

Semi-supervised Learning for Neural Machine Translation

Yong Cheng



清华大学
Tsinghua University



交叉信息研究院
Institute for Interdisciplinary
Information Sciences


joint work with Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, Yang Liu

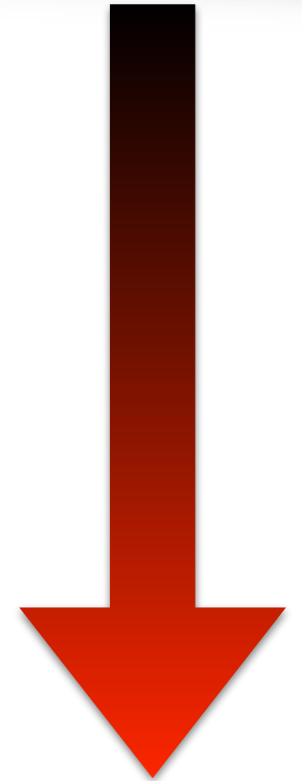
Machine Translation



Automated translation using computer software

Machine Translation

- * Rule-based Machine Translation 1970s
- * Example-based Machine Translation 1984
- * Statical Machine Translation (SMT) 1993
-  **Neural Machine Translation (NMT) 2014**



Trends: learning to translate from **DATA**

Machine Translation

Parallel corpora are usually limited in

quantity

&

quality

&

coverage

Monolingual Corpora



Parallel Corpora



Monolingual Corpora Used in SMT and NMT

- * N-gram language model in SMT
Koehn et al., [2007]
- * Monolingual corpora as decipherment
Ravi and Knight [2011]
- * Integrate a neural language model into NMT.
Gulccehre et al. [2015]
- * Additional pseudo parallel corpus.
Sennrich et al. [2016]

Supervised Training

Parallel Corpus

$$\mathcal{D} = \{ \langle \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \rangle \}_{n=1}^N$$

Objective

$$L(\boldsymbol{\theta}) = \sum_{n=1}^N \log P(\mathbf{y}^{(n)} | \mathbf{x}^{(n)}; \boldsymbol{\theta})$$

