

SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing

نکته: تمامی عناوین و کلمات قرمز رنگ به صورت لینک های قابل کلیک هستند.

چکیده

این مقاله SentencePiece را توصیف می کند؛ یک توکنایزر و دیتوکنایزر زیرواحدی مستقل از زبان که برای پردازش متون مبتنی بر شبکه های عصبی، شامل ترجمه ماشینی عصبی، طراحی شده است. این ابزار پیاده سازی متن باز C++ و Python برای واحد های زیرواحدی ارائه می دهد. در حالی که ابزارهای موجود برای بخش بندی زیرواحدی فرض می کنند ورودی پیش تر به دنباله ای از واژه ها توکنیزه شده است، SentencePiece می تواند مدل های زیرواحدی را مستقیماً از جملات خام آموزش دهد؛ که این امکان را فراهم می آورد تا سیستمی کاملاً انتهای انتها به انتها (end-to-end) و مستقل از زبان ایجاد شود. ما یک آزمایش اعتبار سنجی بر روی ترجمه ماشینی عصبی انگلیسی - زبانی انجام داده ایم و مشاهده کردیم که می توان دقیقی قابل مقایسه با آموزش مستقیم زیرواحدی از جملات خام به دست آورد. همچنین عملکرد آموزش و بخش بندی زیرواحدی را در پیکربندی های مختلف مقایسه می کنیم. SentencePiece تحت مجوز Apache 2 در دسترس است به آدرس:

<https://github.com/google/sentencepiece>

شبکه‌های عصبی عمیق تأثیر بزرگی بر پردازش زبان طبیعی داشته‌اند. ترجمه ماشین عصبی (NMT) (Bahdanau et al., 2014; Luong et al., 2015; Wu et al., 2016; Vaswani et al., 2017) به‌طور خاص محبویت روزافروزی پیدا کرده است، چرا که می‌تواند با استفاده از شبکه‌های عصبی، ترجمه‌ها را به‌طور مستقیم و تنها با یک معناری ساده‌ی انتها به انتها انجام دهد. NMT در چندین وظیفه‌ی مشترک نتایج چشمگیری نشان داده است (Denkowski and Neubig, 2017; Nakazawa et al., 2017) و رویکرد مؤثر آن بر سایر وظایف مرتبط در NLP مانند تولید گفت‌و‌گو (Vinyals and Le, 2015) و خلاصه‌سازی خودکار (Rush et al., 2015) نیز تأثیر گذاشته است.

اگرچه NMT به‌طور بالقوه می‌تواند ترجمه‌ی انتها به انتها انجام دهد، بسیاری از سیستم‌های NMT همچنان به پیش‌پردازش‌ها و پس‌پردازش‌های وابسته به زبان متکی هستند، مشابه آنچه در سیستم‌های سنتی ترجمه ماشینی آماری (SMT) استفاده می‌شد. Moses، یک جعبه‌ابزار استاندارد برای SMT، مجموعه‌ای نسبتاً کارآمد از پیش‌پردازش و پس‌پردازش را پیاده‌سازی می‌کند. با این حال، این ابزار بر اساس قواعد دستی و وابسته به زبان ساخته شده است که کارایی آن‌ها برای NMT هنوز اثبات نشده است. علاوه بر این، این ابزارها عمدتاً برای زبان‌های اروپایی طراحی شده‌اند که در آن‌ها کلمات با فاصله از هم جدا می‌شوند. برای آموزش سیستم‌های NMT در زبان‌های بدون جداساز مانند چینی، کره‌ای و ژاپنی، باید جداکننده‌ی واژه‌ها را به صورت مستقل اجرا کنیم.

