

# Dynamic Programming Encoding for Subword Segmentation in Neural Machine Translation

(He et al., ACL 2020)

---

出口 祥之

✉ `deguchi@ai.cs.ehime-u.ac.jp`

2020/07/30 二宮研 論文輪読会

## **Paper**

<https://www.aclweb.org/anthology/D19-1098/>

## **Source Code**

<https://github.com/xlhex/dpe>

# Introduction

## NMT におけるサブワード分割

**貪欲法:** バイトペア符号化 (BPE)<sup>1</sup>, 最長一致法<sup>2</sup>

**確率的アルゴリズム:** ユニグラムLM<sup>3</sup>, BPE-dropout<sup>4</sup>

原言語側, 目的言語側ともに複数分割候補が得られる. 訓練時に分割候補を確率的にサンプリングすることでモデルの頑健性向上.

**動的計画法:** 提案手法. サブワード分割の周辺化.

---

<sup>1</sup> Sennrich, Haddow, and Birch, "Neural Machine Translation of Rare Words with Subword Units".

<sup>2</sup> Wu et al., *Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation*.

<sup>3</sup> Kudo, "Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates".

<sup>4</sup> Provilkov, Emelianenko, and Voita, "BPE-Dropout: Simple and Effective Subword Regularization".

# Related Work (Greedy Algorithms)

## BPE, Wordpiece

- unconscious → un + conscious, likes → like + s
- 隣り合った頻出サブワードから順に，予め指定した語彙数に到達するまで再帰的に結合 (BPE)
- 語彙数とデコード速度はトレードオフ
  - (語彙数を小さくするだけであれば文字単位でよい)
  - テキスト圧縮の技術を利用
  - 語彙数の上限を制約とし，文長が短くなるような分割を得るアルゴリズム

# Related Work (Stochastic Algorithms)

## ユニグラムLM, BPE-dropout

- unconscious  $\rightarrow$  {un + concious, uncon + scious}
- 複数分割候補を得られる
  - ユニグラムLM: 尤度ベースでサンプリング
  - BPE-dropout: 結合時に確率的に棄却
  - NMT 訓練時に分割を確率的に得ることでデータ拡張 (Data Augmentation) の効果
    - ▶ モデルの頑健性, 汎用化

# Related Work (Dynamic Programming Algorithms)

音声認識<sup>5</sup>

非自己回帰 NMT モデル<sup>67</sup>

---

<sup>5</sup> Wang et al., *Sequence Modeling via Segmentations*.

<sup>6</sup> Chan et al., *Imputer: Sequence Modelling via Imputation and Dynamic Programming*.

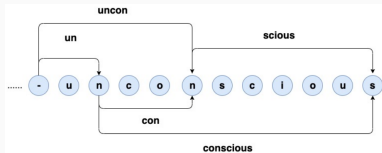
<sup>7</sup> Saharia et al., *Non-Autoregressive Machine Translation with Latent Alignments*.

# Latent Subword Segmentation - Definitions

## 目的言語側の分割を潜在変数とみなす

■  $M$  個のサブワード境界:  $\{y_{z_i, z_{i+1}}\}_{i=1}^M$

- $y = (y_1, \dots, y_T)$ : 目的言語文の文字列
- $z = (z_1(=0), \dots, z_{M+1}(=T))$ : 文字インデックス列
- $y_{a,b}$ :  $(a+1)^{\text{th}}$  から  $b^{\text{th}}$  まで結合したサブワード



例:

- 辞書  $\mathcal{V} = \{'c', 'a', 't', 'ca', 'at'\}$
- 単語: 'cat'

$z$	サブワード列
$(0, 1, 3)$	$(\text{'c'}, \text{'at'})$
$(0, 2, 3)$	$(\text{'ca'}, \text{'t'})$
$(0, 1, 2, 3)$	$(\text{'c'}, \text{'a'}, \text{'t'})$

# Latent Subword Segmentation - Likelihood

## 連鎖律を用いてサブワード列の対数尤度を表現

- 各サブワードにおいて語彙のカテゴリ分布を生成

$$\log p(\mathbf{y}, \mathbf{z} | \mathbf{x}) = \sum_{i=1}^{|\mathbf{z}|} \log p(\mathbf{y}_{z_i, z_{i+1}} | \mathbf{y}_{z_1, z_2}, \dots, \mathbf{y}_{z_{i-1}, z_i}, \mathbf{x})$$

※  $\mathbf{x}$ : 原言語文

- 殆どの NMT は  $\mathbf{z}$  は  $\mathbf{y}$  の決定論的関数とみなされる:

$$\log p(\mathbf{y}, \mathbf{z}) \approx \log p(\mathbf{y})$$



# Latent Subword Segmentation - Latent Variable

$z \in \mathcal{Z}_y$  ( **$y$ の分割集合**) を潜在表現とみなす

■  $p(\mathbf{y}|\mathbf{x}) = \sum_z p(\mathbf{y}, z|\mathbf{x})$  とする

$$\log p(\mathbf{y}|\mathbf{x}) = \log \sum_{z \in \mathcal{Z}_y} \exp \sum_{i=1}^{|z|} \log p(\mathbf{y}_{z_i, z_{i+1}} | \mathbf{y}_{z_1, z_2}, \dots, \mathbf{y}_{z_{i-1}, z_i}, \mathbf{x})$$

