

Evaluating Bilingual Embeddings in Bilingual Dictionary Alignment

Çift Dilli Kelime Temsilleri ile Sözlük Eşlenmesi

Yiğit Sever

Dr. Gönenç Ercan

Supervisor

Submitted to
Graduate School of Science and Engineering of Hacettepe University
as a Partial Fulfilment to the Requirements
for the Award of the Degree of Master of Science
in Computer Engineering

2019

Declaration of Authorship

I, Yiğit Sever, declare that this thesis titled, “Evaluating Bilingual Embeddings in Bilingual Dictionary Alignment” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

Abstract

Yiğit SEVER

Evaluating Bilingual Embeddings in Bilingual Dictionary Alignment

Dictionaries catalog and describe the semantic information of a lexicon. WordNet provides an edge by presenting distinct concepts with the hierarchy information among them. Research in computer science has been using this hand crafted tool in natural language applications such as text summarization and machine translation. Original WordNet has been compiled for English yet counterparts for other languages are not as readily available nor as comprehensive. In order for research on languages other than English to benefit from the power of a WordNet, machine assisted creation and evaluation methods are essential.

Word embeddings can provide a mapping between words and points in a real valued vector space. Using these vectors, representing documents as well as forming geometric relationships between them is a well studied area of research. In this thesis we start by hypothesizing that a dictionary definition captures the semantic basis of the described word. We used word embeddings as building blocks to map dictionary definitions into a multidimensional space. These spaces can be aligned to accommodate two languages, allowing the transfer of information from one language to another. We investigate the success of retrieving and matching discrete senses across languages by employing supervised and unsupervised methods. Our experiments show that dictionary alignment can be evaluated successfully by using both unsupervised and supervised methods but corpora sizes should be taken into consideration. We further argue that some methods are not viable considering their poor performance.

Acknowledgements

The acknowledgments and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	iii
Abstract	v
Acknowledgements	vii
1 Introduction	1
1.1 Dictionaries	1
1.2 WordNet	2
1.3 Multilingual Wordnets	5
1.4 Thesis Goals	6
1.5 Thesis Outline	7
2 Background Information & Related Work	9
2.1 Word Embeddings	9
2.1.1 History of Word Representations	10
Linguistic Background	10
Vector Space Model	11
2.1.2 Latent Semantic Analysis	12
2.1.3 Building Upon Distributional Hypothesis	13
2.1.4 Distributed Vector Representations	13
2.1.5 Pre-trained Embeddings	15
2.1.6 Popularization of Word Embeddings	15
2.1.7 fastText	17
2.1.8 ConceptNet Numberbatch	18
2.2 Bilingual Word Embeddings	19
2.3 Approaches in Wordnet Generation	19
2.3.1 History of Wordnet Generation	19
2.3.2 Examples of Wordnet Generation	22

3 Unsupervised Matching	25
3.1 Linear Assignment Using Sentence Embeddings	25
3.2 Linear Assignment Algorithm	27
3.3 Results	27
4 Dictionary Alignment as Pseudo-Document Retrieval	29
4.1 Monolingual Retrieval	29
4.1.1 Term Document Matrix	31
4.1.2 Term Weighting	31
4.1.3 Similarity Measure	32
4.1.4 Retrieval	33
4.1.5 Evaluation	34
4.1.6 Results	34
4.2 Cross Lingual Document Retrieval	34
4.2.1 Word Movers Distance	34
4.3 Balikas <i>et al.</i> Implementation	36
4.3.1 Optimal Transport	36
4.3.2 Sinkhorn	36
5 Supervised Validation	37
5.0.1 Vanishing Gradient Problem	37
6 Experiments and Evaluation	41
7 Conclusion	43
Bibliography	45

List of Figures

1.1	WordNet result for the query “run”, truncated for brevity.	3
1.2	Example of a WordNet relationship graph. Hyponymy/hypernymy relations are shown from the hierarchy level of	4
2.1	Sense representation using human judgement scores for the concept “Father” adapted from Osgood <i>et al.</i> [23].	11
2.2	Skipgram architecture adapted from Figure 1 of Mikolov <i>et al.</i> [38] . .	16
2.3	Overview of fastText embedding model	17
4.1	Term document matrix to queries and documents	33
5.1	Graphical representation of vanishing gradient problem where the shades of the nodes represent the influence of the input signal, adapted from Figure 4.1 of Graves [89]	38
5.2	Simplified long short-term memory cell architecture	39
5.3	Preserving the input signal through blocking (-) or allowing (O) the input signal, adapted from Figure 4.4 of Graves [89]	40

List of Tables

1.1	Summary of the Wordnets used	6
2.1	Example synsets and their respective glosses	20
3.1	Some definitions from English Princeton WordNet	27
3.2	Linear Assignment Using 2000 Definitions	27
3.3	Linear Assignment Using 3000 Definitions	28
3.4	Summary of Linear Assignment	28
4.1	Example definitions and translations	30
4.2	English Princeton WordNet definitions and the target wordnet definitions we want to match	35
4.3	Experiment results for monolingual retrieval	35
6.1	Accuracy Scores of the Vectors Aligned Using VecMap	41
6.2	Coverage Scores for the Vectors Aligned Using VecMap	41

For/Dedicated to/To my...

Chapter 1

Introduction

1.1 Dictionaries

Dictionaries are living records of a society's language usage. Languages change over time, people adopt new words for new senses while others fall out of use. Concepts appear as a result of technological advancements or social shifts, giving birth to new senses and words to define them. Meanwhile, the term *dictionary* is a broad one to define. On its own, it brings forth the monolingual dictionary into consideration [1]. This type of dictionary presents words alongside their definitions following an alphabetical order. The intention is to inform the user about the words [2]. Other types of dictionaries vary with regard to their use case, target audience, and scope. For instance, bilingual dictionaries present words alongside their translations in the target language, often used by language learners or translators. Domain specific dictionaries list technical terms that target people who are familiar with the terminology.

The term that precedes the entries is called *headword* or *lemma*. Usually, lemmas are the form of a word without inflections. The sense they convey is as comprehensive as possible, reducing the number of otherwise redundant entries that would have been the derivatives of the unmarked form [3].

Dictionaries also inform the user about how senses relate to each other. **Polysemous** words share the same spelling while having related, often derivative meanings. For example; under the entry for the term *bank*, a definition might clarify the meaning *financial institution* while another can define *the building of a financial institution*. In contrast, **homonymous** words have distinct meanings while having identical spellings through coincidence. Formal definition of homonymy separates sound based and spelling based homonymy differently as homophones and homographs but for the

purposes of our text based arguments, we do not delve into the specifics. The *bank of a river* is homonym to the given examples. Homonyms are often shown in discrete blocks of descriptions.

Synonymity is another lexical relation we are interested in. A word is synonymous to another if they share the same meaning but are not spelled alike, such as the terms *right* and *correct*. However, synonymity is seldom shown in dictionaries.

Dictionaries take an immense amount of time and expertise to prepare. We can talk about the examples after narrowing our scope down to the dictionaries that are still available today. A survey by Uzun [4] notes that the first instalment of the modern Turkish dictionary, led by a team of experts, has taken over 6 years to prepare. Kendall [5] talks about how Noah Webster, the writer of the *An American Dictionary of the English Language* had to mortgage off his home in order to finish his project which took over 26 years. The bulk of this effort is collecting documents and other written material in order to establish a *corpus* [4]. This endeavour is necessary since a corpus is crucial to create the vocabulary of a language. Once the corpus is at hand, researchers can extract the lemmas. The resulting wordstock is called the *lexicon* of the language.

The internet radically changed the way researchers aggregate data. The advancements in digital storage technology allowed the data to be persistent. Improvements in networking ensured that people can share the volume of it among themselves. With the popularization of social media, the internet generates everyday conversations at an unprecedented rate that researchers are using for natural language applications. Moreover, efforts on open, collaborative, web based encyclopedias generate structured, multilingual data often used in machine translation and text categorization tasks. Once the cumbersome task of corpus attainment is now akin to web crawling. With the digitized data, it was only natural for dictionaries to go digital as well since it's generally acknowledged that they are no longer viable if they are not electronic [1].

1.2 WordNet

George A. Miller started the WordNet project in the mid-1980s. On its early days, project members studied theories that were aimed towards enabling computers to understand natural language as intrinsically as humans do. While working on then popular semantic networks and sense graphs, they have started something that will evolve into an expansive, influential resource [6].

Traditional dictionaries are rigid, constrained by the nature of the printed form. Today, people can browse WordNet via queries, like an online dictionary or a thesaurus. Behind the scenes, a sprawling lexical database has relationship information for more than 117000 senses. Figure 1.1 shows a brief result for the query string “run”.

Noun

S: (n) **run**, tally (a score in baseball made by a runner touching all four bases safely) *“the Yankees scored 3 runs in the bottom of the 9th”; “their first tally came in the 3rd inning”*

direct hyponym / full hyponym

- S: (n) earned run (a run that was not scored as the result of an error by the other team)
- S: (n) unearned run (a run that was scored as a result of an error by the other team)
- S: (n) run batted in, rbi (a run that is the result of the batter’s performance) *“he had more than 100 rbi last season”*

direct hypernym / inherited hypernym / sister term

- S: (n) score (the act of scoring in a game or sport) *“the winning score came with less than a minute left to play”*

derivationally related form

- W: (v) run [Related to: run] (make without a miss)
- W: (v) tally [Related to: tally] (keep score, as in games)
- W: (v) tally [Related to: tally] (gain points in a game) *“The home team scored many times”; “He hit a home run”; “He hit .300 in the past season”*

S: (n) test, trial, **run** (the act of testing something) *“in the experimental trials the amount of carbon was measured separately”; “he called each flip of the coin a new trial”*

S: (n) footrace, foot race, **run** (a race run on foot) *“she broke the record for the half-mile run”*

Verb

- S: (v) **run** (move fast by using one’s feet, with one foot off the ground at any given time) *“Don’t run—you’ll be out of breath”; “The children ran to the store”*
- S: (v) scat, **run**, scarper, turn tail, lam, run away, hightail it, bunk, head for the hills, take to the woods, escape, fly the coop, break away (flee; take to one’s heels; cut and run) *“If you see this man, run!”; “The burglars escaped before the police showed up”*

FIGURE 1.1: WordNet result for the query “run”, truncated for brevity.

WordNet lists terms, much like a traditional dictionary, alongside its polysemes but also their homonyms. Additionally, there is a horizontal association; for any sense, the lemmas that share the row with the target term are synonyms. This set of synonyms is aptly named *synsets*. A short description is also provided to clarify the meaning. These descriptions, hence the meanings for any synset is unique within the WordNet. During this discussion, we have used sense and synset interchangeably.

WordNet also includes other relationships such as *hypernymy* and *hyponymy*, semantic

relation of senses being type-of one another [7].¹ For instance, the term “building” is a hyponym of “restaurant” since it encompasses a more general sense; the restaurant is type of a building. While coffee shop is a hypernym to the restaurant since it is a more specific sense.

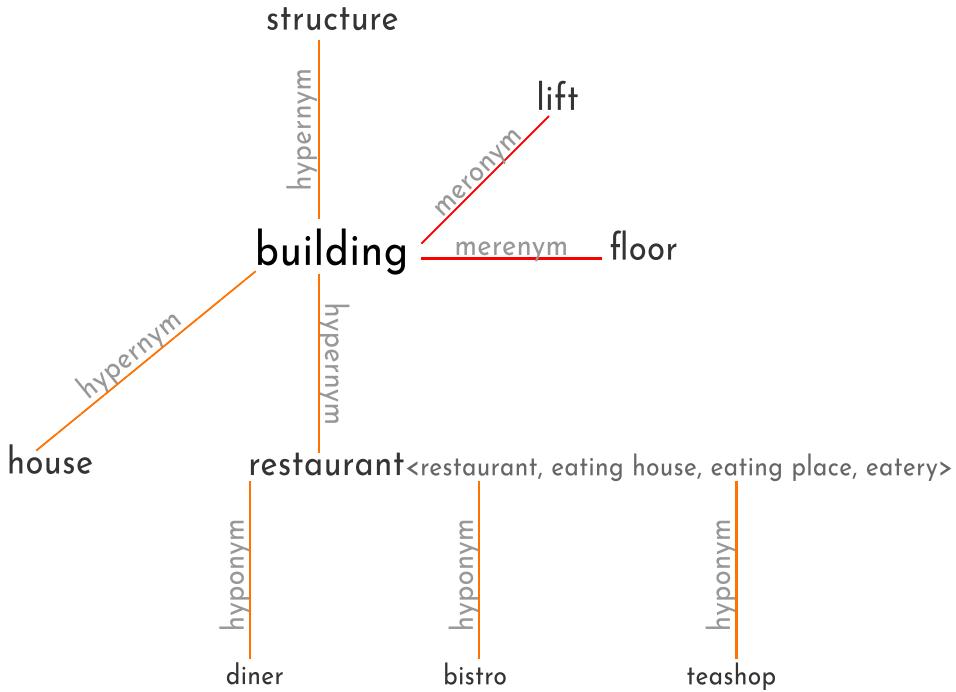


FIGURE 1.2: Example of a WordNet relationship graph. Hyponymy/hypernymy relations are shown from the hierarchy level of

One other relation is the meronymy, defined as a sense being part of or a member of another [8]. Keeping to our building example, windows are meronym to buildings. Other relationships exist but listing them is outside the scope of this thesis. Bottom line is the effort that has gone through to map 117,000 senses according to different semantic relationships. Sagot & Fišer [9] argue that the semantic relationships between senses are not tied to a specific language. With this assumption at hand, we can infer the effort behind the WordNet does not need to be repeated but can be translated to other languages.

Since its inception, other projects built lexical databases, using the same WordNet design. Fellbaum [10] talks about the correct terminology that we abide for the thesis; “As WordNet became synonymous with a particular kind of lexicon design, the proper

¹not to be confused with homonymy

name shed its capital letters and became a common designator for semantic networks of natural languages”. Hence *WordNet* refers to English Princeton WordNet, while *wordnets* created for other languages are not stylized.

1.3 Multilingual Wordnets

Authorities list more than 7000² living languages but only 40³ of them have a sizeable presence on the internet. Among this small fraction, English is the dominant language of the web. English is not the centrepiece for natural language processing research because of any linguistic attribute. It is simply the most abundant language on web, giving researchers data to work with.

Natural language processing library spaCy ⁴ resorts to lemmatizations such as =PRON= to denote pronouns in order to collapse the senses for “I” “you”, “them” etc.. The sense and the accompanying word for the brother of a person’s father or mother differs in Turkish, Danish, Chinese and Swedish among other languages while both collapse to “uncle” in English. Studying other languages can provide insight towards concepts that are not present in English.

Translation, information transfer from foreign languages is a valid way of enriching a language’s corpora; if a term that for a sense does not have a match in the target language, it is a good indication for the linguists of that language to look into their lexicons and work towards expanding it [3]. Further research in the area contributes to languages other than English having access to tools that will incorporate them into the literature.

Open Multilingual WordNet [11] set out to discover the effects related to the choice of license for wordnets. Their criteria for usefulness is the number of citations a publication tied to the wordnet has gotten on literature. They identified two major problems with the current distributions;

- some projects have picked restrictive licenses, effectively barring access to their tools for research purposes.

²<https://www.ethnologue.com/statistics>

³https://w3techs.com/technologies/history_overview/content_language/

⁴<https://spacy.io/>

- the structures of the wordnets are not standardized, creating additional cost for creating programs to parse and use the wordnets.

In order to overcome the standardization issue, Bond & Paik have aligned the wordnets according to their English Princeton WordNet lemma ids and have written individual scripts to parse them. They are currently hosting the results from a single source.⁵

With alignment information at hand, we have created our dataset that we will assume to be perfectly aligned; a golden corpus. Among the 34 wordnets available on Open Multilingual WordNet, only 6 of them have gloss information available. Given this thesis will only investigate the ability to map senses using definitions of the sense, we used the subset of Albanian [12], Bulgarian [13], Greek [14], Italian [15], Slovenian [16] and Romanian [17] wordnets. Table 1.1 shows brief statistics about them. We should note that the languages of the wordnets used in the thesis are all present in the 40 languages that have a significant presence on the internet that we have mentioned before. We have constrained this study to use only the freely available wordnets and not considered wordnets that are gated behind restrictive licenses.

Name of the Project	Language	Number of Definitions
Albanet	Albanian	4681
BulTreeBank WordNet	Bulgarian	4959
Greek Wordnet	Greek	18136
ItalWordnet	Italian	12688
Romanian Wordnet	Romanian	58754
SloWNet	Slovenian	3144

TABLE 1.1: Summary of the Wordnets used.

1.4 Thesis Goals

In this thesis, we will study document matching and document retrieval methods.

We will evaluate existing methods for their performance on cross-lingual document retrieval but our documents are dictionary definitions which are short, descriptive snippets of text. At the end of this study, we will answer the following research questions;

⁵<http://compling.hss.ntu.edu.sg/omw/>

1. Is it possible to create wordnet like lexical databases using unsupervised document matching and retrieval techniques.
2. How well does the studied techniques perform.
3. What attributes need to be considered regarding the available data.

1.5 Thesis Outline

Fill later...

Chapter 2

Background Information & Related Work

Somers puts down modern dictionaries in their article *You're Probably Using the Wrong Dictionary*; “The definitions are these desiccated little husks of technocratic meaningese, as if a word were no more than its coordinates in semantic space.” [18]. From the perspective of an author, the efficiency of dictionaries might be worrisome but we will build this thesis on the presumption that dictionary definitions can indeed be represented in some semantic space.

2.1 Word Embeddings

Word embeddings are real valued dense vectors for words. Recently, language modelling studies have been focusing on explicitly learning word embeddings in order to show words or phrases as points on a low dimensional latent space. On the other hand, earlier research had been obtaining what we can call feature vectors for words while studying natural language processing tasks such as named entity recognition or part of speech tagging [19, 20]. The vector representation allows researchers access to the tools of the broad literature in linear algebra and machine learning, since they are intuitive for humans to interpret, more so for machines. Vectors can be compared as a measure of semantic similarity or composed together to build more expansive sentence, paragraph or document representations.

Induced embeddings can be saved to the disk in matrix notation. Each row is labelled with a token that is denoted by the following n real numbers, as popularized by the

open source package *word2vec*.¹ Researchers have been sharing their models on the internet so that other researchers can simply download and use them in their own applications. Word embeddings available in this manner are often called as *pre-trained* models. Examples of pre-trained word embeddings are *word2vec*, *GloVe*, *Numberbatch* and *fastText*.

In the following section, we will briefly present the history of word embeddings. Word embeddings is a sprawling subject that researchers has been building upon using the ideas from probabilistic, statistical and neural network models. We have omitted crucial contributions that optimized models and brought the literature where it is today due to space constraints, following a path that will lead into the preliminary behind the models we have chosen.

2.1.1 History of Word Representations

In order to talk about how words can be mapped to a multidimensional space, we should first talk about how the idea that *they can* has came about.

Linguistic Background

In his 1954 article, Harris [21] introduced his ideas which later came to known as *distributional hypothesis* in the field of linguistics. He argued that similar words appear within similar contexts. The famous quote by Firth [22] captures the idea as; “You shall know a word by the company it keeps!” For instance, the semantic similarity or relatedness between the terms *jacket* and *coat* can be theoretically shown since they will be accompanied by similar verbs, such as *wear*, *dry clean* or *hang*, and similar adjectives such as *warm* or *leather*.

Early attempts to show words on a semantic space is studied in Osgood *et al.* [23]. Authors suggested representing concepts using orthogonal scales and relying on human judgement to score meanings on the axes. An example concept from their study is given in Figure 2.1.

However, for a researcher to pick appropriate scales or have meaning extracted by hand would be infeasible for natural language processing tasks [24].

¹<https://code.google.com/archive/p/word2vec/>

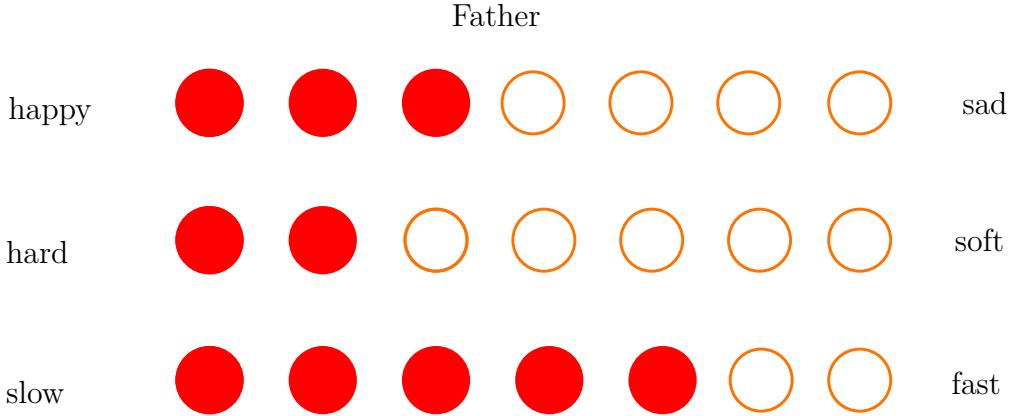


FIGURE 2.1: Sense representation using human judgement scores for the concept “Father” adapted from Osgood *et al.* [23].

Even though Harris argued that “language is not merely a bag of words”, using un-ordered collection of word counts to capture the semantic information will be used in the literature and be known as the *bag-of-words* model.

Vector Space Model

The history of word embeddings is tightly coupled with *vector space models* that initially appeared in the field of information retrieval. The intent was to extract vectors that represented *documents*. First vector space model developed by Salton *et al.* [25] was presented in “A Vector Space Model for Automatic Indexing”. It was the first application of bag-of-words hypothesis on a corpus to extract semantic information [26]. Salton *et al.* presented the novel idea of a *document space*, that consist of fixed sized vectors as the columns of a term document matrix. The dimensions of the row vectors were the whole vocabulary of the corpus.

In this space, a document D_i is represented using t distinct terms as a row vector;

$$D_i = (d_{i1}, d_{i2}, \dots, d_{it})$$

The weights for each of the index terms d_{ix} are calculated using the *tf-idf* measure introduced by Jones [27]. *tf-idf* is the multiplication of two metrics;

tf the number of times a term k occurs in a document

idf the inverse of the number of documents that contain k .

The weighting scheme was selected to “assign the largest weight to those terms which arise with high frequency in individual documents, but are at the same time relatively rare in the collection as a whole” [25].

The vector space model allowed Salton *et al.* to handle the similarity between documents as the angle between two vectors. Above all, they have shown that there is merit to handling documents as real valued vectors.

2.1.2 Latent Semantic Analysis

Deerwester *et al.* [28] introduced latent semantic analysis in order to address a crucial problem with the vector space model. They have identified that the term document matrix approach cannot handle synonyms and homonyms due to the fact that vector space model requires the words to match exactly between the two documents. Synonymy is an issue because the query can have terms that have the same meaning as the target word, without getting matched. On the other hand, homonyms can match with an unrelated word. In order to answer these issues, their model seeks the higher order latent semantic structure in order to learn the *similarity between words*.

Latent semantic analysis starts with a word co-occurrence matrix X . The terms of the matrix is weighted by some weighting scheme. While original study used raw term frequencies, ***tf-idf*** is a possibility while Levy *et al.* [29] reports pointwise mutual information (PMI) [30] as a popular choice. A term document matrix X is then factorized into three matrices using singular value decomposition [31];

$$X = T_0 S_0 D'_0$$

Where the columns of T_0 and D'_0 are orthogonal to each other and S_0 is the diagonal matrix of singular values. The singular values of S_0 can be ordered by size to keep only the k largest elements, setting others to zero [28]. The resulting matrix is shown as S .

Deerwester *et al.* talk about the use for the resulting matrix as follows;

Each term or document is then characterized by a vector of weights indicating its strength of association with each of these underlying concepts.

That is, the “meaning” of a particular term, query, or document can be expressed by k factor values, or equivalently, by the location of a vector in the k -space defined by the factors.

They have used this technique to solve document similarity task and briefly mentioned *word similarity*. Landauer & Dumais [32] later studied word similarity in full using latent semantic analysis.

2.1.3 Building Upon Distributional Hypothesis

While Deerwester *et al.* studied relatedness between words using vectors, their approach used the whole document for the co-occurrence information while the focus was clearly still on the document similarity. Schütze [33] proposed “to represent the semantics of words and contexts in a text as vectors” and built upon word co-occurrence. They theorized a context window of 1000 *characters* in order to consider words that are close to the target word instead of the whole document in co-occurrence calculations. However, they claimed that the computation power available was not suitable yet to fully tackle the task.

Lund & Burgess [24] took the challenge and experimented upon 160 million *words* taken from the Usenet, a precursor to the internet. They used a context window of 10 words and provided a method to obtain feature vectors to represent the meaning of words. However, intricate tuning of word co-occurrence generated associatively similar vectors instead of semantically similar ones. To give an example, *school* is related to *student* while they are not similar.

2.1.4 Distributed Vector Representations

Bengio *et al.* [34] proposed learning word representations using a feedforward neural network. Their model would learn feature vectors for words using a predictive approach instead of counting based approaches we have presented until now. Although neural networks have been proposed to learn a language model by Xu & Rudnicky [35], the main contribution of Bengio *et al.* is to use an *embedding layer*, in order to attack *curse of dimensionality*. For a corpus with vocabulary V , there are $|V|$ dimensions for the language model to learn and taking n -*gram* representations into consideration, the problem grows exponentially. Using m dimensions in the embedding layer allowed

Bengio *et al.* to represent words using manageable, more representative dimensions and the problem scaled linearly.

The setup for the neural network starts with the one hot encoded vector representation of the context for a word w . This context window is similar to those used in statistical models that predicts the word w_t using the words that lead up to w_t . In other words, context window of w_t is T words on the left of w_t .

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1})$$

The input layer is projected into an embedding layer, later to a softmax layer to get a probability distribution in order to minimize the following cost function.

$$\hat{P}(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}} \quad (2.1)$$

However, this formulation is too expensive computationally since all vocabulary needs to be considered for the sum in the denominator. The curse of dimensionality problem were shifted to the final layer of the neural network which will be solved later using hierarchical softmax. Authors reported training times around 3 weeks using 3 to 5 context window sizes and vocabulary sizes around 17000.

Collobert & Weston [20] suggested a deep neural network model in order to learn feature vectors for various natural language processing tasks. Their proposed approach for language model is important for our case since it explicitly learned distributed word representations or simply *word embeddings*. They have introduced two key ideas;

- Instead of using a context window that used words left of the target word to estimate the probability of the target word, they have placed the context window *on* the target window, using n words for left and right of the target word.
- They introduced negative examples, where they randomly changed the middle word with a random one. This allowed them to use the ranking cost;

$$\sum_{s \in S} \sum_{w \in D} \max(0, 1 - f(s) + f(s^w))$$

Collobert & Weston [20] claim that these additions allowed their system to learn the representation better rather than the probability.

2.1.5 Pre-trained Embeddings

P. Turian *et al.* [36] evaluated the performance of different word representations as word features you can include into an existing task. They summarize their contribution as follows;

Word features can be learned in advance in an unsupervised, task-inspecific, and model-agnostic manner. These word features, once learned, are easily disseminated with other researchers, and easily integrated into existing supervised NLP systems.

[...]

With this contribution, word embeddings can now be used off-the-shelf as word features, with no tuning.

2.1.6 Popularization of Word Embeddings

word2vec package [37–39] popularized word embeddings. There are two aspects of the work done by Mikolov *et al.* that contributed to the fact;

- Their model captures the semantic and syntactic attributes of words and phrases on a large scale with good accuracy, trained on billions of words using a shallow neural network, keeping the computational cost down.
- They published their code as well as their pre-trained embeddings as an open source project².

The second point is self explanatory but in order to argue about the first one, we should report the algorithms behind word2vec.

The *skip-gram model* introduced by Mikolov *et al.* [37] differs from the previous methods by predicting the surrounding words given the target word (Figure 2.2).

²<https://code.google.com/archive/p/word2vec/>

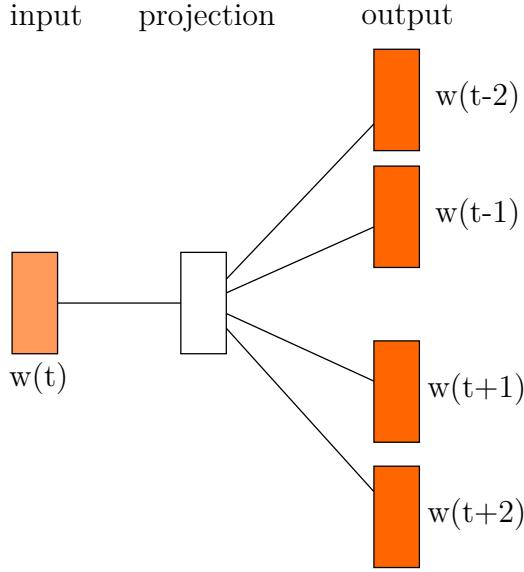


FIGURE 2.2: Skipgram architecture adapted from Figure 1 of Mikolov *et al.* [38]

In “Distributed Representations of Words and Phrases and Their Compositionality”, it is defined as follows;

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t) \quad (2.2)$$

The w_1, w_2, \dots, w_T are the context of the word w_t . We should note that, Levy *et al.* [29] has identified that this window size is *dynamic* in the open source implementation of word2vec, where the actual window size is sampled between 1 and T .

Levy *et al.* [29] compared the performance of count based and prediction based word representation models. Representation algorithms they considered are;

- Positive pointwise mutual information (PPMI) [30, 40]
- Singular Value Decomposition on PPMI Matrix (Latent Semantic Analysis) [28]
- Skip-Gram with Negative Sampling [38]
- Global Vectors for Word Representation [41]

They found out that choice of a particular algorithm played an insignificant role compared to choosing the right approach during training, mainly picking correct *hyperparameters*. Besides, the amount of data the models are trained on were more important

than the model itself. They used this finding to counter the results reported by Baroni *et al.* [42]. Baroni *et al.* claimed that predictive models outperformed count based models. On the other hand, Levy *et al.* noted that Baroni *et al.* used count based models without hyperparameter tuning, denying them from “tricks” developed in the word representation literature.

2.1.7 fastText

Armed with the fact that a good word representation model should tune their hyperparameters and should be trained on a large dataset, we set our sights on *fastText*. On their website, authors define fastText as (a) “Library for efficient text classification and representation learning”. The ideas behind it are presented in Mikolov *et al.* [43]. Overall, it builds upon word2vec [38] by adding position dependent features presented in Mnih & Kavukcuoglu [44] and character n-grams suggested on Bojanowski *et al.* [45].

Let us turn our focus towards “Enriching Word Vectors with Subword Information” [45]. Instead of using a context window to learn the representation of the target word or predicting surrounding words given the centre word like in the skip-gram model, Bojanowski *et al.* learn representations for character n-grams.

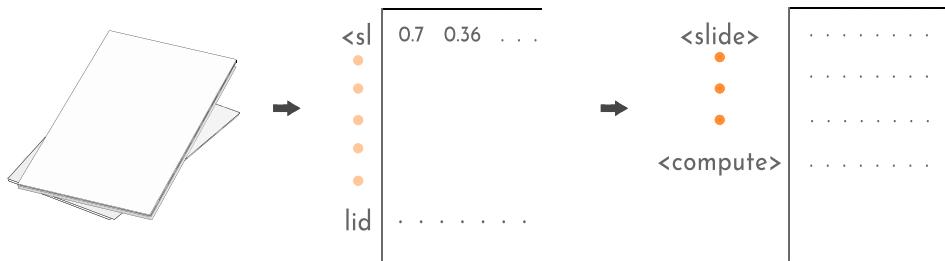


FIGURE 2.3: Overview of fastText embedding model

With the subword vectors z_g for every n-gram g at hand, authors take a dictionary of n-grams of size G and for a given word w , they denote $G_w \subset 1, \dots, G$ as the n-grams of w .

$$s(w, c) = \sum_{g \in G_w} z_g^T v_c$$

Grave *et al.* [46] trained fastText model on different languages using Wikipedia and Common Crawl data. Wikipedia is a curated encyclopaedia that can be used as a multilingual corpora with 28 languages that have over 100 million tokens and 82 languages with 10 million tokens [46]. Considering the original word2vec was trained on 100 billion tokens, Grave *et al.* also used Common Crawl, a non profit project that collocates web pages and publishes the data publicly. This data can be inadvertently noisy which Grave *et al.* addressed using linewise language identification and by removing duplicate lines that often appear as leftover boilerplate on many cites. While Wikipedia provided them with a curated signal, Common Crawl data helped with capturing as many contexts as possible to train their distributed model. They are currently hosting the pre-trained word embeddings on their website ³.

2.1.8 ConceptNet Numberbatch

ConceptNet⁴ is a knowledge graph. Like the WordNet, it presents semantic relationships between concepts but the defined relationships are more fine grained. The relationships are collocated from various resources like Wikimedia, English Princeton WordNet and knowledge databases like DBPedia. ConceptNet stylizes its relations with the /r/Relationship syntax. For instance, being distinct members of a set relation is defined as /r/DistinctFrom which includes concepts like August and September.⁵. Compared to WordNet, more subjective, human centric relations are defined such as /r/MotivatedByGoal, relationship between *compete* and *win* or /r/ObstructedBy, the relationship between *sleep* and *noise*. Total number of relations that are available between two concepts is 36. Moreover, these relationships are *weighted* depending on how present they are. For example, the concept *run* is related to *fast* with a weight of 9.42 but the similarity between *race* and *run* is weighted at 2.54. Finally, the ConceptNet is multilingual, encompassing 304 languages in total with 77 of them reported as having moderate support that can be used in downstream tasks. Details of the knowledge

³<https://fasttext.cc/>

⁴<http://conceptnet.io/>

⁵<https://github.com/commonsense/conceptnet5/wiki/Relations>

graph and the procedure is explained by Speer *et al.* [47] in the paper “ConceptNet 5.5: An Open Multilingual Graph of General Knowledge”.

Speer *et al.* also present their resource in word embedding matrix form for ease of use in the current plug-and-play environment. In order to get word embeddings from the knowledge graph, first they prepare a symmetric term-term matrix X where an element $w_{i,j}$ is the sum of all edge weights between terms i and j . This is similar to a word co-occurrence matrix but instead of constructing the matrix using unstructured text, this approach builds upon semantic connections between senses. Speer *et al.* reports that their approach learns *relatedness* more so than *similarity* which was reported by the word co-occurrence approaches before.

With the term-term matrix at hand, Speer *et al.* weigh the terms using positive pointwise mutual information (PPMI) as suggested by Levy *et al.* [29] and reduce it to 300 dimensions, as is standard set by word2vec, using truncated SVD. Resulting matrix is similar to approaches set by Deerwester *et al.* [28] or Pennington *et al.* [41] but Speer *et al.* enrich it further using *retrofitting* as proposed by Faruqui & Dyer [48]. Pre-trained word2vec and GloVe embeddings are incorporated into the reduced term-term matrix to finally obtain embeddings that include both distributional and semantic relatedness signal. Authors call their finalized model *ConceptNet numberbatch*. Speer *et al.* reports state of the art results compared to word2vec embeddings on word relatedness and models built using their embeddings get scores equivalent to humans on SAT-style analogy tasks.

2.2 Bilingual Word Embeddings

Cross lingual embedding models optimize similar objectives. Only source of variation is due to the data used and the monolingual regularization objectives employed [49].

2.3 Approaches in Wordnet Generation

2.3.1 History of Wordnet Generation

We have mentioned the lexical database WordNet created by Princeton University. To reiterate, WordNet is a lexical database with human annotated collection of senses and relationships among them. The relationships are hierarchical so they can be followed

along to reach new nodes due to the transitive property. The format itself has become the standard for databases that present meanings and concepts [50].

Glosses or the definitions that go along with synsets were not initially part of the WordNet design. Authors believed that “definition by synonymy” would be enough. In other words, definition of a synset can be derived from the lemmas that make up the synset. As the number of items in the WordNet grew, only then short glosses, followed by longer definitions got included in WordNet [6].

Synset	Gloss
{glossary, gloss}	an alphabetical list of technical terms in some specialized field of knowledge; usually published as an appendix to a text on that field
{dog, domestic dog, Canis familiaris}	(a member of the genus <i>Canis</i> (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds)
{university}	(the body of faculty and students at a university)
{depository financial institution, bank, banking concern, banking company}	(a financial institution that accepts deposits and channels the money into lending activities)

TABLE 2.1: Example synsets and their respective glosses

WordNet has been used in various natural language processing applications over the years such as text summarization [51]. Since the original WordNet was prepared for English over many years of work, efforts for creating an equivalent resource for other languages went under way. Arguably, EuroWordNet set the standard for creating wordnets for languages other than English [52, 53].

EuroWordNet⁶ project was initiated to introduce the benefits of English Princeton WordNet for other languages. Additionally, an interlinked semantic network can be a research topic on lexicalization patterns of languages, finding conceptual clusters of vocabularies or cross lingual text retrieval [52, 54]. EuroWordNet project included 7 wordnets for languages other than English and an adapted English wordnet. Due to

⁶<http://projects.illc.uva.nl/EuroWordNet/>

the effort needed to create a wordnet from scratch, Vossen averted the EuroWordNet project from creating full scale semantic lexicons and prioritized the connectivity between wordnets. All in all, Vossen [52] defines the aims as;

1. to create a multilingual database;
2. to maintain language-specific relations in the wordnets;
3. to achieve maximal compatibility across the different resources;
4. to build the wordnets relatively independently (re)-using existing resources;

A sense in one language might not have a direct equivalent in another. Cultural differences or linguistic differences between languages contribute to this fact [55]. This phenomenon is called a lexical gap or untranslatability. EuroWordNet addresses lexical gaps using Inter-Lingual-Index (ILI). ILI is a higher order list of meanings just for wordnet synsets to align themselves to. Elevating alignment from English Princeton WordNet to ILI allowed language specific structures to exist in wordnets while keeping the connections among themselves.

Two approaches for wordnet generation were defined;

Merge Approach where a wordnet structure is formed in the target language with synset selection and relation mapping. Then the connections between the new wordnet and English Princeton WordNet can be established.

Expand Approach where English Princeton WordNet is (machine) translated to target language [56], preserving connection information with a trade-off where the target language wordnet will be biased towards the relationships of the English Princeton WordNet and may not include target language specific lexical connections.

In order to maintain as much language specific lexical connections as possible while having a starting point for evaluation of target wordnets, EuroWordNet project offered “Base Concepts”. This idea later evolved into 1000 and to 5000 core synsets that are compiled from most frequent, connected and representative synsets to be used for evaluating wordnet generation [57].

2.3.2 Examples of Wordnet Generation

Whether an Arabic wordnet is attainable or not using parallel corpora was explored by Diab [58]. Arabic differs from other languages explored so far due to its unusual morphological nature. Using their proposed method, they observed 52.3% of words they processed are sufficient for a future Arabic wordnet. They have also reiterated that semantic relationships in a language are transferable to a target language’s wordnet.

Further approaches using parallel corpora to align the target language with English Princeton WordNet were used by Sagot & Fišer [9] in order to create a production ready French wordnet and by Fiser [59] to create a Slovene wordnet.

Approaches that used machine translation to get potential synsets for the target language was explored by **lam_automatically_2014**. They proposed two approaches for this task. First approach uses a single bilingual dictionary to translate English Princeton WordNet lemmas to target language to form synsets. Second approach translates existing wordnet synsets to English and then translates them to target language.

Our focus is to study approaches that proposed *word embedding* based methods for the task. Sand *et al.* [60] used word embedding models to extend Norwegian wordnet with relationships on existing synsets or by introducing new synsets. Cosine similarity is used to retrieve nearest neighbours. Then a threshold value is used to cut off any potential synsets below a value. They evaluated against accuracy, the percentage of relations that were correct. Attachment was also used similar to recall scoring. It is the percentage of considered synsets that were above the threshold value.

Arguably most relevant to our study, Khodak *et al.* [61] proposed an unsupervised method for automated construction of wordnets. In their paper “Automated WordNet Construction Using Word Embeddings” they present 3 approaches and compare their precision, recall and coverage.

First off, they picked 200 French and Russian adjectives, nouns and verbs. These words are selected based on if they have a translation in the core sense list provided by English Princeton WordNet. Their approach is presented by handling a word w in the selected list.

Initially, the presented method collects the possible translations of w using a bilingual dictionary and uses them as lemmas to query English Princeton WordNet. As we

have mentioned, lemmas return synsets from English Princeton WordNet in the form `<lemma.pos.offset>`. Furthermore, every synset includes possible lemmas that can represent the sense of the synset. These lemmas form the set S . By translating the English Princeton WordNet lemmas to target language, they obtain the set T_S . Elements of T_S and w are both in target language so a monolingual word embeddings can be used; As a baseline, they calculate the average cosine similarity between word embeddings of every element of S and w .

$$\frac{1}{S} \sum_{w' \in S} v_w \cdot v_{w'}$$

In order to improve the discriminative power of their target vectors, they use synsets relations, definitions and example sentences. These are used to get a more representative vector v_L . Given resources can be thought as sentences and sentence embeddings are calculated by getting a weighted average of word vectors of words that make up the sentence. The weighting is called smooth inverse frequency by Arora *et al.* [62].

$$v_L = \sum_{w' \in L} \frac{a}{a + Pw'} v_{w'}$$

We will talk about sentence embeddings in detail for our approach in Section 3.1.

Third approach is word sense induction [63].

1. Expand model, Vossen [52], fixed synsets are translated from English to target language.
2. Link English entries from machine-readable bilingual dictionaries to English Princeton WordNet senses Knight & Luk [56].

Gordeev *et al.* [64] uses unsupervised cross-lingual embeddings to match cross-lingual product classifications. Working on taxonomy matching, they use out of domain pre-trained embeddings due to small size of their corpora and investigate methods using untranslated and translated text.

Lesk [65] represent words using their gloss. Relied upon traditional dictionaries.

Banerjee & Pedersen [66] developed on lesk algorithm and included WordNet definitions.

Metzler *et al.* [67] talked about short text retrieval and lexical matching. They reported that lexical matching is good for finding semantically identical matches.

Chapter 3

Unsupervised Matching

3.1 Linear Assignment Using Sentence Embeddings

From word embeddings that represented single tokens, research also tackled the representation of longer pieces of text like documents, paragraphs and sentences. Implementations like Le & Mikolov [68] extended the skip-gram idea of Mikolov *et al.* [38] to represent *paragraphs* as feature vectors and used them to predict surrounding paragraphs in the text. While approaches like Kiros *et al.* [69] trained an encoder that constructed surrounding sentences to learn to learn sentence representations. However, there is an assumption of a continuous text for the given models. When the text that we would like to show on a latent space is not part of a longer piece of text but *discrete* pieces, that assumption does not hold. With the dictionary definitions, we have such a case. Our dictionary definitions are comprised of 10 to 11 words and there is no relation from one distinct dictionary definition to another. In other words, they are not continuous. One other case where a similar situation occurs is *twitter*. *Tweets* are short pieces of text due to the 280 character constraint imposed by the platform. With such short pieces of text, instead of paragraph embeddings, we can talk about *sentence embeddings*. A sentence embedding model should ideally capture the collective meaning of the short text where every word is potentially informative.

Thankfully, Wieting *et al.* [70] studied sentence embeddings and reported that averaging word embeddings that make up a sentence to get sentence embeddings is a valid and surprisingly effective approach. Arora *et al.* [62] built upon that idea and presented that weighed average of word vectors perform so well, they decided to call it; “A Simple but Tough-to-Beat Baseline for Sentence Embeddings”. In the suggested approach, word embeddings are weighed with a scale called *smooth inverse frequency*. Smooth

inverse frequency is reported as $a/(a + p(w))$ where a is a parameter and $p(w)$ is the word frequency. The authors point out that the metric is similar to *tf-idf* weighting scheme if “one treats a ‘sentence’ as a ‘document’ and make the reasonable assumption that the sentence doesn’t typically contain repeated words”. These assumptions hold for us so we scaled our word embeddings using *tf-idf* weights to get sentence embeddings.

Parallel to Zhao *et al.* [71] used two approaches for SemEval-2015 Task 2: Semantic Textual Similarity ¹. First, for a sentence $S = (w_1, w_2, \dots, w_s)$ where the length of the presumably small sentence is $|S| = s$ and the word embedding of a w_t is v_t ;

- They summed up the word embeddings of the sentence $\sum_{t \in S} v_t$
- Used information content [72] to weigh each word’s LSA vector $\sum_{t \in S} I(w_t)v_t$

Both approaches results in a vector that is in the same dimensions R^d as the original word representations.

Edilson A. Corrêa *et al.* [73] expanded upon this simple yet effective idea to tackle the SemEval-2017 Task 4², Sentiment Analysis in Twitter. In order to acquire embeddings that represented *tweets*, they weighed the word embeddings that made up a tweet; $\text{tweet}_i = (w_{i1}, w_{i2}, \dots, w_{im})$ with the *tf-idf* weights. For the *tf-idf* calculation, they cast individual weights as documents so that term frequency become the term count in a single tweet while document frequency become the number of tweets the term w_t occurs.

We have mentioned that our dictionary definitions are not continuous. Yet, we advocate using *tf-idf* weights to weigh our word embeddings to get sentence embeddings. In order to clarify, let us present Table 3.1.

For the *tf-idf* calculations, we followed a similar approach. The term frequency is the raw count of a term in a dictionary definition. While the document frequency is the number of dictionary definitions where w_t occurs.

Then, with the term-embedding matrix at hand, we have calculated definition embeddings using:

$$S_{\text{emb}}(S) = \sum_{w_i \in S} tf_{w_i, S} \cdot idf_{w_i} \cdot Emb_w(w_i) \quad (3.1)$$

¹<http://alt.qcri.org/semeval2015/task2/>

²<http://alt.qcri.org/semeval2017/task4>

TABLE 3.1: Some definitions from English Princeton WordNet

turn red, as if in embarrassment or shame
a feeling of extreme joy
a person who charms others (usually by personal attractiveness)
so as to appear worn and threadbare or dilapidated
a large indefinite number
distributed in portions (often equal) on the basis of a plan or purpose
a lengthy rebuke

Every word that makes up a definition is scaled by its vector in \mathbb{R}^n , then concatenated to form sentence embeddings on \mathbb{R}^n . Given the N vectors from source and target language, we hypothesize that there exists a matching where every source definition vector is perfectly mapped to one target vector. Given that this problem naively iterates over $N!$ matchings, we have looked into an algorithm.

3.2 Linear Assignment Algorithm

3.3 Results

Language	Percentage of Correctly Matched Definitions		
	fastText 1M	fastText 500k	Numberbatch
bg	0.39	0.41	0.19
el	0.37	0.38	0.14
it	0.28	0.28	0.36
ro	0.39	0.39	0.20
sl	0.15	0.15	0.06
sq	0.55	0.54	0.27

TABLE 3.2: Linear Assignment Using 2000 Definitions

Language	Percentage of Correctly Matched Definitions		
	fastText 1M	fastText 500k	Numberbatch
bg	0.35	0.36	0.18
el	0.36	0.36	0.12
it	0.25	0.25	0.32
ro	0.36	0.37	0.19
sl	0.11	0.11	0.05
sq	0.39	0.40	0.19

TABLE 3.3: Linear Assignment Using 3000 Definitions

	fastText 1M	fastText 500k	Numberbatch
Best	0.55	0.54	0.36
Worst	0.11	0.11	0.05
Average	0.33	0.33	0.19

TABLE 3.4: Summary of Linear Assignment

Chapter 4

Dictionary Alignment as Pseudo-Document Retrieval

Document retrieval is the prototypical information retrieval task. Bush [74] theorized the possibilities of the automatic information retrieval by machines in his essay titled “As We May Think”. On his “Modern Information Retrieval”, Singhal [75] gives due credit to Luhn [76] for the initial suggestion of using word overlap in document retrieval.

Modern information retrieval techniques are far from the scope of this thesis. