



M Ű E G Y E T E M 1 7 8 2

**Budapesti Műszaki és Gazdaságtudományi Egyetem**

Villamosmérnöki és Informatikai Kar

Department of Automation and Applied Informatics

# Universal embedding

DIPLOMATERV

*Készítette*

Eszter Iklódi

*Konzulens*

Gábor Recski

March 27, 2018

# Contents

<b>Kivonat</b>	<b>4</b>
<b>Abstract</b>	<b>5</b>
<b>1 Introduction</b>	<b>6</b>
1.1 Natural Language Processing [17], [1], [2] . . . . .	6
1.1.1 Common tasks of NLP . . . . .	6
1.1.2 Motivation for NLP research . . . . .	7
1.2 Universal Embeddings ( 1.5 page) . . . . .	7
<b>2 Word embeddings</b>	<b>9</b>
2.1 Semantic encoding of words . . . . .	9
2.2 State-of-the-art models for learning word embeddings . . . . .	10
2.3 Multilingual word embeddings . . . . .	11
2.3.1 Motivation . . . . .	12
2.3.2 Tasks . . . . .	12
2.3.3 Under-resourced languages . . . . .	14
2.3.4 State-of-the-art models . . . . .	14
2.4 Multilingual data . . . . .	17
2.4.1 English-Italian setup of Dinu . . . . .	17
2.4.2 The <b>fastText</b> embedding . . . . .	18
2.4.3 Panlex . . . . .	19
<b>3 Proposed model</b>	<b>21</b>
3.1 Description of our method . . . . .	21
3.2 Software architecture . . . . .	21
<b>4 Experiments</b>	<b>22</b>
4.1 Description of different tasks . . . . .	22
4.1.1 Word translation . . . . .	22
4.1.2 Cross-lingual semantic word similarity . . . . .	22
4.2 Summarizing the contributions of the thesis . . . . .	23
<b>5 Future work</b>	<b>24</b>

Köszönetnyilvánítás	25
Irodalomjegyzék	28

## HALLGATÓI NYILATKOZAT

Alulírott *Eszter Iklódi*, szigorló hallgató kijelentem, hogy ezt a szakdolgozatot/ diplomatervet **(nem kívánt törlendő)** meg nem engedett segítség nélkül, saját magam készítettem, csak a megadott forrásokat (szakirodalom, eszközök stb.) használtam fel. Minden olyan részt, melyet szó szerint, vagy azonos értelemben, de átfogalmazva más forrásból átvettem, egyértelműen, a forrás megadásával megjelöltem.

Hozzájárulok, hogy a jelen munkám alapadatait (szerző(k), cím, angol és magyar nyelvű tartalmi kivonat, készítés éve, konzulens(ek) neve) a BME VIK nyilvánosan hozzáférhető elektronikus formában, a munka teljes szövegét pedig az egyetem belső hálózatán keresztül (vagy autentikált felhasználók számára) közzétegye. Kijelentem, hogy a benyújtott munka és annak elektronikus verziója megegyezik. Dékáni engedéllyel titkosított diplomatervek esetén a dolgozat szövege csak 3 év eltelte után válik hozzáférhetővé.

Budapest, March 27, 2018

---

*Eszter Iklódi*  
hallgató

# Kivonat

Mindennapi életünkben egyre fontosabb szerepet tölt be a természetes nyelv számítógép segítségével történő feldolgozása. Digitalizált világunkban egyre inkább alapkövetelmény, hogy a gép és ember közötti kommunikáció természetes nyelven történjen. Ennek a megvalósításához elengedhetetlen az emberi nyelv szemantikai értelmezése.

Manapság a state-of-the-art rendszerekben a szavak szemantikai reprezentációja sokdimenziós vektorokkal, word embeddingek-kel történik. Diplomaterv munkámban már feltanított word embeddingek-hez keresek olyan fordítási mátrixokat, amelyek képesek egy adott nyelvű word embedding univerzális térbe történő leképzésére.

Az így nyert fordítási mátrixokat különböző feladatokon értékelem ki. (eredmények 2 mondatban)

# Abstract

Computer-driven natural language processing plays an increasingly important role in our everyday life. In our digital world using natural language for human-machine communication has become a basic requirement. In order to meet this requirement it is inevitable to analyze human languages semantically.

Nowadays state-of-the-art systems represent word meaning with high dimensional vectors, i.e. with word embeddings. In my thesis work I am searching for translation matrices to already trained word embeddings, such that the translation matrices will be able to map these embeddings into a universal space.

The obtained translation matrices will be evaluated on different tasks. (2 sentences about results)

# Chapter 1

## Introduction

### 1.1 Natural Language Processing [17], [1], [2]

Natural language processing (NLP) is a vibrant interdisciplinary field with many different names, all reflecting a different facet of it. It is often referred to as speech and language processing, human language technology, computational linguistics, or speech recognition and synthesis. The main goal of this field is to get computers to be able to understand and express themselves in human languages.

Natural Language Processing is hard since it deals with what is considered to be one of the most delicate characteristics of human being: with human language. This field is strongly connected with artificial intelligence since the way humans think and feel about the world is mainly happening in terms of human languages.

Being nowhere near as fast as digital channels, expressing ourselves by means of human languages is a very effective way of communication, though. Saying only the minimum message our listeners can fill up the rest with their world and common knowledge, and can easily figure out the missing or misunderstood parts from the context of the situation. This way they are also able to resolve ambiguities, homonyms etc. without even noticing it. Nonetheless, these tasks for a computer are not at all that trivial.

