

# Machine Translation with Weakly Paired Documents



<sup>1</sup>Lijun Wu, <sup>2</sup>Jinhua Zhu, <sup>3</sup>Fei Gao, <sup>4</sup>Di He,  
<sup>5</sup>Tao Qin, <sup>1</sup>Jianhuang Lai and <sup>5</sup>Tie-Yan Liu  
<sup>1</sup>Sun Yat-sen University; <sup>2</sup>University of Science and Technology of China;  
<sup>3</sup>Institute of Computing Technology, Chinese Academy of Sciences;  
<sup>4</sup>Peking University, <sup>5</sup>Microsoft Research Asia



## 1. Introduction

- NMT achieves strong performance in rich-resource language pairs with large amount of parallel data.
- Low-resource** language pairs have much lower translation accuracy due to the lack of bilingual sentence pairs.
- Unsupervised** machine translation has been explored with monolingual data only.
- In reality, large amount of **weakly paired bilingual documents** can be leveraged.
- We propose to boost the unsupervised machine translation with weakly paired documents using two innovated components.*
- We achieve strong performances in various language pairs and reduce the gap between supervised and unsupervised translation up to 50%.*

## 2. Approach

- We propose to leverage weakly paired bilingual documents from **Wikipedia**.
- Notations:
  - $D = \{(d_i^X, d_i^Y)\}$  as the set of weakly paired documents (e.g., two cross-lingual linked Wikipedia pages)
  - $n_i^X, n_i^Y$  are the number of sentences in paired documents  $d_i^X, d_i^Y$ , usually  $n_i^X \neq n_i^Y$
  - $x, y$  are the sentences of language  $X, Y$

### ◆ Mining implicitly aligned sentence pairs

- $e_w$ , cross-lingual word embedding from MUSE
- $p_w$ , the estimated frequency from the document
- $a$ , predefined parameter and  $\hat{e}_s$  is the weighted sentence embedding
- $u_1$ , the first principal component from all sentence embedding

$$\hat{e}_s = \sum_{w \in s} \frac{a}{a + p(w)} e_w,$$

$$u_1 \leftarrow PCA(E),$$

- $e_s = \hat{e}_s - u_1 u_1^T \hat{e}_s$
- Select sentence pairs by  $\text{sim}(s^X, s^Y) = \frac{\langle e_{s^X}, e_{s^Y} \rangle}{\|e_{s^X}\| \|e_{s^Y}\|}$  larger than  $c_1$ , also ensure this pair is larger than others pairs by  $c_2$
- The implicitly aligned sentence training loss of two-sides is

$$L_p(S; \theta) = -\frac{1}{|S|} \sum_{(s^X, s^Y) \in S} \log P_{X \rightarrow Y}(s^Y | s^X; \theta) - \frac{1}{|S|} \sum_{(s^X, s^Y) \in S} \log P_{Y \rightarrow X}(s^X | s^Y; \theta).$$

### ◆ Aligning Topic Distribution

- Translate  $d_i^X$  to  $\hat{d}_i^Y$
- Evaluate the word distribution between  $d_i^Y$  and  $\hat{d}_i^Y$
- Feed pair  $(s_{i,k}^X, \hat{s}_{i,k}^Y)$  into NMT model and calculate  $P(w^Y; d_i^X)$  by

$$P(w_{i,k,t}^Y | s_{i,k}^X, \hat{s}_{i,k}^Y) \sim P_{X \rightarrow Y}(w_{i,k,t}^Y | s_{i,k}^X, \hat{s}_{i,k}^Y; \theta),$$

$$P(w^Y; d_i^X, \theta) \propto \prod_{i,k,t} P(w_{i,k,t}^Y | s_{i,k}^X, \hat{s}_{i,k}^Y; \theta),$$

- The ground-truth document word distribution is

$$P(w^Y; d_i^Y) = \frac{\#w \text{ in } d_i^Y}{\#token \text{ in } d_i^Y}.$$

- The document alignment loss of  $X \rightarrow Y$  is

$$L_d(D; \theta, X \rightarrow Y) = \frac{1}{|D|} \sum_{(d_i^X, d_i^Y) \in D} KL(P(w^Y; d_i^X) || P(w^Y; d_i^Y, \theta)).$$

- The detailed two-sides loss

$$L_d(D; \theta) = L_d(D; \theta, X \rightarrow Y) + L_d(D; \theta, Y \rightarrow X).$$

## 3. Algorithm

- The overall loss function is

$$L = L_m(M; \theta) + \alpha L_p(S; \theta) + \beta L_d(D; \theta).$$

### Algorithm 1 Training Algorithm

**Require:** Initial translation model with parameter  $\theta$ ; monolingual dataset  $M$ , implicitly aligned sentence pairs dataset  $S$ , weakly paired documents dataset  $D$ ; optimizer  $Opt$