چنین پردازش‌های وابسته به زبان همچنین آموزش مدل‌های NMT چندزبانه (Johnson et al., 2016) را دشوار می‌سازد، چرا که نیازمند مدیریت دقیق پیکربندی‌های پیش‌پردازش و پس‌پردازش برای هر زبان است؛ در حالی که معناری‌های درونی شبکه‌های عصبی مستقل از زبان دارند. با استاندارد شدن رویکردهای NMT و حرکت به سمت معناری‌های مستقل‌تر از زبان، برای جامعه‌ی NLP اهمیت بیشتری پیدا می‌کند که یک پیش‌پردازش و پس‌پردازش ساده، کارآمد، قابل بازتولید و مستقل از زبان توسعه یابد که بتواند به راحتی در سیستم‌های NLP مبتنی بر شبکه‌های عصبی، شامل NMT، ادغام شود. در این مقاله، ما SentencePiece را توصیف می‌کنیم؛ یک توکنایزر و دیتوکنایزر متنی ساده و مستقل از زبان که عمدتاً برای سیستم‌های تولید متن مبتنی بر شبکه‌های عصبی طراحی شده است، جایی که اندازه واژگان پیش از آموزش BPE مدل عصبی مشخص می‌شود. SentencePiece دو الگوریتم بخش‌بندی زیرواحدی، یعنی رمزگذاری جفت‌بایتی (Kudo, 2018) و مدل زبانی تک‌واحدی Unigram LM (Sennrich et al., 2016) را با امکان آموزش مستقیم از جملات خام پیاده‌سازی می‌کند. SentencePiece ساخت یک سیستم کاملاً انتها به انتها را ممکن می‌سازد که به هیچ پردازش وابسته به زبان خاصی نیاز ندارد.

```
% spm_train --input=data/input.txt
--model_prefix=spm --vocab_size=1000

% echo "Hello_world." | spm_encode --
model=spm.model
Hello_world.

% echo "Hello_world." | spm_encode --
model=spm.model --output_format=id
151 88 21 887 6

% echo "Hello_world." |
spm_decode --model=spm.model
Hello world.

% echo "151 88 21 887 6" | spm_decode
--model=spm.model
--input_format=id
Hello world.
```

شکل ۱: کاربرد استفاده خط فرمان SentencePiece

۲ مروری بر سیستم

شامل چهار مؤلفه اصلی است: نرم‌الایزر، آموزنده، کدگذار و کدگشا. نرم‌الایزر مأذولی است برای نرم‌السازی کاراکترهای یونیکد هم معنی به فرم‌های استاندارد. آموزنده، مدل بخش‌بندی زیر واحدی را از روی پیکره نرم‌الشده آموزش می‌دهد. نوع مدل زیر واحدی به عنوان پارامتر آموزنده مشخص می‌شود. کدگذار به طور داخلی نرم‌الایزر را اجرا می‌کند تا متن ورودی را نرم‌السازی کند و سپس آن را با استفاده از مدل زیر واحدی آموزش دیده به دنباله‌ای از زیر واحدها توکنیزه می‌کند. کدگشا دنباله زیر واحدی را دوباره به متن نرم‌الشده تبدیل می‌کند. نقش کدگذار و کدگشا به ترتیب معادل postprocessing (detokenization) و preprocessing (tokenization) هستند. با این حال، آن‌ها را «کدگذاری» و «کدگشایی» می‌نامیم، چرا که SentencePiece نگاشت واژگان به شناسه (id) را مدیریت می‌کند و می‌تواند متن را مستقیماً به دنباله‌ای از شناسه‌ها و بالعکس تبدیل کند. کدگذاری و کدگشایی مستقیم به از دنباله‌های شناسه برای اکثر سیستم‌های NMT مفید است، چرا که ورودی و خروجی آن‌ها دنباله‌هایی از شناسه‌ها هستند. شکل؟ یک نمونه‌ی انتهای آموزش (spm_train)، کدگذاری (spm_encode) و کدگشایی (spm_decode) را نشان می‌دهد. می‌بینیم که متن ورودی به صورت برگشت‌پذیر از طریق spm_decode و spm_encode تبدیل می‌شود.

۳ طراحی کتابخانه

این بخش طراحی و جزئیات پیاده‌سازی SentencePiece را همراه با نمونه‌هایی از خط فرمان و کد شرح می‌دهد.

۱.۳ توکنیزاسیون بدون از دست دادن اطلاعات (Lossless Tokenization)

جمله خام و جمله توکنیزه شده زیر نمونه‌ای از پیش‌پردازش وابسته به زبان هستند:

- متن خام: .Hello world

- توکنیزه شده: [.] [world] [Hello]

یک نکته این است که متن خام و دنباله توکنیزه شده به صورت برگشت‌پذیر قابل تبدیل نیستند. برای مثال، این که بین «world» و «.» فاصله‌ای وجود نداشته، در دنباله توکنیزه شده نگهداری نمی‌شود. بنابراین، دیتوکنایزیشن (بازسازی متن خام از دنباله توکن‌ها) باید به صورت وابسته به زبان انجام شود، چرا که این عملیات برگشت‌ناپذیر هستند. برای مثال، در زبان‌های اروپایی دیتوکنایزر معمولاً فاصله‌ها را بین توکن‌ها اضافه می‌کند، در حالی که در زبان‌هایی مثل ژاپنی یا چینی نیازی به فاصله وجود ندارد.

