# Learning Multilingual Word Representations using a Bag-of-Words Autoencoder

**Stanislas Lauly**
Département d'informatique
Université de Sherbrooke
*stanislas.lauly@usherbrooke.ca*

**Alex Boulanger**
Département d'informatique
Université de Sherbrooke
*alex.boulanger@usherbrooke.ca*

**Hugo Larochelle**
Département d'informatique
Université de Sherbrooke
*hugo.larochelle@usherbrooke.ca*

## Abstract

Recent work on learning multilingual word representations usually relies on the use of word-level alignements (e.g. infered with the help of GIZA++) between translated sentences, in order to align the word embeddings in different languages. In this workshop paper, we investigate an autoencoder model for learning multilingual word representations that does without such word-level alignements. The autoencoder is trained to reconstruct the bag-of-word representation of given sentence from an encoded representation extracted from its translation. We evaluate our approach on a multilingual document classification task, where labeled data is available only for one language (e.g. English) while classification must be performed in a different language (e.g. French). In our experiments, we observe that our method compares favorably with a previously proposed method that exploits word-level alignments to learn word representations.

## 1 Introduction

Vectorial word representations have proven useful for multiple NLP tasks [1, 2]. It's been shown that meaningful representations, capturing syntactic and semantic similarity, can be learned from unlabled data. Along with a (usually smaller) set of labeled data, these representations allows to exploit unlabeled data and improve the generalization performance on some given task, even allowing to generalize out of the vocabulary observed in the labeled data only.

While the majority of previous work has concentrated on the monolingual case, recent work has started looking at learning word representations that are aligned across languages [3, 4, 5]. These representations have been applied to a variety of problems, including cross-lingual document classification [3] and phrase-based machine translation [4]. A common property of these approaches is that a word-level alignment of translated sentences is leveraged, either to derive a regularization term relating word embeddings across languages [3, 4].

In this workshop paper, we experiment with a method to learn multilingual word representations that does without word-to-word alignment of bilingual corpora during training. We only require aligned sentences and do not exploit word-level alignments (e.g. extracted using GIZA++, as is usual). To do so, we propose a multilingual autoencoder model, that learns to relate the hidden representation of paired bag-of-words sentences.

We use these representations in the context of cross-lingual document classification where labeled dataset can be available in one language, but not in another one. With the multilingual word representations, we want to learn a classifier with documents in one language and then use it on documents in another language. Our preliminary experiments suggest that our method is competitive with the representations learned by [3], which rely on word-level alignments.

In Section 2, we describe the initial autoencoder model that can learn a representation from which an input bag-of-words can be reconstructed. Then, in Section 3, we extend this autoencoder for the multilingual setting. Related work is discussed in Section 4 and experiments are presented in Section 5.

## 2  Autoencoder for Bags-of-Words

Let $\mathbf{x}$ be the bag-of-words representation of a sentence. Specifically, each $x_i$ is a word index from a fixed vocabulary of $V$ words. As this is a bag-of-words, the order of the words within $\mathbf{x}$ does not correspond to the word order in the original sentence. We wish to learn a $D$-dimensional vectorial representation of our words from a training set of sentence bag-of-words $\{\mathbf{x}^{(t)}\}_{t=1}^T$.

We propose to achieve this by using an autoencoder model that encodes an input bag-of-words $\mathbf{x}$ as the sum of its word representations (embeddings) and, using a non-linear decoder, is trained to reproduce the original bag-of-words.

Specifically, let matrix $\mathbf{W}$ be the $D \times V$ matrix whose columns are the vector representations for each word. The aggregated representation for a given bag-of-words will then be:

$$\phi(\mathbf{x}) = \sum_{i=1}^{|\mathbf{x}|} \mathbf{W}_{\cdot, x_i} \ . \tag{1}$$

To learn meaningful word representations, we wish to encourage $\phi(\mathbf{x})$ to contain information that allows for the reconstruction of the original bag-of-words $\mathbf{x}$. This is done by choosing a reconstruction loss and by designing a parametrized decoder which will be trained jointly with the word representations $\mathbf{W}$ so as to minimize this loss.