※ 対数周辺尤度の下限:  $\log p(\mathbf{y}|\mathbf{x}) \geq \log p(\mathbf{y}, z|\mathbf{x})$

■ 各サブワードの確率が条件部のコンテキストの分割に依存するため、巨大な空間  $\mathcal{Z}_y$  上での厳密な周辺化は組み合わせ爆発を起こす

# A Mixed Character-Subword Transformer

## 文字に基づいてサブワードを生成する Transformer

- 条件部のコンテキストを文字のみに

$$\log p(\mathbf{y}, \mathbf{z} | \mathbf{x}) = \sum_{i=1}^{|\mathbf{z}|} \log p(\mathbf{y}_{z_i, z_{i+1}} | y_{z_1}, \dots, y_{z_i}, \mathbf{x})$$

- $\mathbf{y}$  の各文字位置  $t$  において、次に来るサブワード  $w \in \mathcal{V}$  の分布を以下に基づいて生成

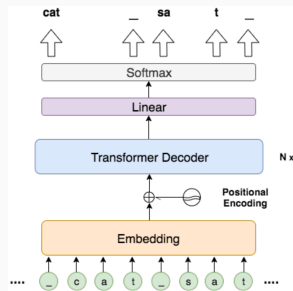
$$p(w | y_1, \dots, y_t, \mathbf{x}) = \frac{\exp(f(y_1, \dots, y_t)^\top e(w))}{\sum_{w' \in \mathcal{V}} \exp(f(y_1, \dots, y_t)^\top e(w'))}$$

- $f(\cdot)$  : Transformer により条件部の計算
- $e(\cdot)$  : ソフトマックス層の重み

# A Mixed Character-Subword Transformer

## $t$ ステップ目のモデル出力

1.  $t$  ステップ目でサブワード  $w$  を生成
2. サブワード  $w$  の文字をデコーダに入力 ( $t + 1$  から  $t + |w|$  まで)
3.  $t + |w|$  ステップ目で次のサブワードを生成



目的関数  $\mathcal{L}(\theta)$  を最大化

$$\mathcal{L}(\theta) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log P(\mathbf{y}|\mathbf{x})$$

- 周辺化と対数周辺尤度の勾配計算が必要

# Exact Marginalization

## 動的計画法を用いて周辺尤度を計算

- サブワードの出力確率が文字のみによって得られるため動的計画法によって対数周辺尤度が計算可能

---

**Algorithm 1** Dynamic Programming (DP) for Exact Marginalization

---

**Input:**  $\mathbf{y}$  is a sequence of  $T$  characters,  $V$  is a subword vocabulary,  $m$  is the maximum subword length

**Output:**  $\log p(\mathbf{y})$  marginalizing out different subword segmentations.

```
1:  $\alpha_0 \leftarrow 0$ 
2: for  $k = 1$  to  $T$  do
3:    $\alpha_k \leftarrow \log \sum_{j=k-m}^{k-1} \mathbb{1}[\mathbf{y}_{j,k} \in V] \exp \left( \alpha_j + \log P_{\theta}(\mathbf{y}_{j,k} | y_1, \dots, y_j) \right)$ 
4: end for
5: return  $\alpha_T$                                  $\triangleright$  the marginal probability  $\log p(\mathbf{y}) = \log \sum_{\mathbf{z} \in \mathcal{Z}_{\mathbf{y}}} p(\mathbf{y}, \mathbf{z})$ 
```

---

- 計算量:  $\mathcal{O}(mT)$

- $m$ : 語彙に含まれる最長の単語の文字数

# Gradient Computation

## 計算量増加への対処

- PyTorch での著者実装で通常の Transformer デコーダより 8 倍遅く，メモリ使用量も増加
  - DP アルゴリズムと文字レベルでの演算による系列長の増加が原因
- Transformer のレイヤ数を 6 から 4 に減らし，16 ステップ勾配蓄積 (Gradient Accumulation) してからパラメタ更新

# Segmenting Target Sentences

## Dynamic Programming Encoding (DPE): 最大事後確率を持つ目的言語文の分割を探索

---

**Algorithm 2** Dynamic Programming Encoding (DPE) for Subword Segmentation

---

**Input:**  $y$  is a sequence of  $T$  characters,  $V$  is a subword vocabulary,  $m$  is the maximum subword length

**Output:** Segmentation  $z$  with highest posterior probability.

**for**  $k = 1$  **to**  $T$  **do**

$\beta_k \leftarrow \max_{\{j \in [k-m, k-1] \mid y_{j,k} \in V\}} \beta_j + \log P_\theta(y_{j,k} | y_1, \dots, y_j)$

$b_k \leftarrow \operatorname{argmax}_{\{j \in [k-m, k-1] \mid y_{j,k} \in V\}} \beta_j + \log P_\theta(y_{j,k} | y_1, \dots, y_j)$

**end for**

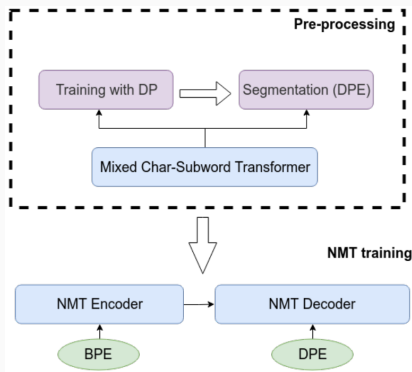
$z \leftarrow \operatorname{backtrace}(b_1, \dots, b_T)$

▷ backtrace the best segmentation using  $b$

---

# Segmenting Target Sentences

- 混合文字サブワード Transformer は訓練データの目的言語文の分割のためのみに使用
- 分割した文で通常のサブワード Transformer を訓練





# Experiments

**データセット** WMT09 En-Hu, WMT14 En-De, WMT15  
En-Fi, WMT16 En-Ro, WMT18 En-Et

## モデル

---

アーキテクチャ	Transformer base
分割 (原言語側)	BPE-dropout ( $p = 0.05$ )
(目的言語側)	DPE

---

# Main Results

Method	BPE	BPE dropout	$\Delta_1$	This paper	$\Delta_2$
Source segmentation	BPE	BPE dropout		BPE dropout	
Target segmentation	BPE	BPE dropout		DPE	
En→De	27.11	27.27	+0.16	27.61	+0.34
En→Ro	27.90	28.07	+0.17	28.66	+0.59
En→Et	17.64	18.20	+0.56	18.80	+0.60
En→Fi	15.88	16.18	+0.30	16.89	+0.71
En→Hu	12.80	12.94	+0.14	13.36	+0.42
De→En	30.82	30.85	+0.03	31.21	+0.36
Ro→En	31.67	32.56	+0.89	32.99	+0.43
Et→En	23.13	23.65	+0.52	24.62	+0.97
Fi→En	19.10	19.34	+0.24	19.87	+0.53
Hu→En	16.14	16.61	+0.47	17.05	+0.44
Average	22.22	22.57	+0.35	23.12	+0.55

# Segmentation Examples

---

BPE source:

Die G@@ le@@ is@@ anlage war so ausgestattet , dass dort elektr@@ isch betrie@@ bene Wagen eingesetzt werden konnten .

DPE target:

The railway system was equipped in such a way that electrical@@ ly powered cart@@ s could be used on it .

BPE target:

The railway system was equipped in such a way that elect@@ r@@ ically powered car@@ ts could be used on it .

---

BPE source:

Normalerweise wird Kok@@ ain in kleineren Mengen und nicht durch Tunnel geschm@@ ug@@ gelt .

DPE target:

Normal@@ ly c@@ oca@@ ine is sm@@ ugg@@ led in smaller quantities and not through tunnel@@ s .

BPE target:

Norm@@ ally co@@ c@@ aine is sm@@ ugg@@ led in smaller quantities and not through tun@@ nels .

---

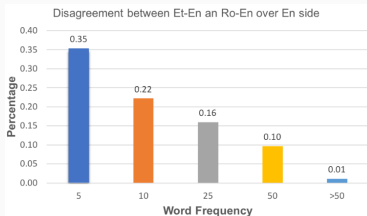
## ■ 他の例は論文参照

# Conditional Subword Segmentation

原言語文を条件部に入れず，LM で分割

Source Target	BPE drop BPE drop	BPE drop LM DPE	BPE drop DPE
En→Ro	28.07	28.07	28.66
En→Hu	12.94	12.87	13.36
Ro→En	32.56	32.57	32.99
Hu→En	16.61	16.41	17.05

同一の目的言語文で原言語側を変えて違いを比較



# Conditional Subword Segmentation

## 原言語文が BPE-dropout によって変化することの有効性

Source Target	BPE drop DPE Fixed	BPE drop DPE On The Fly
En→Ro	28.58	28.66
En→Hu	13.14	13.36
En→Et	18.51	18.80
Ro→En	32.73	32.99
Hu→En	16.82	17.05
Et→En	24.37	24.62

## 目的言語側の分割アルゴリズムを変えて比較

Source Target	BPE drop BPE	BPE drop BPE drop	BPE drop DPE
En→Ro	28.04	28.07	28.66
En→Et	18.09	18.20	18.80
Ro→En	32.40	32.56	32.99
Et→En	23.52	23.65	24.62

## ■ **Dynamic Programming Encoding** を提案

- 訓練時は目的言語側の分割を潜在変数とみなして周辺化
- 推論時は事後確率が最も高くなる分割を出力

## ■ BPE だけでなく BPE-dropout と比較しても翻訳性能が向上