



**HAL**  
open science

# Imputing out-of-vocabulary embeddings with LOVE makes language models robust with little cost

Lihu Chen, Gaël Varoquaux, Fabian Suchanek

► **To cite this version:**

Lihu Chen, Gaël Varoquaux, Fabian Suchanek. Imputing out-of-vocabulary embeddings with LOVE makes language models robust with little cost. ACL 2022 - 60th Annual Meeting of the Association for Computational Linguistics, May 2022, Dublin, Ireland. hal-03613101

**HAL Id: hal-03613101**

**<https://hal.science/hal-03613101>**

Submitted on 19 Mar 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Imputing Out-of-Vocabulary Embeddings with LOVE Makes Language Models Robust with Little Cost

Lihu Chen<sup>1</sup>, Gaël Varoquaux<sup>2</sup>, Fabian M. Suchanek<sup>1</sup>

<sup>1</sup> LTCI & Télécom Paris & Institut Polytechnique de Paris, France

<sup>2</sup> Soda, Inria Saclay & CEA & Université Paris-Saclay, France

{lihu.chen, fabian.suchanek}@telecom-paris.fr

{gael.varoquaux}@inria.fr

## Abstract

State-of-the-art NLP systems represent inputs with word embeddings, but these are brittle when faced with Out-of-Vocabulary (OOV) words. To address this issue, we follow the principle of mimick-like models to generate vectors for unseen words, by learning the behavior of pre-trained embeddings using only the surface form of words. We present a simple contrastive learning framework, LOVE, which extends the word representation of an existing pre-trained language model (such as BERT), and makes it robust to OOV with few additional parameters. Extensive evaluations demonstrate that our lightweight model achieves similar or even better performances than prior competitors, both on original datasets and on corrupted variants. Moreover, it can be used in a plug-and-play fashion with FastText and BERT, where it significantly improves their robustness.

## 1 Introduction

Word embeddings represent words as vectors (Mikolov et al., 2013a,b; Pennington et al., 2014). They have been instrumental in neural network approaches that brought impressive performance gains to many natural language processing (NLP) tasks. These approaches use a fixed-size vocabulary. Thus they can deal only with words that have been seen during training. While this works well on many benchmark datasets, real-world corpora are typically much noisier and contain Out-of-Vocabulary (OOV) words, i.e., rare words, domain-specific words, slang words, and words with typos, which have not been seen during training. Model performance deteriorates a lot with unseen words, and minor character perturbations can flip the prediction of a model (Liang et al., 2018; Belinkov and Bisk, 2018; Sun et al., 2020; Jin et al., 2020). Simple experiments (Figure 1) show that the addition of typos to datasets degrades the performance for text classification models considerably.

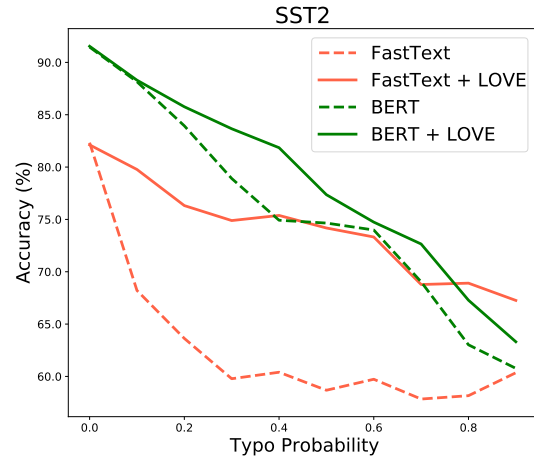


Figure 1: Performances of existing word embeddings as we gradually add typos to the datasets. Using our model, LOVE, to produce vectors for OOV words makes the models more robust.

To alleviate this problem, pioneering work pre-trained word embeddings with morphological features (sub-word tokens) on large-scale datasets (Wieting et al., 2016; Bojanowski et al., 2017; Heizerling and Strube, 2017; Zhang et al., 2019). One of the most prominent works in this direction is FastText (Bojanowski et al., 2017), which incorporates character n-grams into the skip-gram model. With FastText, vectors of unseen words can be imputed by summing up the n-gram vectors. However, these subword-level models come with great costs: the requirements of pre-training from scratch and high memory footprint. Hence, several simpler approaches have been developed, e.g., MIMICK (Pinter et al., 2017), BoS (Zhao et al., 2018) and KVQ-FH (Sasaki et al., 2019). These use only the surface form of words to generate vectors for unseen words, through learning from pre-trained embeddings.

Although MIMICK-like models can efficiently reduce the parameters of pre-trained representa-

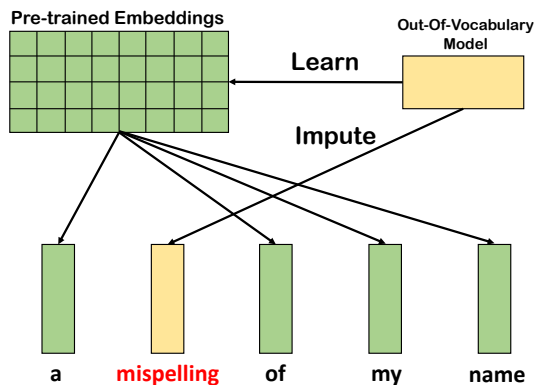


Figure 2: Our lightweight OOV model, LOVE, learns the behavior of pre-trained embeddings (e.g., FastText and BERT), and is then able to impute vectors for unseen words. LOVE can enhance the robustness of existing word representations in a plug-and-play fashion.

tions and alleviate the OOV problem, two main challenges remain. First, the models remain bound in the trade-off between complexity and performance: The original MIMICK is lightweight but does not produce high-quality word vectors consistently. BoS and KVQ-FH obtain better word representations but need more parameters. Second, these models cannot be used with existing pre-trained language models such as BERT. It is these models, however, to which we owe so much progress in the domain (Peters et al., 2018; Devlin et al., 2019; Yang et al., 2019; Liu et al., 2020). To date, these high-performant models are still fragile when dealing with rare words (Schick and Schütze, 2020), misspellings (Sun et al., 2020) and domain-specific words (El Boukkouri et al., 2020).

We address these two challenges head-on: we design a new contrastive learning framework to learn the behavior of pre-trained embeddings, dubbed LOVE, Learning Out-of-Vocabulary Embeddings. Our model builds upon a memory-saving mixed input of character and subwords instead of n-gram characters. It encodes this input by a lightweight Positional Attention Module. During training, LOVE uses novel types of data augmentation and hard negative generation. The model is then able to produce high-quality word representations that are robust to character perturbations, while consuming only a fraction of the cost of existing models. For instance, LOVE with 6.5M parameters can obtain similar representations as the original FastText model with more than 900M parameters. What is more, our model can be used in a plug-and-play

fashion to robustify existing language models. We find that using LOVE to produce vectors for unseen words improves the performance of FastText and BERT by around 1.4-6.8 percentage points on noisy text – without hampering their original capabilities (As shown in Figure 2).

In the following, Section 2 discusses related work, Section 3 introduces preliminaries, Section 4 presents our approach, Section 5 shows our experiments, and Section 6 concludes. The appendix contains additional experiments and analyses. Our code and data is available at <https://github.com/tigerchen52/LOVE>

## 2 Related Work

### 2.1 Character-level Embeddings

To address OOV problems, some approaches inject character-level features into word embeddings during the pre-training (Wieting et al., 2016; Cao and Rei, 2016; Bojanowski et al., 2017; Heinzerling and Strube, 2017; Kim et al., 2018; Li et al., 2018; Üstün et al., 2018; Piktus et al., 2019; Zhu et al., 2019; Zhang et al., 2019; Hu et al., 2019). One drawback of these methods is that they need to pre-train on a large-scale corpus from scratch. Therefore, simpler models have been developed, which directly mimic the well-trained word embeddings to impute vectors for OOV words. Some of these methods use only the surface form of words to generate embeddings for unseen words (Pinter et al., 2017; Zhao et al., 2018; Sasaki et al., 2019; Fukuda et al., 2020; Jinman et al., 2020), while others use both surface and contextual information to create OOV vectors (Schick and Schütze, 2019a,b). In both cases, the models need an excessive number of parameters. FastText, e.g., uses ~2 million n-gram characters to impute vectors for unseen words.