Because words are implicitly high-dimensional objects, care must be taken in the choice of reconstruction loss and decoder for stochastic gradient descent to be efficient. For instance, Dauphin et al. [6] recently designed an efficient algorithm for reconstructing binary bag-of-words representations of documents, in which the input is a fixed size vector where each element is associated with a word and is set to 1 only if the word appears at least once in the document. They use importance sampling to avoid reconstructing the whole $V$-dimensional input vector, which would be expensive.

In this work, we propose a different approach. We assume that, from the decoder, we can obtain a probability distribution over any word $\widehat{x}$ observed at the reconstruction output layer $p(\widehat{x}|\phi(\mathbf{x}))$. Then, we treat the input bag-of-words as a $|\mathbf{x}|$-trials multinomial sample from that distribution and use as the reconstruction loss its negative log-likelihood:

$$\ell(\mathbf{x}) = \sum_{i=1}^{|\mathbf{v}|} - \log p(\widehat{x} = x_i|\phi(\mathbf{x})) \ . \tag{2}$$

We now must ensure that the decoder can compute $p(\widehat{x} = x_i|\phi(\mathbf{x}))$ efficiently from $\phi(\mathbf{x})$. Specifically, we'd like to avoid a procedure scaling linear with the vocabulary size $V$, since $V$ will be very large in practice. This precludes any procedure that would compute the numerator of $p(\widehat{x} = w|\phi(\mathbf{x}))$ for each possible word $w$ separetly and normalize so it sums to one.

We instead opt for an approach borrowed from the work on neural network language models [7, 8]. Specifically, we use a probabilistic tree decomposition of $p(\widehat{x} = x_i|\phi(\mathbf{x}))$. Let's assume each word has been placed at the leaf of a binary tree. We can then treat the sampling of a word as a stochastic path from the root of the tree to one of the leaf.

We note as $\mathbf{l}(x)$ as the sequence of internal nodes in the path from the root to a given word $x$, with $l(x)_1$ always corresponding to the root. We will not as $\boldsymbol{\pi}(x)$ the vector of associated left/right branching choices on that path, where $\pi(x)_k = 0$ means the path branches left at internal node
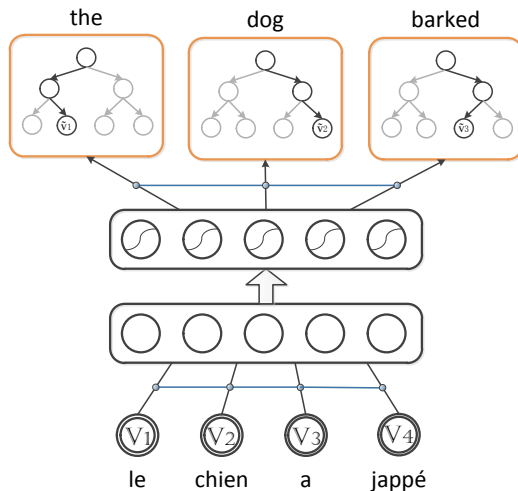
Figure 1: Illustration of a bilingual autoencoder that learns to construct the bag-of-word of the English sentence "*the dog barked*" from its French translation "*le chien a jappé*". The horizontal blue line across the input-to-hidden connections highlights the fact that these connections share the same parameters (similarly for the hidden-to-output connections).

$l(x)_k$ and branches right if $\pi(x)_k = 1$ otherwise. Then, the probability $p(\widehat{x}|\phi(\mathbf{x}))$ of a certain word $x$ observed in the bag-of-words is computed as

$$p(\widehat{x}|\phi(\mathbf{x})) = \prod_{k=1}^{|\boldsymbol{\pi}(\hat{x})|} p(\pi(\widehat{x})_k|\phi(\mathbf{x})) \tag{3}$$

where $p(\pi(\widehat{x})_k|\phi(\mathbf{x}))$ is output by the decoder. By using a full binary tree of words, the number of different decoder outputs required to compute $p(\widehat{x}|\phi(\mathbf{x}))$ will be logarithmic in the vocabulary size $V$. Since there are $|\mathbf{x}|$ words in the bag-of-words, at most $O(|\mathbf{x}| \log V)$ outputs are thus required from the decoder. This is of course a worse case scenario, since words will share internal nodes between their paths, for which the decoder output can be computed just once. As for organizing words into a tree, as in Larochelle and Lauly [9] we used a random assignment of words to the leaves of the full binary tree, which we have found to work well in practice.

Finally, we need to choose of parametrized form for the decoder. We choose the following non-linear form:

$$p(\pi(\widehat{x})_k = 1|\phi(\mathbf{x})) = \mathrm{sigm}(b_{l(\hat{x}_i)_k} + \mathbf{V}_{l(\hat{x}_i)_k,}\mathbf{h}(\mathbf{c} + \phi(\mathbf{x}))) \tag{4}$$

where $\mathbf{h}(\cdot)$ is an element wise non-linearity, $\mathbf{c}$ is a $D$-dimensional bias vector, $\mathbf{b}$ is a $(V-1)$-dimensional bias vector, $\mathbf{V}$ is a $(V-1) \times D$ matrix and $\mathrm{sigm}(a) = 1/(1 + \exp(-a))$ is the Sigmoid non-linearity. Each left/right branching probability is thus modeled with a logistic regression model applied on the non-linearly transformed representation of the input bag-of-words $\phi(\mathbf{x})$[1].

## 3   Multilingual Bag-of-words

Let's now assume that for each sentence bag-of-words $\mathbf{x}$ in some source language $\mathcal{X}$, we have an associated bag-of-words $\mathbf{y}$ for the same sentence translated in some target language $\mathcal{Y}$ by a human expert. Assuming we have a training set of such $(\mathbf{x}, \mathbf{y})$ pairs, we'd like to use it to learn representations in both languages that are aligned, such that pairs of translated words have similar representations.

---

[1]While the literature on autoencoders usually refers to the post-nonlinearity activation vector as the hidden layer, we use a different description here simply to be consistent with the representation we will use for documents in our experiments, where the non-linearity will not be used

To achieve this, we propose to augment the regular autoencoder proposed in Section 2 so that, from the sentence representation in a given language, a reconstruction can be attempted of the original sentence in the other language.

Specifically, we now define language specific word representation matrices $\mathbf{W}^x$ and $\mathbf{W}^y$, corresponding to the languages of the words in $\mathbf{x}$ and $\mathbf{y}$ respectively. Let $V^x$ and $V^y$ also be the number of words in the vocabulary of both languages, which can be different. The word representations however are of the same size $D$ in both languages. The sentence-level representation extracted by the encoder becomes

$$\phi(\mathbf{x}) = \sum_{i=1}^{|\mathbf{x}|} \mathbf{W}^x_{:,x_i} , \qquad \phi(\mathbf{y}) = \sum_{i=1}^{|\mathbf{y}|} \mathbf{W}^y_{:,y_i} . \tag{5}$$

From the sentence in either languages, we want to be able to perform a reconstruction of the original sentence in any of the languages. In particular, given a representation in any language, we'd like a decoder that can perfrom a reconstruction in language $\mathcal{X}$ and another decoder that can reconstruct in language $\mathcal{Y}$. Again, we use decoders of the form proposed in Section 2, but let the decoders of each language have their own parameters $(\mathbf{b}^x, \mathbf{V}^x)$ and $(\mathbf{b}^y, \mathbf{V}^y)$:

$$p(\widehat{x}|\phi(\mathbf{z})) = \prod_{k=1}^{|\boldsymbol{\pi}(\widehat{x})|} p(\pi(\widehat{x})_k|\phi(\mathbf{z})), \quad p(\pi(\widehat{x})_k = 1|\phi(\mathbf{z})) = \mathrm{sigm}(b^x_{l(\hat{x}_i)_k} + \mathbf{V}^x_{l(\hat{x}_i)_k,:}\mathbf{h}(\mathbf{c} + \phi(\mathbf{z})))$$

$$p(\widehat{y}|\phi(\mathbf{z})) = \prod_{k=1}^{|\boldsymbol{\pi}(\widehat{y})|} p(\pi(\widehat{y})_k|\phi(\mathbf{z})), \quad p(\pi(\widehat{y})_k = 1|\phi(\mathbf{z})) = \mathrm{sigm}(b^y_{l(\hat{x}_i)_k} + \mathbf{V}^y_{l(\hat{y}_i)_k,:}\mathbf{h}(\mathbf{c} + \phi(\mathbf{z})))$$

where $\mathbf{z}$ can be either $\mathbf{x}$ or $\mathbf{y}$. Notice that we share the bias $\mathbf{c}$ in the nonlinearity $\mathbf{h}(\cdot)$ across decoders.

This encoder/decoder structure allows us to learn a mapping within each language and across the languages. Specifically, for a given pair $(\mathbf{x}, \mathbf{y})$, we can train the model to (1) construct $\mathbf{y}$ from $\mathbf{x}$, (2) construct $\mathbf{x}$ from $\mathbf{y}$, (3) reconstruct $\mathbf{x}$ from itself and (4) reconstruct $\mathbf{y}$ from itself. In our experiments, performed each of these 4 tasks simultaneously, combining the equally weighting the learning gradient from each. Experiments on various weighting schemes should be investigated and are left for future work. Another promising direction of investigation to the exploit the fact that tasks (3) and (4) could be performed on extra monolingual corpora, which is more plentiful.

## 4  Related work

We mentioned that recent work has considered the problem of learning multilingual representations of words and usually relies on word-level alignments. Klementiev et al. [3] propose to train simultaneously two neural network languages models, along with a regularization term that encourages pairs of frequently aligned words to have similar word embeddings. Zou et al. [4] use a similar approach, with a different form for the regularizor and neural network language models as in [2]. In our work, we specifically investigate whether a method that does not rely on word-level alignments can learn comparably useful multilingual embeddings in the context of document classification.

Looking more generally at neural networks that learn multilingual representations of words or phrases, we mention the work of Gao et al. [10] which showed that a useful linear mapping between *separately training* monolingual skip-gram language models could be learned. They too however rely on the specification of pairs of words in the two languages to align. Mikolov et al. [5] also propose a method for training a neural network to learn useful representations of phrases (i.e. short segments of words), in the context of a phrase-based translation model. In this case, phrase-level alignments (usually extracted from word-level alignments) are required.

## 5  Experiments

To evaluate the quality of the word embeddigns learned by our model, we experiment with a task of cross-lingual document classication. The setup is as follows. A labeled data set of documents

in some language $\mathcal{X}$ is available to train a classifier, however we are interested in classifying documents in a different language $\mathcal{Y}$ at test time. To achieve this, we leverage some bilingual corpora, which importantly is not labeled with any document-level categories. This bilingual corpora is used instead to learn document representations in both languages $\mathcal{X}$ and $\mathcal{Y}$ that are enroucaged to be invariant to translations from one language to another. The hope is thus that we can successfully apply the classifier trained on document representations for language $\mathcal{X}$ directly to the document representations for language $\mathcal{Y}$.

## 5.1 Data

We trained our multilingual autoencoder to learn bilingual word representation between English and French and between English and German. To train the autoencoder, we used the English/French and English/German section pairs of the Europarl-v7 dataset[2]. This data is composed of about two million sentences, where each sentence is translated in all the relevant languages.

For our crosslingual classification problem, we used the English, French and German sections of the Reuters RCV1/RCV2 corpus, as provided by Amini et al. [11][3]. There are 18758, 26648 and 29953 documents (news stories) for English, French and German respectively. Document categories are organized in a hierarchy in this dataset. A 4-category classification problem was thus created by using the 4 top-level categories in the hierarchy (CCAT, ECAT, GCAT and MCAT). The set of documents for each language is split into training, validation and testing sets of size 70%, 15% and 15% respectively. The raw documents are represented in the form of a bag-of-words using a TFIDF-based weighting scheme. Generally, this setup follows the one used by Klementiev et al. [3], but uses the preprocessing pipeline of Amini et al. [11].

## 5.2 Crosslingual classification

As described earlier, crosslingual document classification is performed by training a document classifier on documents in one language and applying that classifier on documents in a different language at test time. Documents in a language are represented as a linear combination of its word embeddings learned for that language. Thus, classification performance relies heavily on the quality of multlingual word embeddigns between languages, and specifically on whether similar words across languages have similar embeddings.

Overall, the experimental procedure is as follows.

1. Train bilingual word representations $\mathbf{W}^x$ and $\mathbf{W}^y$ on sentence pairs extracted from Europarl-v7 for languages $\mathcal{X}$ and $\mathcal{Y}$ (we use a separate validation set to early-stop training).

2. Train document classifier on the Reuters training set for language $\mathcal{X}$, where documents are represented using the word representations $\mathbf{W}^x$ (we use the validation set for the same language to perform model selection).

3. Use the classifier trained in the previous step on the Reuters test set for language $\mathcal{Y}$, using the word representations $\mathbf{W}^y$ to represent the documents.

We used a linear SVM as our classifier.

We compare our representations to those learned by Klementiev et al. [3][4]. This is achieved by simply skipping the first step of training the bilingual word representations and directly using those of Klementiev et al. [3] in step 2 and 3. The provided word embeddings are of size 80 for the English and French language pair, and of size 40 for the English and German pair. The vocabulary used by Klementiev et al. [3] consisted in 43614 words in English, 35891 words in French and 50110 words in German. The same vocabulary was used by our model, to represent the Reuters documents.

In all cases, document representations were obtained by multiplying the word embeddings matrix with either the TFIDF-based bag-of-words feature vector or its binary version (the choice of which

---

[2]`http://www.statmt.org/europarl/`

[3]`http://multilingreuters.iit.nrc.ca/ReutersMultiLingualMultiView.htm`

[4]The trained word embeddings were downloaded from `http://people.mmci.uni-saarland.de/~aklement/data/distrib/`

|  | Train FR / Test EN | Train EN / Test FR |
| --- | --- | --- |
| Klementiev | 34.9% | 49.2% |
| Our embeddings | **27.7%** | **32.4%** |
|  | Train GR / Test EN | Train EN / Test GR |
| Klementiev et al. | 42.7% | 59.5% |
| Our embeddings | **29.8%** | **37.7%** |

| | | |
| --- | --- | --- |
| Soleil | ⟹ | Sun |
| Personne | ⟹ | Person |
| Roi | ⟹ | King |
| Voiture | ⟹ | Car |
| Maison | ⟹ | House |
| Arme | ⟹ | Weapon |
| Vivre | ⟹ | Live |
| Apprendre | ⟹ | Learn |
| Papier | ⟹ | Paper |

Table 1: **Left:** Crosslingual classification error results, for English/French pair **(Top)** and English/German pair **(Bottom)**. **Right:** For each French word, its nearest neighbor in the English word embedding space.

method to use is made based on the validation set performance). We normalized the TFIDF-based weights to sum to one.

Test set classification error results are reported in Table 1. We observe that the word representations learned by our autoencoder are competitive with those provided by Klementiev et al. [3]. One will notice that the results for Klementiev et al. [3] are worse than those reported in the original reference. This difference might be due to the fact that our preprocessing of the Reuters data, which comes from Amini et al. [11], is different from the one in Klementiev et al. [3]. In particular, we note that Klementiev et al. [3] ignored documents that belonged to multiple categories, while Amini et al. [11] included them by assigning them to the category with the least training examples.

Table 1 also shows, for a few French words, what are the nearest neighbor words in the English embedding space. A more complete picture is presented in the t-SNE visualization [12] of Figure 2. It shows a 2D visualization of the French/English word embeddings, for the 600 most frequent words in both languages. Both illustrations confirm that the multilingual autoencoder was able to learn similar embeddings for similar words across the two languages.

# 6 Conclusion and Future Work

We presented evidence that meaningful multilingual word representations could be learned without relying on word-level alignments. Our proposed multilingual autoencoder was able to perform competitively on a crosslingual document classification task, compared to a word representation learning method that exploits word-level alignments.

Encouraged by these preliminary results, our future work will investigate extensions of our bag-of-words multilingual autoencoder to bags-of-ngrams, where the model would also have to learn representations for short phrases. Such a model should be particularly useful in the context of a machine translation system. Thanks to the use of a probabilistic tree in the output layer, our model could efficiently assign scores to pairs of sentences. We thus think it could act as a useful, complementary metric in a standard phrase-based translation system.

# References

[1] Joseph Turian, Lev Ratinov, and Yoshua Bengio. Word representations: A simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the*

*Association for Computational Linguistics (ACL2010)*, pages 384–394. Association for Computational Linguistics, 2010.

[2] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural Language Processing (Almost) from Scratch. *Journal of Machine Learning Research*, 12:2493–2537, 2011.

[3] Alexandre Klementiev, Ivan Titov, and Binod Bhattarai. Inducing Crosslingual Distributed Representations of Words. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, 2012.

[4] Will Y. Zou, Richard Socher, Daniel Cer, and Christopher D. Manning. Bilingual Word Embeddings for Phrase-Based Machine Translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2013)*, 2013.

[5] Tomas Mikolov, Quoc Le, and Ilya Sutskever. Exploiting Similarities among Languages for Machine Translation. Technical report, arXiv, 2013.

[6] Yann Dauphin, Xavier Glorot, and Yoshua Bengio. Large-Scale Learning of Embeddings with Reconstruction Sampling. In *Proceedings of the 28th International Conference on Machine Learning (ICML 2011)*, pages 945–952. Omnipress, 2011.

[7] Frederic Morin and Yoshua Bengio. Hierarchical Probabilistic Neural Network Language Model. In *Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics (AISTATS 2005)*, pages 246–252. Society for Artificial Intelligence and Statistics, 2005.

[8] Andriy Mnih and Geoffrey E Hinton. A Scalable Hierarchical Distributed Language Model. In *Advances in Neural Information Processing Systems 21 (NIPS 2008)*, pages 1081–1088, 2009.

[9] Hugo Larochelle and Stanislas Lauly. A Neural Autoregressive Topic Model. In *Advances in Neural Information Processing Systems 25 (NIPS 25)*, 2012.

[10] Jianfeng Gao, Xiaodong He, Wen-tau Yih, and Li Deng. Learning Semantic Representations for the Phrase Translation Model. Technical report, Microsoft Research, 2013.

[11] Massih-Reza Amini, Nicolas Usunier, and Cyril Goutte. Learning from Multiple Partially Observed Views - an Application to Multilingual Text Categorization. In *Advances in Neural Information Processing Systems 22 (NIPS 2009)*, pages 28–36, 2009.

[12] Laurens van der Maaten and Geoffrey E Hinton. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008. URL `http://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf`.
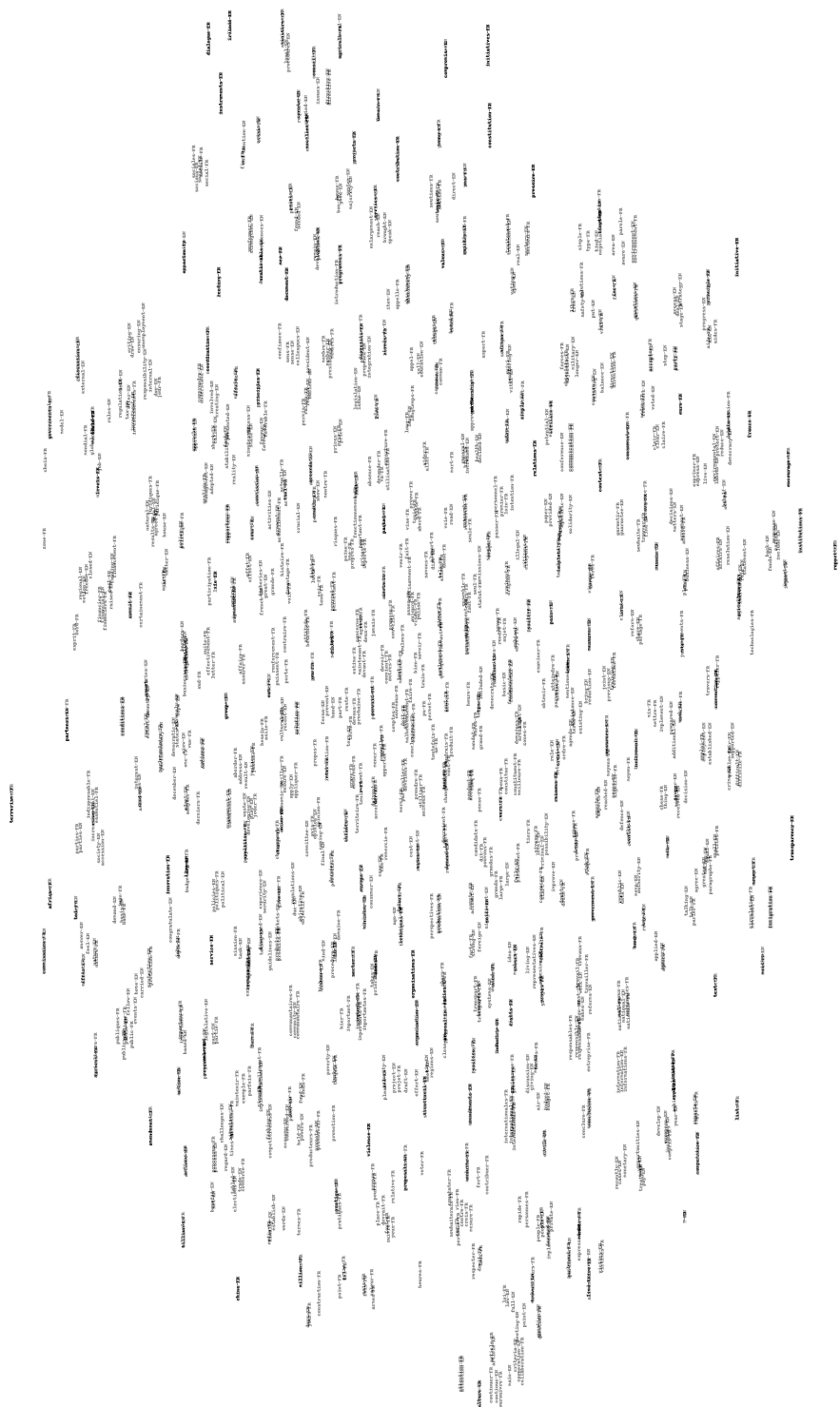
Figure 2: A t-SNE 2D visualization of the learned English/French word representations (better visualized on a computer). Words hyphenated with "EN" and "FR" are English and French words respectively.