- متن خام: (Hello world.) [こんにちちは世界。]

- توکنیزه شده: [。] [世界] [こんにちちは]

چینی پردازش‌های وابسته به زبان معمولاً با قواعد دستی پیاده‌سازی شده‌اند که نوشتمن و نگهداری آن‌ها پرهازینه است. دیتوکنایزر را به عنوان عمل معکوس کدگذار پیاده‌سازی می‌کند، به این صورت:

$$Decode(Encode(Normalize(text))) = Normalize(text)$$

ما این طراحی را «توکنیزاسیون بدون از دست دادن اطلاعات» می‌نامیم، چرا که همه اطلاعات لازم برای بازتولید متن نرمال‌شده در خروجی کدگذار حفظ می‌شود. ایده اصلی این روش آن است که متن ورودی صرفاً به عنوان دنباله‌ای از کاراکترهای یونیکد در نظر گرفته شود. حتی فاصله نیز به عنوان یک نماد عادی پردازش می‌شود. برای شفافیت بیشتر، ابتدا فاصله‌ها را با نماد متأ (U+2581) _جایگزین می‌کند و سپس متن را به دنباله‌ای دلخواه از زیر واحدها توکنیزه می‌کند. برای مثال:

- متن خام: Hello world.

- توکنیزه شده: [.][ld][_wor][Hello]

از آنجا که فاصله در متن توکنیزه شده حفظ می شود، می توان توکن ها را بدون ابهام بازسازی کرد. برای نمونه در پایتون:

```
detok = ''.join(tokens).replace('_', ' ')
```

لازم به ذکر است که ابزار **subword-nmt** نمایش متفاوتی برای واحدهای زیروحادی دارد. این ابزار بر نحوه بخش بندی واژه به زیروحد تمرکز می کند و از علامت **@@** به عنوان مرز درون واژه ای استفاده می کند:

- توکنیزه شده: [@@.][@@ld][wor][Hello]

این نمایش همیشه قادر به انجام توکنیزاسیون بدون از دست دادن اطلاعات نیست، چرا که در برخورد با فاصله ها ابهام باقی می ماند. به طور خاص، امکان رمزگذاری فاصله های پیاپی در این نمایش وجود ندارد.

۲.۳ آموزش و بخش بندی زیروحادی کارآمد

ابزارهای موجود برای بخش بندی زیروحادی، مدل های زیروحادی را از جملات پیش توکنیزه شده آموزش می دهند. این پیش توکنیزاسیون در ابتدا برای سرعت بخشیدن به آموزش زیروحد معرفی شد (Sennrich et al., 2016). با این حال، همیشه نمی توان فرض کرد که پیش توکنیزاسیون در دسترس است، به ویژه برای زبان های بدون جداساز. علاوه بر این، پیش توکنیزاسیون باعث می شود که انجام توکنیزاسیون بدون از دست دادن اطلاعات دشوار گردد. SentencePiece چندین تکنیک افزایش سرعت هم برای آموزش و هم برای بخش بندی به کار می گیرد تا امکان توکنیزاسیون بدون از دست دادن اطلاعات روی داده های خام حجیم فراهم شود. برای مثال، اگر جمله (یا کلمه) ورودی طول N داشته باشد، بخش بندی مبتنی بر BPE در هر تکرار نیازمند هزینه محاسباتی $O(N^2)$ است، چرا که باید جفت نمادها را در هر مرحله جست و جو کرد. SentencePiece الگوریتمی با پیچیدگی $O(N \log N)$ به کار می گیرد که در آن نمادهای ادغام شده توسط یک صف اولویت دار (binary heap) مدیریت می شوند. علاوه بر این، پیچیدگی زمانی آموزش و بخش بندی در مدل های زبانی تک واحدی (unigram language models) مناسب با اندازه داده های ورودی (یعنی خطی) است.

