# Between words and characters: A Brief History of Open-Vocabulary Modeling and Tokenization in NLP

Sabrina J. Mielke  $^{1,2}$  Zaid Alyafeai  $^3$  Elizabeth Salesky  $^1$  Colin Raffel  $^2$  Manan Dey  $^4$  Matthias Gallé  $^5$  Arun Raja  $^6$  Chenglei Si  $^7$  Wilson Y. Lee  $^8$  Benoît Sagot  $^{9*}$  Samson Tan  $^{10*}$ 

BigScience Workshop Tokenization Working Group

#### Abstract

What are the units of text that we want to model? From bytes to multi-word expressions, text can be analyzed and generated at many granularities. Until recently, most natural language processing (NLP) models operated over words, treating those as discrete and atomic tokens, but starting with byte-pair encoding (BPE), subword-based approaches have become dominant in many areas, enabling small vocabularies while still allowing for fast inference. Is the end of the road character-level model or byte-level processing? In this survey, we connect several lines of work from the pre-neural and neural era, by showing how hybrid approaches of words and characters as well as subwordbased approaches based on learned segmentation have been proposed and evaluated. We conclude that there is and likely will never be a silver bullet singular solution for all applications and that thinking seriously about tokenization remains important for many applications.

#### 1 Introduction

"'tokens' are not a real thing. they are a computer generated illusion created by a clever engineer" —@dril\_gpt<sup>1</sup>

When we first introduce people to NLP models, we often take for granted the idea that text is cut up into little pieces that are fed to a computer, eventually as nothing but a sequence of integers. Following Webster and Kit (1992), we call these (usually)

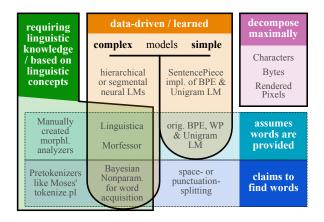


Figure 1: A taxonomy of segmentation and tokenization algorithms and research directions

contiguous substrings *tokens*. In teaching settings and in fact historically in NLP, these tokens are somewhat naturally implied to be *words*—at first, perhaps naively as "space-separated substrings" in English. Understandably, as soon as we try to split punctuation from words, things get a little tricky. Take the word "don't" for example: not splitting it is reasonable, but if we split on all punctuation, we would get three somewhat nonsensical tokens (don't)—and we might argue that the most sensible split actually ought to yield the two units "do" and "n't", as found in the Penn Treebank (Marcus et al., 1993).

This survey deals with such questions of *tokenization* and we will elaborate on fundamental questions and terminology in §2, and show how this important if somewhat unglamorous part of all NLP work has historically been treated (§3). However, since especially in the last five years there has been renewed interest in going *beyond* intuitive definitions of a "token" as a somewhat atomic word-like space-separated unit. One way to do

<sup>&</sup>lt;sup>1</sup>Johns Hopkins University <sup>2</sup>HuggingFace <sup>3</sup>King Fahd University of Petroleum and Minerals <sup>4</sup>SAP <sup>5</sup>Naver Labs Europe <sup>6</sup>Institute for Infocomm Research, A\*STAR Singapore <sup>7</sup>University of Maryland <sup>8</sup>BigScience Workshop <sup>9</sup>Inria Paris <sup>10</sup>Salesforce Research Asia & National University of Singapore sjmesjmielke.com

<sup>\*</sup> Working group chairs

<sup>&</sup>lt;sup>1</sup>A Twitter bot with (human-curated) outputs of a language model based on GPT2 (Radford et al., 2019) and trained on tweets of Twitter poet @dril; https://twitter.com/dril\_gpt2/status/1373596260612067333.

so is to use word-internal information to augment word-like units (§4), which neural modeling has made easier than ever, leading to models that can even *learn* word boundaries with no overt indication (e.g., when spaces are present). That notion of unsupervised word segmentation or discovery has been present on its own in decades of work (§5), but we will find it useful to consider when finally looking at the now prevalent idea of using subword units as atomic tokens (§6). We will close out this survey with a look at some issues of sharing and competition in multilingual vocabularies (§7) and work using the simplest tokenization possible: maximal decomposition into characters, bytes, or even pixels (§8).

Equipped with all this knowledge, we will conclude the survey in §9 by making a case for why "complicated" tokenization is something to practice or even learn about even now in 2021 when the easy solution of bytes seems in reach. We will argue that while recent advances like CANINE (Clark et al., 2021), ByT5 (Xue et al., 2021), or Charformer (Tay et al., 2021) make maximally decomposed processing feasible for certain domains and use-cases, they do not cover the wide variety of NLP scenarios and come with their own drawbacks and biases.

#### 2 Tokens, word-forms, and sub-words

In NLP, textual data has been traditionally segmented into "sentences" (or "utterances", etc.) and "words" due to linguistic motivations and technical constraints. The macroscopic units ("sentences") are often considered independently from one another and themselves segmented into microscopic units. The definition of these microscopic units has always been a matter of approximation and compromise. On the one hand, these units receive linguistic annotations (e.g. part-of-speech tags, morphosyntactic annotation, syntactic dependency information), which would require them to be linguistically motivated units. On the other hand, a large range of phenomena make it highly non-trivial to identify and even to consistently define linguistic units, denoted by the Morphological Annotation Framework (MAF) ISO standard (Clément et al., 2005) as word-forms. Such phenomena include contractions (e.g. English don't, cited above, and French aux 'to thepl'), compounds (e.g. French copier-coller 'copy-paste'<sup>2</sup>), morphological derivatives (e.g. English or French *anti-Trump*), as well as numerous classes of named entities and other sequences following type-specific grammars (e.g. numbers, URLs).

As a result, typographic units, generally called tokens have been used as an approximation for such linguistically motivated units. For instance, MAF defines a token as a "non-empty contiguous sequence of graphemes or phonemes in a document." In the case of writing systems using a typographic separator such as the whitespace, now universally used with the Latin script for instance, tokens have been widely used and broadly defined as either contiguous sequences of non-punctuation non-whitespace marks or punctuation marks. Provided a few arbitrary decisions are made regarding certain punctuation marks (e.g. the hyphen or the apostrophe), such a definition makes it possible to deterministically split a sentence into atomic units, resulting in a segmentation into tokens that are acceptable approximations of word-forms. Crucially, as discussed in detail by Clément et al. (2005), Sagot and Boullier (2008) and elsewhere, there is no one-to-one correspondence between tokens and word-forms; a word-form can be made of several tokens (e.g. French or English sine die) whereas several word-forms can be represented by the same token (e.g. English don't = do + not, Spanish dam'elo = da + me + lo). This is what the Universal Dependencies guidelines<sup>3</sup> refer to as "multitoken words" and "multiword tokens," respectively, a topic further discussed by More et al. (2018). In fact, both phenomena can interfere in non trivial ways (e.g. French à l'instar  $du = a_l$ 'instar\_de +le).<sup>4</sup>

In recent years, the spread of approaches based on neural language models resulted in an evolution in how sentences are split into atomic units, thereby resulting in a redefinition of the notion of tokenization. Indeed, based both on scientific results (e.g. the impact of sub-word segmentation on machine translation performance (Sennrich et al., 2016)) and on technical requirements (e.g. language models such as BERT that require a fixed-size vocabulary), the need for the atomic processing units (still called tokens) to be an approximation of

<sup>&</sup>lt;sup>2</sup>Cf. inflected forms *copié-collé* 'copy-pasted', lit. 'copied-pasted', but *copie-collerai* 'will<sub>1sg</sub> copy-paste'.

<sup>&</sup>lt;sup>3</sup>https://universaldependencies.org/u/ overview/tokenization.html

<sup>&</sup>lt;sup>4</sup>When the writing system at hand does not have a typographic separator, tokens must be defined differently. With scripts like the Chinese or Japanese script, an option for instance is to consider each character as a token on its own.

word-forms has faded out. As a result, in current NLP, the notion of token still perfectly matches its MAF definition, but it no longer corresponds to the traditional definition of a typographic unit. "Tokenization" now denotes the task of segmenting a sentence into such non-typographically (and indeed non-linguistically) motivated units, which are often smaller than classical tokens and wordforms, and therefore often called sub-words. Typographic units (the "old" tokens) are now often called "pre-tokens," and what used to be called "tokenization" is therefore called nowadays "pretokenization." This term is motivated by the fact that the first approaches to the new notion of "tokenization" often involved segmenting sentences into proper typographic units (i.e. the "old" notion of tokenization) before further segmenting (some of) the resulting units (formerly "tokens", now "pretokens") into "sub-words".

# 3 Pre-tokenization yields word-like typographic units

As a compromise between the linguistic irrelevance of purely typographic tokens and the difficulty of automatically splitting a text into linguistically motivated word-forms, units that are halfway between purely typographic tokens and purely linguistic word-forms have been widely used,<sup>5</sup> albeit often (improperly) denoted by the term "token" before the spread of sub-word tokenization, and "pretoken" since then. Many tools, formerly known as "tokenizers" and nowadays as "pre-tokenizers", have been developed and used for a long time. Some of them are relatively simple and remain faithful to the typographic token. Amongst the most widely used, we can cite the venerable Moses (Koehn et al., 2007) tokenizer<sup>6</sup> and the more recent pre-tokenizers package in Hugging Face's Tokenizers package.<sup>7</sup> This practice resulted in the word "tokenization," now "pre-tokenization," ending up denoting the task of segmenting sentences into atomic units in general, i.e "basic units which need not be decomposed in a subsequent processing" (Webster and Kit, 1992), even when such units are closer to word-forms than to typographic units.

Moreover, it is fairly common for tokenizers to not only segment sentences but also modify the raw text, for instance for normalization, spelling correction or named entity detection purposes, thereby departing from the standard definition of token. Thus a string like "some 'quoted' text" might be tokenized into five units: "some " quoted " text." On the other hand, a number of tools have been developed and used that attempt a joint segmentation of sentences into tokens (which remain proper substrings of the input text) and word-forms (which can be the output of normalization, spelling correction and named entity recognition steps). Some of these tools take into account the inherent nondeterminism of the mapping between the two types of units (Sagot and Boullier, 2008). Normalization operations like these or conflation and merging of different whitespace symbols leads to most tokenizers being irreversible, i.e., we cannot recover the raw text definitively from the tokenized output.<sup>8,9</sup>

### 4 Augmenting word-level pretokenizer tokens with character information

While word-level models are conceptually easy to understand and in the neural era (Bengio et al., 2001; Collobert and Weston, 2008; Mikolov et al., 2013) offer features at an interpretable granularity, their central weakness is the inability to deal with rare and novel words, i.e., words that were seen very rarely during training or not even at all (outof-vocabulary, OOV)—they are closed-vocabulary models. In particular, historically rare word types were replaced with a new word type UNK (unknown) at training time; at test time, any token that was not part of the model's vocabulary could then be replaced by UNK. That approach however comes with a number of drawbacks: 1) UNKs are not acceptable when performing natural language generation (NLG), 2) they do not allow us to extract features for novel words that are useful anchors of meaning and not just one-off events (Church, 2000) when used in large-scale models like ELMo (Peters et al., 2018) or BERT (Devlin et al., 2018), and

<sup>&</sup>lt;sup>5</sup>Cf. the Penn TreeBank (Marcus et al., 1993), where don't is split into two units do and n't, as mentioned above.

<sup>6</sup>https://github.com/moses-smt/
mosesdecoder/blob/master/scripts/
tokenizer/tokenizer.perl

<sup>7</sup>https://github.com/huggingface/ tokenizers

<sup>&</sup>lt;sup>8</sup>While this is not a major issue for most applications, it means that we no longer model the original text, but a string that may correspond to many original strings, inflating probability estimates of language models; this issue is also highlighted in the context of ambiguous tokenization (see §6.4.3) by Cao and Rimell (2021).