Unsupervised Training

Monolingual Corpus $\mathcal{T} = \{\mathbf{y}^{(t)}\}_{t=1}^T$



Our Approach — Autoencoders

bushi yu shalong juxing le huitan

X

Our Approach — Autoencoders



$$P(\mathbf{y} | \mathbf{x}; \vec{\theta})$$

bushi yu shalong juxing le huitan

X

Our Approach — Autoencoders

latent

Bush held a talk with sharon

y

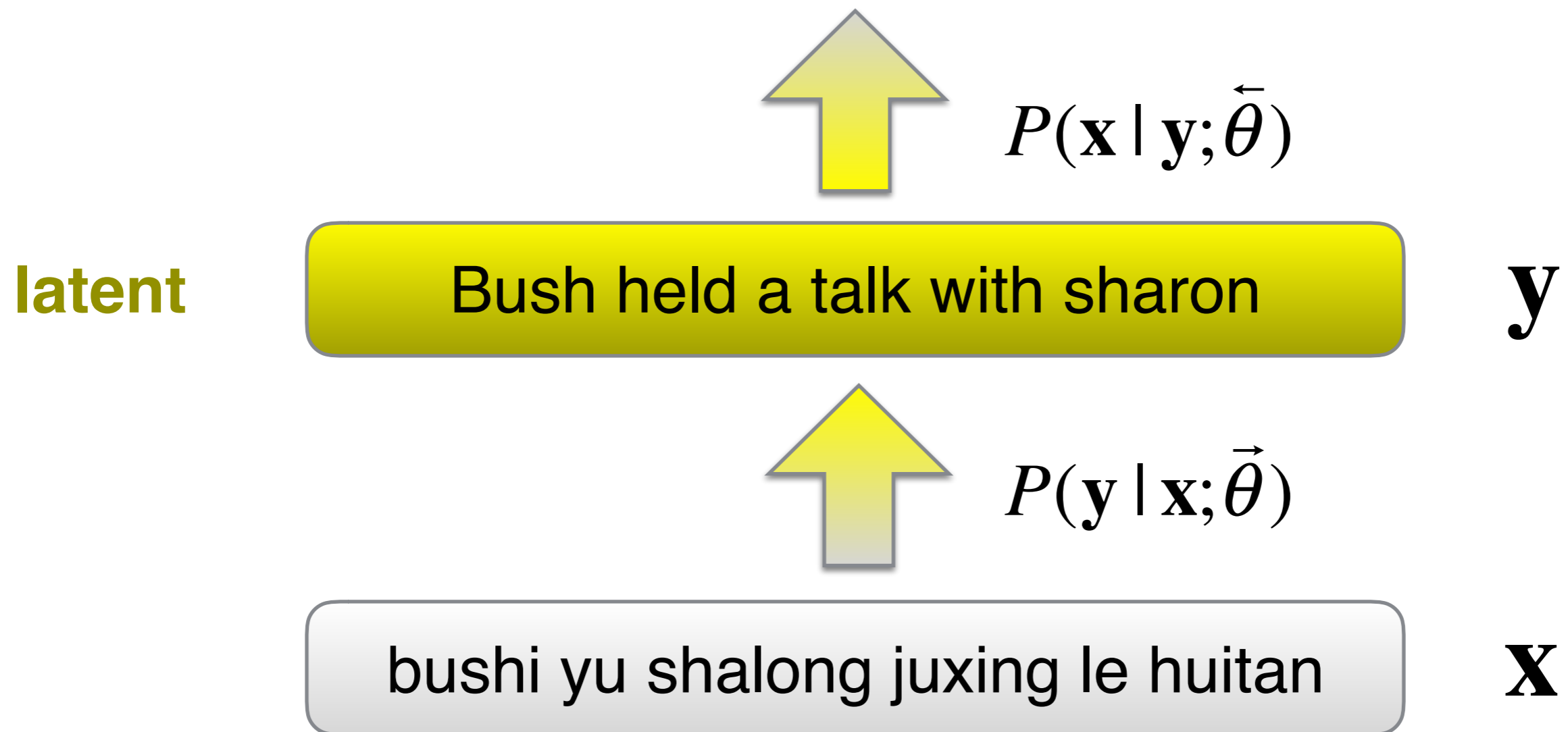


$P(\mathbf{y} | \mathbf{x}; \vec{\theta})$

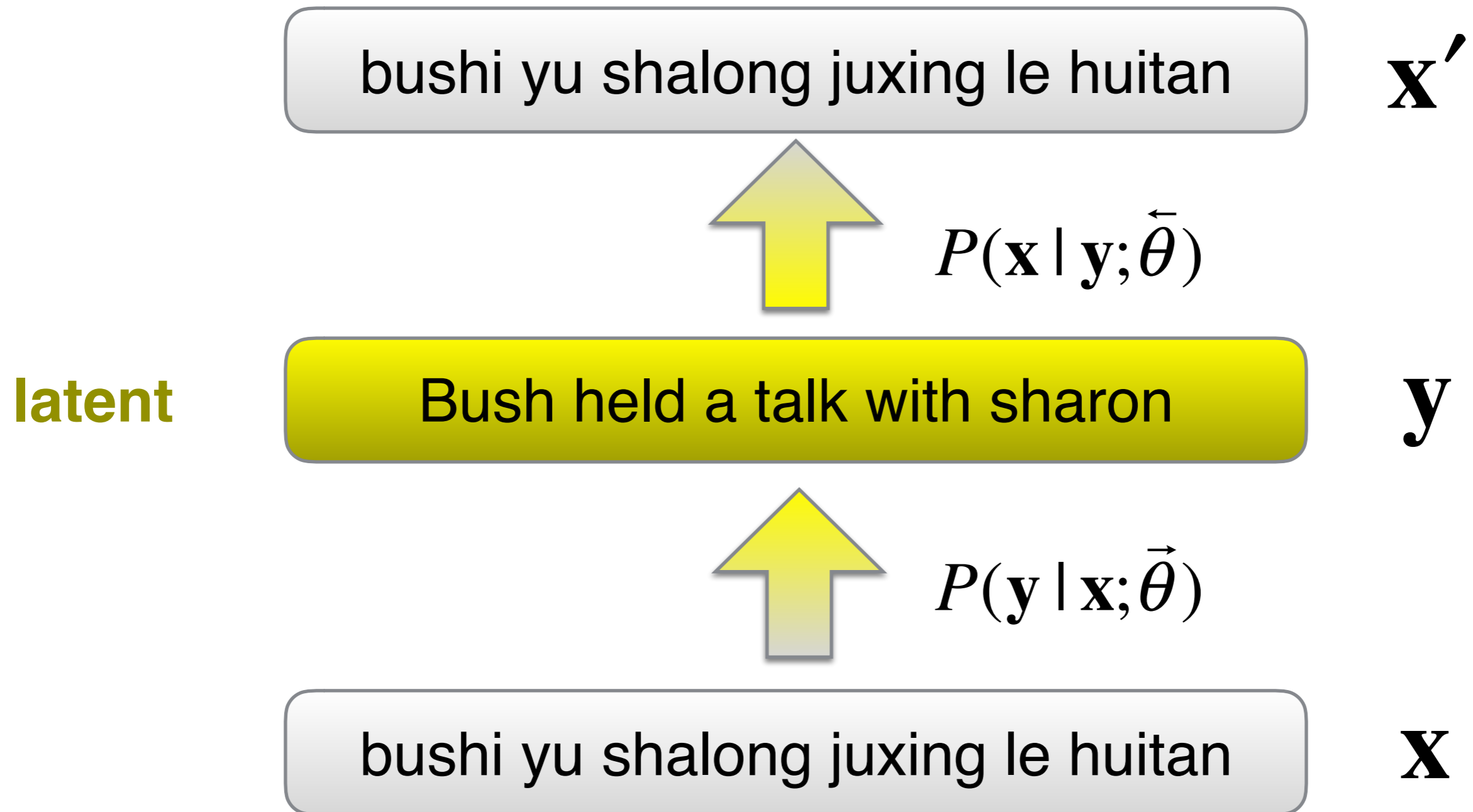
bushi yu shalong juxing le huitan

X

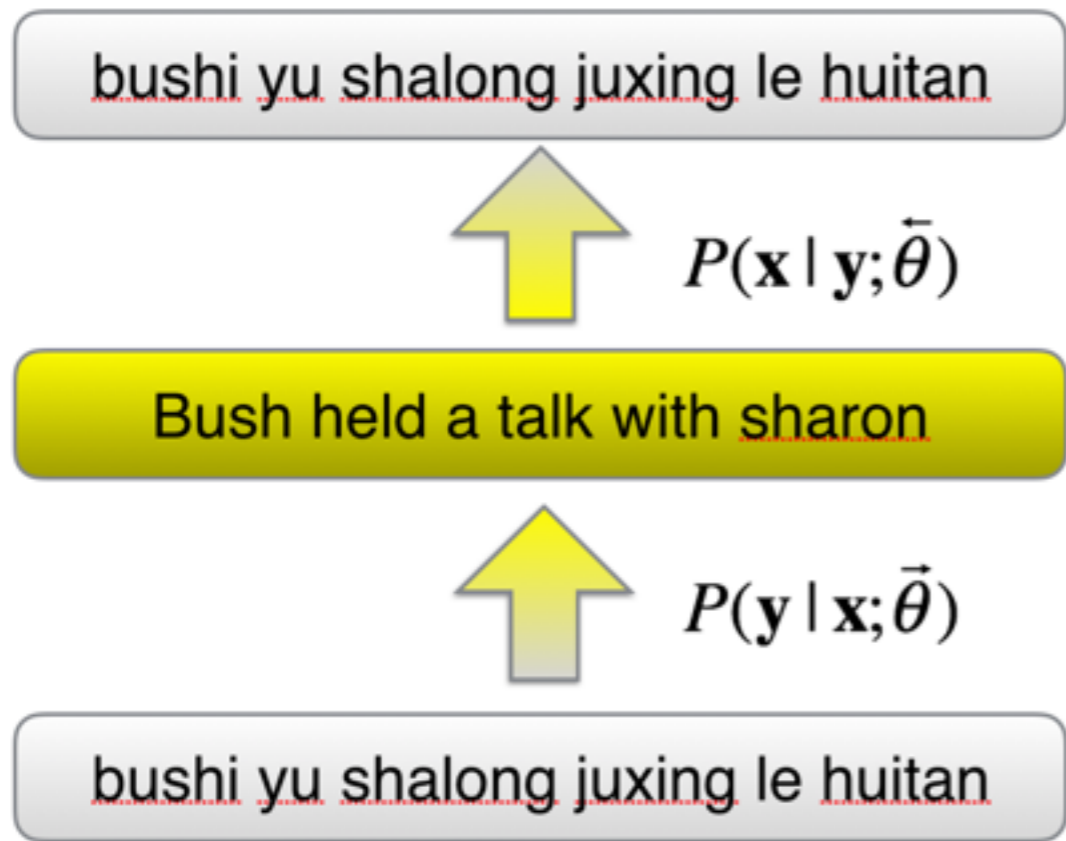
Our Approach — Autoencoders



Our Approach — Autoencoders

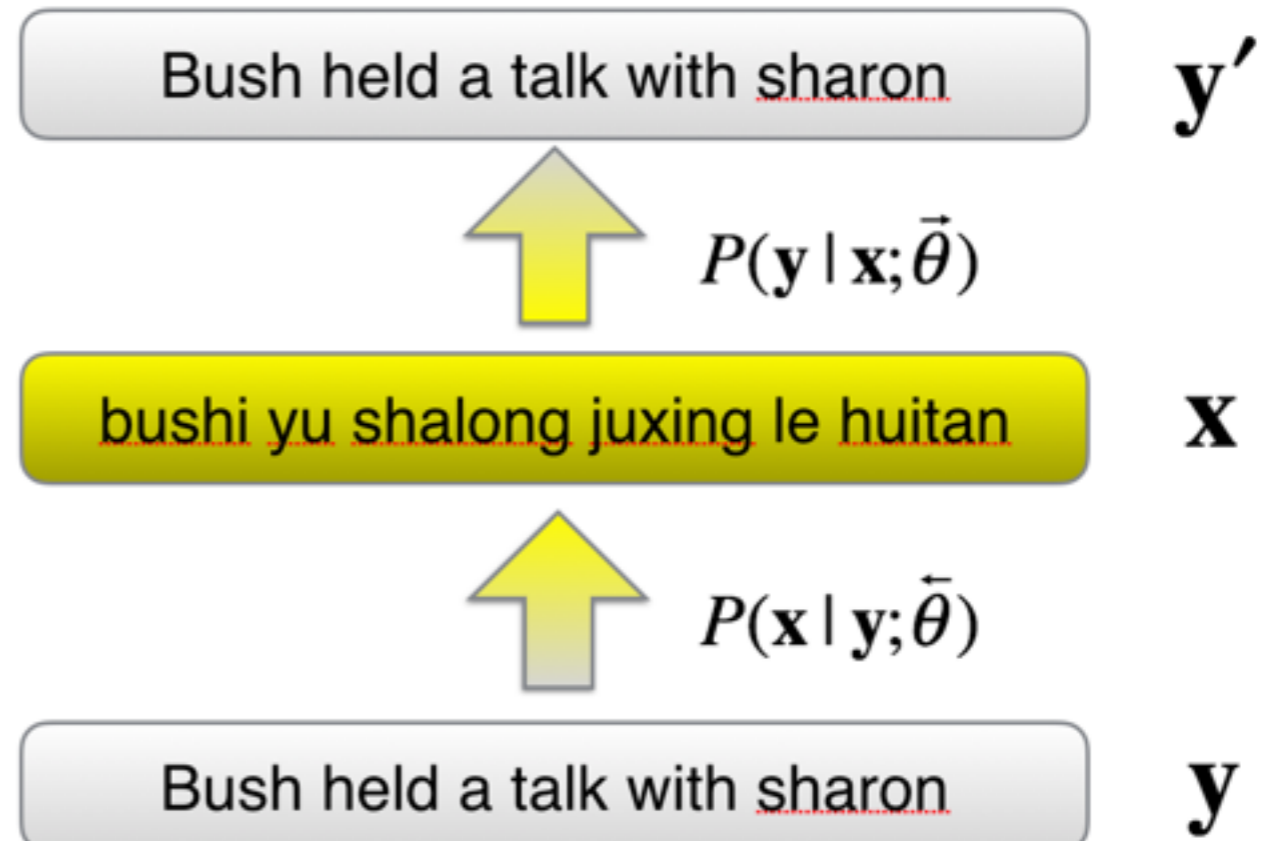


Our Approach — Autoencoders