۳.۳ مدیریت شناسه واژگان

SentencePiece نگاشت واژگان به شناسه (id) را مدیریت می کند تا بتوان متن ورودی را مستقیماً به دنباله ای از شناسه ها تبدیل کرد و برعکس. اندازه واژگان با استفاده از پرچم `-vocab_size=<size>` در دستور `spm_train` مشخص می شود. در حالی که ابزار **subword-nmt** تعداد عملیات ادغام را تعیین می کند، SentencePiece اندازه نهایی واژگان را

تعیین می کند، چرا که تعداد عملیات ادغام یک پارامتر اختصاصی برای BPE است و نمی تواند برای الگوریتم های بخش بندی دیگر مثل مدل زبانی تک واحدی (Kudo, 2018) به کار گرفته شود. SentencePiece تعدادی شناسه ویژه برای نمادهای متا رزرو می کند، از جمله: نماد ناشناخته <unk>, نماد آغاز جمله <s>, نماد پایان جمله </s> و نماد پر کننده <pad>. شناسه دقیق این نمادها با استفاده از پرچم های خط فرمان تنظیم می شود. همچنین می توان نمادهای متای سفارشی تعریف کرد تا اطلاعات متی (contextual information) را به صورت توکن های مجازی کد گذاری کنند. برای مثال، شاخص های زبانی مانند <2ja> و <2de> برای مدل های چند زبانه (Johnson et al., 2016) به کار می روند.

U+41 U+302 U+300 <tab> U+1EA6
U+41 U+302 U+301 <tab> U+1EA4
...

شکل ۲: قانون نرمال سازی سفارشی در قالب TSV

۴.۳ نرمال سازی قابل سفارشی سازی کاراکترها

نرمال سازی کاراکترها یک گام پیش پردازشی مهم برای پردازش متون واقعی است؛ متونی که شامل کاراکترهای یونیکد هم معنی از نظر معنایی هستند. برای مثال، در زبان ژاپنی، کاراکترهای لاتین تمام عرض (fullwidth) می توانند به کاراکترهای لاتین استاندارد (ASCII) نرمال شوند. همچنین تبدیل حروف بزرگ به کوچک (lowercasing) نیز بسته به کاربرد می تواند نوعی نرمال سازی مؤثر باشد. به طور سنتی، نرمال سازی کاراکترها با مجموعه ای از قواعد دستی پیاده سازی می شد. در سال های اخیر، NFKC و NFC (مانند Unicode Normalization Forms) به طور گسترده ای در برنامه های NLP مورد استفاده قرار گرفته اند، چرا که تکرار پذیری بالاتری دارند و به عنوان بخشی از استاندارد یونیکد به خوبی پشتیبانی می شوند. به صورت پیش فرض، SentencePiece متن ورودی را با استفاده از نرمال سازی NFKC یونیکد نرمال می کند. قوانین نرمال سازی با استفاده از پارامتر خط فرمان normalization_rule_name=nfkc به صورت نگاشت رشته به رشته (string-to-string mapping) و با قاعده هی «بیشترین تطابق از چپ» (leftmost longest matching) پیاده سازی شده است. این قوانین نرمال سازی در قالب یک ترانسیدیوسر حالت متناهی (finite state transducer) از نوع Aho–Corasick کامپایل می شوند تا نرمال سازی به شکل کارآمدی انجام شود. همچنین از قوانین نرمال سازی automaton سفارشی پشتیبانی می کند که می توان آن ها را به صورت فایل TSV تعریف کرد. شکل ۲ نمونه ای از چنین فایل TSV را نشان می دهد. در این مثال، دنباله یونیکد [U+41 U+302 U+300] به کاراکتر U+1EA6 تبدیل می شود. زمانی که در تبدیل، ابهامی وجود داشته باشد، طولانی ترین قاعده اعمال می شود. فایل های TSV تعریف شده توسط کاربر با استفاده از پرچم spm_train مشخص می شوند. قواعد خاص هر وظیفه normalization_rule_tsv=<file>- در دستور spm_train تعیین می شوند. قواعد task-specific rules (NFKC) را می توان با گسترش قواعد پیش فرض TSV تعیین کرد؛ این قواعد به صورت فایل

در بسته SentencePiece ارائه می‌شوند.