<sup>&</sup>lt;sup>9</sup>The "reversible" language-agnostic tokenizer of Mielke and Eisner (2018) attempts to remedy some of these issues, but still conflates whitespace.

3) in languages other than English, in particular those with more productive morphology and thus higher type-token-ratio, removing rare words is infeasible (Cotterell et al., 2018; Mielke et al., 2019). Nevertheless, since the word *is* a fundamental unit of language, a number of approaches emerged to improve handling of rare and novel words under a fundamentally word-based framework by basing their handling on the characters that make up a word. We will present some of these approaches in this section, for a more comprehensive treatment of word representation, Pinter (2021) surveys linguistic background and multiple approaches.

### 4.1 Augmenting word-level models with spelling information

The idea of somehow using information about the spellings of a word to inform the word's representations of course is decades old. In neural models for language, research in the 90s and 2000s often forewent the focus on words altogether and processed strings of characters instead (see §8 and §4.2), but as soon as neural models became important in NLP, combinations of word- and character-level information for use in neural networks emerged there, too.

Dos Santos and Zadrozny (2014) first proposed to use information about the words themselves to aid word embedding estimation. Soon thereafter Ling et al. (2015), Kim et al. (2016), and Jozefowicz et al. (2016) popularized the idea of deterministically constructing a word's embedding from its spelling, <sup>10</sup> both for textual input as well as for generative language modeling, that is, prediction of strings. However, even when replacing embedding matrices with convolutional neural network (CNN) layers, their generative models are still closedvocabulary, meaning they can only predict words that were seen (often enough) in training data, so the CNN construction only helps with rare words, not novel words. Furthermore, constructing embeddings from spellings for each token (as opposed to every type like Mielke and Eisner (2018), see §4.2) implicitly trains the CNN-powered embedding function to "get frequent words right" instead of anticipating novel words, an issue discussed in Mielke and Eisner (2018). Similar constructions

led to advances in other classic NLP tasks like POS tagging (Plank et al., 2016) and ultimately powered the first big *contextual word embedding* model ELMo (Peters et al., 2018).

The popular fastText embeddings (Bojanowski et al., 2017) propose constructing word embeddings not from characters, but from overlapping *n*-grams, allowing one to obtain embeddings for novel words (making it "open-vocabulary" in that sense, though not in the generative sense). Ataman and Federico (2018a) likewise obtain better performance on machine translation by using (overlapping) *n*-grams instead of characters (also beating BPE on morphologically rich languages).

In more recent times, El Boukkouri et al. (2020, CharacterBERT) and Ma et al. (2020, CharBERT) use the same CNN construction as in Kim et al. (2016) on a modern BERT-style model, this time enhancing the BPE units' embedding with their constituent characters' embedding, motivated by better handling noisy texts with spelling errors or transferring to new domains like medical text; concurrently, Aguilar et al. (2021) do almost the same, but using a small Transformer instead of CNNs.

Finally, construction-based approaches have also been integrated into pretrained word-level input models. Specifically, Pinter et al. (2017) learn a model that is trained to *mimic* the embedding of a word given its spelling using a helper RNN model that is called whenever an unknown word appears during test time.

# 4.2 Open-vocabulary language modeling with (tokenizer-defined) words made of characters

Extending closed-vocabulary generative models to open-vocabulary models, i.e., those that can predict and generate novel words at test time, is somewhat more difficult than being open-vocabulary on the *input* side because it must be possible to hold out probability mass for the infinite set of sentences that contain completely novel words.

Inspired by Luong and Manning (2016), Mielke and Eisner (2018) propose a probabilistic two-stage model that essentially augments the ordinary closed-vocab word-level recurrent neural network language model (RNNLM) setup by regularizing word embeddings to be predictive of their spellings using a smaller character-level RNNLM and using that smaller model to generate novel words on the fly whenever the word-level RNNLM pre-

<sup>&</sup>lt;sup>10</sup>It should be noted that Jozefowicz et al. (2016) also propose a variant in which output tokens are not scored through a softmax, but generated character by character, anticipating the advancements described in §4.2, but still staying in a closed-vocabulary setup.

dicts UNK, yielding an open-vocabulary model motivated by linguistic notions and intuitive modeling and proven successful qualitatively and quantitatively.

Independently developed, the model of Kawakami et al. (2017) follows a similar two-level setup of word- and character-level RNN, but where each word has to be spelled out using a character-level RNN if it cannot be directly copied from the recent past using a cache model (Grave et al., 2016). 11 Their analysis shows clearly that the cache model not only copies "bursty" unknown words like Noriega (Church, 2000), but also extremely common function words like the in an attempt to keep itself from forgetting them. The idea is picked up by Ataman et al. (2019) for a machine translation decoder (creating word embeddings on the encoder side from character-level BiRNNs as in ELMo (Peters et al., 2018, see §4.1)) and later extended by Ataman et al. (2020) with some additional stochasticity that is intended to pick up on lemmata and inflections unsupervisedly.

A different approach is having higher layers of multi-layer RNNs run at lower speed (skipping updates to the hidden state) This is an old idea, first present in El Hihi and Bengio (1995) (building on Schmidhuber (1991, 1992)'s "neural sequence chunker") and revived in Koutnik et al. (2014) for fixed-frequency skipping and Hwang and Sung (2017) for skipping on word boundaries (which are assumed to be observed). This approach leads to the first of a number of ways in which we can actually *learn* word boundaries and thus segmentations.

# 5 Learning segmentations to find concatenative word-like pretokenizer tokens

So far we have relied on having a predefined notion of word (or pretokenization output) despite the conceptual struggles outlined in §2. But what if such a definition is not given, not obtainable, or simply not desirable (for reasons of robustness and in languages other than English etc.)? Is there a way to let our data-driven machine learning approach also *learn* the tokenization? Most approaches described in this section propose to tackle tokenization by treating the implied *segmentation* as a latent variable (with an exponentially-sized domain) on which we can perform approximate or (using more assumptions) exact inference to find segments and boundaries that hopefully correspond to meaningful units. The various techniques described in this section yield units of varying size and quality.

### 5.1 Character-level neural models that learn to skip steps at higher levels

Already in the 90s, Elman (1990) manually analyzed character-level RNNs and correlated their prediction surprisal with word boundaries. This idea that was then expanded on in Schmidhuber (1991, 1992)'s "neural sequence chunker". More recently, surprisal was applied to not only character-level neural models but also n-gram models under a beam search framework by Doval and Gómez-Rodríguez (2019) to split microblog texts in which spaces are deleted.

Instead of using post-hoc surprisal thresholding, the HM-RNN (Chung et al., 2017) takes the idea of multiple timescales motivated in §4.2, but learns the binary decision to skip or update (thereby providing a sense of word boundaries), optimizing with approximate gradient descent using the straight-through estimator (Bengio et al., 2013). In their model, communication between layers happens bidirectionally: the lower network reports its final state to the higher one; that higher network reports its new state to the lower layer that then proceeds to run by itself and so on. While they "recover" word boundaries when including spaces in their data, Kawakami et al. (2019) claim to get unusable segments with this model when not including spaces. Furthermore, when trying to use the HM-RNN for NMT, Cherry et al. (2018) report that it took a lot of fixing to get it to train at all; its performance on the task was competitive but not superior. This finding corroborates that of Kádár et al. (2018), who dedicate a paper to trying to get the HM-RNN to train well, ablating it, and also showing subpar segmentations on text data (as well as the worrying inability to reach the original re-

<sup>&</sup>lt;sup>11</sup>As mentioned before, the idea of spelling out words in isolation from hidden states had previously proven unsuccessful in Jozefowicz et al. (2016)'s comparison, but this was in a closed-vocab setup and without the caching mechanism Kawakami et al. (2017) employ.

<sup>&</sup>lt;sup>12</sup>Specifically, Hwang and Sung (2017) describe an architecture in which character-level and word-level models run in parallel from left to right and send vector-valued messages to each other. The word model sends its hidden state to the character model, which generates the next word, one character at a time, and then sends its hidden state back to update the state of the word model.

ported numbers). Kreutzer and Sokolov (2018) try to use a similar paradigm of skipping steps and generating summaries with lower layers for NMT and find (similarly to Kádár et al. (2018)) that skipping is rarely used and thus seems to be unnecessary for good performance. Nevertheless, the model is extended to phrase- and sentence-level boundaries by Luo and Zhu (2021).

It is worth pointing out that despite having coarser layers of computation, these models still have to "spell out" a word every time it is generated, i.e., they cannot *memoize* tokens as reusable units.

### 5.2 Marginalization over all possible segmentations

Finally, a conceptually straightforward approach is to treat the segmentation of a string as a latent variable that needs to be marginalized over both at training and test time. This essentially means having a vocabulary that contains strings of differing lengths that overlap, i.e., it may contain "cat," "at," and "foster cat," such that the string "my foster cat" can be decomposed a number of ways corresponding to different sequences of latent units. As the number of segmentations is exponential in the sequence or context length, we need to either resort to approximations for marginalizing over latent decompositions (§5.2.1) or simplify the model with independence assumptions e.g. by using an *n*-gram model (§5.2.2).

### 5.2.1 Approximate marginalization

Chan et al. (2017) propose an estimator to approximate the marginal probability of observations using approximate MAP inference through beam search. They find that the model is very hard to train, but manage to obtain promising results. Buckman and Neubig (2018) confirm this model's instability and propose some approximate inference schemes based on averaging RNN hidden states that produce better results in terms of LM perplexity. Hiraoka et al. (2020) implement a similar model, based on a Unigram LM tokenization proposal distribution (see  $\S6.4.3$ ), whose *n*-best tokenizations of a sentence are fed into any sentence encoder model independently and whose resulting sentence embeddings are averaged in line with their a priori tokenization likelihood. Hiraoka et al. (2021) extend this model to sequence-to-sequence settings by training a tokenizer and downstream model with separate losses, the former by rewarding tokenizations that produced a low downstream loss, and the latter using just one tokenization sampled from the conditioned (and tempered) LM.

# 5.2.2 Exact marginalization using additional independence assumptions: segmental neural language models

The more popular solution of segmental neural language models was pioneered by Kong et al. (2016), who cast the problem of segmentation as a monotonic<sup>13</sup> seq2seq task, going from characters to a covering sequence of substrings, i.e., a segmentation. By conditioning segment prediction on the entire raw string, processed and embedded using a BiRNN, segmentation decisions/scores can use context, but by scoring every individual possible substring independently as a segment using these embeddings and then adding up individual scores to score entire segmentations, they can find a covering of the entire input string with segments efficiently using dynamic programming. The reason for this ability is the central independence assumption: the model does not depend on any other segments when scoring a segment, but merely on surrounding characters. Wang et al. (2017) extend this by also having a per-segment RNN over characters for the outputs that can run without knowing the segmentation and whose past representations can thus be used by the individual segment generation processes, allowing for left-to-right sharing of information about segments without breaking dynamic programming.