The importance of computer integrated human language communication has gone as far as assigning truly intelligent machines to the ability of being capable of processing language as skillfully as humans do. This idea was first introduced by Alan Turing (1950) who proposed what has come to be known as the Turing test.

#### 1.1.1 Common tasks of NLP

Natural language processing comprises a wide variety of tasks. Some of them like spam detection, POS-tagging, or named entity recognition are considered to be mostly solved problems. Applications for these tasks are now out in the market and are usually integrated to our smart devices even by default.

With some other tasks great progress has been made recently which implies the existence of already fair enough applications but means that research work is still has to be done. Among them there are tasks like sentiment analysis, words sense disambiguation, syntactic

parsing, machine translation etc.

What is still considered to be very hard is to understand the meaning of a text. There are numerous interesting tasks for example question answering, dialogues, summarization, paraphrases, or text inference, for which in order to make relevant progress dealing with the semantics is inevitable.

### 1.1.2 Motivation for NLP research

We are living in an era when natural language processing is becoming more and more integrated into our every life. With the advent of smart phones the importance of language has gone even further. These devices are having small and rather inconvenient keyboards thus speech-driven communication seems very appealing. Big companies like Amazon, Apple, Facebook, Google, etc. are all putting out of products that use natural language to communicate with users. Since the contributions of this thesis are aiming the research field of word meaning and universal semantic representations, below I will only list applications that can directly take advantages of these contributions.

Speech-driven assistance applications can really make our everyday life more comfortable and more convenient. Let's just consider the fact how much help they could provide to people living with disabilities. These systems include speech input for which at first automatic speech recognition technologies should be applied. But after that, in order to understand the goal of the user, we must run a semantic analysis as well.

An early version of conversational agents and certain strongly domain-based chatbots are already out on the market, providing 24 hour, immediate assistance for customers. By letting computers do the monotone and not at all creative tasks employees could have more interesting jobs, jobs that only humans can do, or they could have less working hours in a week. Any of these would be a great progress for the society.

Advances in machine translation has already created a world where non-English speakers can also enjoy the benefits of the English-based web services. Generally, we can say that for widespread languages machine translation has already reached a fairly usable state, for rare languages, however, it is still facing difficulties.

There are also numerous Web related tasks that are strongly relying on the semantic analysis of the text. One promising task is for example the Web-based question answering task which is basically an extended version of the classical Web search, whereby instead of searching just for key words it would also be possible to ask complete questions and thus communicate with the search engine just like human beings do. For all these applications, however, it is inevitable to look way behind the syntactic surface and dig deep into the underlying semantics.

## 1.2 Universal Embeddings ( 1.5 page)

Given the need for robust representations for many languages, the question of whether human conceptual structure is universal has recently gained interest not only among cognitive scientists ([29], [20], [14]), but among computational linguists as well. Youn et al. [36]



has shown that human conceptual structure is independent of certain non-linguistic factors such as geography, climate, topology or literary traditions. Based on such findings we propose a procedure to construct a universal semantic representation in form of translation matrices that serve to map each language to the universal space. We use the pretrained fastText word embeddings [8], which are available for 294 languages.

**TODO: summary of experiments and results.**

**Github links, self-references, contributions**

The thesis is structured as follows:

- **Chapter 1** explains briefly the goals and motivations of the research field of natural language processing. It also summarizes the main contributions and the results of the accomplished work.
- **Chapter 2** discusses the state-of-the-art semantic word representations, i.e. the word embedding. It briefly describes the standard learning procedure for monolingual word-embeddings (word2vec) and it introduces the concept and resources for multilingual word-embeddings.
- **Chapter 3** describes the proposed model. It details the learning procedure and discusses the basic infrastructural and architectural features.
- **Chapter 4** presents all the experiments. It summarizes our results and compares them with the performance of other systems.
- **Chapter 5** is devoted to the description of the future work. This chapter suggests modifications, follow-ups for which we didn't have time to accomplish, or which are beyond the scope of this thesis work.

## Chapter 2

# Word embeddings

### 2.1 Semantic encoding of words

Within the field of natural language processing a more specific area concentrates on semantic representations which are being leveraged both by classical semantic tasks such as question answering or chatbots and by other NLP tasks which in the strict sense of the word are not considered semantic tasks such as machine translation or syntactic parsing. A crucial part of all semantic tasks is to have a proper word representation which is capable of encoding the meaning as well.

One way to build a semantic representation is to use a distributional model. The idea is based upon the observation that synonyms or words with similar meanings tend to occur in similar contexts, or as it was articulated by Firth in 1957: *You shall know a word by the company it keeps* [13]. For example in the following two sentences “**The cat is walking in the bedroom**” and “**A dog was running in a room**” words like “dog” and “cat” have exactly the same semantic and grammatical roles therefore we could easily imagine the two sentences in the following variations: “**The dog is walking in the bedroom**” and “**A cat was running in a room**” [7]. Based upon this intuition, what distributional models are aiming to do is to compute the meaning of a word from the distribution of words around it. [17]. The obtained meaning representations are usually high dimensional vectors, called as word embeddings, which refers to their characteristic feature that they model a word by embedding it into a vector space.

