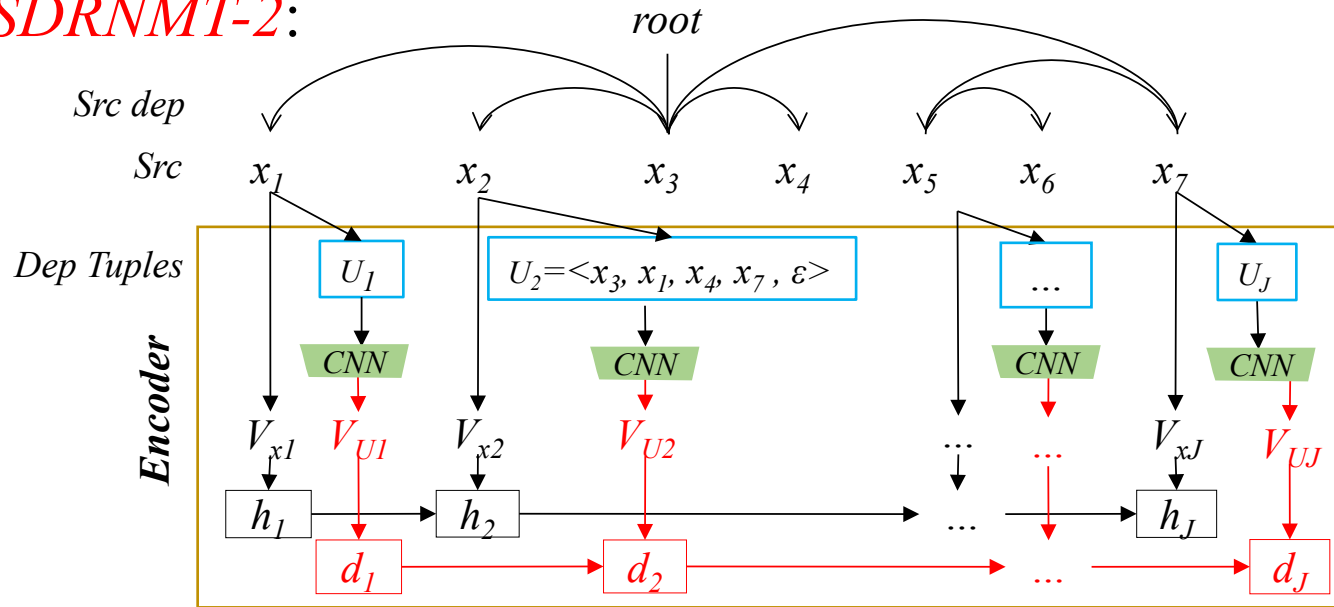
# Neural Machine Translation with SDR

*SDRNMT-2:*



# Neural Machine Translation with SDR

*SDRNMT-2:*



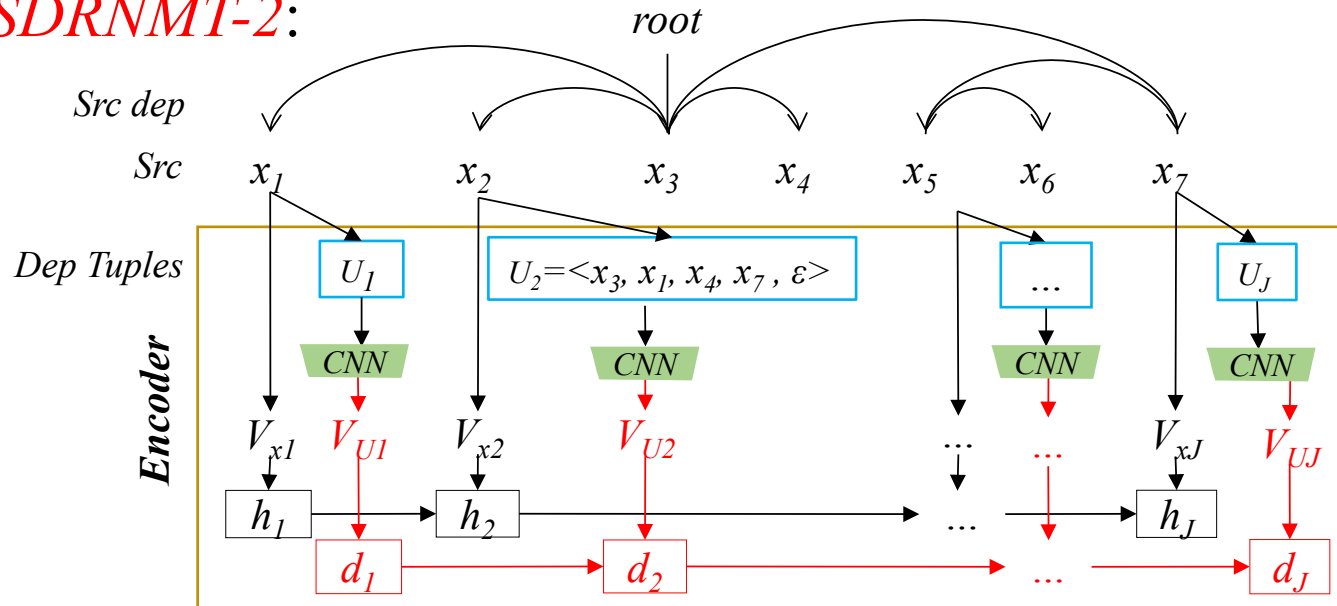
**Encoder:**

$$h_j = f_{enc}(V_{x_j}, h_{j-1}),$$

$$d_j = f_{enc}(V_{U_j}, d_{j-1})$$

# Neural Machine Translation with SDR

*SDRNMT-2:*



*Attention  $\tilde{\alpha}$*

**Encoder:**  $h_j = f_{enc}(V_{x_j}, h_{j-1}),$

$d_j = f_{enc}(V_{U_j}, d_{j-1})$

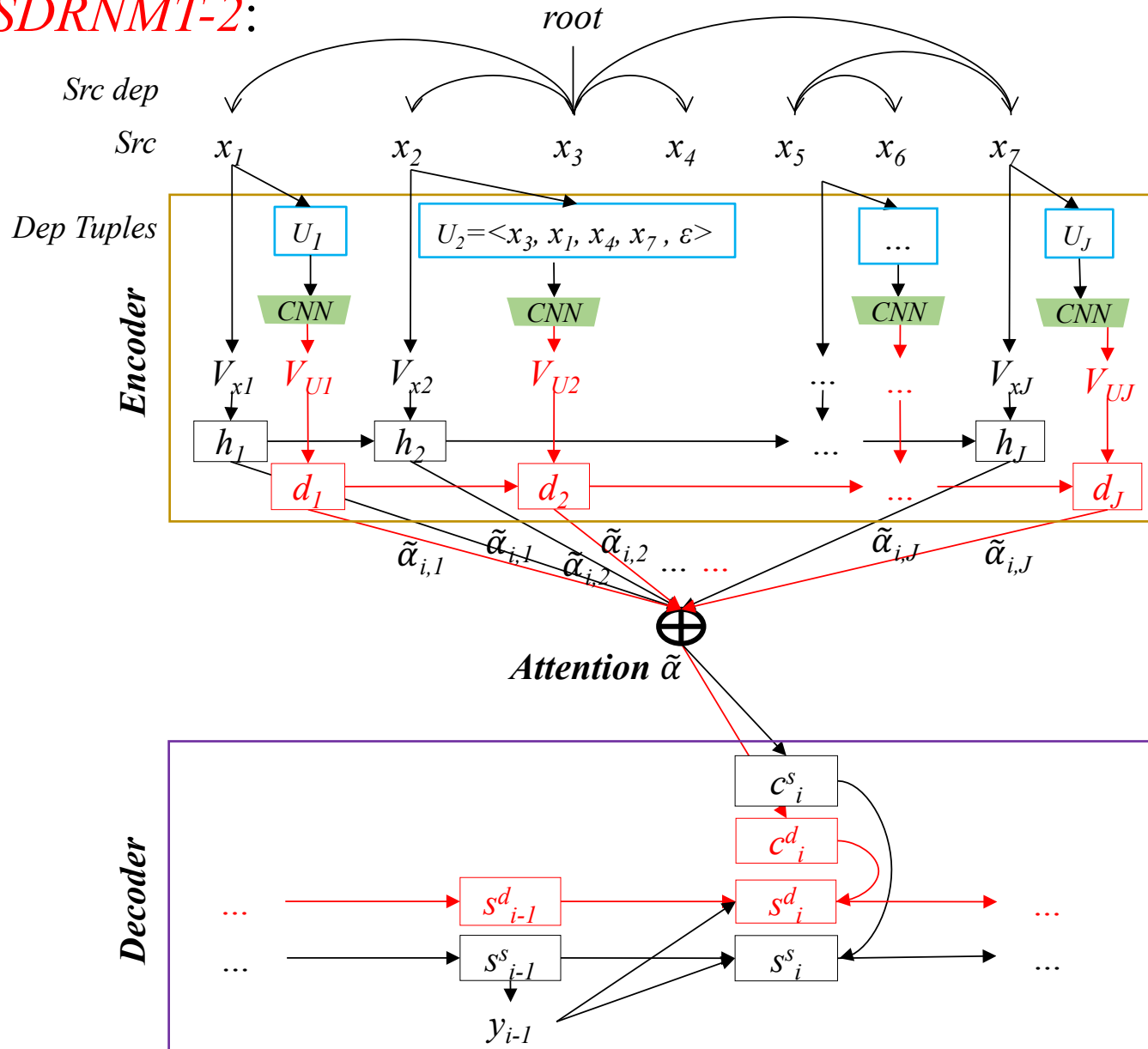
**Attention:**  $e_{i,j}^s = f(s_{i-1}^s + h_j),$

$e_{i,j}^d = f(s_{i-1}^d + d_j).$

$$\alpha_{i,j} = \frac{\exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}{\sum_{j=1}^J \exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}$$

# Neural Machine Translation with SDR

## SDRNMT-2:



**Encoder:**  $h_j = f_{enc}(V_{x_j}, h_{j-1}),$

$d_j = f_{enc}(V_{U_j}, d_{j-1})$

**Attention:**  $e_{i,j}^s = f(s_{i-1}^s + h_j),$

$e_{i,j}^d = f(s_{i-1}^d + d_j).$

$$\alpha_{i,j} = \frac{\exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}{\sum_{j=1}^J \exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}$$