The jump to LMing is now made simply by omitting the conditioning on an input, yielding the model of Sun and Deng (2018), who coin the term segmental language model, training on Chinese characters and using the unsupervisedly learned segments to compete on Chinese Word Segmentation. To keep the computation of the quadratic number of segments feasible, they restrict the segments to a maximum length of 4 characters (a sensible prior for Chinese). Grave et al. (2019) make the same jump independently, using Transformers as the independent character-level global backbone. When evaluating on English open-vocabulary language modeling, Grave et al. (2019) notice improved perplexity, but not using or evaluating the obtained segmentation, most likely because they, too, only use 4-grams that appear at least 400

<sup>&</sup>lt;sup>13</sup>Interestingly, with some reordering of the input one can break monotonicity between input and output, making the model similar to phrase-based MT (Huang et al., 2018).

times. Contemporaneously, Kawakami et al. (2019) use the same independence idea, but have emissions of string segments come from a context-dependent *mixture* of a character-level model like in Kawakami et al. (2017) (see §4.2) and a large set of substrings (up to 10-grams that appear often enough in training data) with learned embeddings. They evaluate not only on perplexity, but also on word segmentation performance, where they do beat some baselines (see §5.3), but still perform much worse than some previous models, <sup>14</sup> which they argue tuned their hyperparameters on segmentation performance instead of marginal likelihood and thus have an unfair advantage.

Interestingly, when training on image captions, Kawakami et al. (2019) find that both perplexity and segmentation performance improve when the model also has access to the image that is being described, showing that learning segmentation only from monolingual and unimodal text may be harder than when other modalities or languages are present. This observation is shared by He et al. (2020), who build a similar segmental model (in their case, a Transformer-based version that still follows the character backbone idea to allow for dynamic programming) as the target-side generator in an NMT system and use it not as the final model, but merely as a learned tokenizer. This is easy to achieve by changing the dynamic program from marginalization to maximization and thus obtaining a new segmentation, called DPE, that can be used in place of BPE or unigram LM (see §6.4). He et al. (2020) proceed to show that learning to tokenize with a small Transformer-based NMT model<sup>15</sup> produces better segmentations than BPE for use in a bigger model; in particular, training the tokenizing model on the translation task produces different segmentations depending on the source language, and, more importantly, better segmentations (as measured through downstream translation performance) than training on target-side language modeling alone.

The idea of conditioning on characters and pre-

dicting segments is extended to the adirectional masked language modeling setting found in Transformers and left-to-right autoregressive Transformers by Downey et al. (2021), though results do not outperform RNN-based SNLMs consistently.

Note that many of these models can also be seen as relatives of models based on UnigramLM, which we will cover in §6.4.3.

# 5.3 Finding words through Bayesian non-parametrics

In the era of *n*-gram and word-based language models, MacKay and Peto (1995) noticed that a Bayesian view of autoregressive language models may prove beneficial, reinterpreting smoothing and backoff in n-gram models as inference in a hierarchical model where higher-order distributions are drawn from a Dirichlet distribution whose mean is a lower-order distributions. Teh (2006) extends this thinking, proposing a hierarchical PYP language model where we again have *n*-gram distributions of arbitarily large orders, drawn through a hierarchy of PYP distributions that lead to a model that still bears resemblance to *n*-gram language model smoothing, but offers a principled way to forego the choice of n. The pinnacle of this idea of modeling was reached in Wood et al. (2011)'s sequence memoizer, which boasted great compression performance for arbitrary binary data and still performed very well on language modeling tasks, although neural models at this time already proved to be strong competitors.

At the same time, Goldwater et al. (2006b) extended this Bayesian perspective to also explain how new words are first coined and how they are then used in running text: a process they call twostage language modeling (see §4.2), with the two stages being referred to as generator (which creates new lexemes) and adaptor (which governs reuse; here, a Pitman-Yor Process (PYP)), relating the resulting interaction between types and tokens to interpolated Kneser-Ney smoothing as presented in Chen and Goodman (1999).<sup>16</sup> Given such a twostage model to explain text and the use of Bayesian nonparametrics that can assign positive probability to an infinite number of possible lexemes, it becomes possible to also try to infer word boundaries, that is to perform unsupervised word segmentation.

<sup>&</sup>lt;sup>14</sup>On English they cite the pre-neural models of Johnson and Goldwater (2009) and Berg-Kirkpatrick et al. (2010) as significantly better; on Chinese, they are beaten by pre-neural models like the one of Mochihashi et al. (2009) and the neural model of Sun and Deng (2018). More information about some pre-neural models is given in §5.3.

<sup>&</sup>lt;sup>15</sup>Unlike previously mentioned papers, they however restrict the vocabulary to units of an input BPE vocabulary instead of using length and frequency heuristics.

<sup>&</sup>lt;sup>16</sup>The formalism of generators and adaptors is extended and formally specified under the name *adaptor grammars* in Johnson et al. (2007) and used very successfully for state-of-the-art word segmentation in Johnson and Goldwater (2009).

Motivated more by trying to explain and model cognitive processes and in particular child language acquisition, Goldwater et al. (2009)<sup>17</sup> summarize Unigram and Bigram Dirichlet Processes (DPs) for segmenting transcribed infant-directed speech, showing superiority over older non-Bayesian approaches. Mochihashi et al. (2009) extend the idea from bigram DPs to ∞-gram nested/hierarchical PYPs to improve word segmentation for English written text; Elsner et al. (2013) additionally model phonotactic processes that convert a sequence of segments into observed realizations.

### 5.4 Related task: Unsupervised Chinese Word Segmentation

Word segmentation for languages without whitespace delimiters such as Chinese, Japanese and Vietnamese (Shao et al., 2018) is an important area of research and can be notoriously tricky.

In Chinese word segmentation (CWS), there is growing interest in exploring unsupervised word segmentation involving neural language models. Traditionally, popular unsupervised approaches take on two primary paths: 1) discriminative models and 2) generative models. Discriminative models rely on goodness measures for candidate segmentation. These statistical measures incude Mutual Information (MI), normalized Variation of Branching Entropy (nVBE) and Minimum Description Length (MDL), etc., see §6.3. Generative models focus on designing statistical models to find the optimal segmentation of the highest generative probability. These models include Hidden Markov Model (HMM), Hierarchical Dirichlet Process (HDP), Nested Pitman-Yor Process (NPY), etc., see §5.3. It is trivial to extend discriminative approaches by replacing *n*-gram language model with neural language models. For generative approaches, previous work has shown that constructing a neural language model with a context encoder and a segment decoder achieves competitive performance to its statistical counterparts (Sun and Deng, 2018, see previous subsection §5.2.2).

# 6 Learning subword vocabularies and segmentations

As teased in §3, *subword* units allow for a smooth transition between a word-level model and a

character-level model: split the word-like tokens obtained by pre-tokenization into smaller units: the set of all possible subword units is finite and can be determined from training data, but it is assumed to include all characters (or bytes, see §8.3) that appear at test time, making it possible to explain any novel word in held-out data.

While thinking about subword information may have more tradition for processing morphologically rich languages, Mikolov et al. (2012) already proposed using *subword units*<sup>18</sup> instead of words for language modeling English to avoid out-of-vocabulary (*OOV*) issues. Since Sennrich et al. (2016), however, it has become customary in many if not most current-day NLP models to combat large and infinite vocabulary size.

What then should we use as subword units? One option are manually constructed, linguistically informed rule-based systems (§6.1), another is given by data-driven segmentation learners, which traditionally have been motivated and evaluated either linguistically (§6.3) or given by simple heuristics to be fast and easy and improve downstream performance (§6.4).

It is important to point out that despite reasonable motivation, segmentation may be a bad idea in for example Semitic languages like Arabic and Hebrew (Shapiro and Duh, 2018) or other languages with non-concatenative morphological phenomena, which Amrhein and Sennrich (2021) claim are better served by character-level models or those with very small subword inventories.

#### 6.1 Manually constructed linguistic analyzers

Morphological analysis is of great importance for morphologically rich languages and various tools have been developed to analyze word forms into their lemmata and inflections, earliest and most famous of them the Porter stemmer (Porter, 1980) for English. These tools are often constructed manually by linguists using finite-state tools (FST; Beesley and Karttunen, 2003) as often morphological processes lend themselves to systematic description, making it faster and cheaper to manually construct finite-state analyzers than to try to learn complicated data-driven models (Beemer et al., 2020).

Interestingly, such finite-state tools can not only be used for overt segmentation or discovery of lemmata and other subword units, but even a morpho-

<sup>&</sup>lt;sup>17</sup>The idea and partial results are already presented in Goldwater et al. (2006a), but the authors request citing the updated 2009 paper. Goldwater et al. (2011) summarized this thread of research.

<sup>&</sup>lt;sup>18</sup>They don't supply many details, but claim that the units they use are *syllables*—and that they help.

logical tagger can be used to induce segmentations, as finite-state machines allow for easy tracking of what parts of an input string resulted in output choices, allowing one to identify for example affixes.

It is worth pointing out that such analysis, yes, relies on manual annotation, but beyond that is often considered slow and needlessly complicated. Nevertheless combinations of lemmatization of tagging have been used successfully to tackle large and potentially noisy vocabularies, for example by Tan et al. (2020)'s BITE, which converts inflected forms into lemma and tag to protect against noise and improve on dialectal data.

An important difference to the rest of this survey is that such an approach has the potential to be stronger even, as foregoing purely concatenative segmentation allows one to "segment" for example the word "hoping" as "hope V.PTCP;PRS" or "ate" as "eat PST," allowing sharing of information with other forms in the respective paradigm. The benefit of such an approach is also shown by Hofmann et al. (2021), who observe that undoing derivational processes by splitting words into morphemes before tokenizing can improve sentiment and topicality classification results.

#### **6.2** Other Language-Specific Methods

German, where compounds are never separated with spaces, has prompted research into compound splitting (Koehn and Knight, 2003; Braschler and Ripplinger, 2004; Macherey et al., 2011). Another tricky example is Sanskrit, where segmentation is complicated by the fact that there are processes that occur at / cross word boundaries (Krishna et al., 2017; Huet, 2003). More recently, in the era of neural and subword-based models, questions of tokenization have most recently been researched for Arabic, where Alyafeai et al. (2021) examined various language-agnostic and language-specific tokenizers and find that the performance varies depending on factors like the size and morphological complexity of the datasets. For Chinese, Si et al. (2021) converted characters into stroke orders or romanization sequences before applying BPE in order to capture potential sub-character information based on glyph or pronunciation. Park et al. (2020) shows that a hybrid approach of morphological segmentation followed by BPE (§6.4) works best for most Korean datasets.

### 6.3 Unsupervised morphological segmentation

While subword representations are now commonly evaluated based on their use for a downstream application, initial work in this space often directly assessed the *linguistic* validity of subword<sup>19</sup> segmentations by means of databases such as CELEX (Baayen et al., 1995) or annotations from linguistic grammars and analyzers.

In both computational models of corpora and speakers, it has been found that "distributional regularity and phonotactic constraints are useful for segmentation" (Brent and Cartwright, 1996). de Marcken (1996) proposed to deconstruct text recursively from sentences over words and morphs into characters through "composition and perturbation," presenting results towards recovering linguistically plausible words. Brent et al. (1995) proposed an essentially minimum description length (MDL; Rissanen, 1989) based approach to morphological segmentation, in which the sum of potential vocabulary units' length and the length of text encoded with this vocabulary is minimized. MDL-based approaches were prolifically adapted and expanded for unsupervised morphological segmentation, as in Linguistica (Goldsmith, 2001), and found to generate segmentations with high correlations to morpheme boundaries on English and Romance languages (Baroni, 2000; Goldsmith, 2001). Initially, these approaches were only lightly guided by additional information about possible morphological structure or paradigms—partitioning word types into sets of stems with either suffixes<sup>20</sup> or prefixes, they could not recursively split morphs into additional subwords or account for conditioned character changes—and so with only one morpheme boundary they were most appropriate only for the languages on which they were initially tested.

The use of morphological 'categories' and additional structure within segmentation models expanded their recall and applicability. The *Morfessor* family<sup>21</sup> comprises several unsupervised and semi-supervised segmentation models which aimed to incorporate linguistic biases to improve initial naïve MDL models. *Morfessor 1.0* (Creutz and Lagus, 2002), later called the *Morfessor Baseline*,

<sup>&</sup>lt;sup>19</sup>The resulting units are often termed 'morphs' in such settings, representing the surface forms of morphemes.