### 2.2 Pre-trained Language Models

Currently, the state-of-the-art word representations are pre-trained language models, such as ELMo (Peters et al., 2018), BERT (Devlin et al., 2019) and XLnet (Yang et al., 2019), which adopt subwords to avoid OOV problems. However, BERT is brittle when faced with rare words (Schick and Schütze, 2020) and misspellings (Sun et al., 2020). To make BERT more robust, CharacterBERT (El Boukkouri et al., 2020) and CharBERT (Ma et al., 2020) infuse character-level features into BERT and pre-train the variant from

scratch. This method can significantly improve the performance and robustness of BERT, but requires pre-training an adapted transformer on a large amount of data. Another work on combating spelling mistakes recommends placing a word corrector before downstream models (Pruthi et al., 2019), which is effective and reusable. The main weakness of this method is that an error generated by the word corrector propagates to downstream tasks. For example, converting “*aleph*” to “*alpha*” may break the meaning of a mathematical statement. And indeed, using the word corrector consistently leads to a drop (0.5-2.0 percentage points) in BERT’s performance on the SST dataset (Socher et al., 2013).

### 2.3 Contrastive Learning

The origin of contrastive learning can be traced back to the work by Becker and Hinton (1992) and Bromley et al. (1993). This method has achieved outstanding success in self-supervised representation learning for images (Oord et al., 2018; Hjelm et al., 2018; He et al., 2020; Chen et al., 2020; Grill et al., 2020). The contrastive learning framework learns representations from unlabeled data by pulling positive pairs together and pushing negative pairs apart. For training, the positive pairs are often obtained by taking two randomly augmented versions of the same sample and treating the other augmented examples within a mini-batch as negative examples (Chen et al., 2017, 2020). The most widely used loss is the infoNCE loss (or contrastive loss) (Hjelm et al., 2018; Logeswaran and Lee, 2018; Chen et al., 2020; He et al., 2020). Although many approaches adopt contrastive learning to represent sentences (Giorgi et al., 2020; Wu et al., 2020; Gao et al., 2021), it has so far not been applied to word representations.

	Input	Encoder	Loss
MIMICK (2017)	character sequence {s, p, e, l, l}	RNNs	$\mathcal{L}_{\text{dis}}$
BoS (2018)	n-gram subword {spe, pel, ell}	SUM	$\mathcal{L}_{\text{dis}}$
KVQ-FH (2019)	adapted n-gram subword {spe, pel, ell}	Attention	$\mathcal{L}_{\text{dis}}$

Table 1: Details of different mimick-like models, with the word `spell` as an example.

## 3 Preliminaries

### 3.1 Mimick-like Model

Given pre-trained word embeddings, and given an OOV word, the core idea of MIMICK (Pinter et al., 2017) is to impute an embedding for the OOV word using the surface form of the word, so as to mimic the behavior of the known embeddings. BoS (Zhao et al., 2018), KVQ-FH (Sasaki et al., 2019), Robust Backed-off Estimation (Fukuda et al., 2020), and PBoS (Jinman et al., 2020) work similarly, and we refer to them as mimick-like models.

Formally, we have a fixed-size vocabulary set  $\mathcal{V}$ , with an embedding matrix  $\mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times m}$ , in which each row is a word vector  $\mathbf{u}_w \in \mathbb{R}^m$  for the word  $w$ . A mimick-like model aims to impute a vector  $\mathbf{v}_w$  for an arbitrary word  $w \notin \mathcal{V}$ . The training objective of mimick-like models is to minimize the expected distance between  $\mathbf{u}_w$  and  $\mathbf{v}_w$  pairs:

$$\mathcal{L}_{\text{dis}} = \frac{1}{|\mathcal{V}|} \sum_{w \in \mathcal{V}} \psi(\mathbf{u}_w, \mathbf{v}_w) \quad (1)$$

Here,  $\psi(\cdot)$  is a distance function, e.g., the Euclidean distance  $\psi = \|\mathbf{u}_w - \mathbf{v}_w\|_2^2$  or the cosine distance  $\psi = 1 - \cos(\mathbf{u}_w, \mathbf{v}_w)$ . The vector  $\mathbf{v}_w$  is generated by the following equation:

$$\mathbf{v}_w = \phi(\zeta(w)), \text{ for } w \in \mathcal{V} \text{ or } w \notin \mathcal{V} \quad (2)$$

Here,  $\zeta(\cdot)$  is a function that maps  $w$  to a list of subunits based on the surface form of the word (e.g., a character or subword sequence). After that, the sequence is fed into the function  $\phi(\cdot)$  to produce vectors, and the inside structure can be CNNs, RNNs, or a simple summation function. After training, the model can impute vectors for arbitrary words. Table 1 shows some features of three mimick-like models.

### 3.2 Contrastive Learning

Contrastive learning methods have achieved significant success for image representation (Oord et al., 2018; Chen et al., 2020). The core idea of these methods is to encourage learned representations for positive pairs to be close, while pushing representations from sampled negative pairs apart. The widely used contrastive loss (Hjelm et al., 2018; Logeswaran and Lee, 2018; Chen et al., 2020; He et al., 2020) is defined as:

$$\ell_{\text{cl}} = -\log \frac{e^{\text{sim}(\mathbf{u}_i^+ \mathbf{u}^+)/\tau}}{e^{\text{sim}(\mathbf{u}_i^+ \mathbf{u}^+)/\tau} + \sum e^{\text{sim}(\mathbf{u}_i^+ \mathbf{u}^-)/\tau}} \quad (3)$$

Here,  $\tau$  is a temperature parameter,  $\text{sim}(\cdot)$  is a similarity function such as cosine similarity, and  $(u_i, u^+)$ ,  $(u_i, u^-)$  are positive pairs and negative pairs, respectively (assuming that all vectors are normalized). During training, positive pairs are usually obtained by augmentation for the same sample, and negative examples are the other samples in the mini-batch. This process learns representations that are invariant against noisy factors to some extent.

## 4 Our Approach: LOVE

LOVE (Learning Out-of-Vocabulary Embeddings) draws on the principles of contrastive learning to maximize the similarity between target and generated vectors, and to push apart negative pairs. An overview of our framework is shown in Figure 3. It is inspired by work in visual representation learning (Chen et al., 2020), but differs in that one of the positive pairs is obtained from pre-trained embeddings instead of using two augmented versions. We adopt five novel types of word-level augmentations and a lightweight Positional Attention Module in this framework. Moreover, we find that adding hard negatives during training can effectively yield better representations. We removed the nonlinear projection head after the encoder layer, because its improvements are specific to the representation quality in the visual field. Furthermore, our approach is not an unsupervised contrastive learning framework, but a supervised learning approach.

Our framework takes a word from the original vocabulary and uses data augmentation to produce a corruption of it. For example, "misspelling" becomes "mispelling" after dropping one letter "s". Next, we obtain a target vector from the pre-trained embeddings for the original word, and we generate a vector for the corrupted word. These two vectors are a pair of positive samples, and we maximize the similarity between them while making the distance of negative pairs (other samples in the same mini-batch) as large as possible. As mentioned before, we use the contrastive loss as an objective function (Eq 3). There are five key ingredients in the framework that we will detail in the following (similar to the ones in Table 1): the Input Method, the Encoder, the Loss Function, our Data Augmentation, and the choice of Hard Negatives.