- while** not converged **do**
- Randomly sample a mini-batch monolingual sentences from  $M$ , implicitly aligned sentence pairs from  $S$  and weakly paired documents from  $D$
- Calculate loss  $L_m, L_p$  and  $L_d$
- Update  $\theta$  by minimizing the overall objective  $L$  using optimizer  $Opt$
- end while**

- $L_m$  is the original unsupervised NMT training loss

## 4. Experiments

### • Data Statistics

Language	#Wiki Documents
English	5,684,240
German	2,201,782
Spanish	1,389,469
Romanian	387,627

Task	#Document Pairs
English-German	948,631
English-Spanish	836,564
English-Romanian	87,289

### • Overall Results

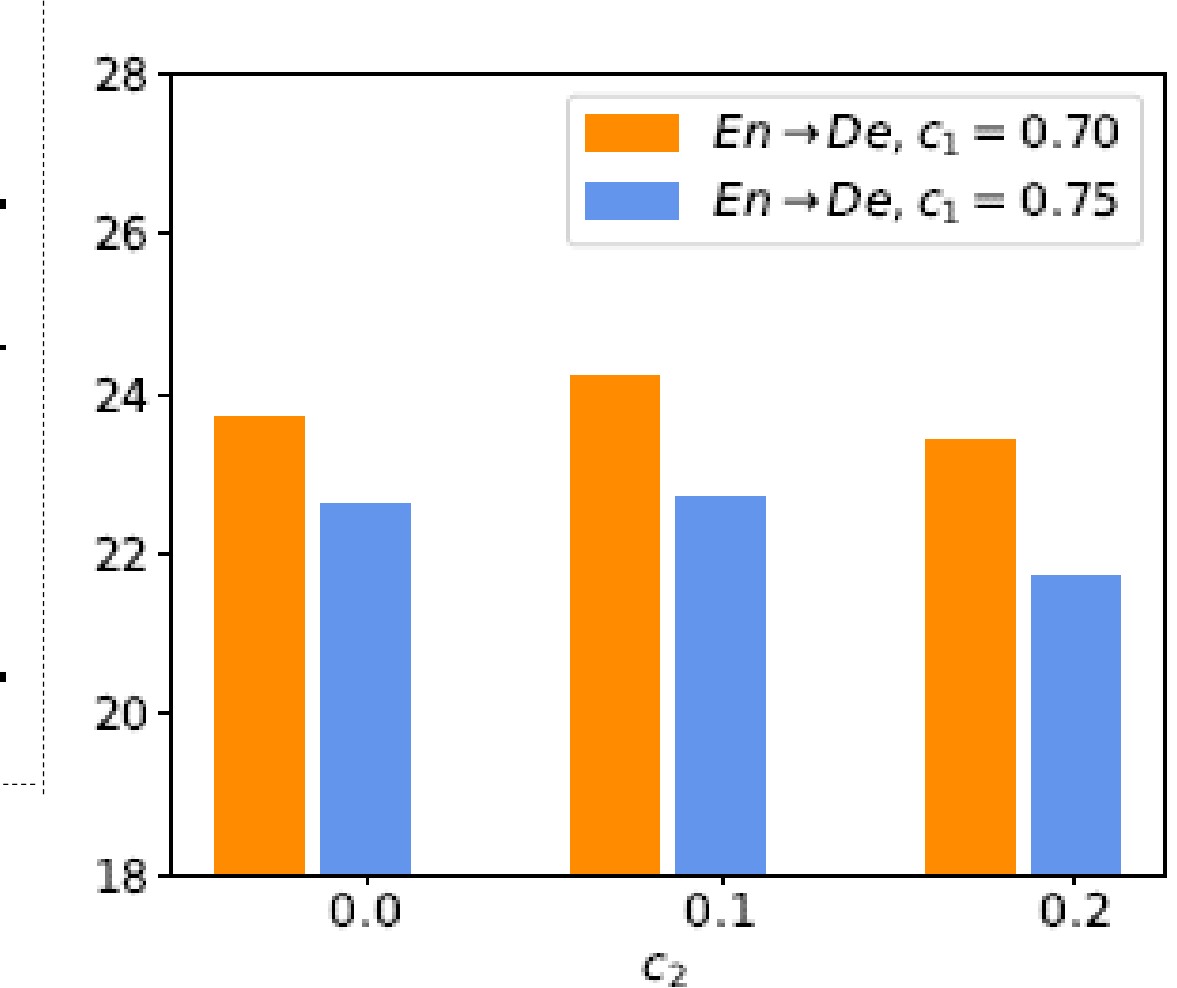
Unsupervised Method	En→De	De→En	En→Es	Es→En	En→Ro	Ro→En
Lample et al. (2017)	9.6	13.3	-	-	-	-
Yang et al. (2018)	10.9	14.6	-	-	-	-
NMT (Lample et al., 2018)	17.2	21.0	19.7	20.0	21.2	19.5
PBSMT (Lample et al., 2018)	17.9	22.9	-	-	22.0	23.7
PBSMT + NMT (Lample et al., 2018)	20.2	25.2	-	-	25.1	23.9
NMT + First Wiki Sentence	16.3	19.3	17.3	18.3	19.4	18.1
NMT + Document Translation	12.0	14.9	14.5	15.3	16.8	15.7
<b>Ours</b>	<b>24.2</b>	<b>30.3</b>	<b>28.1</b>	<b>27.6</b>	<b>30.1</b>	<b>27.6</b>
Supervised NMT	33.6	38.2	33.2	32.9	32.8	35.4