<sup>&</sup>lt;sup>20</sup>Termed 'signatures' by Goldsmith (2001)

<sup>&</sup>lt;sup>21</sup>A Python implementation of the Morfessor algorithms is provided by Smit et al. (2014).

is a recursive MDL model based on unigram morph frequencies and lengths; without additional structure it has a clear tendency to over-segment and create spurious splits such as 's + plit.' Morfessor CatMAP (Creutz and Lagus, 2005), or categories maximum-a-posteriori, added a hierarchical HMM-based model of the sequential nature of loose morphological categories (prefixes, stems, and suffixes), where priors could be learned in a semi-supervised way from wordlists; this model remained ideal for the concatenative morphology found in such languages as evaluated in the Morpho Challenge<sup>22</sup>—English, Finnish, Turkish, German (Kurimo et al., 2010). The FlatCat model (Grönroos et al., 2014) flattened this structure, which reduced accuracy under unsupervised conditions but simplified and improved semi-supervised learning. The most recent Morfessor model, EM+Prune (Grönroos et al., 2020) merges the above tradition with recent techniques, leading to a model that is a strict generalization of Unigram LM (Kudo, 2018, see §6.4.3).

Many of the approaches to morphological segmentation implicitly treated the task as *paradigm* learning by incorporating the notion of morphological paradigms and inflection classes. Within this perspective, one research arch focused on expanding the limited paradigmatic structure in early MDL models either through explicit rules, clustering, or 'chains' (Snover and Brent, 2002; Creutz and Lagus, 2005; Monson et al., 2007, 2009; Lignos, 2010; Narasimhan et al., 2015). Another focused on improving segmentation by discovering forms with shared paradigms, by inducing morphological relatedness across surface forms and further allowing for e.g., spelling changes to improve discovery of shared structure across irregular forms (Schone and Jurafsky, 2001; Snover and Brent, 2001; Yarowsky and Wicentowski, 2000; Bergmanis and Goldwater, 2017). While segmentation from this perspective can result in more compact corpus representations, with higher segmentation recall and greater corpus-level consistency, their precision is often lower than with e.g., Mor*fessor* or frequency-driven techniques. As briefly discussed in §6.5, use of morphological segmentation as tokenization for downstream tasks provides only inconsistent improvements compared to the

lighter-weight techniques of §6.4, with recent work predominantly electing for the simplicity of these approaches.

### **6.4** Modern fast subword segmentation algorithms

As explained earlier, the breakthrough for subword tokenization nowadays considered central was the use of Byte-Pair-Encoding (BPE; Gage, 1994) by Sennrich et al. (2016) for machine translation.<sup>23</sup>

### 6.4.1 BPE (Gage, 1994; Sennrich et al., 2016)

BPE is a compression algorithm from a family termed "macro-schemas" (Storer and Szymanski, 1982) in which substrings are replaced with references to them. The name was coined in Gage (1994), although equivalent algorithms have been applied for pattern discovery in natural language (Wolff, 1975) and complexity of genetic sequences (Ángel Jiménez-Montaño, 1984) earlier.<sup>24</sup> When learning a tokenization, BPE replaces pairs of adjacent symbols with a new symbol representing that pair, iteratively merging all occurrences of the pair that occurs most often at any given time. At test time, the same procedure of merging can be performed by executing all recorded merges in the order in which they were conducted during training of the tokenization model.

Byte-Level BPE (Wang et al., 2019) applies BPE not to characters, but raw bytes (see §8.3); it is used in GPT-2 (Radford et al., 2019) and other models. BPE-dropout (Provilkov et al., 2020) is an extension allowing for subword regularization (see §6.4.3.

# 6.4.2 WordPiece (Schuster and Nakajima, 2012)

A very similar technique had been proposed under the name "WordPiece" by Schuster and Nakajima (2012) for Japanese and Korean text (where reliance on space-separated tokens is impossible as text is written without spaces), though it is also used in BERT (Devlin et al., 2018) and other models. Unlike BPE, WordPiece doesn't merge the

<sup>&</sup>lt;sup>22</sup>An evaluation campaign from 2005-10 which focused on unsupervised morphological segmentation; see Virpioja et al. (2011) for evaluation methodology.

<sup>&</sup>lt;sup>23</sup>In language modeling, van Merriënboer et al. (2017) were the first to apply BPE to language modeling and Mielke and Eisner (2018) show that a BPE-based baseline beat all state-of-the-art and even their proposed model on some languages, but the idea didn't really take off until really put to the test by state-of-the-art models like the GPT models (Radford et al., 2018) and BERT (Devlin et al., 2018).

<sup>&</sup>lt;sup>24</sup>See Gallé (2019) for more historical connection and corresponding analyses, e.g., its linear-time implementation by Larsson and Moffat (2000).

most often co-occuring pair but pairs that increase the likelihood that an *n*-gram based language model trained with this updated vocabulary reaches on data (the fact that only some counts need to be updated in such a model and the use of frequency-based heuristics for selecting and batching of multiple merge candidates keep the process computationally feasible). To segment text, WordPiece follows a per-word left-to-right longest-match-first strategy, allowing for very fast linear-time processing (Song et al., 2021).

### **6.4.3** UnigramLM (Kudo, 2018)

Kudo (2018) picks up on the idea of judging subword candidates by evaluating their use in a language model, but it uses a simple unigram language model (hence calling the algorithm *unigram LM*) and iteratively *removes* subword units from a starting vocabulary that contains far more subword units than are desired: on every iteration, the unigram LM is trained using EM and then the lowest-probability items are pruned from the vocabulary—the process is repeated a few times until a desired vocabulary size has been reached.

Interestingly, this probabilistic setup also cleanly models the fact that there are many possible segmentations that are consistent with a given string (in the extreme case, one could always fall back to characters). They report that training with *sampled* segmentation (termed "subword regularization") instead of just using one deterministic segmentation indeed improves machine translation performance. The same motivation led Provilkov et al. (2020) to propose BPE-dropout where the skipping of individual merges in the BPE construction leads to variety in segmentations. Subword regularization not only has been shown to help in monolingual in-domain tasks, but also for improving transfer in multilingual models, see §7.

The observation that sampling segmentation helps is confirmed by Hiraoka et al. (2019), who employ a Bayesian nonparametric model (see §5.3) as the LM that defines the tokenization.

Wang et al. (2021a) build a similar model where the unigram LM is based on character-level BiL-STM encodings of the input and apply it to unsupervised Chinese Word Segmentation (see §5.4).<sup>25</sup>

### 6.4.4 SentencePiece (Kudo and Richardson, 2018)

Not itself an algorithm as often assumed, but actually a software package, SentencePiece (Kudo and Richardson, 2018) offers both BPE and Unigram LM algorithms (so specifying "Sentence-Piece" is certainly not informative enough). Importantly, unlike their other implementations it does not treat spaces as special guaranteed word boundaries, allowing learned units to cross these boundaries and obviating the need for pre-tokenization in languages without whitespace-tokenized words like Chinese and Japanese.

### 6.5 Comparing morphological segmentation to BPE and friends

Several studies have compared linguistically motivated segmentation with data-driven ones, without conclusive results (to say the least).

Bostrom and Durrett (2020) claim that UnigramLM obtains better segmentation than BPE, both qualitatively (they tend to better correspond to morphemes and Morfessor (Creutz and Lagus, 2007) morphs) and quantitatively (they improve BERT-style models' performance on modern NLP tasks a little in English and a lot in Japanese).

When using (manually analyzed or gold) morphological analysis, Matthews et al. (2018) show that language modeling can be improved for agglutinative languages like Turkish. In Schwartz et al. (2020)'s low-resource study shows Morfessorbased language models (and character-based ones, see §8) outperform BPE-based ones. Pan et al. (2020) likewise improve NMT on Turkish and Uyghur by using morphological analyzers before applying BPE.

Using unsupervisedly obtained "morphological" subwords on the other hand, only Ataman and Federico (2018b) find that a model based on Morfessor FlatCat can outperform BPE; Zhou (2018), Domingo et al. (2018), Macháček et al. (2018), and Saleva and Lignos (2021) find no reliable improvement over BPE for translation. Banerjee and Bhattacharyya (2018) analyze translations obtained segmenting with Morfessor and BPE, and conclude that a possible improvement depends on the similarity of the languages. Huck et al. (2017) propose thus to combine both approaches.

As a possible explanation for the good performance of BPE, Gallé (2019) claims that the performance of BPE is related to its compression capac-

<sup>&</sup>lt;sup>25</sup>Note that this thus character-conditioned model can also be seen as an example of segmental neural language models (§5.1).

ity: with respect to members from the same class of compression algorithm, BPE performs close to the top in data compression benchmarks.

#### 6.6 How many units do we need?

Another open question is the question of how many merges (or what other prior) one should select for optimal performance. The best-performing number may depend on the task and the domain and differ by language (Mielke et al., 2019; Domingo et al., 2018). More and thus larger subword units allow for and lead to more memorization (Kharitonov et al., 2021), which may or may not be desired depending on the application.

Predicting the number of merges that works best without having to try different sizes would be desirable and Gowda and May (2020) claim to have found one such heuristic: merge as much as possible to shorten the overall sequence lengths while making sure that 95% of subword units appear at least 100 times (presumably in training data). Their motivation is that neural machine translation models are frequency-biased in their outputs and thus maintaining a more uniform frequency distribution is better.<sup>27</sup> A similar study undertaken by Ding et al. (2019) reiterate how contradictory suggestions for number of merges in past work are and add that in low-resource scenarios far fewer merges seem to be better, a trend with Transformers which differs from that with LSTMs, leading to an interesting question: should smaller corpora mean you can't afford characters or is it rather that you can't afford words?

A simple online answer to the question of how to select merges is presented by Salesky et al. (2018): while training an NMT model using BPE segments, gradually *increase* the vocabulary size by merging BPE vocabulary items, adding new, bigger BPE segments until they obtain diminishing returns. Embeddings for the newly introduced subwords are initialized by merging the embeddings of the two merged BPE segments with an autoencoder. Formulating the vocabulary selection prob-

lem as a search for the set of tokens with the highest entropy, Xu et al. (2021) proposes an optimal transport driven selection from BPE units that obtains vocabulary merges that often outperform a language-independent standard setting for translation. Another recent method that comes with a stopping criteria (and therefore dispenses with an additional hyperparameter) is Vilar and Federico (2021) which defines the likelihood of a vocabulary with respect to a sequence, and improves that likelihood greedily.

### 7 Shared vocabularies in multilingual models

Many NLP applications process text in more than one language at a time, the most obvious example perhaps being a machine translation system. In such cases, one could either use (and potentially learn) a tokenizer per language or a single tokenizer for both languages (also allowing sharing of embeddings if desired). Building a highly multilingual system that translates between more than two languages, Johnson et al. (2017) perform the former and first encounter questions like whether to oversample low-resource languages for learning a data-driven tokenizer and if so to which degree. These questions are addressed differently in the now more common highly multilingual pre-trained Transformers like mBERT (Devlin et al., 2018) and XLM (CONNEAU and Lample, 2019), and XLM-R (Conneau et al., 2020). In these models the sharing of learned representations is hypothesized to help transfer between languages by Pires et al. (2019), though Wu and Dredze (2019) provide inclonclusive results for this claim. It is worth pointing out that K et al. (2020) disagree and claim that subword overlap is not as important for transfer.

