

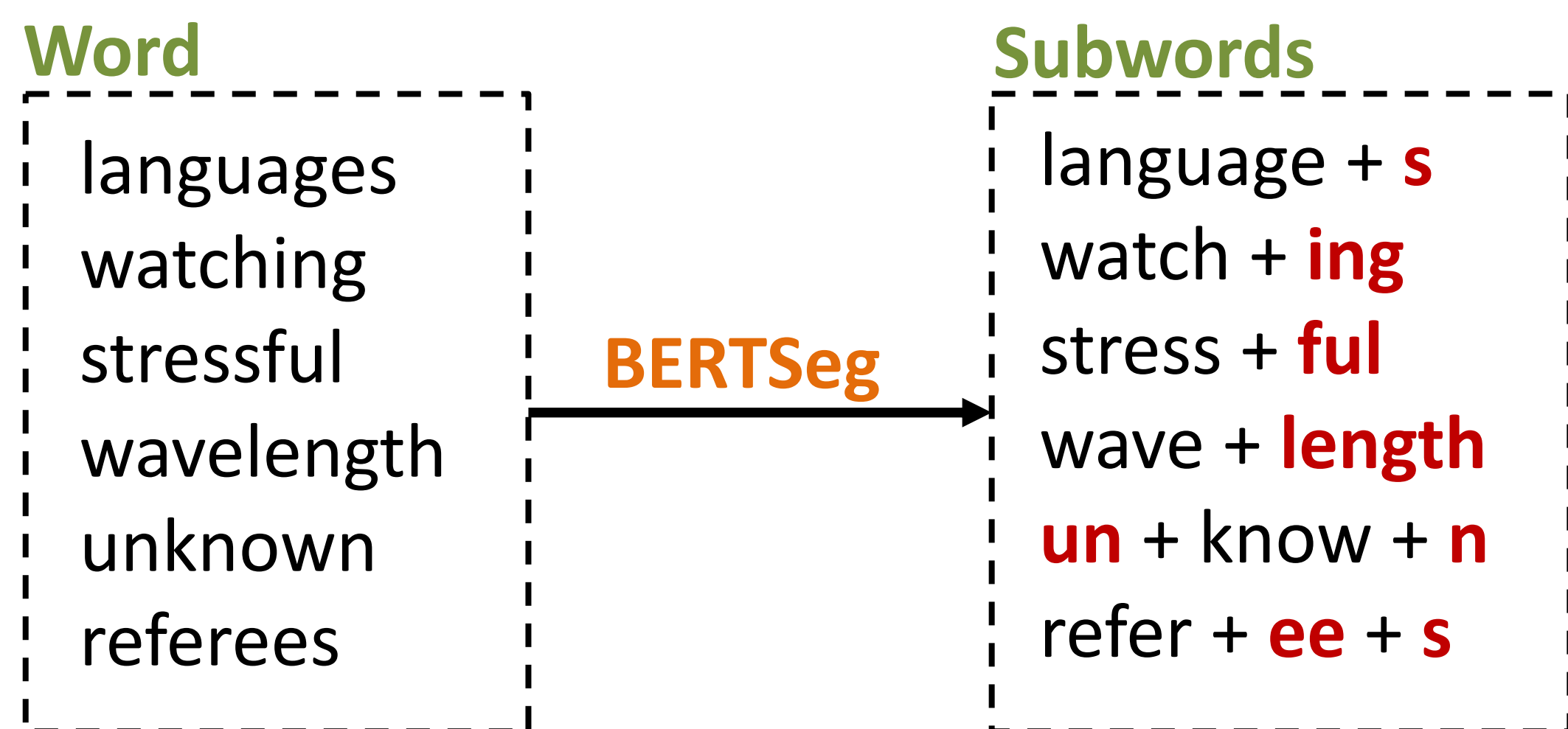
# BERTSeg: BERT Based Subword Segmentation for Neural Machine Translation