## 5. Studies

### • Analysis

#### Ablation Study

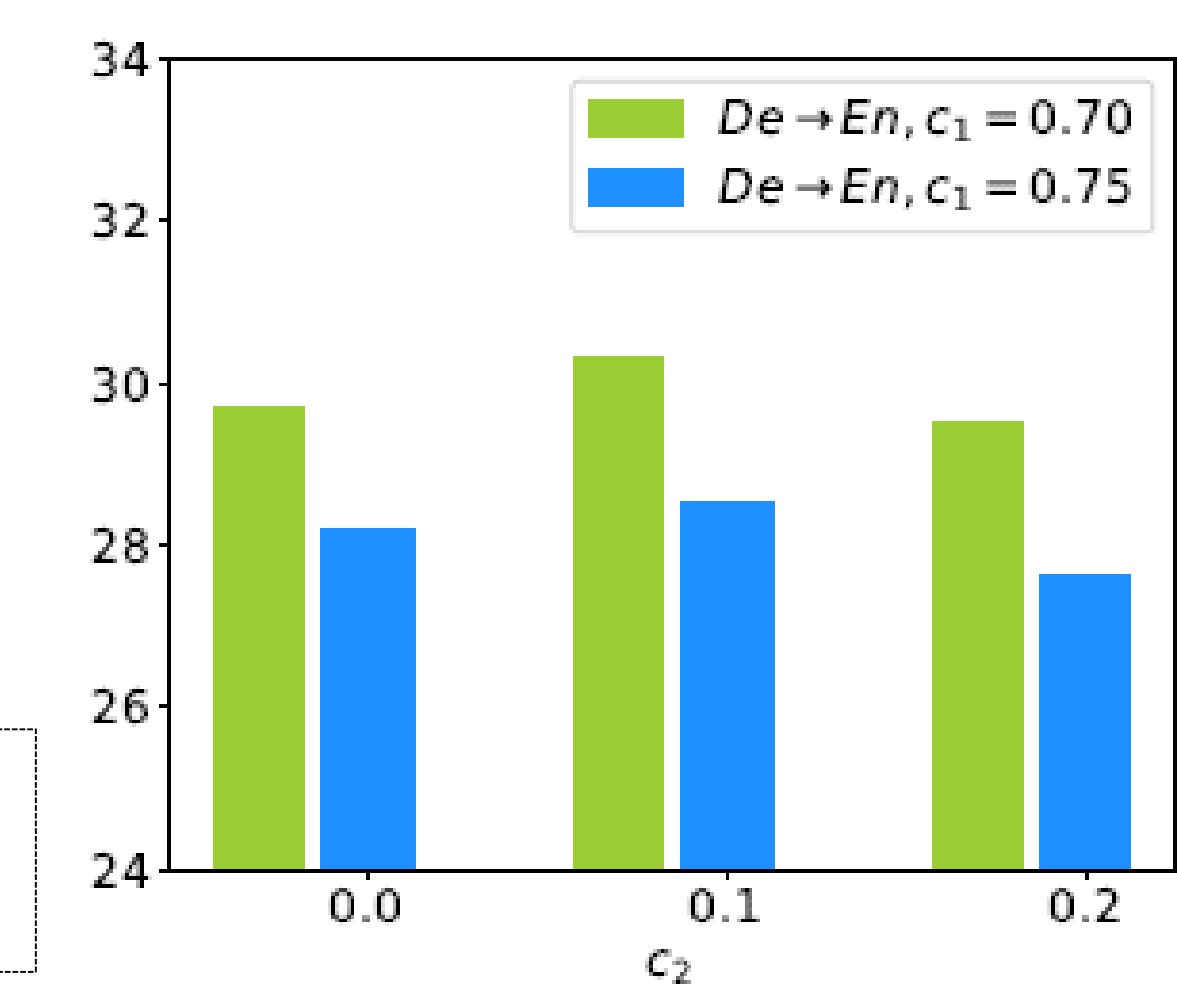
Our Method	En→De	De→En
with $L_p$ and $L_d$	24.2	30.3
without $L_d$	22.9	28.7
without $L_p$	18.5	23.3



(a) En→De BLEU performance.

#### Impact of Sentence Quality

English-German			
$c_1/c_2$	0.0	0.1	0.2
0.70	257,947	199,965	132,403
0.75	100,497	84,271	58,814



(b) De→En BLEU performance.

### Contact

- wulijun3@mail2.sysu.edu.cn