Even though all these models settle make sure to oversample low-resource languages at least some amount, Ács (2019) and Rust et al. (2021) show that tokenization in BERT-based Transformers is still biased towards high-resource languages. This bias is visible in a word's "fertility," i.e., the number of subwords a word is split into on average (which for example is much lower for English than it is for, say, Finnish), but they also find it affecting results in controlled (comparing monolingual to multilingual tokenization) downstream tasks. Maronikolakis et al. (2021) find that these granularity differences in tokenization between languages also

<sup>&</sup>lt;sup>26</sup>Relatedly, Novotný et al. (2021) show that the subword size matters and differs somewhat systematically between languages in the *n*-gram based fastText embeddings (Bojanowski et al., 2017).

<sup>&</sup>lt;sup>27</sup>A somewhat similar approach is taken by Gutierrez-Vasques et al. (2021), who look at the entropy (and transformations) of the distribution over subword types to identify a "turning point" at which one of the transformed quantities is minimal—but this turning point happens with far fewer merges than are generally required to reach good performance.

greatly affect sharing of semantic representations.

For selecting appropriate tokenization in a multilingual setting, Chung et al. (2020) offer an approach for retraining models from scratch, selecting subword vocabularies for language clusters to explicitly control for allocation and sharing. If on the other hand retraining from scratch is not feasible, one option is to add new subword units for the underresourced/oversegmented languages. Wang et al. (2020b) and Chau et al. (2020) both propose such additions with randomly initialized embeddings, but these approaches did not perform well when studied by Ebrahimi and Kann (2021); extending the idea, Liu et al. (2021) propose to use information about existing subword units to estimate embeddings instead of initializing newly added units randomly (similar to Salesky et al. (2018)). A different option is proposed by Wang et al. (2021b), who instead force the model to use (already existing) smaller subword units in highresource languages like English to make the segmentations across languages more similar and thus aid transfer—thus avoiding the complete retraining that comes with changing the segmentation method or allocation. In particular, they fine-tune the model to perform the target task 1) well with the deterministic segmentation (that undersegments highresource languages relative to low-resource ones), 2) well with sampled segmentations (so even highresource languages' words are segmented more), and 3) equally (in terms of a low divergence between the respective output distributions) between the two.

# 8 "Tokenization-free" character- and byte-level modeling

In sections §4.1 and §4.2 we discussed augmenting word models with character-level information in closed- and open-vocabulary settings. The idea of pure character- or byte-level modeling seems like an obvious simplification. Indeed, Sutskever et al. (2011) successfully modeled strings one character at a time using multiplicative RNNs, Chrupała (2013) suggest using character-level RNN modeling for agglutinative languages and tasks that require character information like character-level text segmentation (even anticipating contextualized embeddings!), and Conneau et al. (2017) successfully perform text classification from raw characters. The big breakthrough for generative character/byte-level models however came only with Al-Rfou et al.

(2019), who showed that sufficiently deep Transformers (64 layers in their case) can *greatly* outperform previous subword-based and hybrid (§4.2) open-vocabulary language models. This finding was updated by Choe et al. (2019), who again manage to match previous word- and subword-based state-of-the-art language modeling results.

#### 8.1 Characters?

A major factor limiting the adoption of characterlevel models is the fact that character sequences tend to be much longer than their word- or subwordlevel counterparts, making training and inference slower. To improve training speed and efficiency, Libovický and Fraser (2020) propose to start with a subword-based model and then fine-tune that model to instead work over characters, though they find improvements only in one of two evaluated language pairs. The more common approach to both training and inference however are various architectures for subsampling sequences, particularly in applications to machine translation: Chung et al. (2016) introduce a bi-scale recurrent neural network to enable the decoder of an encoder-decoder model to produce character sequences; they demonstrate improved performance over a subword-level decoder. Lee et al. (2017) (and later Gao et al. (2020)) advocate for the use of convolution and pooling layers at the input of the encoder. Cherry et al. (2018) evaluate various temporal pooling strategies including the HM-RNN of Chung et al. (2017) (discussed in §5.1) and conclude that none of them offered a suitable trade-off of performance and speed.

It has been argued that character-level models are more robust to noise and out-of-distribution data (Gupta et al., 2019), possibly because a wordor subword-level token sequence will change dramatically in the presence of noise. Libovický et al. (2021) however survey multiple character-level MT systems and conclude that they "show neither better domain robustness, nor better morphological generalization, despite being often so motivated," a finding corroborated by Rosales Núñez et al. (2021) for noisy user-generated text. Specifically under the lens of gender bias, Gaido et al. (2021) argue that character-level processing can lead to less gender bias in models: data-driven BPE vocabularies are biased towards male forms, for one because of frequency, but also because in languages like French and Italian, female forms are often constructed by affixing male forms, leading to a disparity in tokenization. They show that character-level processing can ameliorate these disparities to some degree.

#### 8.2 Character hashes?

In massively multilingual settings, naïve use of a character-level model can result in a very large vocabulary. For example, Unicode 13.0 specifies 143,859 codepoints, most of which will never be seen in training. This can lead to "out-ofvocabulary" problems just like the ones encountered with words (discussed in §4). One possible solution is used in Clark et al. (2021)'s CANINE, where instead of giving each character its own embedding, they hash all possible codepoints with multiple hash functions into smaller sets to share parameters across codepoints (similar to Svenstrup et al. (2017)). CANINE further includes downsampling and upsampling operations to attain reasonable computational efficiency, and focuses on non-generative sequence labeling tasks.

### **8.3** Bytes?

An alternative way to avoid the huge-vocabulary problem of character-level models is to instead represent text using the byte sequence resulting from a standard encoding like UTF-8. This can be seen as using a fixed tokenization scheme that was created by a standards body (the Unicode consortium) rather than a linguistically- or statisticallyinformed process. The main advantage of using Unicode bytes specifically is their incredible scope - the Unicode consortium states that their goal is to cover "all the characters for all the writing systems of the world, modern and ancient" in addition to "technical symbols, punctuations, and many other characters used in writing text". 28 Separately, byte-level modeling is presented as a natural choice by Graves (2013) when following the "principle of modelling the smallest meaningful units in the data". As shown recently by ByT5 (Xue et al., 2021), byte-level models confer similar benefits to character-level models (better robustness to noise, avoidance of out-of-vocabulary issues, etc.) with similar drawbacks (longer training and inference time due to longer sequences). As such, similar results and models have been proposed for byte-level models as for character-level. For example, the recent Charformer model (Tay et al., 2021) downsamples input sequences to avoid increased computational costs, and Al-Rfou et al. (2019) demonstrate strong language modeling performance with a deep byte-level Transformer.

### 8.4 So are these maximal decompositions the solution then?

While byte-level modeling is often presented as an unbiased, token-free approach, we argue it is more constructive to consider byte-level modeling as simply using a fixed, predefined, and standardized tokenization scheme (that incidentally is often the same as the way that underlying the text data is stored on disk). This tokenization scheme is by no means the "best" or most fundamental – indeed, the Unicode standard was not created with any linguistically-motivated goal in mind, apart from being able to represent a huge variety of content. Indeed, using a Unicode-based byte-level tokenization scheme will result in significantly different trade-offs for different languages by virtue of the way the standard was created: Latin (i.e. ASCII) characters are represented as a single byte, whereas characters in other languages are represented as multiple bytes. An immediate consequence is that a UTF-8-based byte-level tokenization scheme can result in dramatically longer sequences when representing the same underlying semantic content in different languages. This could unfairly increase computational costs for downstream users of a model who do not communicate in ASCII symbols.

### 8.5 Visual featurization: Pixels?

Another approach to "tokenization-free" modeling utilizes *visual* rather than byte-based representations. Where byte-level models aim to cover the full underlying 'vocabulary' as represented *on disk*, visual representations aim to encode similarities among characters which *human readers* may use when processing text. One motivation is robustness: byte-level models are reliant on consistent observed sequences, which leaves them susceptible to variation and noise which may render similarly visually.<sup>29</sup>

The initial motivation for much work on using visual features were to create embeddings that reflected the shared character components found in e.g., Chinese characters (radicals) or Korean

<sup>28</sup>https://unicode.org/faq/basic\_q.html

<sup>&</sup>lt;sup>29</sup>For languages without a digital orthographic standard (e.g., Pashto), multiple byte sequences render similarly and are in free variation among speakers.

syllable blocks, and so were better able to generalize to rare or unseen characters. The first work in this space initialized embeddings for Chinese by linearizing bitmaps of characters or words (Aldón Mínguez et al., 2016; Costa-jussà et al., 2017). Subsequent work focused on character segmentation, either through pre-computed embeddings based on linearized images (Wang et al., 2020a) or learned with the downstream task via convolutions (Liu et al., 2017; Dai and Cai, 2017; Broscheit, 2018; Meng et al., 2019), improving results for rare characters. Character segmentation was in part motivated by application to Chinese and the focus on sub-character components, but also enabled the use of fixed-width images, a problem which Sun et al. (2018, 2019) instead tackled by rendering words into square images, wrapping characters as necessary.

Recent work has proposed "tokenization-free" visual representations (Mansimov et al., 2020; Salesky et al., 2021). Rather than rendering images for a given segmentation into words, characters, or bytes, and thus incorporating their previously discussed challenges, they render each sentence as an image and translate directly from pixel decompositions. Salesky et al. (2021) show that such models can perform competitively across a range of languages and scripts for machine translation and are significantly more robust to induced noise, including but not limited to unicode errors. Mansimov et al. (2020) explore pixel-to-pixel models; a challenging setting that is not yet competitive but neatly sidesteps many concerns with text-based models.

### 9 Discussion and Conclusion

Tokenization, both the process and the term itself, has evolved significantly since the early days of NLP. In this paper, we traced their evolution and highlighted major changes over the years, connecting disparate intuitions.

Despite significant advancement, we have seen that there is no perfect method of handling tokenization. All options—from whitespace-delimited pre-tokens to learned subwords and down to bytes—have drawbacks and strengths. Many desiderata may be fundamentally at odds with each other, e.g., the desire to decompose maximally for simple and robust processing with a desire to be computationally efficient in a way that is fair across languages—a question that particularly pertinent as the field

turns its attention to greener NLP. While there are applications for which characters and bytes or even pixels may be the tool of choice, there are a myriad of applications in which larger discrete tokens are desirable, both for interpretability and efficiency. Recent work gives us new examples of this: Zhang et al. (2021) improve BERT by feeding inputs tokenized multiple ways and Dossou and Emezue (2021) use human-annotated larger-than-word vocabulary units to improve low-resource MT. Itzhak and Levy (2021) show that using subword units need not mean giving up on spelling-information, showing that spellings can be recovered from subword-based pretrained models.

In conclusion, despite all the promises of neural networks of allowing "end-to-endness", tokenization is a clear testimony that that promised land is still far away. Like often, this drawback is evident for practitioners who have to verify that the pre-trained models matches the tokenizer that is being used (a separated model in most NLP frameworks). The fact that tokenization is handled completely independent of the down-stream tasks adds to this separation. As Henderson (2020) puts it: "It remains to find effective neural architectures for learning the set of entities jointly with the rest of the neural model, and for generalising such methods from the level of character strings to higher levels of representation".

#### Acknowledgements

We would like to thank Shiyue Zhang, Vassilina Nikoulina, and Kamalkumar R for their notes, Kaustubh Dhole and Mike Tian-Jian Jiang for paper suggestions, and Jason Eisner for feedback on a draft. Samson is supported by Salesforce and Singapore's Economic Development Board under the Industrial Postgraduate Programme.

#### References

Judit Ács. 2019. Exploring BERT's vocabulary.

Gustavo Aguilar, Bryan McCann, Tong Niu, Nazneen Rajani, Nitish Keskar, and Thamar Solorio. 2021. Char2subword: Extending the subword embedding space using robust character compositionality.

Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. 2019. Characterlevel language modeling with deeper self-

- attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3159–3166.
- David Aldón Mínguez, Marta Ruiz Costa-Jussà, and José Adrián Rodríguez Fonollosa. 2016. Neural machine translation using bitmap fonts. In *Proceedings of the EAMT 2016 Fifth Workshop on Hybrid Approaches to Translation (Hy-Tra)*, pages 1–9.
- Zaid Alyafeai, Maged S. Al-shaibani, Mustafa Ghaleb, and Irfan Ahmad. 2021. Evaluating various tokenizers for arabic text classification. *ArXiv*, abs/2106.07540.
- Chantal Amrhein and Rico Sennrich. 2021. How suitable are subword segmentation strategies for translating non-concatenative morphology? In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 689–705, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Duygu Ataman, Wilker Aziz, and Alexandra Birch. 2020. A latent morphology model for open-vocabulary neural machine translation. In *International Conference on Learning Representations*.
- Duygu Ataman and Marcello Federico. 2018a. Compositional representation of morphologically-rich input for neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 305–311, Melbourne, Australia. Association for Computational Linguistics.
- Duygu Ataman and Marcello Federico. 2018b. An evaluation of two vocabulary reduction methods for neural machine translation. In *Proceedings* of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track), pages 97–110.
- Duygu Ataman, Orhan Firat, Mattia A. Di Gangi, Marcello Federico, and Alexandra Birch. 2019. On the importance of word boundaries in character-level neural machine translation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 187–193, Hong Kong. Association for Computational Linguistics.

- R.H. Baayen, R Piepenbrock, and L Gulikers. 1995. CELEX2. In *LDC96L14*. Linguistic Data Consortium, University of Pennsylvania.
- Tamali Banerjee and Pushpak Bhattacharyya. 2018. Meaningless yet meaningful: Morphology grounded subword-level nmt. In *Proceedings of the second workshop on subword/character level models*, pages 55–60.
- Marco Baroni. 2000. Distributional Cues in Morpheme Discovery: A Computational Model and Empirical Evidence. Ph.D. thesis, University of California Los Angeles.
- Sarah Beemer, Zak Boston, April Bukoski, Daniel Chen, Princess Dickens, Andrew Gerlach, Torin Hopkins, Parth Anand Jawale, Chris Koski, Akanksha Malhotra, Piyush Mishra, Saliha Muradoglu, Lan Sang, Tyler Short, Sagarika Shreevastava, Elizabeth Spaulding, Testumichi Umada, Beilei Xiang, Changbing Yang, and Mans Hulden. 2020. Linguist vs. machine: Rapid development of finite-state morphological grammars. In *Proceedings of the 17th SIGMOR-PHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 162–170, Online. Association for Computational Linguistics.
- Kenneth R Beesley and Lauri Karttunen. 2003. Finite-state morphology: Xerox tools and techniques. *CSLI*, *Stanford*.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2001. A neural probabilistic language model. In *Advances in Neural Information Processing Systems*, volume 13. MIT Press.
- Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432.
- Taylor Berg-Kirkpatrick, Alexandre Bouchard-Côté, John DeNero, and Dan Klein. 2010. Painless unsupervised learning with features. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 582–590, Los Angeles, California. Association for Computational Linguistics.
- Toms Bergmanis and Sharon Goldwater. 2017. From segmentation to analyses: a probabilistic

- model for unsupervised morphology induction. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 337–346, Valencia, Spain. Association for Computational Linguistics.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5(0):135–146.
- Kaj Bostrom and Greg Durrett. 2020. Byte pair encoding is suboptimal for language model pretraining. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4617–4624, Online. Association for Computational Linguistics.
- Martin Braschler and Bärbel Ripplinger. 2004. How effective is stemming and decompounding for german text retrieval? *Information Retrieval*, 7(3):291–316.
- Michael R Brent and Timothy A Cartwright. 1996. Distributional regularity and phonotactic constraints are useful for segmentation. *Cognition*, 61(1-2):93–125.
- Michael R Brent, Sreerama K Murthy, and Andrew Lundberg. 1995. Discovering morphemic suffixes: a case study in MDL induction. In *In Fifth International Workshop on AI and Statistics*, Ft. Citeseer.
- Samuel Broscheit. 2018. Learning distributional token representations from visual features. In *Proceedings of The Third Workshop on Representation Learning for NLP*, pages 187–194, Melbourne, Australia. Association for Computational Linguistics.
- Jacob Buckman and Graham Neubig. 2018. Neural lattice language models. *Transactions of the Association for Computational Linguistics (TACL)*, 6.
- Kris Cao and Laura Rimell. 2021. You should evaluate your language model on marginal likelihood over tokenisations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2104–2114, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- William Chan, Yu Zhang, Quoc Le, and Navdeep Jaitly. 2017. Latent sequence decompositions. In *International Conference on Learning Repre*sentations 2017.
- Ethan C. Chau, Lucy H. Lin, and Noah A. Smith. 2020. Parsing with multilingual BERT, a small corpus, and a small treebank. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1324–1334, Online. Association for Computational Linguistics.
- Stanley F. Chen and Joshua Goodman. 1999. An empirical study of smoothing techniques for language modeling. *Computer Speech & Language*, 13(4):359–394.
- Colin Cherry, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. 2018. Revisiting character-based neural machine translation with capacity and compression. In *EMNLP*.
- Dokook Choe, Rami Al-Rfou, Mandy Guo, Heeyoung Lee, and Noah Constant. 2019. Bridging the gap for tokenizer-free language models.
- Grzegorz Chrupała. 2013. Text segmentation with character-level text embeddings. Workshop on Deep Learning for Audio, Speech and Language Processing, ICML 2013; Workshop on Deep Learning for Audio, Speech and Language Processing, ICML 2013; Conference date: 16-06-2013.
- Hyung Won Chung, Dan Garrette, Kiat Chuan Tan, and Jason Riesa. 2020. Improving multilingual models with language-clustered vocabularies. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4536–4546, Online. Association for Computational Linguistics.
- Junyoung Chung, Sungjin Ahn, and Yoshua Bengio. 2017. Hierarchical multiscale recurrent neural networks. In *International Conference on Learning Representations* 2017.
- Junyoung Chung, Kyunghyun Cho, and Yoshua Bengio. 2016. A character-level decoder without explicit segmentation for neural machine translation. *arXiv preprint arXiv:1603.06147*.
- Kenneth W. Church. 2000. Empirical estimates of adaptation: The chance of two Noriegas is closer

- to p/2 than  $p^2$ . In Proceedings of the 18th Conference on Computational Linguistics Volume 1, pages 180–186. Association for Computational Linguistics.
- Jonathan H Clark, Dan Garrette, Iulia Turc, and John Wieting. 2021. Canine: Pre-training an efficient tokenization-free encoder for language representation. *arXiv preprint arXiv:2103.06874*.
- Lionel Clément, Eric De la Clergerie, and Lionel Net. 2005. Maf: a morphosyntactic annotation framework.
- Ronan Collobert and Jason Weston. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th International Conference on Machine Learning*, pages 160–167.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Alexis CONNEAU and Guillaume Lample. 2019. Cross-lingual language model pretraining. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Alexis Conneau, Holger Schwenk, Loïc Barrault, and Yann Lecun. 2017. Very deep convolutional networks for text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1107–1116, Valencia, Spain. Association for Computational Linguistics.
- Marta R. Costa-jussà, David Aldón, and José A. R. Fonollosa. 2017. Chinese–spanish neural machine translation enhanced with character and word bitmap fonts. *Machine Translation*, 31(1):35–47.
- Ryan Cotterell, Sabrina J. Mielke, Jason Eisner, and Brian Roark. 2018. Are all languages equally hard to language-model? In *Proceedings of the 2018 Conference of the North Amer-*

- ican Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 536–541, New Orleans, Louisiana. Association for Computational Linguistics.
- Mathias Creutz and Krista Lagus. 2002. Unsupervised discovery of morphemes. In *Proceedings* of the ACL-02 Workshop on Morphological and Phonological Learning, pages 21–30.
- Mathias Creutz and Krista Lagus. 2005. Inducing the morphological lexicon of a natural language from unannotated text. In *Proceedings of the International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning (AKRR'05)*, volume 1, pages 51–59.
- Mathias Creutz and Krista Lagus. 2007. Unsupervised models for morpheme segmentation and morphology learning. *ACM Trans. Speech Lang. Process.*, 4(1).
- Falcon Dai and Zheng Cai. 2017. Glyph-aware embedding of Chinese characters. In *Proceedings* of the First Workshop on Subword and Character Level Models in NLP, pages 64–69, Copenhagen, Denmark. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Shuoyang Ding, Adithya Renduchintala, and Kevin Duh. 2019. A call for prudent choice of subword merge operations in neural machine translation. In *Proceedings of Machine Translation Summit XVII Volume 1: Research Track*, pages 204–213, Dublin, Ireland. European Association for Machine Translation.
- M. Domingo, M. Garcia-Martinez, A. Helle,
   F. Casacuberta, and M. Herranz. 2018. How
   Much Does Tokenization Affect Neural Machine
   Translation? arXiv e-prints.
- Miguel Domingo, Mercedes Garcia-Martinez, Alexandre Helle, Francisco Casacuberta, and Manuel Herranz. 2018. How much does to-kenization affect neural machine translation? arXiv preprint arXiv:1812.08621.

- Bonaventure F. P. Dossou and Chris C. Emezue. 2021. Crowdsourced phrase-based tokenization for low-resourced neural machine translation: The case of fon language.
- Yerai Doval and Carlos Gómez-Rodríguez. 2019. Comparing neural- and n-gram-based language models for word segmentation. *Journal of the Association for Information Science and Technology*, 70(2):187–197.
- C. M. Downey, Fei Xia, Gina-Anne Levow, and Shane Steinert-Threlkeld. 2021. A masked segmental language model for unsupervised natural language segmentation.
- Abteen Ebrahimi and Katharina Kann. 2021. How to adapt your pretrained multilingual model to 1600 languages.
- Hicham El Boukkouri, Olivier Ferret, Thomas Lavergne, Hiroshi Noji, Pierre Zweigenbaum, and Jun'ichi Tsujii. 2020. CharacterBERT: Reconciling ELMo and BERT for word-level openvocabulary representations from characters. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6903–6915, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Salah El Hihi and Yoshua Bengio. 1995. Hierarchical recurrent neural networks for long-term dependencies. In *Proceedings of the 8th International Conference on Neural Information Processing Systems*, NIPS'95, page 493–499, Cambridge, MA, USA. MIT Press.
- Jeffrey L. Elman. 1990. Finding structure in time. *Cognitive Science*, 14(2):179–211.
- Micha Elsner, Sharon Goldwater, Naomi Feldman, and Frank Wood. 2013. A joint learning model of word segmentation, lexical acquisition, and phonetic variability. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 42–54, Seattle, Washington, USA. Association for Computational Linguistics.
- Philip Gage. 1994. A new algorithm for data compression. *C Users J.*, 12(2):23–38.
- Marco Gaido, Beatrice Savoldi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. How to

- split: the effect of word segmentation on gender bias in speech translation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3576–3589, Online. Association for Computational Linguistics.
- Matthias Gallé. 2019. Investigating the effectiveness of BPE: The power of shorter sequences. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1375–1381, Hong Kong, China. Association for Computational Linguistics.
- Yingqiang Gao, Nikola I Nikolov, Yuhuang Hu, and Richard HR Hahnloser. 2020. Character-level translation with self-attention. *arXiv* preprint arXiv:2004.14788.
- John Goldsmith. 2001. Unsupervised learning of the morphology of a natural language. *Computational linguistics*, 27(2):153–198.
- S. Goldwater, T. L. Griffiths, and M. Johnson. 2006a. Contextual dependencies in unsupervised word segmentation. In *Proc. of COLING-*ACL.
- Sharon Goldwater, Thomas L. Griffiths, and Mark Johnson. 2009. A Bayesian framework for word segmentation: Exploring the effects of context. *Cognition*, 112(1):21–54.
- Sharon Goldwater, Thomas L. Griffiths, and Mark Johnson. 2011. Producing power-law distributions and damping word frequencies with two-stage language models. *Journal of Machine Learning Research*, 12(Jul):2335–2382.
- Sharon Goldwater, Mark Johnson, and Thomas Griffiths. 2006b. Interpolating between types and tokens by estimating power-law generators. In *Advances in Neural Information Processing Systems*, volume 18. MIT Press.
- Thamme Gowda and Jonathan May. 2020. Finding the optimal vocabulary size for neural machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3955–3964, Online. Association for Computational Linguistics.