Embeddings are not only proved to be a better alternative of n-grams used for language modelling [7], but they are doing quite well on semantic tasks too. Mikolov et al. [24] has shown that the characteristics of word embeddings go way beyond the simple syntactic regularities. They showed that applying simple vector operations (e.g. vector addition and subtraction) can often produce meaningful results. For example it was shown that  $vector("King") - vector("Man") + vector("Woman")$  results in a vector that is closest to the vector representation of the word *Queen* [25]. Moreover, these days' state-of-the-art results on word similarity tasks are all held by word embeddings, where the similarity of two words are calculated as the normalized dot product of the corresponding word vectors. This measure is the so-called cosine similarity of word embeddings.

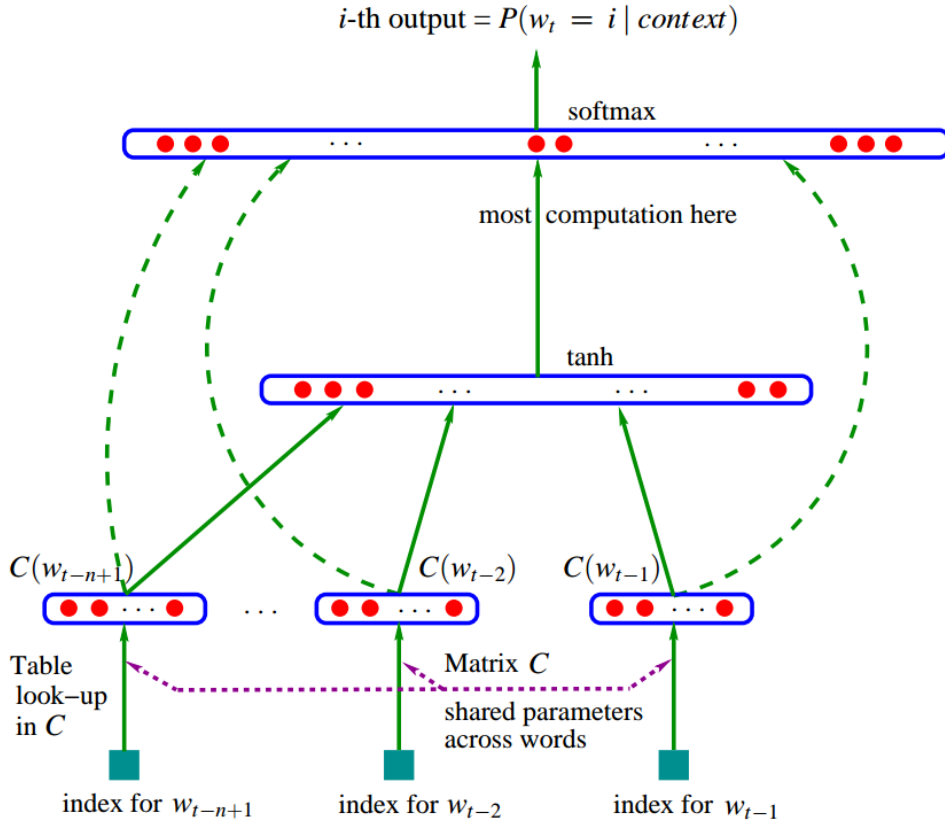
Another way to build semantic representations is to utilize lexical databases. In our previous work we created a hybrid system leveraging both the 4lang orthological model described in [18], [19] and [3] and various distributional models, i.e. various word embeddings. With this system in 2016 we reached a state-of-the-art score on the SimLex-999 [16] benchmark data [28].

In the following sections I will describe the basic procedure of training word embeddings and following that I will focus on a more specific field of computational semantics, on multilingual word embeddings.

## 2.2 State-of-the-art models for learning word embeddings

In 2003 Bengio et al.[7] suggested a probabilistic feedforward neural network language model (NNLM) for learning a distributed representation for words. The network consists of input, projection, hidden and output layers, where at the input layer the  $N$  previous words are encoded using 1-of- $V$  coding, where  $V$  is the size of the vocabulary. Figure 2.1 shows an overview of the architecture. This procedure was evaluated in terms of perplexity at which in comparison with the best of the  $n$ -grams it proved to be much better. The drawback of this architecture is that it becomes complex for computation between the projection and the hidden layer, as values in the projection layer are dense.

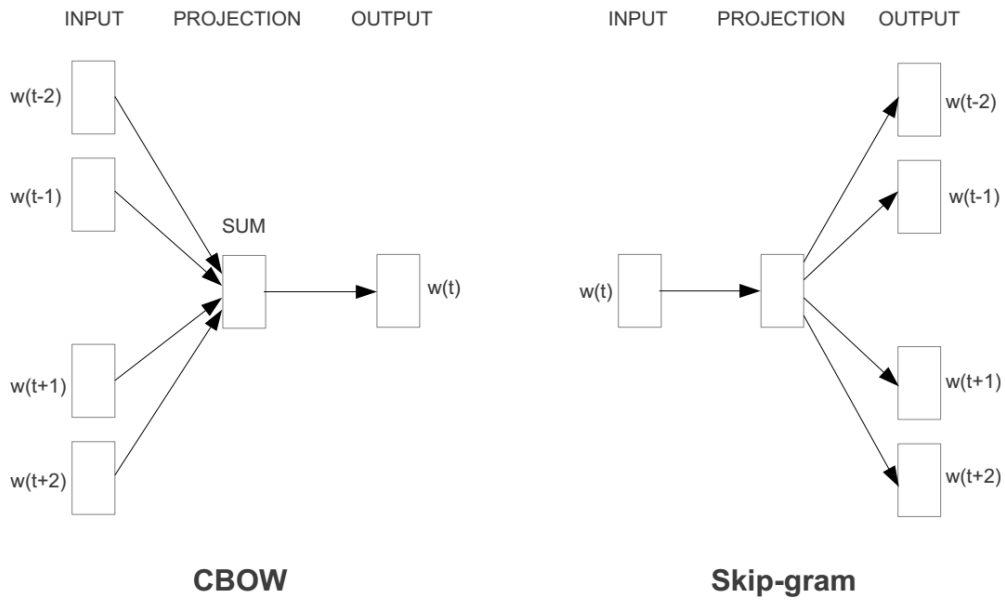
**Figure 2.1:** Network architecture proposed by Bengio et al.[7]