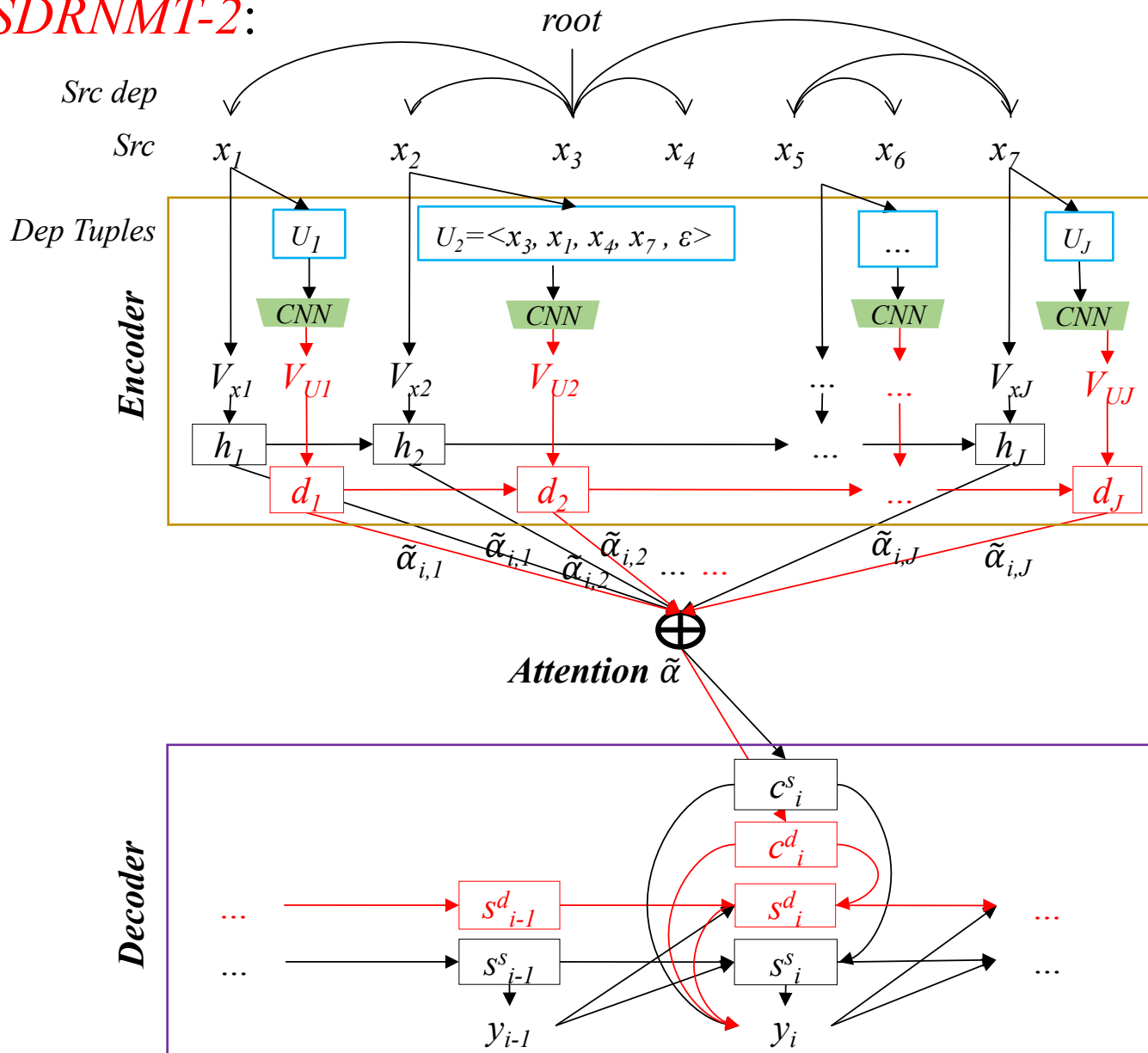
**Decoder:**  $c_{i,j}^s = \sum_{j=1}^J \alpha_{i,j} h_j, c_{i,j}^d = \sum_{j=1}^J \alpha_{i,j} d_j$

$s_i^s = \varphi(s_{i-1}^s, y_{i-1}, c_i^s),$

$s_i^d = \varphi(s_{i-1}^d, y_{i-1}, c_i^d).$

# Neural Machine Translation with SDR

*SDRNMT-2:*



**Encoder:**  $h_j = f_{enc}(V_{x_j}, h_{j-1}),$

$d_j = f_{enc}(V_{U_j}, d_{j-1})$

**Attention:**  $e_{i,j}^s = f(s_{i-1}^s + h_j),$

$e_{i,j}^d = f(s_{i-1}^d + d_j).$

$$\alpha_{i,j} = \frac{\exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}{\sum_{j=1}^J \exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}$$

**Decoder:**  $c_{i,j}^s = \sum_{j=1}^J \alpha_{i,j} h_j, c_{i,j}^d = \sum_{j=1}^J \alpha_{i,j} d_j$

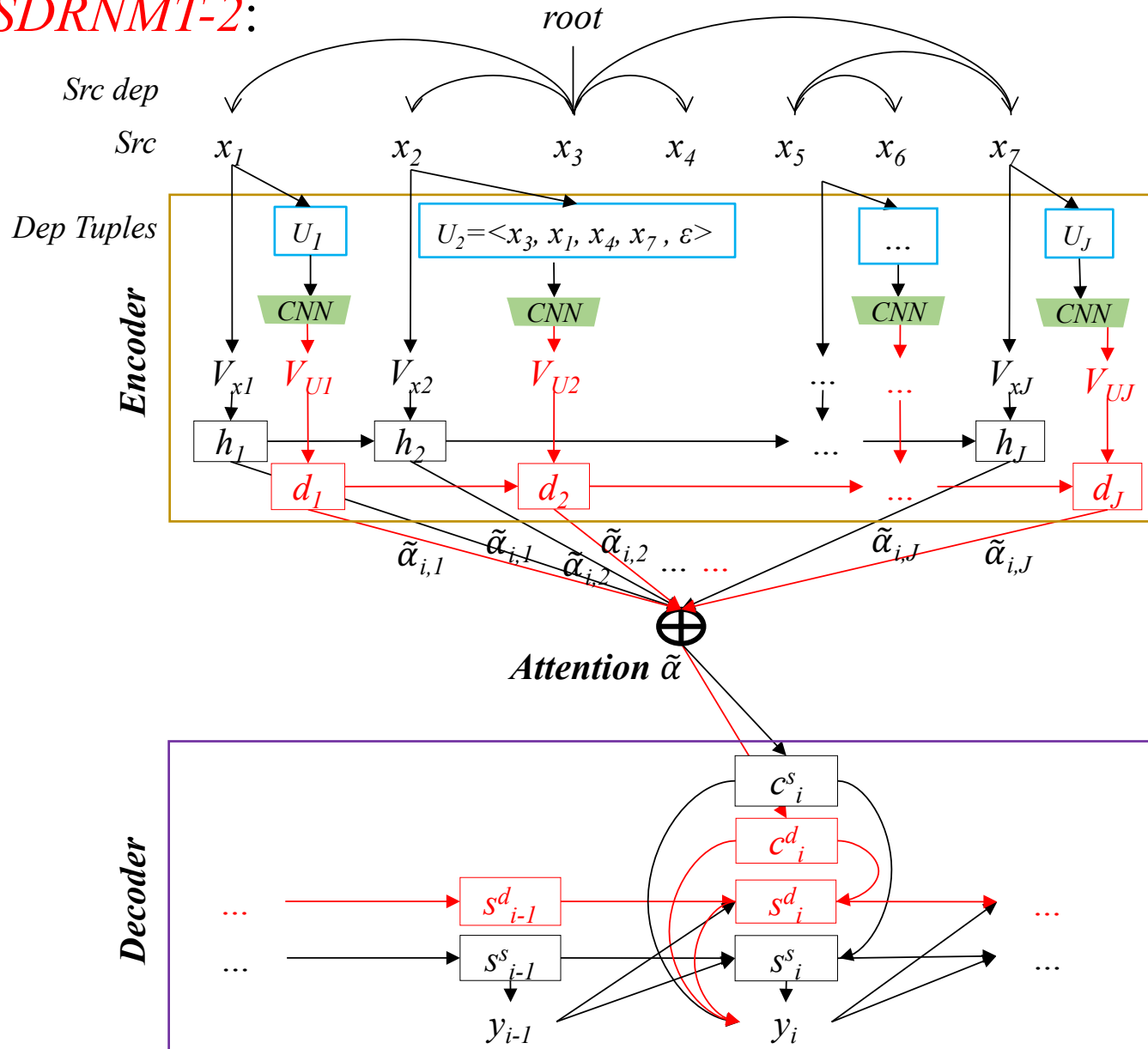
$s_i^s = \varphi(s_{i-1}^s, y_{i-1}, c_i^s),$