- Edouard Grave, Armand Joulin, and Nicolas Usunier. 2016. Improving neural language models with a continuous cache. In *International Conference on Learning Representations* 2016.
- Edouard Grave, Sainbayar Sukhbaatar, Piotr Bojanowski, and Armand Joulin. 2019. Training hybrid language models by marginalizing over segmentations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1477–1482, Florence, Italy. Association for Computational Linguistics.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*.
- Stig-Arne Grönroos, Sami Virpioja, and Mikko Kurimo. 2020. Morfessor EM+Prune: Improved subword segmentation with expectation maximization and pruning. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3944–3953, Marseille, France. European Language Resources Association.
- Stig-Arne Grönroos, Sami Virpioja, Peter Smit, and Mikko Kurimo. 2014. Morfessor FlatCat: An HMM-based method for unsupervised and semi-supervised learning of morphology. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1177–1185.
- Rohit Gupta, Laurent Besacier, Marc Dymetman, and Matthias Gallé. 2019. Characterbased nmt with transformer. *arXiv preprint arXiv:1911.04997*.
- Ximena Gutierrez-Vasques, Christian Bentz, Olga Sozinova, and Tanja Samardzic. 2021. From characters to words: the turning point of BPE merges. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3454–3468, Online. Association for Computational Linguistics.
- Xuanli He, Gholamreza Haffari, and Mohammad Norouzi. 2020. Dynamic programming encoding for subword segmentation in neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3042–3051, Online. Association for Computational Linguistics.

- James Henderson. 2020. The unstoppable rise of computational linguistics in deep learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6294–6306, Online. Association for Computational Linguistics.
- Tatsuya Hiraoka, Hiroyuki Shindo, and Yuji Matsumoto. 2019. Stochastic tokenization with a language model for neural text classification. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1620–1629, Florence, Italy. Association for Computational Linguistics.
- Tatsuya Hiraoka, Sho Takase, Kei Uchiumi, Atsushi Keyaki, and Naoaki Okazaki. 2020. Optimizing word segmentation for downstream task. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1341–1351, Online. Association for Computational Linguistics.
- Tatsuya Hiraoka, Sho Takase, Kei Uchiumi, Atsushi Keyaki, and Naoaki Okazaki. 2021. Joint optimization of tokenization and downstream model. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 244–255, Online. Association for Computational Linguistics.
- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb: Derivational morphology improves BERT's interpretation of complex words. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3594–3608, Online. Association for Computational Linguistics.
- Po-Sen Huang, Chong Wang, Sitao Huang, Dengyong Zhou, and Li Deng. 2018. Towards neural phrase-based machine translation. In *International Conference on Learning Representations*.
- Matthias Huck, Simon Riess, and Alexander Fraser. 2017. Target-side word segmentation strategies for neural machine translation. In *Proceedings* of the Second Conference on Machine Translation, pages 56–67.

- Gérard Huet. 2003. Lexicon-directed segmentation and tagging of sanskrit. In *in « XIIth World Sanskrit Conference*, pages 307–325.
- Kyuyeon Hwang and Wonyong Sung. 2017. Character-level language modeling with hierarchical recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pages 5720–5724. IEEE.
- Itay Itzhak and Omer Levy. 2021. Models in a spelling bee: Language models implicitly learn the character composition of tokens.
- Miguel Ángel Jiménez-Montaño. 1984. On the syntactic structure of protein sequences and the concept of grammar complexity. *Bulletin of Mathematical Biology*, 46(4):641–659.
- Mark Johnson and Sharon Goldwater. 2009. Improving nonparameteric Bayesian inference: experiments on unsupervised word segmentation with adaptor grammars. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 317–325, Boulder, Colorado. Association for Computational Linguistics.
- Mark Johnson, Thomas L. Griffiths, and Sharon Goldwater. 2007. Adaptor grammars: A framework for specifying compositional nonparametric Bayesian models. In *Advances in Neural Information Processing Systems* 19, pages 641–648. MIT Press.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling.
- Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. 2020. Cross-lingual ability of multilingual bert: An empirical study. In *Interna*tional Conference on Learning Representations.

- Ákos Kádár, Marc-Alexandre Côté, Grzegorz Chrupała, and Afra Alishahi. 2018. Revisiting the hierarchical multiscale lstm. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3215–3227. Association for Computational Linguistics.
- Kazuya Kawakami, Chris Dyer, and Phil Blunsom. 2017. Learning to create and reuse words in open-vocabulary neural language modeling. In *Proceedings of the 55th Annual Meeting of the ACL (Volume 1: Long Papers)*, pages 1492–1502, Vancouver, Canada. Association for Computational Linguistics.
- Kazuya Kawakami, Chris Dyer, and Phil Blunsom. 2019. Learning to discover, ground and use words with segmental neural language models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6429–6441, Florence, Italy. Association for Computational Linguistics.
- Eugene Kharitonov, Marco Baroni, and Dieuwke Hupkes. 2021. How bpe affects memorization in transformers.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander Rush. 2016. Character-aware neural language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1).
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.
- Philipp Koehn and Kevin Knight. 2003. Empirical methods for compound splitting. In 10th Conference of the European Chapter of the Association for Computational Linguistics, Budapest, Hungary. Association for Computational Linguistics.
- Lingpeng Kong, Chris Dyer, and Noah A. Smith. 2016. Segmental recurrent neural networks. In *International Conference on Learning Representations*.

- Jan Koutnik, Klaus Greff, Faustino Gomez, and Juergen Schmidhuber. 2014. A clockwork rnn. In *Proceedings of Machine Learning Research*, volume 32, pages 1863–1871, Bejing, China. PMLR.
- Julia Kreutzer and Artem Sokolov. 2018. Learning to Segment Inputs for NMT Favors Character-Level Processing. In *Proceedings of the 15th International Workshop on Spoken Language Translation (IWSLT)*.
- Amrith Krishna, Pavan Kumar Satuluri, and Pawan Goyal. 2017. A dataset for Sanskrit word segmentation. In *Proceedings of the Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 105–114, Vancouver, Canada. Association for Computational Linguistics.
- Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 66–75, Melbourne, Australia. Association for Computational Linguistics.
- Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Mikko Kurimo, Sami Virpioja, Ville Turunen, and Krista Lagus. 2010. Morpho Challenge 2005-2010: Evaluations and Results. In *Proceedings of the 11th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology*, pages 87–95, Uppsala, Sweden. Association for Computational Linguistics.
- N Jesper Larsson and Alistair Moffat. 2000. Off-line dictionary-based compression. *Proceedings* of the IEEE, 88(11):1722–1732.
- Jason Lee, Kyunghyun Cho, and Thomas Hofmann. 2017. Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, 5(0):365–378.

- Jindřich Libovický and Alexander Fraser. 2020. Towards reasonably-sized character-level transformer NMT by finetuning subword systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2572–2579, Online. Association for Computational Linguistics.
- Jindřich Libovický, Helmut Schmid, and Alexander Fraser. 2021. Why don't people use character-level machine translation?
- Constantine Lignos. 2010. Learning from unseen data. In *Proceedings of the Morpho Challenge* 2010 Workshop, pages 35–38, Helsinki, Finland. Aalto University School of Science and Technology.
- Wang Ling, Chris Dyer, Alan W. Black, Isabel Trancoso, Ramon Fermandez, Silvio Amir, Luis Marujo, and Tiago Luis. 2015. Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1520–1530, Lisbon, Portugal. Association for Computational Linguistics.
- Frederick Liu, Han Lu, Chieh Lo, and Graham Neubig. 2017. Learning character-level compositionality with visual features. In *Proceedings* of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2059–2068, Vancouver, Canada. Association for Computational Linguistics.
- Xin Liu, Baosong Yang, Dayiheng Liu, Haibo Zhang, Weihua Luo, Min Zhang, Haiying Zhang, and Jinsong Su. 2021. Bridging subword gaps in pretrain-finetune paradigm for natural language generation.
- Zhaoxin Luo and Michael Zhu. 2021. Recurrent neural networks with mixed hierarchical structures for natural language processing. *CoRR*, abs/2106.02562.
- Minh-Thang Luong and Christopher D. Manning. 2016. Achieving open vocabulary neural machine translation with hybrid word-character models. In *Proceedings of the 54th Annual Meeting of the ACL (Volume 1: Long Papers)*, pages 1054–1063, Berlin, Germany. Association for Computational Linguistics.

- Wentao Ma, Yiming Cui, Chenglei Si, Ting Liu, Shijin Wang, and Guoping Hu. 2020. Char-BERT: Character-aware Pre-trained Language Model. In *COLING*.
- Dominik Macháček, Jonáš Vidra, and Ondřej Bojar. 2018. Morphological and language-agnostic word segmentation for nmt. In *International Conference on Text*, *Speech*, *and Dialogue*, pages 277–284. Springer.
- Klaus Macherey, Andrew Dai, David Talbot, Ashok Popat, and Franz Och. 2011. Language-independent compound splitting with morphological operations. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1395–1404, Portland, Oregon, USA. Association for Computational Linguistics.
- David J. C. MacKay and Linda C. Bauman Peto. 1995. A hierarchical Dirichlet language model. *Journal of Natural Language Engineering*, 1(3):289–308.
- Elman Mansimov, Mitchell Stern, Mia Chen, Orhan Firat, Jakob Uszkoreit, and Puneet Jain. 2020. Towards end-to-end in-image neural machine translation. In *Proceedings of the First International Workshop on Natural Language Processing Beyond Text*, pages 70–74, Online. Association for Computational Linguistics.
- Carl de Marcken. 1996. Linguistic structure as composition and perturbation. In *Proceedings* of the 34th Meeting of the Association for Computational Linguistics.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- Antonis Maronikolakis, Philipp Dufter, and Hinrich Schütze. 2021. Wine is not vin. on the compatibility of tokenizations across languages. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2382–2399, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Austin Matthews, Graham Neubig, and Chris Dyer. 2018. Using morphological knowledge in open-

- vocabulary neural language models. In *HLT-NAACL*.
- Yuxian Meng, Wei Wu, Fei Wang, Xiaoya Li, Ping Nie, Fan Yin, Muyu Li, Qinghong Han, Xiaofei Sun, and Jiwei Li. 2019. Glyce: Glyph-vectors for chinese character representations. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Bart van Merriënboer, Amartya Sanyal, Hugo Larochelle, and Yoshua Bengio. 2017. Multiscale sequence modeling with a learned dictionary. In *MLSLP*.
- Sabrina J. Mielke, Ryan Cotterell, Kyle Gorman, Brian Roark, and Jason Eisner. 2019. What kind of language is hard to language-model? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4975–4989, Florence, Italy. Association for Computational Linguistics.
- Sabrina J. Mielke and Jason Eisner. 2018. Spell once, summon anywhere: A two-level open-vocabulary language model. *CoRR*, abs/1804.08205.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume* 2, pages 3111–3119.
- Tomáš Mikolov, Ilya Sutskever, Anoop Deoras, Hai-Son Le, Stefan Kombrink, and Jan Černocký. 2012. Subword language modeling with neural networks. *preprint* (*rejected from ICASSP 2012*).
- Daichi Mochihashi, Takeshi Yamada, and Naonori Ueda. 2009. Bayesian unsupervised word segmentation with nested pitman-yor language modeling. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 100–108, Suntec, Singapore. Association for Computational Linguistics.
- Christian Monson, Jaime Carbonell, Alon Lavie, and Lori Levin. 2007. ParaMor: Minimally supervised induction of paradigm structure and