۵.۳ مدل‌های خودکفا

در سال‌های اخیر، بسیاری از پژوهشگران مدل‌های ترجمه ماشینی عصبی (NMT) از پیش آموزش دیده را برای بهبود قابلیت بازتولید (reproducibility) نتایج آزمایشی خود ارائه کرده‌اند. با این حال، در بیشتر موارد به‌طور دقیق مشخص نمی‌شود که داده‌ها چگونه پیش‌پردازش شده‌اند. (Post, 2018) گزارش کرده است که تفاوت‌های جزئی در طرح‌های پیش‌پردازش می‌توانند باعث تغییرات چشمگیر در امتیاز BLEU شوند. حتی در صورتی که از ابزار Moses استفاده شود، باز هم تضمینی برای بازتولید تنظیمات دقیق وجود ندارد، مگر این که پیکربندی‌های مربوط به Moses (مانند نسخه و پارامترهای خط فرمان) به‌طور کامل مشخص شده باشند. به‌طور دقیق‌تر، حتی نرمال‌سازی NFKC ممکن است نتایج متفاوتی تولید کند، بسته به نسخه Unicode که استفاده می‌شود. در حالت ایده‌آل، تمام قوانین و پارامترهای مربوط به پیش‌پردازش باید درون فایل مدل به صورت یکپارچه و خودکفا (self-contained) گنجانده شوند تا بتوان شرایط آزمایشی یکسانی را بازتولید کرد، مشروط بر اینکه از همان فایل مدل استفاده شود. مدل SentencePiece به‌گونه‌ای طراحی شده است که کاملاً خودکفا (purely self-contained) باشد. فایل مدل نه تنها شامل واژگان (vocabulary) و پارامترهای بخش‌بندی (segmentation parameters) است، بلکه شامل ترنسدیوسر حالت متناهی از پیش کامپایل شده- (pre-) SentencePiece compiled finite state transducer) برای نرمال‌سازی کاراکترها نیز می‌باشد. رفnar فقط توسط فایل مدل تعیین می‌شود و هیچ وابستگی خارجی‌ای ندارد. این طراحی تضمین می‌کند که بتوان بازتولید کامل (perfect) را به دست آورد، و همچنین اجازه می‌دهد تا فایل مدل SentencePiece به عنوان بخشی از مدل reproducibility منتشر شود. علاوه بر این، توسعه‌دهندگان SentencePiece می‌توانند قوانین نرمال‌سازی پیش‌فرض را بهبود دهنند بدون اینکه نگران از بین رفتن سازگاری (breaking existing preprocessing behaviors) با نسخه‌های قبلی باشند. فایل مدل SentencePiece به صورت فرمت دودویی (binary wire format) از نوع Protocol Buffer ذخیره می‌شود که یک سازوکار پلتفرم‌خنثی (platform-neutral) و قابل توسعه برای سریال‌سازی داده‌های ساختاریافته است. Protocol Buffers کمک می‌کنند تا داده‌های ساختاریافته به شکلی امن و قابل اطمینان سریال‌سازی شوند، در حالی که سازگاری با نسخه‌های قبلی (backward compatibility) و قابلیت گسترش (extensibility) نیز حفظ می‌شود.

۶.۳ رابط کتابخانه‌ای برای پردازش در لحظه (On-the-fly Processing)

پیش‌پردازش متن معمولاً به عنوان یک فرایند آفلاین در نظر گرفته می‌شود. پیش از آغاز آموزش اصلی NMT، ورودی خام مورد پیش‌پردازش قرار گرفته و توسط یک پیش‌پردازنده مستقل به یک دنباله از شناسه‌ها (id sequence) تبدیل

می شود. چنین پیش‌پردازش آفلاینی با دو مشکل همراه است. نخست آنکه ابزارهای مستقل مستقیماً در برنامه‌های NMT که با کاربر در تعامل هستند، ادغام نمی‌شوند؛ این در حالی است که چنین برنامه‌هایی نیاز دارند ورودی کاربر را به صورت بلادرنگ (on-the-fly) پیش‌پردازش کنند. دوم آنکه پیش‌پردازش آفلاین، به کارگیری روش‌هایی مانند افزایش داده در سطح زیر جمله‌ای (subsentence level data augmentation) و تزریق نویز (noise injection) (Kudo, 2018) را دشوار می‌کند؛ روش‌هایی که هدف آنها بهبود دقت و مقاوم سازی مدل‌های NMT است. مطالعات متعددی برای تزریق نویز به جملات ورودی از طریق تغییر تصادفی بازنمایی درونی جملات انجام شده است. (Lample et al., 2017) روش تنظیم زیر واژه‌ای (subword regularization) را پیشنهاد می‌کند که در آن بخش‌بندی زیر واژه‌ای به صورت تصادفی در طول آموخته NMT تغییر می‌کند. به طور مستقل، (Artetxe et al., 2017) و (Lample et al., 2017) یک خود رمزگذار حذف نویز (denoising autoencoder) را در زمینه‌ی یادگیری دنباله‌به‌دنباله (sequence-to-sequence learning) پیشنهاد کردند، که در آن ترتیب کلمات جمله‌ی ورودی به صورت تصادفی تغییر می‌یابد و مدل آموخته می‌بیند تا جمله‌ی اصلی را بازسازی کند. اجرای چنین نمونه‌گیری پویا و تزریق نویزی تنها با پردازش آفلاین دشوار است.

```
#include <sentencepiece_processor.h>
#include <sentencepiece_trainer.h>

SentencePieceTrainer::Train(
    "--input=input.txt",
    "--model_prefix=spm",
    "--vocab_size=1000");

SentencePieceProcessor sp;
sp.Load("spm.model");

std::vector<std::string> pieces;
sp.Encode("Hello_world.", &pieces);

std::vector<int> ids;
sp.Encode("Hello_world.", &ids);

std::string text;
sp.Decode({151, 88, 21, 887, 6}, &text);
```