### 4.1 Input Method

Our goal is to use the surface form to impute vectors for words. The question is thus how to design the function  $\zeta(\cdot)$  mentioned in Section 3.1 to represent each input word. MIMICK (Pinter et al., 2017) straightforwardly uses the character sequence (see Table 1). This, however, loses the information of morphemes, i.e., sequences of characters that together contribute a meaning. Hence, FastText (Bojanowski et al., 2017) adopts character n-grams. Such n-grams, however, are highly redundant. For example, if we use substrings of length 3 to 5 to represent the word `misspelling`, we obtain a list with 24 n-gram characters – while only the substrings `{mis, spell, ing}` are the three crucial units to understand the word. Hence, like BERT, we use WordPiece (Wu et al., 2016) with a vocabulary size of around 30000 to obtain meaningful subwords of the input word. For the word `misspelling`, this yields `{miss, ##pel, ##ling}`. However, if we just swap two letters (as by a typo), then the sequence becomes completely different: `{mi, ##sp, ##sell, ##ing}`. Therefore, we use both the character sequence and subwords (Figure A1).

We shrink our vocabulary by stemming all words and keeping only the base form of each word, and by removing words with numerals. This decreases the size of vocabulary from 30 000 to 21 257 without degrading performance too much (Section A.1).

### 4.2 Encoder

Let us now design the function  $\phi(\cdot)$  mentioned in Section 3.1. We are looking for a function that can encode both local features and global features. Local features are character n-grams, which provide robustness against minor variations such as character swaps or omissions. Global features combine local features regardless of their distance. For the word `misspelling`, a pattern of prefix and suffix `mis+ing` can be obtained by combining the local information at the beginning and the end of the word. Conventional CNNs, RNNs, and self-attention cannot extract such local and global information at the same time. Therefore, we design a new **Positional Attention Module**. Suppose we have an aforementioned mixed input sequence and a corresponding embedding matrix  $\mathbf{V} \in \mathbb{R}^{|\mathcal{V}| \times d}$  where  $d$  is the dimension of vectors. Then the input can be represented by a list of vectors:  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \in \mathbb{R}^{n \times d}$  where  $n$  is the

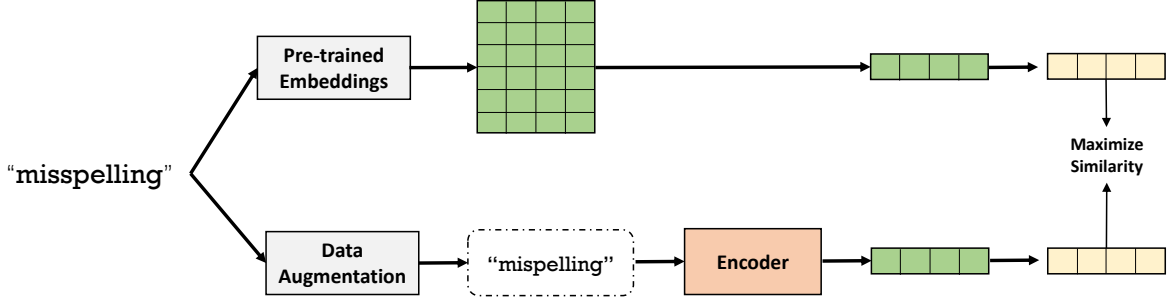


Figure 3: The framework of LOVE with an example of the word `misspelling`.

length of the input. To extract local information, we first adopt positional attention to obtain n-gram features, and then feed them into a conventional self-attention layer to combine them in a global way. This can be written as:

$$\bar{\mathbf{X}} = \text{SA}(\text{PA}(\mathbf{X})) \mathbf{W}^O \quad (4)$$

Here, SA is a standard multi-head self-attention and PA is a positional attention.  $\mathbf{W}^O \in \mathbb{R}^{d_V \times d_O}$  is a trainable parameter matrix, where  $d_V$  are the dimensions of values in SA and PA, and  $d_O$  is that of  $\bar{\mathbf{X}}$ . As for the Positional Attention, we adopt absolute sinusoidal embeddings (Vaswani et al., 2017) to compute positional correlations:

$$\text{PA}(\mathbf{X}) = \text{Softmax}\left(\frac{\mathbf{P}\mathbf{P}^T}{\sqrt{d}}\right)(\mathbf{X}\mathbf{W}^V) \quad (5)$$

Here,  $\mathbf{P} \in \mathbb{R}^{n \times d}$  are the position embeddings, and  $\mathbf{W}^V \in \mathbb{R}^{d \times d_V}$  are the corresponding parameters. More details about the encoder are in Appendix C.4.

### 4.3 Loss Function

In this section, we focus on the loss function  $\mathcal{L}(\cdot)$ . Mimick-like models often adopt the mean squared error (MSE), which tries to give words with the same surface forms similar embeddings. However, the MSE only pulls positive word pairs closer, and does not push negative word pairs apart. Therefore, we use the contrastive loss instead (Equation 3). Wang and Isola (2020) found that the contrastive loss optimizes two key properties: *Alignment* and *Uniformity*. The Alignment describes the expected

distance (closeness) between positive pairs:

$$\ell_{\text{align}} \triangleq \mathbb{E}_{(x,y) \sim p_{\text{pos}}} \psi(\mathbf{u}_x, \mathbf{u}_y) \quad (6)$$

Here,  $p_{\text{pos}}$  is the distribution of positive pairs. The Uniformity measures whether the learned representations are uniformly distributed in the hypersphere:

$$\ell_{\text{uniform}} \triangleq \log \mathbb{E}_{(x,y) \stackrel{i.i.d.}{\sim} p_{\text{data}}} e^{-t \cdot \psi(\mathbf{u}_x, \mathbf{u}_y)} \quad (7)$$

Here,  $p_{\text{data}}$  is the data distribution and  $t > 0$  is a parameter. The two properties are consistent with our expected word representations: positive word pairs should be kept close and negative word pairs should be far from each other, finally scattered over the hypersphere.

### 4.4 Data Augmentation and Hard Negatives

Our positive word pairs are generated by data augmentation, which can increase the amount of training samples by using existing data. We use various strategies (Figure 4) to increase the diversity of our training samples: (1) Swap two adjacent characters, (2) Drop a character, (3) Insert a new character, (4) Replace a character according to keyboard distance, (5) Replace the original word by a synonymous word. The first four augmentations are originally designed to protect against adversarial attacks (Pruthi et al., 2019). We add the synonym replacement strategy to keep semantically similar words close in the embedding space – something that cannot be achieved by the surface form alone. Specifically, a set of synonyms is obtained by retrieving the nearest neighbors from pre-trained embeddings like FastText.

<b>Swap</b>	misspelling -> misspelling
<b>Drop</b>	misspelling -> misspelling
<b>Insert</b>	misspelling -> misspelling
<b>Keyboard</b>	misspelling -> mosspelling
<b>Synonym</b>	misspelling -> heterography

Figure 4: Illustrations of different augmentations for the word misspelling.

Negative word pairs are usually chosen randomly from the mini-batch. However, we train our model to be specifically resilient to *hard negatives* (or *difficult negatives*), i.e., words with similar surface forms but different meanings (e.g., *misspelling* and *dispelling*). To this end, we add a certain number of hard negative samples (currently 3 of them) to the mini-batch, by selecting word pairs that are not synonyms and have a small edit distance.

#### 4.5 Mimicking Dynamical Embeddings

Pre-trained Language Models (e.g., ELMo (Peters et al., 2018) and BERT (Devlin et al., 2019)) dynamically generate word representations based on specific contexts, which cannot be mimicked directly. To this end, we have two options: We can either learn the behavior of the *input embeddings* in BERT before the multi-layer attentions or mimic the static *distilled embeddings* (Bommasani et al., 2020; Gupta and Jaggi, 2021).

We use the BERT as an example to explain these two methods. Suppose we have a subword sequence after applying WordPiece to a sentence:  $W = \{w_1, w_2, \dots, w_n\}$ . For the subword sequence  $W$ , BERT first represents it as a list of subword embeddings:  $\mathbf{E}^{in} = \{e_1^{sub}, e_2^{sub}, \dots, e_n^{sub}\}$ . We refer to this static representation as the Input Embedding of BERT, and we can use our model to mimic the behavior of this part. We call this method *mimicking input embeddings*. For ease of implementation, we learn only from the words that are not separated into pieces. After that step, BERT applies a multi-layer multi-head attention to the input embeddings  $\mathbf{E}^{in}$ , which yields a contextual representation for

each subword:  $\mathbf{E}^{out} = \{e_1^{out}, e_2^{out}, \dots, e_n^{out}\}$ . However, these contextual representations vary according to the input sentence and we cannot learn from them directly. Instead, we choose to mimic the distilled static embeddings from BERT, which are obtained by pooling (max or average) the contextual embeddings of the word in different sentences. We call this method *mimicking distilled embeddings*. The latter allows for better word representations, while the former does not require training on a large-scale corpus. Our empirical studies show that mimicking distilled embeddings performs only marginally better. Therefore, we decided to rather learn the input embeddings of BERT, which is simple yet effective

#### 4.6 Plug and Play

One of the key advantages of our model is that it can be used as a plug-in for other models. For models with static word embeddings like FastText, one can simply use our model to generate vectors for unseen words. For models with dynamic word embeddings like BERT, if a single word is tokenized into several parts, e.g.  $misspelling = \{miss, \#pel, \#ling\}$ , we regard it as an OOV word. Then, we replace the embeddings of the subwords by a single embedding produced by our model before the attention layer. Our final enhanced BERT model has 768 dimensions and 16M parameters. Note that the BERT-base model has ~110M parameters and its distilled one has ~550M parameters.

### 5 Experiments

#### 5.1 Evaluation Datasets

There are two main methods to evaluate word representations: Intrinsic and Extrinsic. Intrinsic evaluations measure syntactic or semantic relationships between words directly, e.g., word similarity in word clusters. Extrinsic evaluations measure the performance of word embeddings as input features to a downstream task, e.g., named entity recognition (NER) and text classification. Several studies have shown that there is no consistent correlation between intrinsic and extrinsic evaluation results (Chiu et al., 2016; Faruqui et al., 2016; Wang et al., 2019). Hence, we evaluate our representation by both intrinsic and extrinsic metrics. Specifically, we use 8 intrinsic datasets (6 word similarity and 2 word cluster tasks): RareWord (Luong et al., 2013), SimLex (Hill et al., 2015), MTurk (Halawi et al., 2012), MEN (Bruni et al., 2014), WordSim (Agirre

	parameters		RareWord	SimLex	Word Similarity			SimVerb	Word Cluster		Avg
	embedding	others			MTurk	MEN	WordSim		AP	BLESS	
FastText (2017)	969M	-	48.1	30.4	66.9	78.1	68.2	25.7	58.0	71.5	55.9
MIMICK (2017)	9M	517K	27.1	15.9	32.5	36.5	15.0	7.5	<b>59.3</b>	<b>72.0</b>	33.2
BoS (2018)	500M	-	<b>44.2</b>	<u>27.4</u>	<u>55.8</u>	<u>65.5</u>	<u>53.8</u>	<u>22.1</u>	41.8	39.0	<u>43.7</u>
KVQ-FH (2019)	12M	-	42.4	20.4	55.2	63.4	53.1	16.4	39.1	42.5	41.6
LOVE	6.3M	200K	42.2	<b>35.0</b>	<b>62.0</b>	<b>68.8</b>	<b>55.1</b>	<b>29.4</b>	<u>53.2</u>	<u>51.5</u>	<b>49.7</b>

Table 2: Performance on the intrinsic tasks, measured as Spearman’s  $\rho$  and purity for word similarity and clustering. Best performance among the mimick-like models in bold, second-best underlined.

	parameters		SST2		MR		CoNLL-03		BC2GM		Avg
	embedding	others	original	+typo	original	+typo	original	+typo	original	+typo	
FastText (2017)	969M	-	82.3	60.5	73.3	62.2	86.4	66.3	71.8	53.4	69.5
Edit Distance	969M	-	-	67.4	-	68.3	-	76.2	-	66.6	-
MIMICK (2018)	9M	517K	69.7	62.3	<u>73.6</u>	61.4	68.0	65.2	56.6	56.7	64.2
BoS (2018)	500M	-	79.7	<u>72.6</u>	<u>73.6</u>	<b>69.5</b>	<b>79.5</b>	68.6	<b>66.4</b>	<u>61.5</u>	<u>71.5</u>
KVQ-FH (2019)	12M	-	77.8	71.4	72.9	66.5	73.1	<b>70.4</b>	46.2	53.5	66.5
LOVE	6.3M	200K	<b>81.4</b>	<b>73.2</b>	<b>74.4</b>	<u>66.7</u>	<u>78.6</u>	<u>69.7</u>	<u>64.7</u>	<b>63.8</b>	<b>71.6</b>

Table 3: Performance on the extrinsic tasks, measured as accuracy and F1 (five runs of different learning rates) for text classification and NER, respectively. Typos are generated by simulated errors of an OCR engine (Ma, 2019). The speed of producing word vectors with Edit Distance and LOVE is 380s/10K words and 0.9s/10K words, respectively.

et al., 2009), Simverb (Agirre et al., 2009), AP (Almuhareb, 2006) and BLESS (Baroni and Lenci, 2011). We use four extrinsic datasets (2 text classification and 2 NER tasks): SST2 (Socher et al., 2013), MR (Pang and Lee, 2005), CoNLL-03 (Sang and De Meulder, 2003) and BC2GM (Smith et al., 2008). It is worth noting that the RareWord dataset contains many long-tail words and the BC2GM is a domain-specific NER dataset. All data augmentations and typo simulations are implemented by NLPAUG<sup>1</sup>. Appendix B provides more details on our datasets and experimental settings.

## 5.2 Results on Intrinsic Tasks

Table 2 shows the experimental results on 8 intrinsic tasks. Compared to other mimick-like models, our model achieves the best average score across 8 datasets while using the least number of parameters. Specifically, our model performs best on 5 word similarity tasks, and second-best on the word cluster tasks. Although there is a gap between our model and the original FastText, we find our performance acceptable, given that our model is 100x times smaller.

## 5.3 Results on Extrinsic Tasks

Table 3 shows the results on four downstream datasets and their corrupted version. In this experiment, we introduce another non-trivial baseline: Edit Distance. For each corrupted word, we find

the most similar word from a vocabulary using edit distance and then use the pre-trained vectors of the retrieved word. Considering the time cost, we only use the first 20K words appearing in FastText (2M words) as reference vocabulary.

The typo words are generated by simulating post-OCR errors. For the original datasets, our model obtains the best results across 2 datasets and the second-best on NER datasets compared to other mimick-like models. For the corrupted datasets, the performance of the FastText model decreases a lot and our model is the second best but has very close scores with BoS consistently. Compared to other mimick-like models, our model with 6.5M achieves the best average score. Although Edit Distance can effectively restore the original meaning of word, it is 400x times more time-consuming than our model.

## 5.4 Robustness Evaluation

In this experiment, we evaluate the robustness of our model by gradually adding simulated post-OCR typos (Ma, 2019). Table 4 shows the performances on SST2 and CoNLL-03 datasets. We observe that our model can improve the robustness of the original embeddings without degrading their performance. Moreover, we find our model can make FastText more robust compared to other commonly used methods against unseen words: a generic UNK token or character-level representation of neural networks. Figure 5 shows the robust-

<sup>1</sup><https://github.com/makcedward/nlpaug>



Typo Probability	SST2						CoNLL-03						Avg
	original	10%	30%	50%	70%	90%	original	10%	30%	50%	70%	90%	
Static Embeddings													
FastText	<b>82.3</b>	68.2	59.8	56.7	57.8	60.3	<b>86.4</b>	81.6	78.9	73.9	70.2	63.4	70.0
FastText + LOVE	82.1	<b>79.8</b>	<b>74.9</b>	<b>74.2</b>	<b>68.8</b>	<b>67.2</b>	86.3	<b>84.7</b>	<b>81.8</b>	<b>77.5</b>	<b>73.1</b>	<b>71.3</b>	<b>76.8</b>
Dynamical Embeddings													
BERT	<b>91.5</b>	88.2	78.9	74.7	69.0	60.1	<b>91.2</b>	<b>89.8</b>	<b>86.2</b>	83.4	79.9	76.5	80.7
BERT + LOVE	<b>91.5</b>	<b>88.3</b>	<b>83.7</b>	<b>77.4</b>	<b>72.7</b>	<b>63.3</b>	89.9	88.3	86.1	<b>84.3</b>	<b>80.8</b>	<b>78.3</b>	<b>82.1</b>

Table 4: Robust evaluation (five runs of different learning rates) on text classification and NER under simulated post-OCR typos. We use uncased and cased BERT-base model for SST2 and CoNLL-03, respectively.

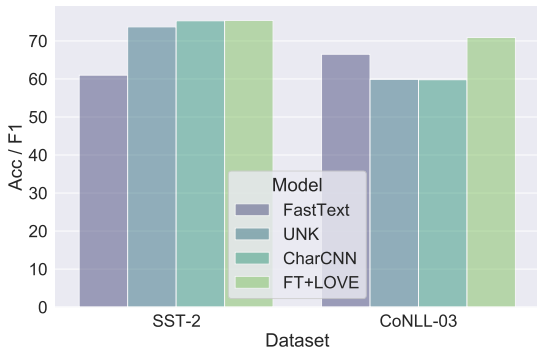


Figure 5: Evaluation of different methods based on FastText under typos.

ness check of different strategies. FastText+LOVE has a consistent improvement on both SST2 and CoNLL-03 datasets. At the same time, LOVE degrades the performance on the original datasets only marginally if at all.

## 5.5 Ablation Study

We now vary the components in our architecture (input method, encoder and loss function) to demonstrate the effectiveness of our architecture.

**Input Method.** To validate the effect of our Mixed Input strategy, we compare it with two other methods: using only the character sequence or only the subword sequence. Table 5 shows that the Mixed method achieves better representations, and any removal of char or subword information can decrease the performance.

**Encoder.** To encode the input sequence, we developed the Positional Attention Module (PAM), which first extracts ngram-like local features and then uses self-attention combine them without distance restrictions. Table 5 shows that PAM performs the best, which validates our strategy of incorporating both local and global parts inside a word. At the same time, the number of parameters

	parameters embedding	others	RareWord	SST2
The original LOVE	6.3M	200K	42.2	81.4
Varying the input method				
only use Char	299K	200K	17.7	71.5
only use Subword	6.0M	200K	25.3	76.0
Varying the encoder				
replace PAM with CNN	6.3M	270K	28.4	61.1
replace PAM with RNN	6.3M	517K	27.2	67.2
replace PAM with SA	6.3M	-	36.9	78.7
Varying the loss function				
use MSE	6.3M	200K	34.5	76.0
use $\ell_{au}(\lambda = 1.0)$	6.3M	200K	40.8	80.8
Ablation of data augmentation and hard negatives				
w/out hard negatives	6.3M	200K	37.7	78.6
w/out hard negatives and augmentation	6.3M	200K	37.8	78.2

Table 5: **Ablation studies for the architecture of LOVE**, measured as Spearman’s  $\rho$  and accuracy, respectively.

of PAM is acceptable in comparison. We visualize the attention weights of PAM in Appendix C.4, to show how the encoder extracts local and global morphological features of a word.

**Loss Function.** LOVE uses the contrastive loss, which increases alignment and uniformity. Wang and Isola (2020) proves that optimizing directly these two metrics leads to comparable or better performance than the original contrastive loss. Such a loss function can be written as:

$$\ell_{au} = \ell_{align} + \lambda \cdot \ell_{uniform} \quad (8)$$

Here,  $\lambda$  is a hyperparameter that controls the impact of  $\ell_{uniform}$ . We set this value to 1.0 because it achieves the best average score on RareWord and SST2. An alternative is to use the Mean Squared Error (MSE), as in mimick-like models. Table 5 compares the performances of these different loss functions. The contrastive loss significantly outperforms the MSE, and there is no obvious improve-

typos per sentence	SST2			
	typo-0	typo-1	typo-2	typo-3
BERT	<b>91.5</b>	77.2	73.2	69.4
Mimicking Input Embeddings				
BERT + Add	91.3	77.2	73.5	70.7
BERT + Linear (2020)	91.4	79.6	77.2	72.8
BERT + Replacement	<b>91.5</b>	81.4	78.7	73.6
Mimicking Distilled Embeddings				
BERT + Add	91.3	78.8	75.6	72.3
BERT + Linear (2020)	91.3	81.4	78.7	73.6
BERT + Replacement	91.4	<b>81.5</b>	<b>78.9</b>	<b>73.8</b>

Table 6: Performances of different strategies that work with BERT together, measured as the accuracy among five different learning rates.

ment by directly using alignment and uniformity. We also tried various temperatures  $\tau$  for the contrastive loss, and the results are shown in Table A3 in the appendix. In the end, a value of  $\tau = 0.07$  provides a good performance.

**Data Augmentation and Hard Negatives.** In Table 5, we observe that the removal of our hard negatives decreases the performance, which demonstrates the importance of semantically different words with similar surface forms.

LOVE uses five types of word augmentation. We find that removing this augmentation does not deteriorate performance too much on the word similarity task, while it causes a 0.4 point drop in the text classification task (the last row in Table 5), where data augmentations prove helpful in dealing with misspellings. We further analyze the performance of single and composite augmentations on RareWord and SST2 in the appendix in Figure A3 and Figure A4. We find that a combination of all five types yields the best results.

## 5.6 The performance of mimicking BERT

As described in Section 4.5, we can mimic the input or distilled embeddings of BERT. After learning from BERT, we use the vectors generated by LOVE to replace the embeddings of OOV subwords. Finally, these new representations are fed into the multi-layer attentions. We call this method the *Replacement* strategy. To valid its effectiveness, we compare it with two other baselines: **(1) Linear Combination** (Fukuda et al., 2020). For each subword  $e^{sub}$ , the generated vectors of word  $e^{word}$  containing the subwords are added to the subword

vectors of BERT:

$$e^{new} = (1 - \alpha) e^{sub} + \alpha e^{word} \quad (9)$$

$$\alpha = \text{sigmoid}(\mathbf{W} \cdot e^{sub})$$

where  $e^{sub} \in \mathbb{R}^d$  is a subword vector of BERT, and  $e^{word} \in \mathbb{R}^d$  is a generated vector of our model.  $\mathbf{W} \in \mathbb{R}^d$  are trainable parameters.

**(2) Add.** A generated word vector is directly added to a corresponding subword vector of BERT:

$$e^{new} = e^{sub} + e^{word} \quad (10)$$

Table 6 shows the result of these strategies. All of them can bring a certain degree of robustness to BERT without decreasing the original capability, which demonstrates the effectiveness of our framework. Second, the replacement strategy consistently performs best. We conjecture that BERT cannot restore a reasonable meaning for those rare and misspelling words that are tokenized into subwords, and our generated vectors can be located nearby the original word in the space. Third, we find mimicking distilled embeddings performs the best while mimicking input embeddings comes close. Considering that the first method needs training on large-scale data, mimicking the input embeddings is our method of choice.

## 6 Conclusion

We have presented a lightweight contrastive-learning framework, LOVE, to learn word representations that are robust even in the face of out-of-vocabulary words. Through a series of empirical studies, we have shown that our model (with only 6.5M parameters) can achieve similar or even better word embeddings on both intrinsic and extrinsic evaluations compared to other mimic-like models. Moreover, our model can be added to models with static embeddings (such as FastText) or dynamical embeddings (such as BERT) in a plug-and-play fashion, and bring significant improvements there. For future work, we aim to extend our model to languages other than English.

## 7 Acknowledgements

We sincerely thank all the reviewers for their insightful comments and helpful suggestions. This work was partially funded by ANR-20-CHIA-0012-01 (“NoRDF”).

## References

- Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Pasca, and Aitor Soroa. 2009. A study on similarity and relatedness using distributional and wordnet-based approaches.
- Abdulrahman Almuhareb. 2006. *Attributes in lexical acquisition*. Ph.D. thesis, University of Essex.
- Marco Baroni and Alessandro Lenci. 2011. How we blessed distributional semantic evaluation. In *Proceedings of the GEMS 2011 Workshop on GEometrical Models of Natural Language Semantics*, pages 1–10.
- Suzanna Becker and Geoffrey E Hinton. 1992. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163.
- Yonatan Belinkov and Yonatan Bisk. 2018. [Synthetic and natural noise both break neural machine translation](#). In *International Conference on Learning Representations*.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Rishi Bommasani, Kelly Davis, and Claire Cardie. 2020. Interpreting pretrained contextualized representations via reductions to static embeddings. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4758–4781.
- Jane Bromley, James W Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. 1993. Signature verification using a “siamese” time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(04):669–688.
- Elia Bruni, Nam-Khanh Tran, and Marco Baroni. 2014. Multimodal distributional semantics. *Journal of artificial intelligence research*, 49:1–47.
- Kris Cao and Marek Rei. 2016. A joint model for word embedding and word morphology. In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 18–26.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- Ting Chen, Yizhou Sun, Yue Shi, and Liangjie Hong. 2017. On sampling strategies for neural network-based collaborative filtering. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 767–776.
- Billy Chiu, Anna Korhonen, and Sampo Pyysalo. 2016. Intrinsic evaluation of word vectors fails to predict extrinsic performance. In *Proceedings of the 1st workshop on evaluating vector-space representations for NLP*, pages 1–6.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Hicham El Boukkouri, Olivier Ferret, Thomas Lavergne, Hiroshi Noji, Pierre Zweigenbaum, and Jun’ichi Tsujii. 2020. Characterbert: Reconciling elmo and bert for word-level open-vocabulary representations from characters. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6903–6915.
- Manaal Faruqi, Yulia Tsvetkov, Pushpendre Rastogi, and Chris Dyer. 2016. Problems with evaluation of word embeddings using word similarity tasks. *arXiv preprint arXiv:1605.02276*.
- Nobukazu Fukuda, Naoki Yoshinaga, and Masaru Kit-suregawa. 2020. Robust backed-off estimation of out-of-vocabulary embeddings. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 4827–4838.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*.
- John M Giorgi, Osvald Nitski, Gary D Bader, and Bo Wang. 2020. Declutr: Deep contrastive learning for unsupervised textual representations. *arXiv preprint arXiv:2006.03659*.
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. 2020. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*.
- Prakhar Gupta and Martin Jaggi. 2021. Obtaining better static word embeddings using contextual embedding models. *arXiv preprint arXiv:2106.04302*.
- Guy Halawi, Gideon Dror, Evgeniy Gabrilovich, and Yehuda Koren. 2012. Large-scale learning of word relatedness with constraints. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1406–1414.

- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738.
- Benjamin Heinzerling and Michael Strube. 2017. Bpemb: Tokenization-free pre-trained subword embeddings in 275 languages. *arXiv preprint arXiv:1710.02187*.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. 2018. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*.
- Ziniu Hu, Ting Chen, Kai-Wei Chang, and Yizhou Sun. 2019. Few-shot representation learning for out-of-vocabulary words. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4102–4112.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025.
- Zhao Jinman, Shawn Zhong, Xiaomin Zhang, and Yingyu Liang. 2020. Pbos: Probabilistic bag-of-subwords for generalizing word embedding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 596–611.
- Yeanchan Kim, Kang-Min Kim, Ji-Min Lee, and SangKeun Lee. 2018. Learning to generate word representations using subword information. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2551–2561.
- Bofang Li, Aleksandr Drozd, Tao Liu, and Xiaoyong Du. 2018. Subword-level composition functions for learning word embeddings. In *Proceedings of the second workshop on subword/character level models*, pages 38–48.
- Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. 2018. Deep text classification can be fooled. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 4208–4215.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Ro{bert}a: A robustly optimized {bert} pretraining approach.
- Lajanugen Logeswaran and Honglak Lee. 2018. An efficient framework for learning sentence representations. In *International Conference on Learning Representations*.
- Minh-Thang Luong, Richard Socher, and Christopher D Manning. 2013. Better word representations with recursive neural networks for morphology. In *Proceedings of the seventeenth conference on computational natural language learning*, pages 104–113.
- Edward Ma. 2019. Nlp augmentation. <https://github.com/makcedward/nlpaug>.
- Wentao Ma, Yiming Cui, Chenglei Si, Ting Liu, Shijin Wang, and Guoping Hu. 2020. Charbert: Character-aware pre-trained language model. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 39–50.
- James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv preprint cs/0506075*.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Matthew E. 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.

- Aleksandra Piktus, Necati Bora Edizel, Piotr Bojanowski, Édouard Grave, Rui Ferreira, and Fabrizio Silvestri. 2019. Misspelling oblivious word embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3226–3234.
- 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.
- Danish Pruthi, Bhuwan Dhingra, and Zachary C Lipton. 2019. Combating adversarial misspellings with robust word recognition. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5582–5591.
- Erik Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.
- Shota Sasaki, Jun Suzuki, and Kentaro Inui. 2019. Subword-based compact reconstruction of word embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3498–3508.
- Timo Schick and Hinrich Schütze. 2019a. Attentive mimicking: Better word embeddings by attending to informative contexts. *arXiv preprint arXiv:1904.01617*.
- Timo Schick and Hinrich Schütze. 2019b. Learning semantic representations for novel words: Leveraging both form and context. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6965–6973.
- Timo Schick and Hinrich Schütze. 2020. Rare words: A major problem for contextualized embeddings and how to fix it by attentive mimicking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8766–8774.
- Larry Smith, Lorraine K Tanabe, Rie Johnson nee Ando, Cheng-Ju Kuo, I-Fang Chung, Chun-Nan Hsu, Yu-Shi Lin, Roman Klinger, Christoph M Friedrich, Kuzman Ganchev, et al. 2008. Overview of biocreative ii gene mention recognition. *Genome biology*, 9(2):1–19.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip Yu, and Caiming Xiong. 2020. Adv-bert: Bert is not robust on misspellings! generating nature adversarial samples on bert. *arXiv preprint arXiv:2003.04985*.
- Ahmet Üstün, Murathan Kurfalı, and Burcu Can. 2018. Characters or morphemes: How to represent words? Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Bin Wang, Angela Wang, Fenxiao Chen, Yuncheng Wang, and C-C Jay Kuo. 2019. Evaluating word embedding models: Methods and experimental results. *APSIPA transactions on signal and information processing*, 8.
- Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016. Charagram: Embedding words and sentences via character n-grams. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1504–1515.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa, Fei Sun, and Hao Ma. 2020. Clear: Contrastive learning for sentence representation. *arXiv preprint arXiv:2012.15466*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.
- Ye Zhang and Byron Wallace. 2015. A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification. *arXiv preprint arXiv:1510.03820*.

- Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. 2019. Biowordvec, improving biomedical word embeddings with subword information and mesh. *Scientific data*, 6(1):1–9.
- Jinman Zhao, Sidharth Mudgal, and Yingyu Liang. 2018. Generalizing word embeddings using bag of subwords. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 601–606.
- Yi Zhu, Ivan Vulić, and Anna Korhonen. 2019. A systematic study of leveraging subword information for learning word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 912–932.

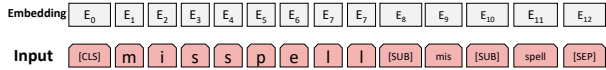


Figure A1: An illustration of our Mixed input for the word `misspell`.

## A Details of Our Approach

### A.1 Shrinking Our Model

We consider the following four methods to reduce the total parameters of our model:

- (1) **Matrix Decomposition.** The original matrix can be decomposed into two smaller matrices  $\mathbf{V} = \mathbf{U} \times \mathbf{M}$ ,  $\mathbf{U} \in \mathbb{R}^{|\mathcal{V}| \times h}$ ,  $\mathbf{M} \in \mathbb{R}^{h \times m}$  and  $h < m$ . Here, we set  $m = 300$  and  $h = 200$  respectively.
- (2) **Top Subword.** We use only the top- $k$  frequent subwords, using the word frequencies from a corpus. We set the parameter  $k = 20000$ .
- (3) **Hashing.** We use a hashing strategy to share memory for subwords aiming to reduce the parameters. We use a bucket size of 20000.
- (4) **Preprocessing.** The original vocabulary contains plurals and conjugations, therefore we stem all complete words and remove words with numerals, obtaining a new vocabulary of 21257 words.

Table A1 shows that the preprocessing method can reduce parameters very effectively while obtaining a very competitive performance.

	parameters		RareWord	SST2
	embedding	non-embedding		
Original	9M	200K	43.5	80.7
Decomposition	5.6M	200K	38.1	80.3
Top-K	6M	200K	39.2	80.1
Hashing	6M	200K	40.5	80.4
Preprocessing	6.3M	200K	<b>42.4</b>	<b>80.7</b>

Table A1: Performance of different shrinkage strategies, measured as Spearman’s  $\rho$  and accuracy, respectively. The target vectors are from `fasttext-crawl-300d-2M`.

## B Details of Our Experiments

### B.1 Training of Mimick-like Models

Our target pre-trained embeddings are those from FastText `fasttext-crawl-300d-2M`, because they provide a strong baseline. They sum up subword-level information to produce vectors for arbitrary words. We also compare to MIMICK, BoS, and KVQ-FH, which do not train on contextual words. We do not compare to Robust Backed-off Estimation (Fukuda et al., 2020) and PBoS (Jin-

Hyperparam	SST2	MR	CONLL-03	BC2GM
model	CNN	CNN	BiLSTM+CRF	BiLSTM+CRF
layer	1	1	1	1
kernel	[3,4,5]	[3,4,5]	-	-
filter	100	100	-	-
hidden size	300	300	300	300
optimizer	Adam	Adam	SGD	SGD
dropout	0.5	0.5	0.5	0.5
batch size	50	50	768	768
epoch	5	5	100	100

Table A2: Hyperparameters for extrinsic datasets.

man et al., 2020), because they need larger and more complex models. Robust Backed-off Estimation uses string matching to find the top- $k$  similar words from the entire vocabulary when imputing. Using the same target vectors, the number of parameter of BoS and PBoS are 163M and 316M, respectively. We re-train MIMICK, BoS, and KVQ-FH as baselines according to the published settings. In order to compare at the same parameter level, we use subwords for MIMICK instead of pure characters and adjust the hashing size  $H = 40K$  for KVQ-FH.

### B.2 Robustness Evaluations

As for our model, we first lower-case and tokenize each word by using WordPiece (Wu et al., 2016) with a vocabulary of 30K subwords and preprocess them by stemming and removing subwords with numerals. This yields a vocabulary of 21257 words. Each subword is represented by corresponding vectors from our model and we adopt a modified attention model to encode the subword sequence. Specifically, the layer number of this encoder is just 1 for efficiency and the hidden dimension is 300. In each block, the number of attention heads is 1 and we use fixed sinusoidal position embeddings (Vaswani et al., 2017) for positional information. To train the contrastive learning framework, we use the open-source tool (Ma, 2019) to augment a word, and use the probabilities  $\{0.07, 0.07, 0.07, 0.07, 0.36, 0.36\}$  for six augmentations: swap, drop, insert, keyboard, synonym, no-operation. Hard negatives are generated by edit distance. For each target word, we store the top-100 similar words and insert 3 of them into a mini-batch as hard negatives. The loss function is a standard contrastive loss with temperature  $\tau = 0.07$ . The optimizer is Adam and the learning rate is 0.002. The dropout rate is 0.2 and we train the model for 20 epochs in total.

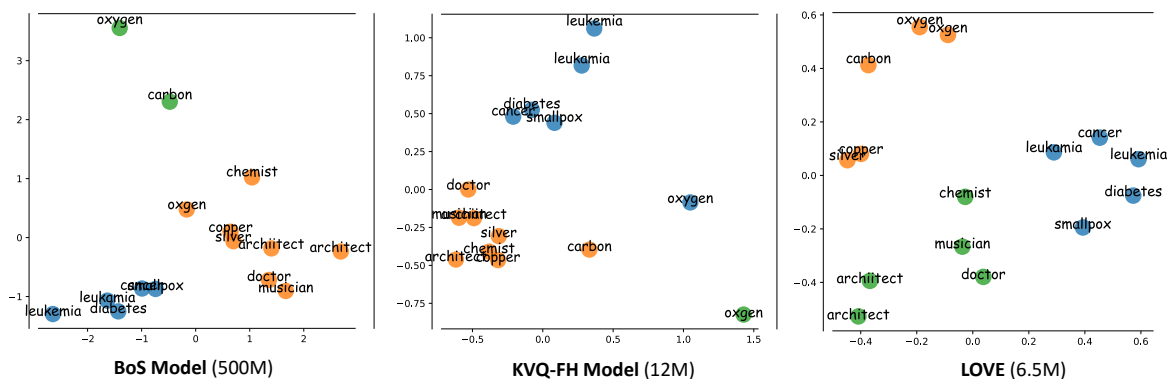


Figure A2: PCA visualizations of word vectors generated by LOVE, BoS, and KVQ-FH. Different colors mean different clusters, as predicted by K-means. There are three OOV words: oxygen, architect and leukemia.



Figure A3: Performances of different augmentations on RareWord, measured as Spearman’s  $\rho$ . Diagonal entries correspond to individual augmentation and off-diagonal entries correspond to composite augmentation.

### B.3 Intrinsic and Extrinsic Evaluations

We choose the setting discussed in Section 4 to train our model for 20 epochs, and evaluate each intrinsic task based on the vectors that the models produce. As for the extrinsic tasks, we feed word vectors into each neural network and fix them during training. We use CNNs for text classification (Zhang and Wallace, 2015) and BiLSTM+CRF for NER (Huang et al., 2015). We compare different embeddings on both intrinsic and extrinsic datasets by using generated vectors. For the word cluster tasks, the produced vectors are clustered by K-Means and then measured by Purity. The hyper-parameters of the extrinsic tasks are shown in Table A2. For each dataset, our model is trained with five learning rates  $\{5e-3, 3e-3, 1e-3, 8e-4, 5e-4\}$ . We select the best one on the develop-

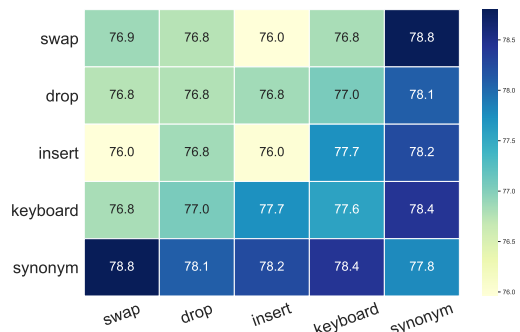


Figure A4: Performances of different augmentations on SST2, measured as accuracy. Diagonal entries correspond to individual augmentation and off-diagonal entries correspond to composite augmentation.

ment set to report its score on the test set.

To generate a corrupted dataset, we simulate post-OCR errors. We adopt the augmentation tool developed by Ma (2019) to corrupt 70% of the original words. To check the robustness of BERT, we directly finetune a BERT-base model using Huggingface (Wolf et al., 2020). During finetuning, the batch size is 16 and we train 5 epochs. We select the best model among five learning rates  $\{9e-5, 7e-5, 5e-5, 3e-5, 1e-5\}$  on the development set and report the score of the model on the test set.

### B.4 Datasets

**Intrinsic Datasets.** We use six word similarity datasets: (1) RareWord (Luong et al., 2013) (2) SimLex (Hill et al., 2015) (3) MTurk (Halawi et al., 2012) (4) MEN (Bruni et al., 2014) (5) Word-



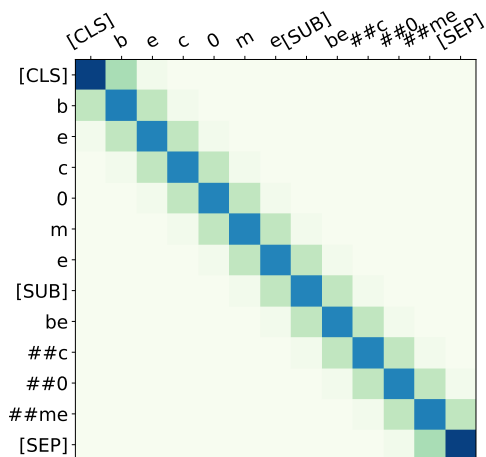


Figure A5: Visualization of positional weights for the post-OCR word `bec0me` (the correct one is `become`).

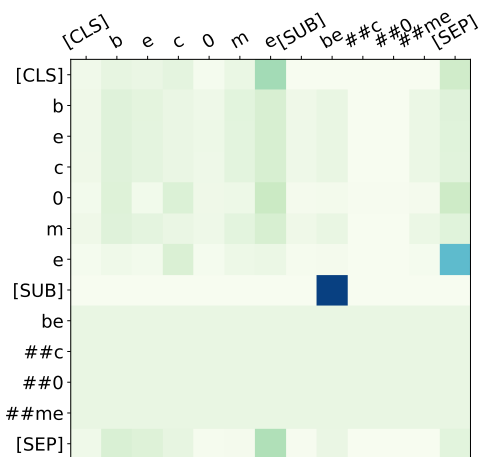


Figure A6: Visualization of self-attention weights for the post-OCR word `bec0me`.

Sim (Agirre et al., 2009), and (6) Simverb (Agirre et al., 2009). The task is scored by Spearman’s  $\rho$ , which computes the correlation between gold similarity and the similarity obtained from generated vectors. For the word cluster task, we use (1) AP (Almuhareb, 2006) and (2) BLESS (Baroni and Lenci, 2011). The generated word vectors are first clustered by K-means (MacQueen et al., 1967) and then scored by the cluster purity.

**Extrinsic Datasets.** We use both sentence-level and token-level downstream datasets to evaluate the quality of word representations. For the sentence level, we use SST2 (Socher et al., 2013) and MR (Pang and Lee, 2005), and the metric is accuracy. For the token level, we use two NER datasets: general CoNLL-03 (Sang and De Meulder, 2003) and biomedical BC2GM (Smith et al., 2008). The metric is the entity-level F1 score. As before, we select the best model among five different learning rates  $\{5e-3, 3e-3, 1e-3, 8e-4, 5e-4\}$  on the development set and then report the model score on the test set.

## C Additional Analyses

### C.1 Qualitative Analysis

To better understand the clusterings produced by LOVE, we chose 15 words from the AP dataset (Almuhareb, 2006), covering three topics (Chemical Substance, Illness, and Occupation). We added 3 corrupted words, `oxgen`, `archiitect` and `leukamia`. Figure A2 shows how LOVE, BoS, and KVQ-FH cluster these words (using a PCA

projection and K-means). All approaches space out the clusters to some degree. In particular, BoS and KVQ-FH have trouble separating professions and chemical substances. For the corrupted words, only LOVE is able to embed them close enough to their original form, so that they appear in the correct cluster.

### C.2 Effect of Augmentation for Text Classification

Figure A4 shows the performance of five augmentation strategies on the text classification task SST2. We observe that synonym is the most effective methods. The first four methods have a weaker effect, but the keyboard replacement can bring a certain degree of improvement. The results on RareWord are similar (Figure A3).

### C.3 Effect of $\tau$ in Contrastive Loss

As discussed in Chen et al. (2020), a proper temperature can yield better representation in the visual area because  $\tau$  is able to weigh the negatives by their relative hardness. As shown in Table A3, we attempt different values of temperature and find that there is no consistent  $\tau$  that makes a model work well both on intrinsic and extrinsic datasets. Hence, we choose the best performer on average, i.e.,  $\tau = 0.07$ .

### C.4 Visualization of Encoder

As mentioned before, we combine two types of attention heads (self-attention and positional attention) to encode a subword sequence. Here, we vi-

	parameters		RareWord	SST2
	embedding	non-embedding		
$\ell_{cl} (\tau = 0.03)$	9M	200K	35.0	<b>81.6</b>
$\ell_{cl} (\tau = 0.07)$	9M	200K	39.8	81.3
$\ell_{cl} (\tau = 0.12)$	9M	200K	<b>39.9</b>	81.1
$\ell_{cl} (\tau = 0.20)$	9M	200K	37.6	81.5
$\ell_{cl} (\tau = 0.50)$	9M	200K	38.3	80.6

Table A3: Performances of contrastive loss with various temperature  $\tau$ , measured as Spearman’s  $\rho$  and accuracy respectively.

sualize the attention weights on each side and show how they work. Figure A5 shows the position-dependent weights. We use sinusoidal functions to generate positional embeddings, and the weights are the dot product between these embeddings. We observe the positional weights tend to the left and right subwords in addition to themselves, which yields trigram representations.

Figure A6 shows the self-attention weights which are computed from the trigram subwords of positional attention. Hence, each subword in this figure is a trigram representation instead of a single subword representation. As we see, self-attention can capture global features regardless of distance. We take the first token [CLS] as an example, and this self-attention assigns high weights for the token e and [SEP], which constructs a representation like this: [CLS]b + me [SUB] + ##me [SEP]. This segment tells us this word starts with b and ends with me.