10 years later, in 2013 Mikolov suggested a bag-of-words neural networks, more specif-

ically two different architecture [22]. The first one, denoted as the CBOW (Continuous Bag-of-Words Model) tries to predict the current word based on the context, whereas the second one, denoted as the continuous skip-gram model tries to maximize the classification of a word based on another word in the same sentence. Both models worked better than the NNLM suggested by Bengio [7] both on semantic and syntactic tasks, while between the two models of Mikolov the CBOW turned out to be slightly better on syntactic tasks and the skip-gram on semantic tasks. Mikolov’s procedure has become known as the **word2vec** procedure and the source code is available on github <http://deeplearning4j.org/word2vec>. The architecture of the CBOW and the skip-gram models are shown in figure 2.2.

**Figure 2.2:** Bag-of-words neural networks suggested by Mikolov et al.[24]



Embeddings are usually evaluated on word similarity and word analogy tasks. Besides providing quite promising results on them, they have also been applied to many downstream tasks, from **named entity recognition** and **chunking** [34] to **dependency parsing** [6]. It has furthermore been shown that weakly supervised embedding algorithms can also lead to huge improvements for tasks like **sentiment analysis** [32].

### 2.3 Multilingual word embeddings

In this section I describe the importance of multilingual word embeddings. I also explain how it is possible to incorporate word embeddings trained on monolingual text corpora into a multilingual context. After that I present a brief summary about the previous attempts on constructing cross-lingual word vector representations. Finally, I describe the available parallel datasets that we were using for our experiments.

### 2.3.1 Motivation

The question how to model representations is a highly interdisciplinary issue to discuss. Within cognitive science, traditionally there are two dominating approaches to this problem. The first one is the *symbolic* which states that cognitive systems can be described as Turing machines. The second one, denoted as *associationism*, says that representations are associations among different kinds of information elements. In his book, *Conceptual Spaces: The Geometry of Thought* [14], Gärdenfors advocates a third approach, which he calls *conceptual* from. This representation is based on using geometrical structures rather than symbols or connections among neurons.

To go a step further one could ask whether these structures are universal among all human beings. Regarding this question with the eyes of a computer scientist we might form this problem as whether it is possible to model meaning universally, i.e. language independently. Current meaning representations are learned from monolingual corpora, therefore they infer language dependency. But is there a way to find one single representation instead of a different one for each and every human language?

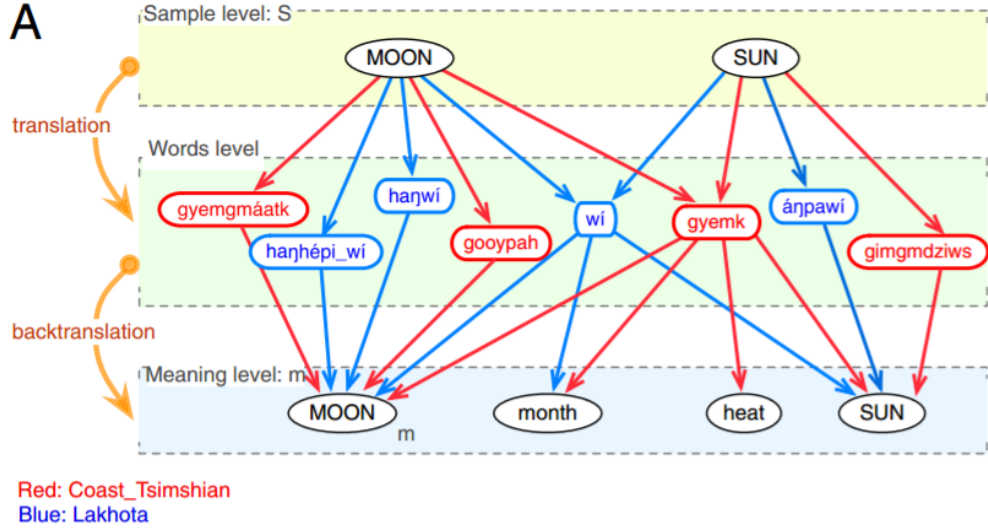
Youn et al. [36] suggested that human brain may reflect distinct features of cultural, historical, and environmental background in addition to properties universal to human cognition. They provided an empirical measure of semantic proximity between concepts using entries of the Swadesh list [31]. The Swadesh list is cross-linguistic dictionary which includes a 110 and a 207 long list of basic concepts in approximately 2000 languages. Youn et al. took 22 concepts of this list that refer to material entities (e.g. STONE, EARTH, SAND, ASHES), celestial objects (e.g., SUN, MOON, STAR), natural settings (e.g., DAY, NIGHT), and geographic features (e.g., LAKE, MOUNTAIN). Then, they applied translation and back-translation through various languages. As a result of numbers of polysemies in the resulting graph originally distinct concepts become connected. For example the Spanish word *cielo* in English both means *heaven* and *sky*. Thus by applying English-Spanish-English translation and back-translation the two English words *heaven* and *sky* become connected. The more such polysemous words we find, the stronger this connection becomes. For example if besides Spanish, we also apply the translation and back-translation through German, the same polysemy appears: the German word *Himmel* both means *heaven* and *sky* in English. The procedure is shown on figures 2.3 and 2.4.