$s_i^d = \varphi(s_{i-1}^d, y_{i-1}, c_i^d).$

$p(y_i | y_{i-1}, x, T) = g(y_{i-1}, s_i^s, s_i^d, c_i^s, c_i^d)$

# Neural Machine Translation with SDR

**SDRNMT-2:**



**Encoder:**  $h_j = f_{enc}(V_{x_j}, h_{j-1}),$

$d_j = f_{enc}(V_{U_j}, d_{j-1})$

**Attention:**  $e_{i,j}^s = f(s_{i-1}^s + h_j),$

$e_{i,j}^d = f(s_{i-1}^d + d_j).$

$$\alpha_{i,j} = \frac{\exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}{\sum_{j=1}^J \exp(\lambda e_{i,j}^s + (1-\lambda)e_{i,j}^d)}$$

**Decoder:**  $c_{i,j}^s = \sum_{j=1}^J \alpha_{i,j} h_j, c_{i,j}^d = \sum_{j=1}^J \alpha_{i,j} d_j$

$s_i^s = \varphi(s_{i-1}^s, y_{i-1}, c_i^s),$

$s_i^d = \varphi(s_{i-1}^d, y_{i-1}, c_i^d).$

$p(y_i | y_{i-1}, x, T) = g(y_{i-1}, s_i^s, s_i^d, c_i^s, c_i^d)$

**Double Context NMT**



# Experimental

- Experiments on Chinese-to-English translation task, 1.42M *LDC corpus*
- Parse source sentences of training data by Stanford Parser ([Chang et al., 2009](#))
- For the *SDRNMT-1* and *SDRNMT-2*, the dimension of  $V_{x_j}$  is 360 and the dimension of  $V_{U_j}$  is 260, and input embedding of the baseline is 620
- The baselines include Phrase-Based Statistical Machine Translation (PBSMT) ([Koehn et al., 2007](#)), standard Attentional NMT (AttNMT) ([Bahdanau et al., 2014](#)), NMT with dependency labels ([Sennrich and Haddow, 2016](#))