Haiyue Song<sup>†,‡</sup>, Raj Dabre<sup>‡</sup>, Zhuoyuan Mao<sup>†</sup>, Chenhui Chu<sup>†</sup>, Sadao Kurohashi<sup>†</sup>  
<sup>†</sup>Kyoto University <sup>‡</sup>NICT

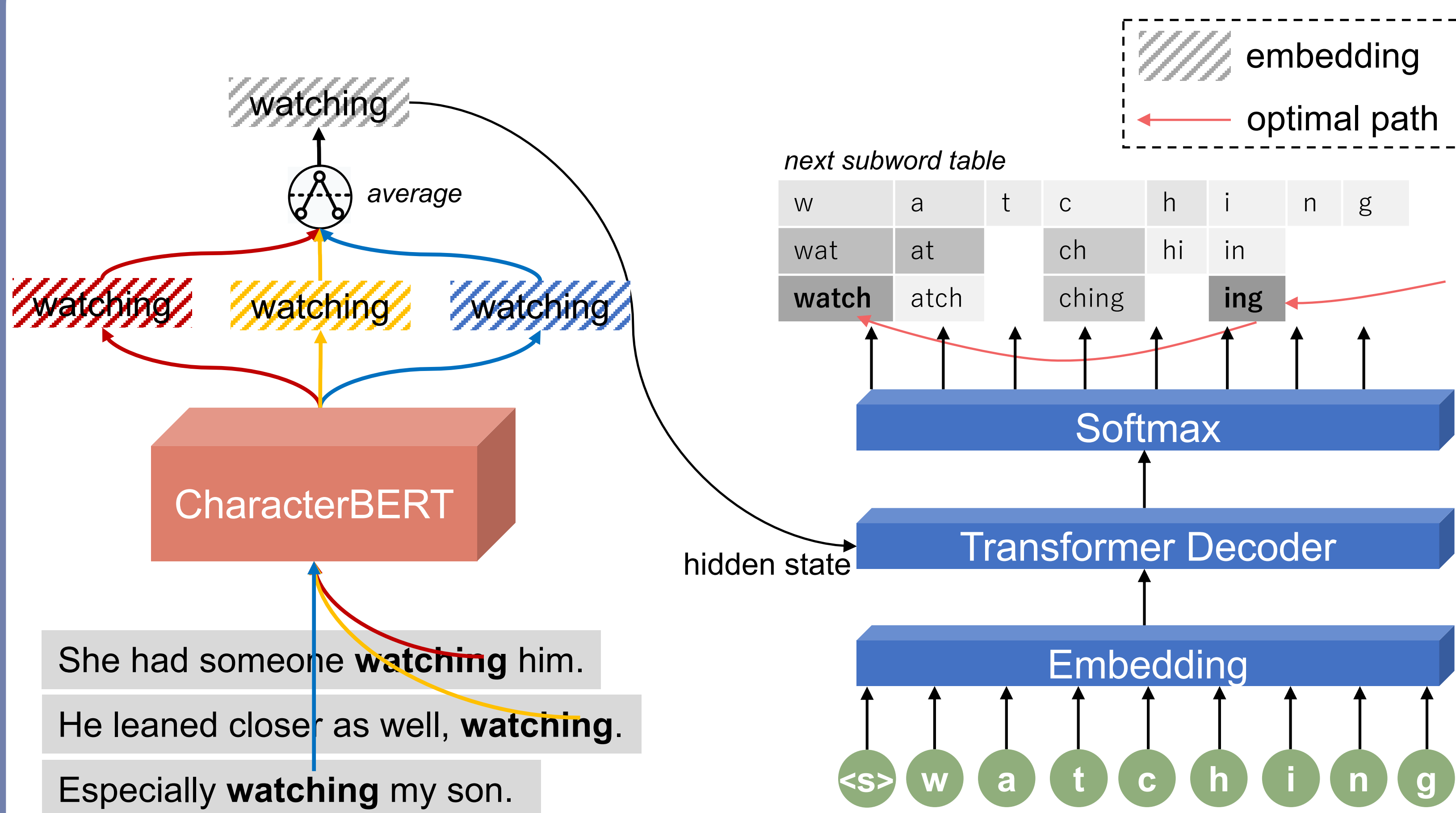
## Introduction

- Subword segmentation methods<sup>1,2</sup> rely on statistical information
- We propose:
  - **BERTSeg**, a BERT-based neural method that relies on semantic information of the target word
  - A **regularization** method for BERTSeg