- morphological analysis. In *Proceedings of Ninth Meeting of the ACL Special Interest Group in Computational Morphology and Phonology*, pages 117–125, Prague, Czech Republic. Association for Computational Linguistics.
- Christian Monson, Kristy Hollingshead, and Brian Roark. 2009. Probabilistic paramor. In *CLEF* (*Working Notes*). Citeseer.
- Amir More, Özlem Çetinoğlu, Çağrı Çöltekin, Nizar Habash, Benoît Sagot, Djamé Seddah, Dima Taji, and Reut Tsarfaty. 2018. CoNLL-UL: Universal morphological lattices for Universal Dependency parsing. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Karthik Narasimhan, Regina Barzilay, and Tommi Jaakkola. 2015. An unsupervised method for uncovering morphological chains. *Transactions of the Association for Computational Linguistics*, 3:157–167.
- Vít Novotný, Eniafe Festus Ayetiran, Dávid Lupták, Michal Stefánik, and Petr Sojka. 2021. One size does not fit all: Finding the optimal n-gram sizes for fasttext models across languages. *CoRR*, abs/2102.02585.
- Yirong Pan, Xiao Li, Yating Yang, and Rui Dong. 2020. Morphological word segmentation on agglutinative languages for neural machine translation. *CoRR*, abs/2001.01589.
- Kyubyong Park, Joohong Lee, Seongbo Jang, and Dawoon Jung. 2020. An empirical study of tokenization strategies for various korean nlp tasks. In *AACL/IJCNLP*.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Yuval Pinter. 2021. Integrating approaches to word representation.

- Yuval Pinter, Robert Guthrie, and Jacob Eisenstein. 2017. Mimicking word embeddings using subword RNNs. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 102–112, Copenhagen, Denmark. Association for Computational Linguistics.
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Barbara Plank, Anders Søgaard, and Yoav Goldberg. 2016. Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418, Berlin, Germany. Association for Computational Linguistics.
- M. F. Porter. 1980. An algorithm for suffix stripping. *Program*, 14(3):130–137.
- Ivan Provilkov, Dmitrii Emelianenko, and Elena Voita. 2020. BPE-dropout: Simple and effective subword regularization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1882–1892, Online. Association for Computational Linguistics.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. *preprint*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Jorma Rissanen. 1989. *Stochastic Complexity in Statistical Inquiry Theory*. World Scientific Publishing Co., Inc., USA.
- José Carlos Rosales Núñez, Guillaume Wisniewski, and Djamé Seddah. 2021. Noisy UGC translation at the character level: Revisiting openvocabulary capabilities and robustness of charbased models. In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT*

- 2021), pages 199–211, Online. Association for Computational Linguistics.
- Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. 2021. How good is your tokenizer? on the monolingual performance of multilingual language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3118–3135, Online. Association for Computational Linguistics.
- Benoît Sagot and Pierre Boullier. 2008. SxPipe 2: architecture pour le traitement pré-syntaxique de corpus bruts. *Revue TAL*, 49(2):155–188.
- Elizabeth Salesky, David Etter, and Matt Post. 2021. Robust open-vocabulary translation from visual text representations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7235–7252, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Elizabeth Salesky, Andrew Runge, Alex Coda, Jan Niehues, and Graham Neubig. 2018. Optimizing segmentation granularity for neural machine translation. *arXiv preprint arXiv:1810.08641*.
- Jonne Saleva and Constantine Lignos. 2021. The effectiveness of morphology-aware segmentation in low-resource neural machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 164–174, Online. Association for Computational Linguistics.
- Cicero dos Santos and Bianca Zadrozny. 2014. Learning character-level representations for partof-speech tagging. In *Proceedings of the 31st International Conference on Machine Learning* (*ICML-14*), pages 1818–1826.
- Jürgen Schmidhuber. 1991. Neural sequence chunkers.
- Jürgen Schmidhuber. 1992. Learning complex, extended sequences using the principle of history compression. *Neural Computation*, 4(2):234–242.

- Patrick Schone and Daniel Jurafsky. 2001. Knowledge-free induction of inflectional morphologies. In Second Meeting of the North American Chapter of the Association for Computational Linguistics.
- Mike Schuster and Kaisuke Nakajima. 2012. Japanese and korean voice search. In 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5149–5152.
- Lane Schwartz, Francis Tyers, Lori Levin, Christo Kirov, Patrick Littell, Chi kiu Lo, Emily Prud'hommeaux, Hyunji Hayley Park, Kenneth Steimel, Rebecca Knowles, Jeffrey Micher, Lonny Strunk, Han Liu, Coleman Haley, Katherine J. Zhang, Robbie Jimmerson, Vasilisa Andriyanets, Aldrian Obaja Muis, Naoki Otani, Jong Hyuk Park, and Zhisong Zhang. 2020. Neural polysynthetic language modelling.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the ACL (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Yan Shao, Christian Hardmeier, and Joakim Nivre. 2018. Universal word segmentation: Implementation and interpretation.
- Pamela Shapiro and Kevin Duh. 2018. BPE and charcnns for translation of morphology: A cross-lingual comparison and analysis. *CoRR*, abs/1809.01301.
- Chenglei Si, Zhengyan Zhang, Yingfa Chen, Fanchao Qi, Xiaozhi Wang, Zhiyuan Liu, and Maosong Sun. 2021. ShuoWen-JieZi: Linguistically Informed Tokenizers For Chinese Language Model Pretraining. *arXiv*, abs/2106.00400.
- Peter Smit, Sami Virpioja, Stig-Arne Grönroos, and Mikko Kurimo. 2014. Morfessor 2.0: Toolkit for statistical morphological segmentation. In Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics, pages 21–24, Gothenburg, Sweden. Association for Computational Linguistics.

- Matthew Snover and Michael Brent. 2002. A probabilistic model for learning concatenative morphology. *Advances in Neural Information Processing Systems*, 15:1537–1544.
- Matthew G. Snover and Michael R. Brent. 2001. A Bayesian model for morpheme and paradigm identification. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 490–498, Toulouse, France. Association for Computational Linguistics.
- Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2021. Fast Word-Piece tokenization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2089–2103, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- James A. Storer and Thomas G. Szymanski. 1982. Data compression via textual substitution. *J. ACM*, 29(4):928–951.
- Baohua Sun, Lin Yang, Catherine Chi, Wenhan Zhang, and Michael Lin. 2019. Squared english word: A method of generating glyph to use super characters for sentiment analysis. In Proceedings of the 2nd Workshop on Affective Content Analysis (AffCon 2019) co-located with the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI 2019), volume 2328, pages 140–151.
- Baohua Sun, Lin Yang, Patrick Dong, Wenhan Zhang, Jason Dong, and Charles Young. 2018. Super characters: A conversion from sentiment classification to image classification. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 309–315, Brussels, Belgium. Association for Computational Linguistics.
- Zhiqing Sun and Zhi-Hong Deng. 2018. Unsupervised neural word segmentation for Chinese via segmental language modeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4915–4920, Brussels, Belgium. Association for Computational Linguistics.
- Ilya Sutskever, James Martens, and Geoffrey Hinton. 2011. Generating text with recurrent neural

- networks. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, ICML '11, pages 1017–1024, New York, NY, USA. ACM.
- Dan Tito Svenstrup, Jonas Hansen, and Ole Winther. 2017. Hash embeddings for efficient word representations. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Samson Tan, Shafiq Joty, Lav Varshney, and Min-Yen Kan. 2020. Mind your inflections! Improving NLP for non-standard Englishes with Base-Inflection Encoding. In *Proceedings of* the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5647–5663, Online. Association for Computational Linguistics.
- Yi Tay, Vinh Q. Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. 2021. Charformer: Fast character transformers via gradient-based subword tokenization.
- Yee Whye Teh. 2006. A hierarchical bayesian language model based on Pitman-Yor processes. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL*, pages 985–992, Sydney, Australia. Association for Computational Linguistics.
- David Vilar and Marcello Federico. 2021. A statistical extension of byte-pair encoding. In *Proceedings of the 18th International Conference on Spoken Language Translation (IWSLT 2021)*, pages 263–275, Bangkok, Thailand (online). Association for Computational Linguistics.
- Sami Virpioja, Ville T. Turunen, Sebastian Spiegler, Oskar Kohonen, and Mikko Kurimo. 2011. Empirical comparison of evaluation methods for unsupervised learning of morphology. *Traitement Automatique des Langues*, 52(2):45–90.
- Changhan Wang, Kyunghyun Cho, and Jiatao Gu. 2019. Neural machine translation with bytelevel subwords.
- Chong Wang, Yining Wang, Po-Sen Huang, Abdelrahman Mohamed, Dengyong Zhou, and

- Li Deng. 2017. Sequence modeling via segmentations. In *International Conference on Machine Learning*, pages 3674–3683.
- Haohan Wang, Peiyan Zhang, and Eric P Xing. 2020a. Word shape matters: Robust machine translation with visual embedding. *arXiv* preprint arXiv:2010.09997.
- Lihao Wang, Zongyi Li, and Xiaoqing Zheng. 2021a. Unsupervised word segmentation with bi-directional neural language model. *CoRR*, abs/2103.01421.
- Xinyi Wang, Sebastian Ruder, and Graham Neubig. 2021b. Multi-view subword regularization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 473–482, Online. Association for Computational Linguistics.
- Zihan Wang, Karthikeyan K, Stephen Mayhew, and Dan Roth. 2020b. Extending multilingual BERT to low-resource languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2649–2656, Online. Association for Computational Linguistics.
- Jonathan J. Webster and Chunyu Kit. 1992. Tokenization as the initial phase in NLP. In COLING 1992 Volume 4: The 14th International Conference on Computational Linguistics.
- J Gerard Wolff. 1975. An algorithm for the segmentation of an artificial language analogue. *British Journal of PsychologyJ*, 66:79–90.
- F. Wood, J. Gasthaus, C. Archambeau, L. James, and Y.W. Teh. 2011. The sequence memoizer. *Communications of the ACM*, 54(2):91–98.
- Shijie Wu and Mark Dredze. 2019. Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 833–844, Hong Kong, China. Association for Computational Linguistics.
- Jingjing Xu, Hao Zhou, Chun Gan, Zaixiang Zheng, and Lei Li. 2021. Vocabulary learning via optimal transport for neural machine translation. In *Proceedings of the 59th Annual Meeting*

- of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7361–7373, Online. Association for Computational Linguistics.
- Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2021. Byt5: Towards a token-free future with pre-trained byte-to-byte models.
- David Yarowsky and Richard Wicentowski. 2000. Minimally supervised morphological analysis by multimodal alignment. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 207–216, Hong Kong. Association for Computational Linguistics
- Xinsong Zhang, Pengshuai Li, and Hang Li. 2021. AMBERT: A pre-trained language model with multi-grained tokenization. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 421–435, Online. Association for Computational Linguistics.
- Giulio Zhou. 2018. Morphological zero-shot neural machine translation. Master's thesis, University of Edinburgh.