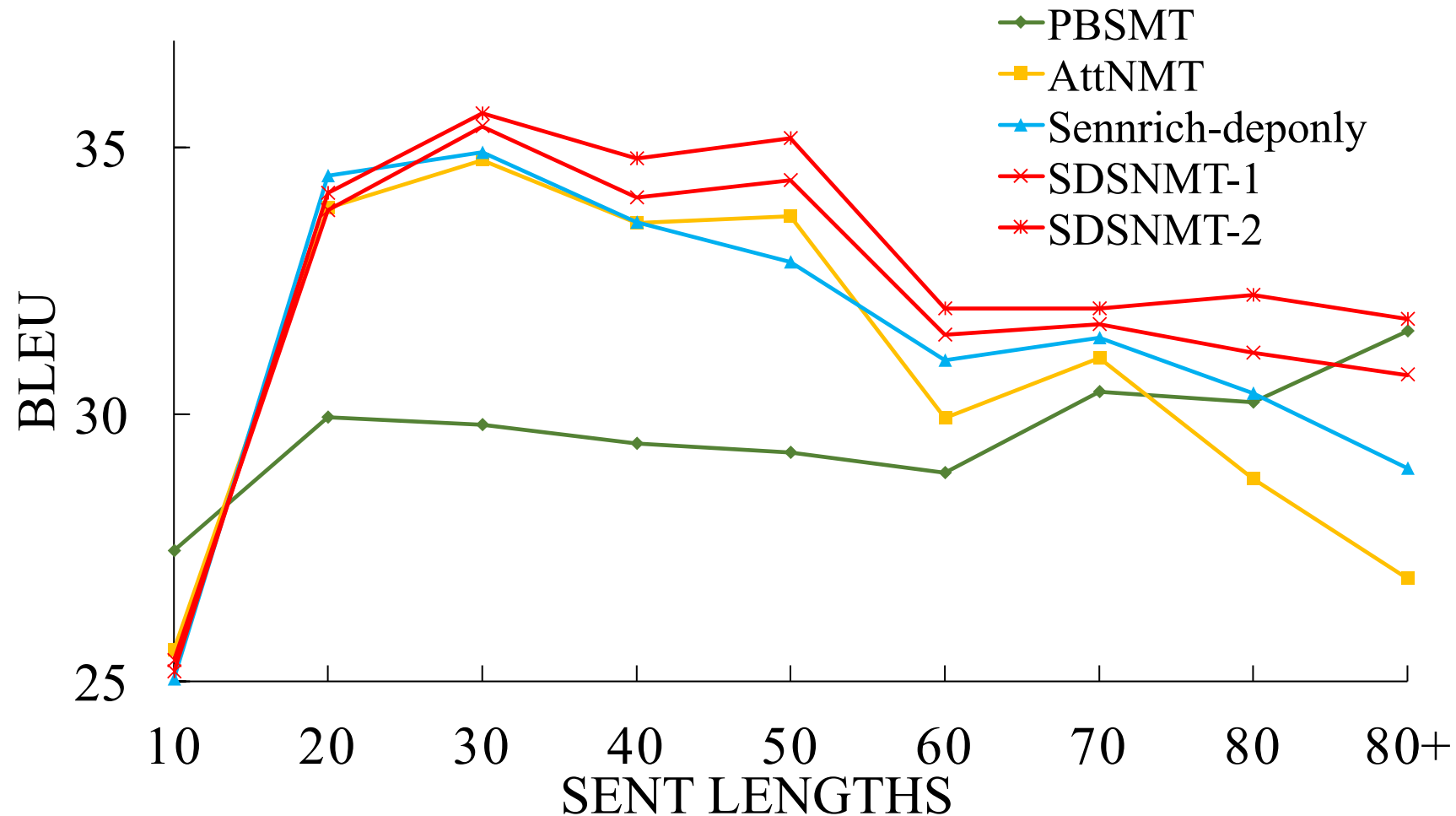
# Experimental

| System           | Dev(NIST02)  | NIST03         | NIST04        | NIST05        | NIST06         | NIST08        | AVG          |
|------------------|--------------|----------------|---------------|---------------|----------------|---------------|--------------|
| PBSMT            | 33.15        | 31.02          | 33.78         | 30.33         | 29.62          | 23.53         | 29.66        |
| AttNMT           | 36.31        | 34.02          | 37.11         | 32.86         | 32.54          | 25.44         | 32.40        |
| Sennrich-deponly | 36.68        | 34.51          | 38.09         | 33.37         | 32.96          | 26.96         | 32.98        |
| SDRNMT-1         | 36.88        | 34.98*         | 38.14         | 34.61**       | 33.58*         | 27.06         | 33.32        |
| SDRNMT-2         | <b>37.34</b> | <b>35.91**</b> | <b>38.73*</b> | <b>34.18*</b> | <b>33.76**</b> | <b>27.64*</b> | <b>34.04</b> |

“\*” indicates statistically significant better than “Sennrich-deponly” at  $p$ -value  $< 0.05$  and “\*\*” at  $p$ -value  $< 0.01$  by bootsrap resampling (Koehn, 2004)

# Experimental Results

- Translation qualities for different sentence lengths



# Conclusion

- Source dependency unit to capture source long-distance dependency constraint
- The proposed *SDRNMT-1* and *SDRNMT-2* consist of NMT and CNN, which are jointly trained to learn SDR and translation instead of separately trained
- Double-Context approach to further utilize source dependency representation