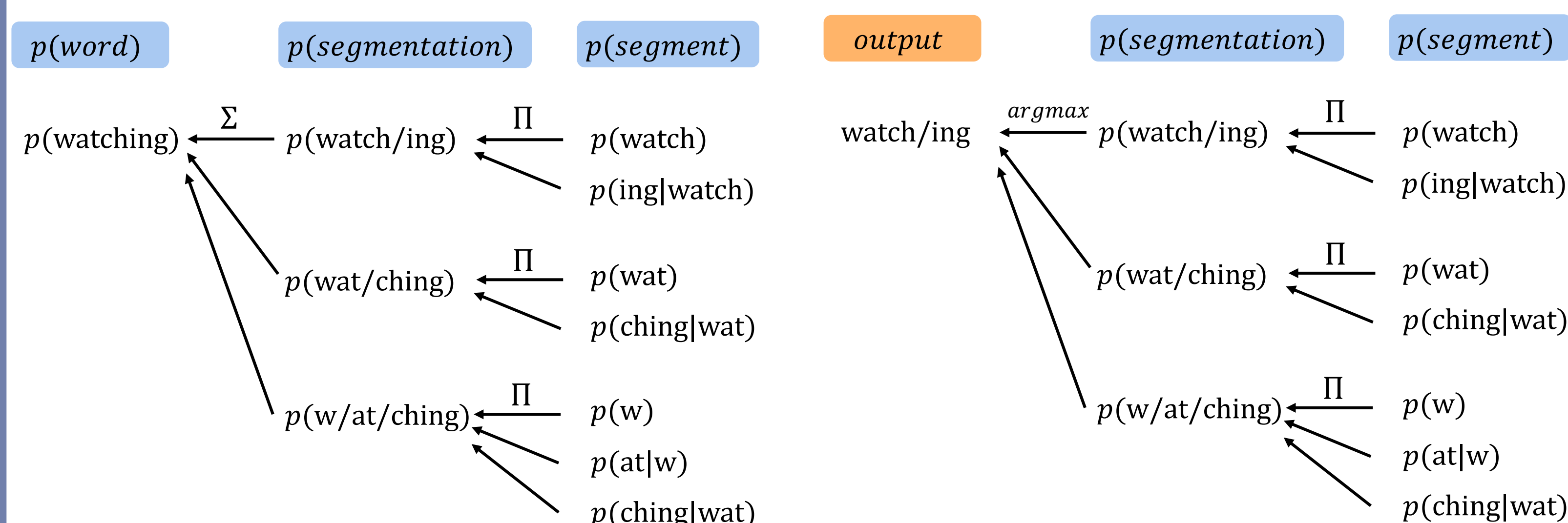


## Method

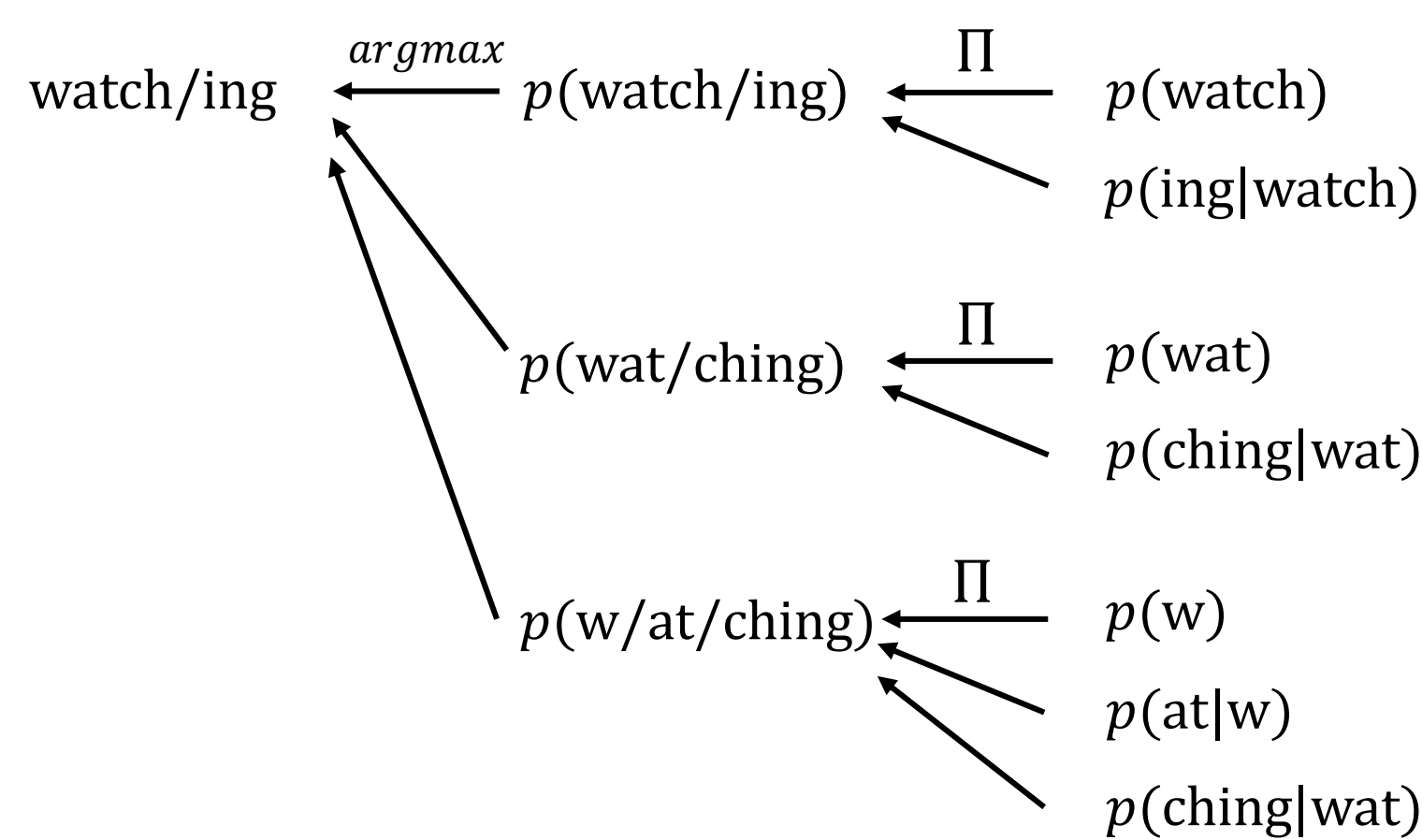
### 1 Architecture



### 2 Train



### 3 Inference



$$\log p(\mathbf{x}_{1:T} | \mathbf{e}_x) = \log \sum_{\mathbf{a}_{1:\tau_a} \in \mathcal{S}(\mathbf{x})} \prod_{i=1}^{\tau_a} p(a_i | \mathbf{e}_x; x_1, \dots, x_j)$$
$$\mathbf{a} = \arg \max_{\mathbf{a}_{1:\tau_a} \in \mathcal{S}(\mathbf{x})} \prod_{i=1}^{\tau_a} p(a_i | \mathbf{e}_x; x_1, \dots, x_j)$$
$$p_{\text{sample}}(\mathbf{a}_i) = \frac{e^{\log p(\mathbf{a}_i)/t}}{\sum_{\mathbf{a}_i \in \mathcal{S}(\mathbf{x})} e^{\log p(\mathbf{a}_i)/t}}$$