Statistical analysis of the obtained graphs constructed over the polysemies observed in the above-mentioned 22-word-long subset of basic vocabulary showed that the structural properties of these graphs are consistent across different language groups, and largely independent of geography, environment, and the presence or absence of a literary tradition. Based upon these findings we assume that meaning, at least to a certain extent, is universal, thus representing semantics at universal level is reasonable.

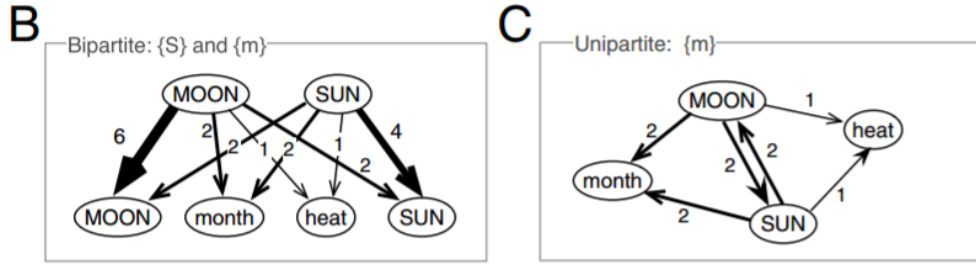
### 2.3.2 Tasks

Beyond the theoretical level of whether meaning is universal there are numerous practical problems for which cross-lingual embeddings come in handy. In this section I write about

**Figure 2.3:** Translating MOON and SUN through polysemous words.



**Figure 2.4:** Making links between English concepts through eliminating the internal nodes.



the different tasks where solutions can be facilitated by utilizing multi-lingual embeddings.

### Cross-language part-of-speech tagging

Part-of-speech tagging, a.k.a. POS-tagging is the task for annotating a text with part-of-speech tags. The fundamental idea behind multilingual learning to part-of-speech tagging is that the patterns of ambiguity inherent in part-of-speech tag assignments differ across languages. A word with part-of-speech tag ambiguity in one language may correspond to an unambiguous word in the other language. For example, the word “can” in English may function as an auxiliary verb, a noun, or a regular verb, however, translating the sentence into other languages the different meanings of "can" are likely to be expressed with different lexemes. By combining natural cues from multiple languages, the structure of each POS-tagger becomes more apparent [26].

Todo: ref to applications of multiling embeddings, is there any ??

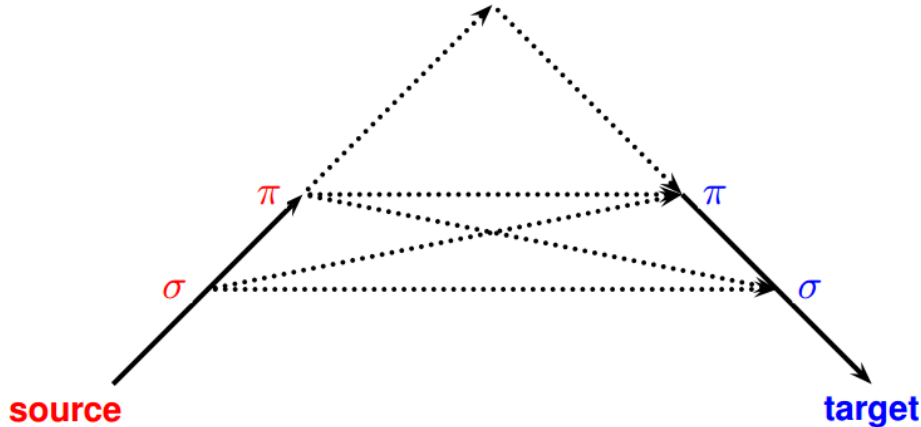
### Cross-language super sense tagging

Todo: [27] WordNet [12], [15]

## Machine translation

Machine translation is the task of translating a text automatically with a computer from a source language to a target language. Current research works are focusing on finding the appropriate level of representations when performing the translations. Basic approaches are: **tree-to-string**, **string-to-string**, and **string-to-tree**, as shown in figure 2.5.

**Figure 2.5:** *Levels of representation in Machine Translation.  $\pi \rightarrow \sigma$  : tree-to-string;  $\sigma \rightarrow \sigma$  : string-to-string;  $\sigma \rightarrow \pi$  : string-to-tree.*



Translation models, however, often fail to generate good translations for infrequent words or phrases. Previous works attacked this problem by inducing new translation rules from monolingual data with a semi-supervised algorithm. Nevertheless, this approach does not scale very well since computationally it is quite expensive. Zhao et al. [37] proposed a much faster and simpler method that creates translation rules for infrequent phrases based on phrases with similar continuous representations, i.e. with similar word vectors, for which a translation is known. Their method improved a phrase-based baseline by up to 1.6 BLEU on Arabic-English translation, and it was three-orders of magnitudes faster than existing semi-supervised methods and 0.5 BLEU more accurate.