شکل ۳: C++ API usage

```
import sentencepiece as spm

params = ('--input=input.txt',
          '--model_prefix=spm',
          '--vocab_size=1000')
spm.SentencePieceTrainer.Train(params)

sp = spm.SentencePieceProcessor()
sp.Load('spm.model')

print(sp.EncodeAsPieces('Hello_world.'))
print(sp.EncodeAsIds('Hello_world.'))
print(sp.DecodeIds([151, 88, 21, 887, 6]))
```

شکل ۴: Python API usage

نہ تنها یک ابزار خط فرمان مستقل برای پیش‌پردازش آف‌لاین ارائه می‌دهد، بلکه از رابطه‌های برنامه‌نویسی کتابخانه‌ای (API) در زبان‌های C++، Python و TensorFlow برای پردازش بالدرنگ (on-the-fly) نیز پشتیبانی می‌کند که به‌آسانی می‌توان آن‌ها را در چارچوب‌های موجود NMT ادغام کرد. شکل‌های ۳، ۴ و ۵ نمونه‌هایی از نحوه‌ی استفاده از رابطه‌های TensorFlow و Python را نشان می‌دهند. شکل ۶ نمونه‌ای از کد Python برای تنظیم زیرواژه‌ای (subword regularization) را نمایش می‌دهد، که در آن یک دنباله‌ی زیرواژه‌ای بر اساس مدل زبان تک‌واژه‌ای (unigram language model) نمونه‌برداری می‌شود. می‌توان مشاهده کرد که متن «New York» به صورت‌های متفاوتی توکنیزه می‌شود.

```
import tensorflow as tf
import tf_sentencepiece as tfs

model = tf.gfile.GFile('spm.model', 'rb').read()

input_text = tf.placeholder(tf.string, [None])
ids, lens, lens = tfs.encode(input_text, model_proto=model,
                             out_type=tf.int32)
output_text = tfs.decode(ids, lens, model_proto=model)

with tf.Session() as sess:
    text = ['Hello_world.', 'New_York']
    ids_, lens_, output_text_ = sess.run([ids, lens, output_text],
                                         feed_dict={input_text:text})
```

شکل ۵: TensorFlow API usage

مدلproto (فایل مدل model) به عنوان یک ویژگی از عملگر SentencePiece تعریف شده و درون گراف TensorFlow جاسازی می‌شود؛ با این ترتیب مدل و گراف به صورت کاملاً مستقل و خودکفا در می‌آیند.

```
>>> sp.Load('spm.model')
>>> for n in range(5):
...     sp.SampleEncodeAsPieces('New_York', -1, 0.1)
['_', 'N', 'e', 'w', '_York']
['_', 'New', '_York']
['_', 'New', '_Y', 'o', 'r', 'k']
['_', 'New', '_York']
['_', 'New', '_York']
```

شکل ۶: Subword sampling with Python API

در هر فراخوانی از تابع SampleEncodeAsPieces، توکنیزه‌سازی به صورت متفاوت انجام می‌شود. برای جزئیات بیشتر درباره‌ی تنظیم زیرواژه‌ای (subword regularization) و ابرپارامترهای نمونه‌برداری آن، به (Kudo, 2018) مراجعه کنید.