## Results

### 1 Machine Translation

- Higher BLEU scores** on low-resource to high-resource datasets compared with BPE<sup>1</sup>, VOLT<sup>2</sup>, DPE<sup>3</sup>, BPE-dropout<sup>4</sup>

	ALT Asian Langs→En	IWSLT15 Vi→En	WMT16 Ro→En	WMT15 Fi→En
<i>w/o Regularization</i>				
BPE <sup>1</sup>	19.76	27.09	<b>32.54</b>	17.45
VOLT <sup>2</sup>	19.91	27.16	31.89	17.25
DPE <sup>3</sup>	19.88	27.40	29.95	16.14
BERTSeg	<b>20.71</b>	<b>27.80</b>	32.33	<b>17.54</b>
<i>With Regularization</i>				
BPE-dropout <sup>4</sup>	23.04	28.76	33.59	<b>18.50</b>
BERTSeg-Regularization	<b>24.68</b>	<b>30.09</b>	<b>33.82</b>	18.46

### 2 Speed

- BERTSeg** requires about 400 seconds on large corpus during training, which is much faster than previous neural method DPE<sup>3</sup>

	ALT	WMT16 Ro-En
BPE <sup>1</sup>	4	13
VOLT <sup>2</sup>	960	1,747
DPE <sup>3</sup>	3,477	68,334
BERTSeg	58	391

### 3 Segmentation Examples

- High generalization ability on **rare** or **unseen** words compared with BPE<sup>1</sup>

BERTSeg		BPE	
<i>Frequent words</i>		<i>Rare words</i>	
official/s	officials	inter/face/s	inter/f/aces
edit/ion	edition	sea/side	se/as/ide
use/d	used	ab/normal/ly	ab/n/orm/ally
farm/er/s	far/mers	b/y/stand/er	by/st/ander
contribute/d	contrib/uted	dis/comfort	disc/om/fort
normal/ly	norm/ally	un/warrant/ed	un/w/arr/anted
seven/th	sevent/h	in/definitely	ind/ef/in/itely
		<i>Unseen words</i>	
		stable/d	st/ab/led
		save/r/s	sa/vers
		M/illion/s	Mill/ions
		Free/way	Fre/ew/ay
		M/i/s/behavior	M/is/be/hav/ior
		m/o/u/r/n/ed	m/our/ned
		M/a/d/a/m/e	Mad/ame

- Subword segmentation with regularization

- Global best N segmentations are obtained through Dynamic Programming algorithm in  $O(N \log N * T^2)$

BERTSeg-Regularization Segmentation	
represent/ed	represented
represent/e/d	re/represented
re/presented	re/present/e/d

## Conclusion

- We proposed **BERTSeg**, an unsupervised neural subword segmenter for NMT, together with a regularization algorithm
- MT results** showed significant improvement over frequency-based and neural network-based methods
- The training is efficient** even compared with non-neural methods

### Future Work

- A **multilingual segmenter** using embeddings from BERT, mBERT, or character-level mBERT
- Remove the dependency on the BPE vocabulary

<sup>1</sup>Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural Machine Translation of Rare Words with Subword Units. ACL.

<sup>2</sup>Jingjing Xu, Hao Zhou, Chun Gan, Zaixiang Zheng, and Lei Li. 2021. Vocabulary Learning via Optimal Transport for Neural Machine Translation. ACL.

<sup>3</sup>Xuanli He, Gholamreza Haffari, and Mohammad Norouzi. 2020. Dynamic Programming Encoding for Subword Segmentation in Neural Machine Translation. ACL.

<sup>4</sup>Ivan Provilkov, Dmitrii Emelianenko, and Elena Voita. 2020. BPE-Dropout: Simple and Effective Subword Regularization. ACL.