source autoencoder

\mathbf{x}'



\mathbf{y}'

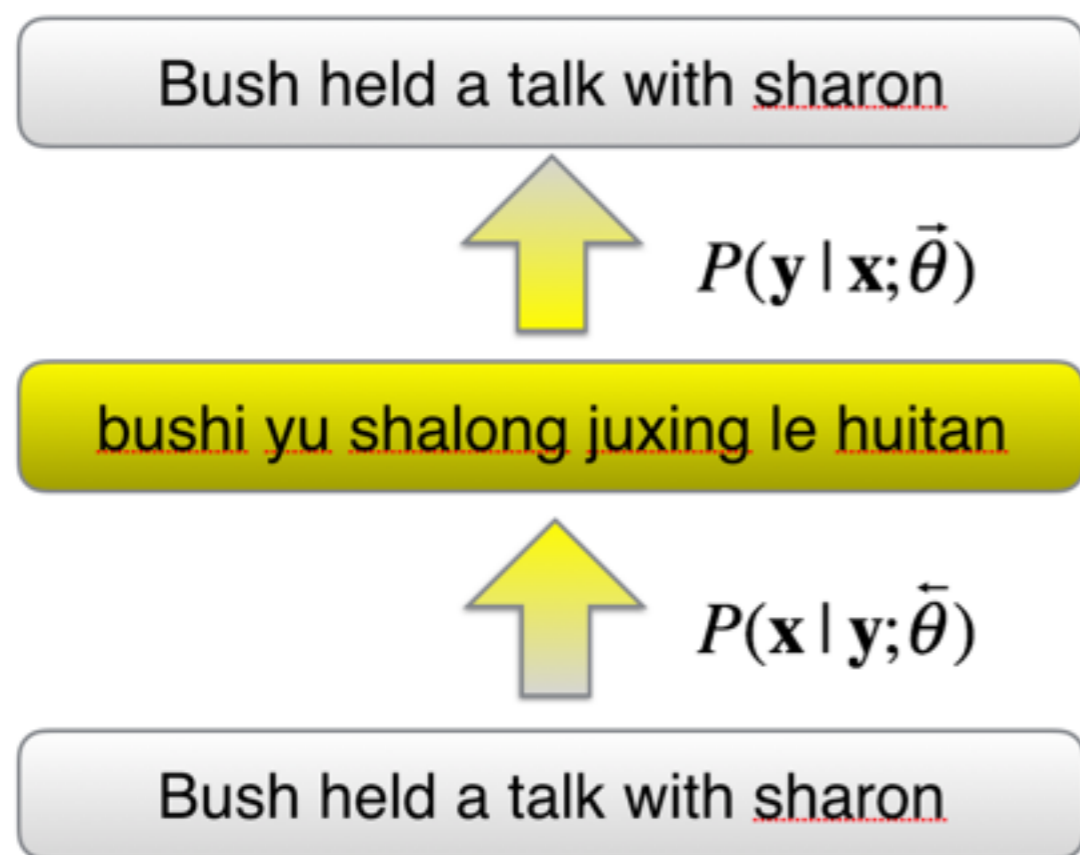
\mathbf{x}

\mathbf{y}

target autoencoder

Unsupervised Training (Autoencoders)

Monolingual Corpus $\mathcal{T} = \{\mathbf{y}^{(t)}\}_{t=1}^T$



\mathbf{y}'

\mathbf{x}

\mathbf{y}

$$\begin{aligned}
 & P(\mathbf{y}' | \mathbf{y}; \vec{\theta}, \overleftarrow{\theta}) \\
 &= \sum_{\mathbf{x}} P(\mathbf{y}', \mathbf{x} | \mathbf{y}; \vec{\theta}, \overleftarrow{\theta}) \\
 &= \sum_{\mathbf{x}} \underbrace{P(\mathbf{x} | \mathbf{y}; \overleftarrow{\theta})}_{\text{encoder}} \underbrace{P(\mathbf{y}' | \mathbf{x}; \vec{\theta})}_{\text{decoder}}
 \end{aligned}$$

target autoencoder

Semi-supervised Training

Training Objective

$$\begin{aligned} & J(\vec{\theta}, \overleftarrow{\theta}) \\ = & \underbrace{\sum_{n=1}^N \log P(\mathbf{y}^{(n)} | \mathbf{x}^{(n)}; \vec{\theta})}_{\text{source-to-target likelihood}} + \underbrace{\sum_{n=1}^N \log P(\mathbf{x}^{(n)} | \mathbf{y}^{(n)}; \overleftarrow{\theta})}_{\text{target-to-source likelihood}} \\ & + \lambda_1 \underbrace{\sum_{t=1}^T \log P(\mathbf{y}' | \mathbf{y}^{(t)}; \vec{\theta}, \overleftarrow{\theta})}_{\text{target autoencoder}} + \lambda_2 \underbrace{\sum_{s=1}^S \log P(\mathbf{x}' | \mathbf{x}^{(s)}; \vec{\theta}, \overleftarrow{\theta})}_{\text{source autoencoder}}, \end{aligned}$$

Translation Results

Compared with Moses (SMT) and RNNSearch (NMT)

System	Training Data			Direction	NIST06	NIST02
	CE	C	E			
MOSES	✓	×	×	C → E	32.48	32.69
				E → C	14.27	18.28
RNNSEARCH	✓	×	×	C → E	30.74	35.16
				E → C	15.71	20.76

Translation Results

Compared with Moses (SMT) and RNNSearch (NMT)

System	Training Data			Direction	NIST06	NIST02
	CE	C	E			
MOSES	✓	×	×	C → E	32.48	32.69
				E → C	14.27	18.28
	✓	×	✓	C → E	34.59	35.21
	✓	✓	×	E → C	20.69	25.85
RNNSEARCH	✓	×	×	C → E	30.74	35.16
				E → C	15.71	20.76

Translation Results

Compared with Moses (SMT) and RNNSearch (NMT)

System	Training Data			Direction	NIST06	NIST02
	CE	C	E			
MOSES	✓	×	×	C → E	32.48	32.69
				E → C	14.27	18.28
	✓	×	✓	C → E	34.59	35.21
	✓	✓	×	E → C	20.69	25.85
RNNSEARCH	✓	×	×	C → E	30.74	35.16
				E → C	15.71	20.76
	✓	×	✓	C → E	35.61 ^{**++}	38.78 ^{**++}
				E → C	17.59 ⁺⁺	23.99 ⁺⁺

Translation Results

Compared with Moses (SMT) and RNNSearch (NMT)

System	Training Data			Direction	NIST06	NIST02
	CE	C	E			
MOSES	✓	×	×	C → E	32.48	32.69
				E → C	14.27	18.28
	✓	×	✓	C → E	34.59	35.21
	✓	✓	×	E → C	20.69	25.85
RNNSEARCH	✓	×	×	C → E	30.74	35.16
				E → C	15.71	20.76
	✓	×	✓	C → E	35.61 ^{**++}	38.78 ^{**++}
				E → C	17.59 ⁺⁺	23.99 ⁺⁺
	✓	✓	×	C → E	35.01 ⁺⁺	38.20 ^{**++}
				E → C	21.12 ^{*++}	29.52 ^{**++}

Translation Results

Compared with Moses (SMT) and RNNSearch (NMT)

System	Training Data			Direction	NIST06	NIST02	NIST03	NIST04	NIST05
	CE	C	E						
MOSES	✓	×	×	C → E	32.48	32.69	32.39	33.62	30.23
				E → C	14.27	18.28	15.36	13.96	14.11
	✓	×	✓	C → E	34.59	35.21	35.71	35.56	33.74
				E → C	20.69	25.85	19.76	18.77	19.74
RNNSEARCH	✓	×	×	C → E	30.74	35.16	33.75	34.63	31.74
				E → C	15.71	20.76	16.56	16.85	15.14
	✓	×	✓	C → E	35.61 ^{**++}	38.78 ^{**++}	38.32 ^{**++}	38.49 ^{**++}	36.45 ^{**++}
				E → C	17.59 ⁺⁺	23.99 ⁺⁺	18.95 ⁺⁺	18.85 ⁺⁺	17.91 ⁺⁺
	✓	✓	×	C → E	35.01 ⁺⁺	38.20 ^{**++}	37.99 ^{**++}	38.16 ^{**++}	36.07 ^{**++}
				E → C	21.12 ^{*++}	29.52 ^{**++}	20.49 ^{**++}	21.59 ^{**++}	19.97 ⁺⁺

Translation Results

Compared with *Sennrich et al. [2015a]*

Method	Training Data			Direction	NIST06	NIST02	NIST03	NIST04	NIST05
	CE	C	E						
<i>Sennrich et al. [2015a]</i>	✓	×	✓	C → E	34.10	36.95	36.80	37.99	35.33
	✓	✓	×	E → C	19.85	28.83	20.61	20.54	19.17
<i>this work</i>	✓	×	✓	C → E	35.61**	38.78**	38.32**	38.49*	36.45**
				E → C	17.59	23.99	18.95	18.85	17.91
	✓	✓	×	C → E	35.01**	38.20**	37.99**	38.16	36.07**
				E → C	21.12**	29.52**	20.49	21.59**	19.97**

Example Translation of Monolingual Corpus

Monolingual	hongsen shuo , ruguo you na jia famu gongsi dangan yishenshifa , name tamen jiang zihui qiancheng .
Reference	hongsen said, if any <i>logging companies</i> dare to defy the law, then they will <i>destroy their own future</i> .
Translation	hun sen said , if any of <i>those companies</i> dare defy the law , then they will <i>have their own fate</i> . [iteration 0]
	hun sen said if any <i>tree felling company</i> dared to break the law , then they would <i>kill themselves</i> . [iteration 40K]
	hun sen said if any <i>logging companies</i> dare to defy the law , they would <i>destroy the future themselves</i> . [iteration 240K]

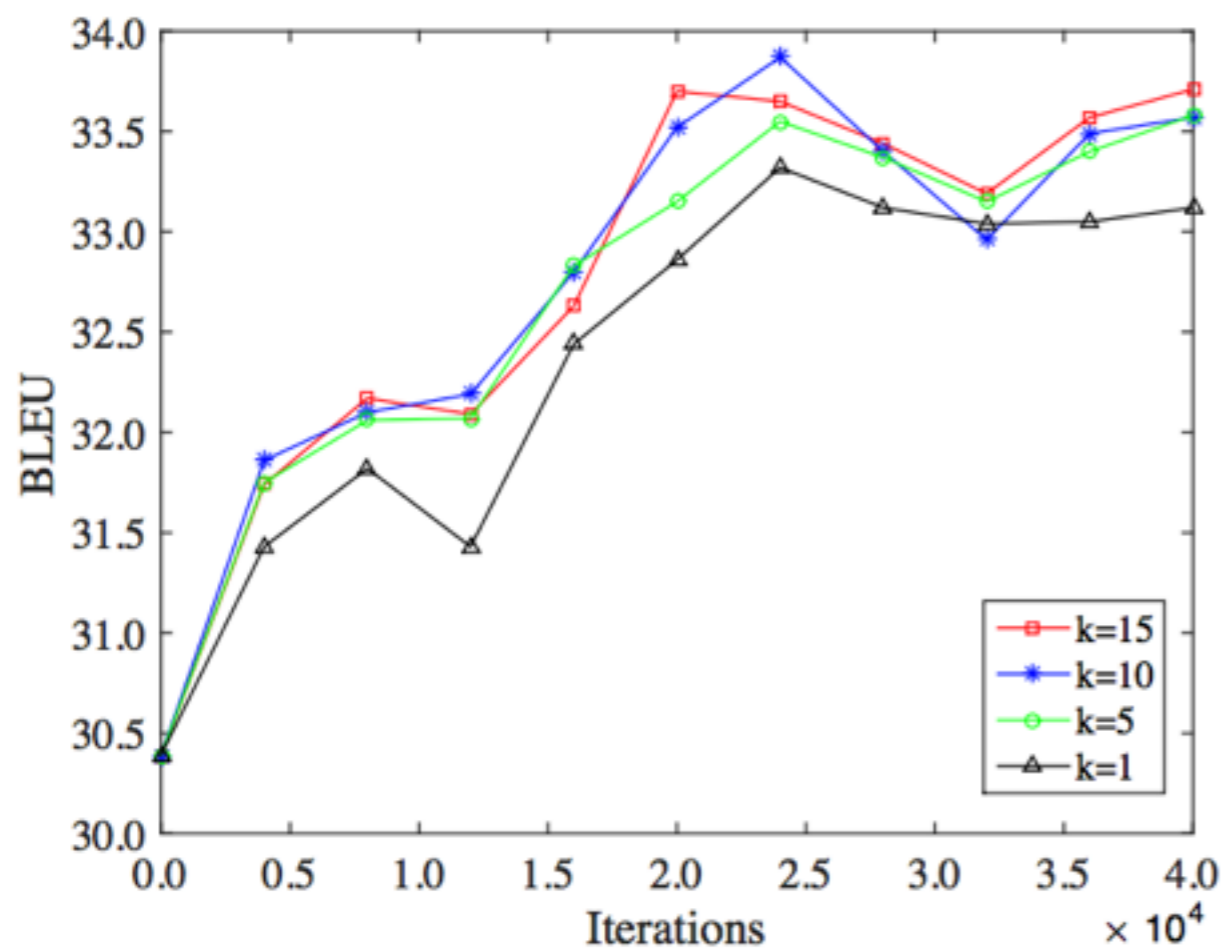
$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{x}} \left\{ P(\mathbf{y}|\mathbf{x}; \vec{\theta}) P(\mathbf{x}'|\mathbf{y}; \overleftarrow{\theta}) \right\}$$

Conclusion

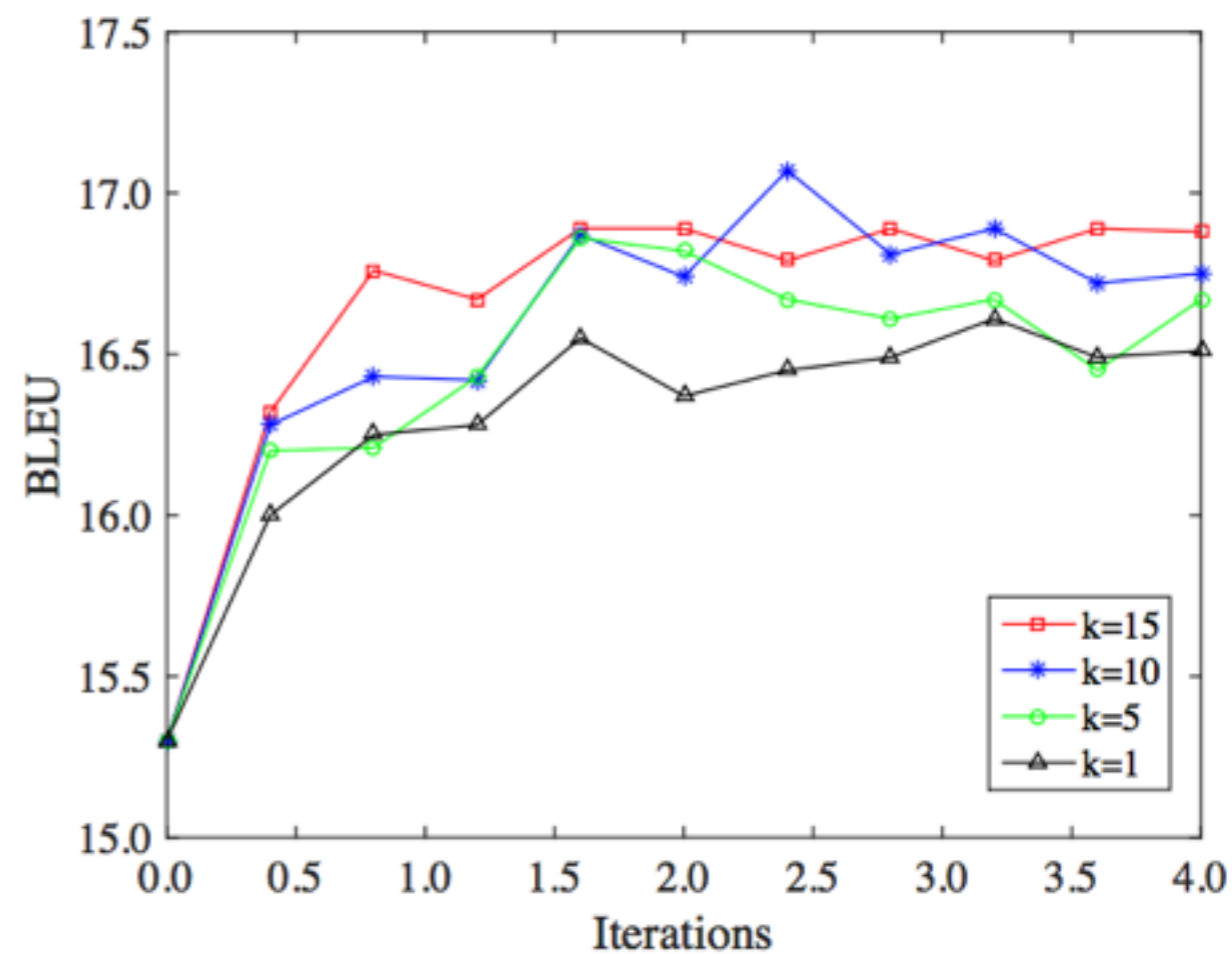
- * Monolingual corpora is an important resource for neural machine translation.
- * We have proposed a semi-supervised approach to training bidirectional neural machine translation models for exploiting monolingual corpora.
- * As our method is sensitive to the OOVs present in monolingual corpora, we plan to integrate Jean et al. (2015)'s technique on using very large vocabulary into our approach.

Thank You !

Effect of Sample Size

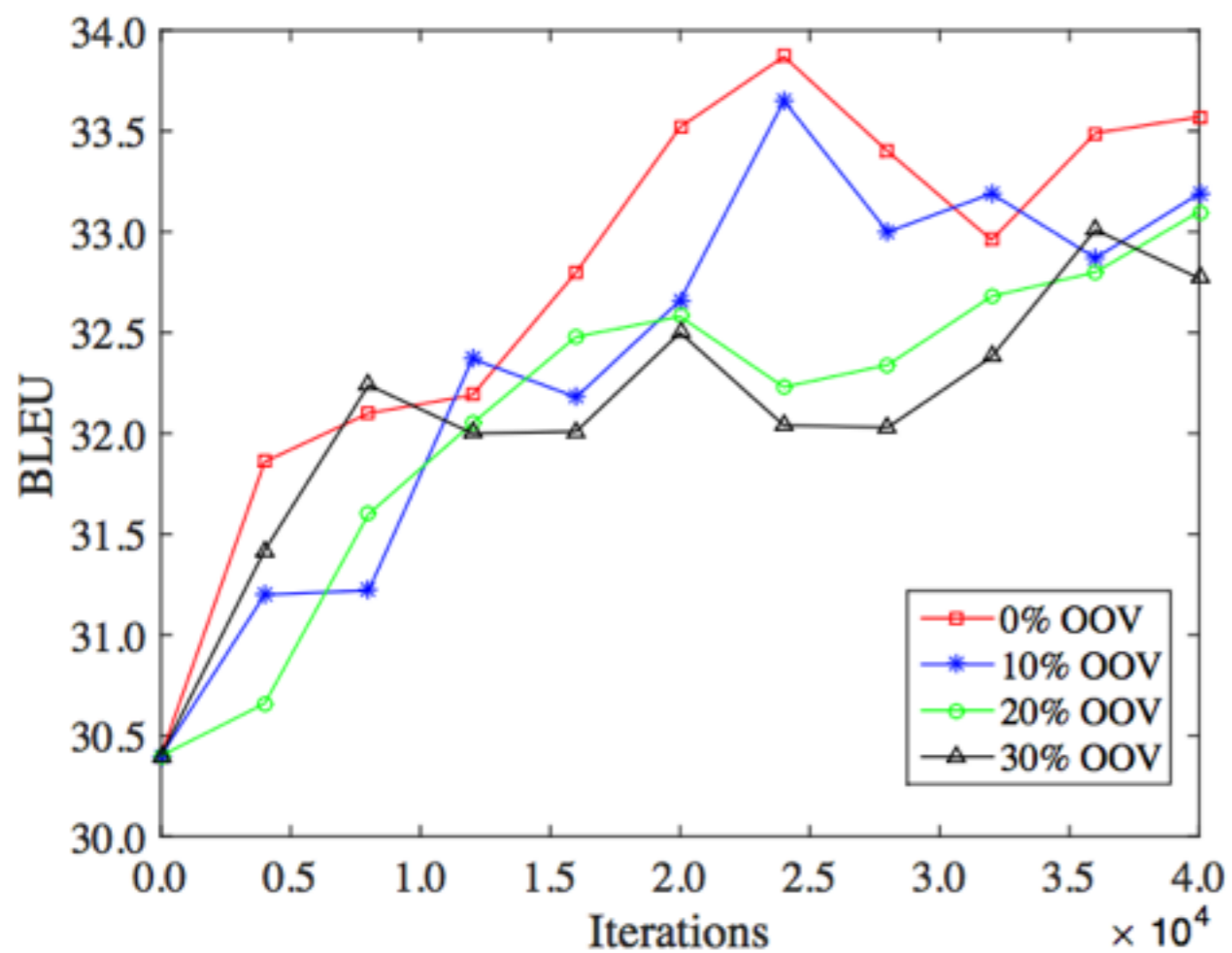


ZH-EN

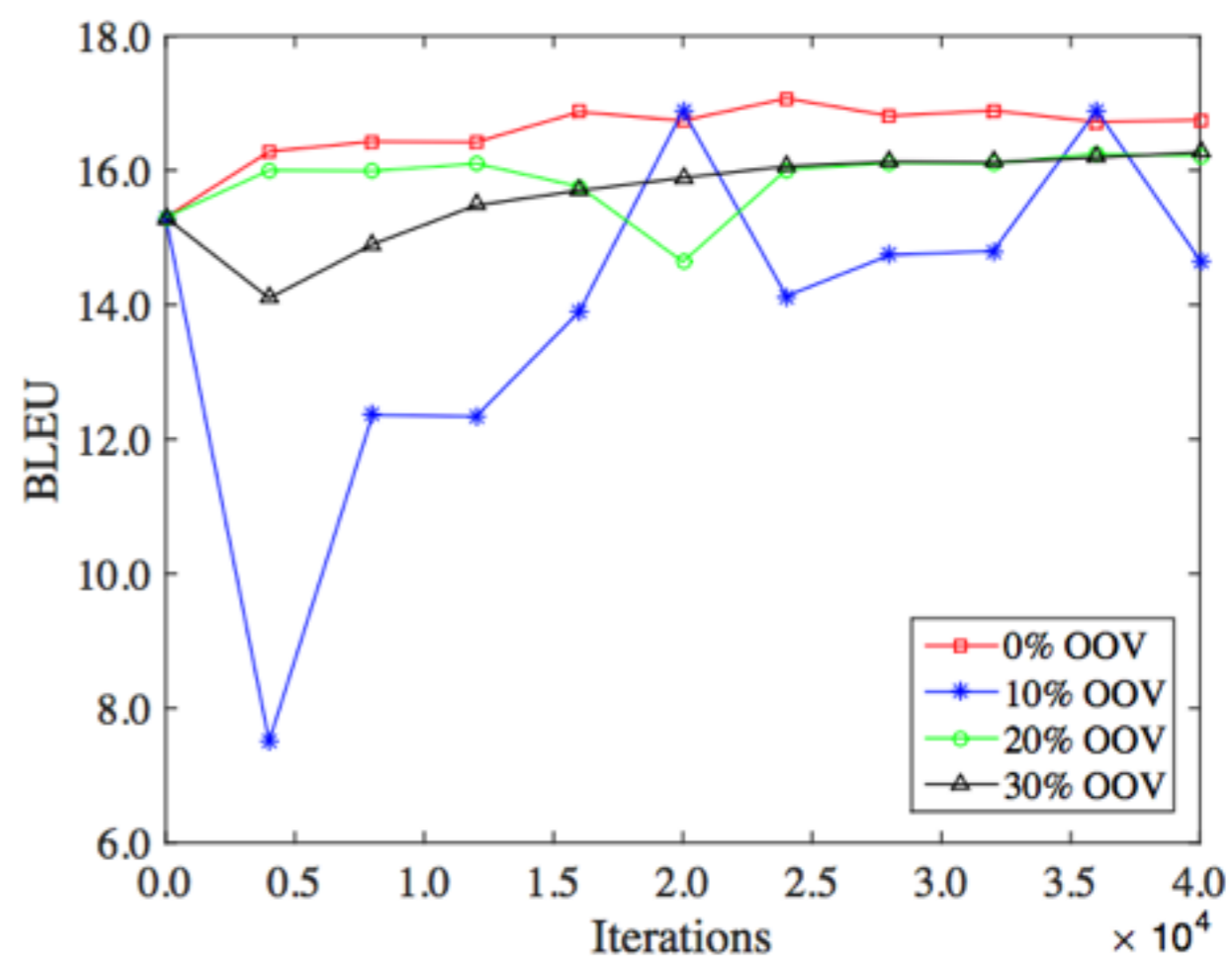


EN-ZH

Effect of OOV ratio



ZH-EN



EN-ZH