۴ آزمایش‌ها

۱.۴ مقایسه روش‌های مختلف پیش‌پردازش

ما کارایی روش‌های مختلف پیش‌پردازش را در ترجمه‌ی انگلیسی–ژاپنی مقالات Wikipedia، مطابق با وظیفه‌ی ترجمه‌ی آزاد کیوتو (KFTT) Kyoto Free Translation Task - KFTT ارزیابی کردیم. داده‌های آموزش، توسعه و آزمون در مجموعه‌داده‌ی KFTT به ترتیب شامل ۴۴۰ هزار، ۱۱۶۰ و ۱۱۶۰ جمله هستند. ما از GNMT (Wu et al., 2016) بهترین نتیجه را در آزمایش‌های خود استفاده کردیم. به طور کلی، تنظیمات و روش آموزشی توصیف شده در عنوان پیاده‌سازی سیستم NMT در آزمایش‌های خود استفاده کردیم. در مجموعه‌داده‌ی KFTT (Wu et al., 2016) را بهترین LSTM را به ترتیب (با برابر با ۵۱۲ و ۶ در نظر گرفتیم. یک مدل واژه‌ای به عنوان سیستم مبنای مورد استفاده قرار گرفت. ما SentencePiece (با مدل زبان تک‌واژه‌ای یا unigram language model) را در دو حالت با پیش‌تکنیزه‌سازی (pre-tokenization) و بدون آن مقایسه کردیم. نسخه‌ی SentencePiece با پیش‌تکنیزه‌سازی در واقع مشابه پیکربندی رایج NMT با استفاده از subword-nmt است. نسخه‌ی SentencePiece بدون پیش‌تکنیزه‌سازی، مدل زیرواژه‌ای را مستقیماً از جملات خام آموزش می‌دهد و از هیچ منع خارجی استفاده نمی‌کند. برای پیش‌تکنیزه‌سازی زبان انگلیسی و ژاپنی به ترتیب از توکنیزه‌کننده‌ی Moses و KyTea استفاده کردیم. همان توکنیزه‌کننده‌ها نیز برای مدل واژه‌ای به کار گرفته شدند. ما از امتیاز BLEU حساس به حروف بزرگ و کوچک (Papineni et al., 2002) به عنوان معیار ارزیابی استفاده کردیم. از آنجا که جملات خروجی در زبان ژاپنی به صورت بخش‌بندی شده نیستند، پیش از محاسبه امتیاز BLEU آن‌ها را با استفاده از KyTea بخش‌بندی کردیم.

جدول ۱ نتایج آزمایش‌ها را نشان می‌دهد. نخست، همان‌طور که در جدول مشاهده می‌شود، بخش‌بندی زیرواژه‌ای با استفاده از SentencePiece به طور پیوسته امتیازهای BLEU را در مقایسه با مدل واژه‌ای بهبود می‌دهد. این نتیجه با یافته‌های پژوهش‌های پیشین (Sennrich et al., 2016) سازگار است. دوم، می‌توان مشاهده کرد که پیش‌تکنیزه‌سازی (pre-tokenization) همواره برای افزایش امتیاز BLEU ضروری نیست. در جهت ترجمه از ژاپنی به انگلیسی، میزان بهبود ناچیز بوده و تفاوت معناداری وجود ندارد. در مقابل، در ترجمه از انگلیسی به ژاپنی، امتیاز BLEU با وجود پیش‌تکنیزه‌سازی کاهش یافته است.

بهبودهای چشمگیرتری در امتیاز BLEU زمانی مشاهده می‌شود که:

(۱) SentencePiece بر روی زبان ژاپنی اعمال شود

(۲) جمله‌ی هدف به زبان ژاپنی باشد.

از آنجا که زبان ژاپنی زبانی بدون بخش‌بندی است، پیش‌تکنیزه‌سازی (pre-tokenization) به عنوان یک محدودیت قوی برای تعیین واژگان نهایی عمل می‌کند. می‌توان چنین در نظر گرفت که اثرات مثبت بخش‌بندی بدون نظارت از ورودی

خام، در شناسایی واژگان خاص حوزه در زبان ژاپنی به طور مؤثری عمل کرده‌اند.

جفتزبان	تنظیم (مبدأ/مقصد)	تعداد واژگان	امتیاز BLEU
ja→en	مدل واژه‌ای (مبنا مقایسه)	80k/80k	۲۸.۲۴
SentencePiece		8k (shared)	۲۹.۵۵
SentencePiece w/ pre-tok.		8k (shared)	۲۹.۸۵
Word/SentencePiece		80k/8k	۲۷.۲۴
SentencePiece/Word		8k/80k	۲۹.۱۴
en→ja	مدل واژه‌ای (مبنا مقایسه)	80k/80k	۲۰.۰۶
SentencePiece		8k (shared)	۲۱.۶۲
SentencePiece w/ pre-tok.		8k (shared)	۲۰.۸۶
Word/SentencePiece		80k/8k	۲۱.۴۱
SentencePiece/Word		8k/80k	۱۹.۹۴

جدول ۱: نتایج ترجمه (%) (BLEU%)

۲.۴ علمکرد بخش بندی

جدول ۲ خلاصه‌ای از عملکرد آموزش و بخش‌بندی در پیکربندی‌های مختلف را نشان می‌دهد. می‌توان مشاهده کرد که سرعت آموزش و بخش‌بندی هر دو روش subword-nmt و SentencePiece در مجموعه‌داده‌ی انگلیسی تقریباً قابل مقایسه است، صرف نظر از انتخاب پیش‌توكینیزه‌سازی (pre-tokenization). این نتیجه قابل انتظار است، زیرا زبان انگلیسی زبانی بخش‌بندی شده است و فضای جستجو برای استخراج واژگان در آن تا حد زیادی محدود است. در مقابل، بهبودهای قابل توجهی در عملکرد هنگام اعمال بر داده‌های خام زبان ژاپنی (بدون پیش‌توكینیزه‌سازی، SentencePiece (w/o pre-tok) نشان می‌دهد. سرعت بخش‌بندی در SentencePiece در این تنظیمات حدود ۳۸۰ برابر سریع‌تر از subword-nmt است. این نتیجه به طور قوی از این ادعا پشتیبانی می‌کند که SentencePiece به اندازه‌ای سریع است که می‌تواند مستقیماً بر داده‌های خام اعمال شود و در نتیجه، پیش‌توكینیزه‌سازی همیشه ضروری نیست. در نتیجه، SentencePiece به ساخت سامانه‌ای کاملاً داده محور و مستقل از زبان کمک می‌کند. سرعت بخش‌بندی SentencePiece حدود ۲۱ هزار جمله در ثانیه برای زبان انگلیسی و ۷۴ هزار جمله در ثانیه برای زبان ژاپنی است، که برای اجرای بلاذرنگ (on-the-fly) به اندازه‌ی کافی سریع است.

وظیفه	ابزار	پیش توکنیزه‌سازی	Japanese (زمان، ثانیه)	English (زمان، ثانیه)
آموزش	subword-nmt	بله	۵۶.۹	۵۴.۱
بخشنده	SentencePiece	بله	۱۰.۱	۱۶.۸
آموزش	subword-nmt	خیر	۵۲۸.۰	۹۴.۷
آموزش	SentencePiece	خیر	۲۱۷.۳	۲۱.۸
آموزش	subword-nmt	بله	۲۳.۷	۲۸.۶
آموزش	SentencePiece	بله	۸.۲	۲۰.۳
آموزش	subword-nmt	خیر	۲۱۶.۲	۳۶.۱
آموزش	SentencePiece	خیر	۵.۹	۲۰.۳
پیش توکنیزه‌سازی	KyTea(ja)/Moses(en)		۲۴.۶	۱۵.۸

جدول ۲: برای ارزیابی از پیکره‌ی KFTT شامل ۴۴۰ هزار جمله استفاده شده است. آزمایش‌ها در محیط Linux با پردازنده‌های Xeon 3.5GHz انجام شدند. اندازه‌ی واژگان برابر با ۱۶ هزار در نظر گرفته شده است. برای زبان انگلیسی از توکنیزه کننده‌ی Moses و برای زبان ژاپنی از KyTea استفاده شده است. توجه داشته باشید که برای مقایسه‌ی منصفانه بین حالت‌های دارای پیش توکنیزه‌سازی و بدون آن، باید زمان صرف شده برای پیش توکنیزه‌سازی نیز در نظر گرفته شود. از آنجا که subword-nmt بروایه‌ی روش BPE است، ما در نیز از مدل BPE استفاده کردیم. مشاهده کردیم که مدل‌های SentencePiece و مدل زبان تک‌واژه‌ای (unigram) عملکردی تقریباً مشابه دارند.