### 2.3.3 Under-resourced languages

Dictionaries and phrase tables are the basis of modern statistical machine translation systems. Mikolov et al. [23] showed a method that can automate the process of generating and extending dictionaries and phrase tables. They could translate missing word and phrase entries by learning language structures based on large monolingual data and mapping between languages from small bilingual data. This is a powerful opportunity for rare languages to join the mostly English-based world of the Web and for non-English speakers to enjoy its benefits without having to speak English.

### 2.3.4 State-of-the-art models

In this section I present a brief history on cross-lingual word vector representations. As a baseline approach I describe the procedure of Mikolov et al. [23] from 2013 and next I study

various attempts made since then to improve this baseline system and to alleviate its errors. Finally, I summarize some recent attempts for obtaining multilingual word embeddings without using any parallel data.

### First attempt: Mikolov et al.

Right after publishing their **word2vec** procedure, Mikolov et al. [23] went even further by noticing that continuous word embedding spaces exhibit similar structures across languages. They applied a simple two-step procedure:

- firstly, monolingual models of languages using huge corpora are built, e.g. by using the **word2vec** method
- secondly, a small bilingual dictionary was used to learn linear projection between the languages. These words are often referred to as anchor points.
- finally, at test time, the translation of any word from the source language is possible by projecting its vector representation from the source language space to the target language space. Once the vector in the target language space is obtained, the most similar word vector can serve as the output of the translation.

This method seemed to be working even for distant language pairs like English and Vietnamese.

In this work for the first step I used a pretrained word embedding described in 2.4.2. The work itself concentrates on finding the linear projections and on the evaluation of these projections.

### Various improvements

#### Faruqui and Dyer

Since Mikolov's experiments various attempts were made to improve the cross-lingual embeddings. Faruqui and Dyer [11] proposed a simple technique based on **canonical correlation analysis (CCA)**.

#### Xing et al.

Xing et al. [35] has shown that bilingual translation can be largely improved by **normalizing** the embeddings and by restricting the transformation matrices into **orthogonal** ones.

#### Dinu et al.

Dinu et al. [10] studied the phenomenon of **hubs**. He showed that the neighbourhoods of the mapped vectors are strongly polluted by hubs. These vectors tend to be near a high proportion of items, and thus their correct labels will be pushed down in the neighbour list when looking up for word translations. They proposed a simple method to alleviate this problem with which they achieved consistent improvements.



#### Lazaridou et al.

Lazaridou et al. [21] studied some theoretical and empirical properties of general cross-space mapping function, and tested them on cross-linguistic (word translation) and cross-modal (image labelling) tasks. By introducing **negative samples** during the learning process they could reach state-of-the-art results on the English-Italian word translation task described in 2.4.1. Settings for the negatives examples were studied both by choosing them random and by choosing "intruders" which are near the mapped vector, but far from the actual gold target space vector. The "intruder" approach achieved better results, and furthermore, it gave better results after just few training epochs.

#### Ammar et al.

Ammar et al. [4] proposed methods for estimating and evaluating embeddings of words in **more than fifty languages** in a **single shared embedding space**. The so-called multiCluster and multiCCA methods were tested on 59 languages, while the multiSkip and translation-invariance methods only on 12 languages for which high-quality parallel data are available. For the 12 languages the bilingual dictionaries were extracted from the Europarl parallel corpora, while for the remaining 47 languages, dictionaries were formed by translating the 20k most common words in the English monolingual corpus with Google Translate.

#### Artetxe et al.

Artetxe et al. [5] built a **generic framework** that generalizes previous works made on cross-linguistic embeddings. For evaluating the methods they used the same English-Italian dataset by Dinu, discussed in 2.4.1. As a conclusion they published that from the proposed methods the ones with orthogonality constraint and a global preprocessing with length normalization and dimension-wise mean centering achieved the best overall results.

#### Smith et al.

Smith et al. [30] also proves that translation matrices should be **orthogonal**. They apply **singular value decomposition** to achieve this. Besides, they introduce a novel "**inverted softmax**" method for identifying translation pairs, with which they improve the precision of Mikolov. Orthogonal transformations also turned out to be more robust to noise which makes it possible to learn the transformation without expert bilingual resource by constructing a "pseudo-dictionary" from the identical character strings.

#### Without parallel data

While all the above mentioned methods rely on biligual word lexicons, most recent studies are aiming to eliminate the need for any parallel data at all. Smith et al. [30] has already made attempts for the alleviation of parallel data supervision by introducing **character-level information**, but the results were not on par with their supervised counterparts, on the one hand, and on the other hand, these methods are strictly limited to pairs of languages sharing a common alphabet.

Conneau et al. [9] introduces an **unsupervised** way for aligning monolingual word embedding spaces between two languages **without using any parallel corpora**. Their experiments show that this method can be applied even for distant language pairs like English-Russian or English-Chinese.

## 2.4 Multilingual data

In this section I briefly describe the data resources that I used. These involve the pretrained embeddings I took for the experiments and the gold bilingual dictionaries I used for the evaluation.

### 2.4.1 English-Italian setup of Dinu

Dinu et al. [10] has constructed an English-Italian dictionary split into a train and a test that are now being used as a benchmark data for evaluating word translation tasks.

Both train and test translation pairs are extracted from a dictionary built from Europarl, available at <http://opus.lingfil.uu.se/> (Europarl, en-it) [33]. For the test set they used 1,500 English words split into 5 frequency bins, 300 randomly chosen in each bin. The bins are defined in terms of rank in the frequency-sorted lexicon: [1-5K], [5K-20K], [20K-50K], [50K-100K] and [100K-200K].

For the training translation pairs they also sampled by frequency, using the top 1K, 5K, 10K and 20K most frequent translation pairs from the Europarl dictionary sorted by English frequency. There is no overlap with test elements on the English side, however, when checking for overlaps with the 5K train data on the Italian side we found 113 Italian words that can be found both in train and test sets. These overlaps can be sorted into the following categories:

- **Singular-plural correspondence:** in Italian when the last vowel of a substantive is accented the plural form is the same as the singular. For example *comunità* and *attività* in table 2.1.
- **Italian word is mistaken for English word:** the English translation is the same as the original Italian word. For example in the test set the Italian word *segnì* is not translated and the same happens with *vecchi*. See table 2.2.
- **Different verb forms:** the same Italian word can be translated into different English verb tenses. For example *sostenere* in table 2.3.
- **Synonyms and homonyms:** one Italian word can be translated into several English words that might be synonyms or might not in case of homonyms. This phenomenon is actually fairly understandable and acceptable under all circumstances. See table 2.4.
- **Errors in the translation:** wrong translations. For example plural form Italian words *gatti* and *passeggeri* are translated both as the correct plural form and the incorrect singular form. See examples in table 2.5.

**Table 2.1:** *Singular-plural correspondence*

Italian	English - train	English - test
comunità	communities	community
attività	activities	activity

**Table 2.2:** *Italian word is mistaken for English word*

Italian	English - train	English - test
signi	signs	<b>signi</b>
vecchi	old	<b>vecchi</b>

**Table 2.3:** *Different verb forms:*

Italian	English - train	English - test
sostenere	support	supporting

A summary of word counts can be seen in table 2.6.

Smith et al. [30] reported results on this English-Italian dataset both in English-Italian and Italian-English direction. They reproduced the methods of Mikolov [23], Faruqui [11] and Dinu . A summary of the English-Italian results can be found in table 2.7 and the Italian-English in table 2.8, respectively. All the methods turned to be more accurate when translating from English to Italian. This is not surprising at all, given the fact that many English words can be translated to either the male or female form of the Italian word.

#### 2.4.2 The fastText embedding

Usual techniques for obtaining continuous word representation, i.e. word embeddings, represent each word of the vocabulary by a distinct vector, without parameter sharing. They ignore completely the morphology of words which is a significant limitation especially for agglutinating languages, e.g. Hungarian. In these languages new words are formed by stringing together morphemes which leads to large vocabularies and many rare words.

In 2017 the Facebook AI Research group proposed a new approach based on the skipgram model [22], but this time, contrary to the previously mentioned methods, parameter sharing was applied, since words are represented as a bag of character n-grams [8]. A vector representation is associated to each character n-gram, and the words are being represented as the sum of these representations. The method turned out to be fast and it allows us to compute word representations for words that did not appear in the training data. The model was evaluated both on word similarity and analogy tasks. The results show that this model outperforms Mikolov’s CBOW and skipgram baseline systems that do not take into account subword information. It also does better than methods relying on morphological analysis.

The pre-trained word vectors for 294 languages, trained on Wikipedia using fastText are available on the following github link:

<https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md>.

**Table 2.4:** *Synonyms and homonyms*

Italian	English - train	English - test
risposte	answers	responses
sufficiente	sufficient	enough

**Table 2.5:** *Errors in the translation*

Italian	English - train	English - test	Explanation
gatti	cat	cats	it only means cats
passengeri	passengers	passenger	it only means passengers

### 2.4.3 Panlex

**Table 2.6:** *Statistics of word counts. The train set contains 5000 and the test set 1500 word pairs, respectively.*

Language	Set	# words
eng	train	3442
	test	1500
ita	train	4549
	test	1849

**Table 2.7:** *English to Italian results published by Smith*

Precision	@1	@5	@10
Mikolov et al. (2013b)	0.338	0.483	0.539
Faruqui et al. (2014)	0.361	0.527	0.581
Dinu et al. (2015)	0.385	0.564	0.639
Smith et al. (2017)	<b>0.431</b>	<b>0.607</b>	<b>0.664</b>

**Table 2.8:** *Italian to English results published by Smith*

Precision	@1	@5	@10
Mikolov et al. (2013b)	0.249	0.410	0.474
Faruqui et al. (2014)	0.310	0.499	0.570
Dinu et al. (2015)	0.246	0.454	0.541
Smith et al. (2017)	<b>0.380</b>	<b>0.585</b>	<b>0.636</b>

## Chapter 3

# Proposed model

### 3.1 Description of our method

In our work we are searching for linear mappings between pretrained word embeddings through a universal vector space. Therefore we

### 3.2 Software architecture

Used packages (tensorflow ...)

Artetxe writes about SW arch, take a look at it!

## Chapter 4

# Experiments

### 4.1 Description of different tasks

#### 4.1.1 Word translation

Mikolov: 90 percent precision@5 for translation of words between English and Spanish

#### 4.1.2 Cross-lingual semantic word similarity

Facebook (MUSE): SemEval 2017

## 4.2 Summarizing the contributions of the thesis



## Chapter 5

### Future work

# Köszönetnyilvánítás

Ez nem kötelező, akár törölhető is. Ha a szerző szükségét érzi, itt lehet köszönetet nyilvánítani azoknak, akik hozzájárultak munkájukkal ahhoz, hogy a hallgató a szakdolgozatban vagy diplomamunkában leírt feladatokat sikeresen elvégezze. A konzulensnek való köszönetnyilvánítás sem kötelező, a konzulensnek hivatalosan is dolga, hogy a hallgatót konzultálja.

# Bibliography

- [1] <https://www.youtube.com/playlist?list=PL3FW7Lu3i5Jsnh1rnUwqTcylNr7EkRe6>.
- [2] URL: <https://www.economist.com/technology-quarterly/2017-05-01/language>.
- [3] Judit Acs, Katalin Pajkossy, and András Kornai. Building basic vocabulary across 40 languages. In *Proceedings of the Sixth Workshop on Building and Using Comparable Corpora*, pages 52–58, 2013.
- [4] Waleed Ammar, George Mulcaire, Yulia Tsvetkov, Guillaume Lample, Chris Dyer, and Noah A Smith. Massively multilingual word embeddings. *arXiv preprint arXiv:1602.01925*, 2016.
- [5] Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2289–2294, 2016.
- [6] Mohit Bansal, Kevin Gimpel, and Karen Livescu. Tailoring continuous word representations for dependency parsing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 809–815, 2014.
- [7] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.
- [8] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*, 2016.
- [9] Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*, 2017.
- [10] Georgiana Dinu, Angeliki Lazaridou, and Marco Baroni. Improving zero-shot learning by mitigating the hubness problem. *arXiv preprint arXiv:1412.6568*, 2014.
- [11] Manaal Faruqui and Chris Dyer. Improving vector space word representations using multilingual correlation. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 462–471, 2014.
- [12] Christiane Fellbaum. *WordNet*. Wiley Online Library, 1998.
- [13] John R Firth. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*, 1957.

- [14] Peter Gärdenfors. *Conceptual spaces: The geometry of thought*. MIT press, 2004.
- [15] Stephan Gouws and Anders Søgaard. Simple task-specific bilingual word embeddings. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1386–1390, 2015.
- [16] Felix Hill, Roi Reichart, and Anna Korhonen. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695, 2015.
- [17] Dan Jurafsky and James H Martin. *Speech and language processing*. Pearson London:, 2017.
- [18] András Kornai. The algebra of lexical semantics. In *the Mathematics of Language*, pages 174–199. Springer, 2010.
- [19] András Kornai. Eliminating ditransitives. In *Formal Grammar*, pages 243–261. Springer, 2012.
- [20] George Lakoff. *Women, fire, and dangerous things*. University of Chicago press, 2008.
- [21] Angeliki Lazaridou, Georgiana Dinu, and Marco Baroni. Hubness and pollution: Delving into cross-space mapping for zero-shot learning. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, volume 1, pages 270–280, 2015.
- [22] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [23] Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*, 2013.
- [24] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [25] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751, 2013.
- [26] Tahira Naseem, Benjamin Snyder, Jacob Eisenstein, and Regina Barzilay. Multilingual part-of-speech tagging: Two unsupervised approaches. *Journal of Artificial Intelligence Research*, 2009.
- [27] Davide Picca, Alfio Massimiliano Gliozzo, and Simone Campora. Bridging languages by supersense entity tagging. In *Proceedings of the 2009 Named Entities Workshop: Shared Task on Transliteration*, pages 136–142. Association for Computational Linguistics, 2009.

- [28] Gábor Recski, Eszter Iklódi, Katalin Pajkossy, and Andras Kornai. Measuring semantic similarity of words using concept networks. In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 193–200, 2016.
- [29] Eleanor Rosch and Carolyn B Mervis. Family resemblances: Studies in the internal structure of categories. *Cognitive psychology*, 7(4):573–605, 1975.
- [30] Samuel L Smith, David HP Turban, Steven Hamblin, and Nils Y Hammerla. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *arXiv preprint arXiv:1702.03859*, 2017.
- [31] Morris Swadesh. Lexico-statistic dating of prehistoric ethnic contacts: with special reference to north american indians and eskimos. *Proceedings of the American philosophical society*, 96(4):452–463, 1952.
- [32] Duyu Tang, Furu Wei, Nan Yang, Ming Zhou, Ting Liu, and Bing Qin. Learning sentiment-specific word embedding for twitter sentiment classification. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1555–1565, 2014.
- [33] Jörg Tiedemann. Parallel data, tools and interfaces in opus. In *LREC*, volume 2012, pages 2214–2218, 2012.
- [34] 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*, pages 384–394. Association for Computational Linguistics, 2010.
- [35] Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. Normalized word embedding and orthogonal transform for bilingual word translation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1006–1011, 2015.
- [36] Hyejin Youn, Logan Sutton, Eric Smith, Cristopher Moore, Jon F Wilkins, Ian Maddieson, William Croft, and Tanmoy Bhattacharya. On the universal structure of human lexical semantics. *Proceedings of the National Academy of Sciences*, 113(7):1766–1771, 2016.
- [37] Kai Zhao, Hany Hassan, and Michael Auli. Learning translation models from monolingual continuous representations. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1527–1536, 2015.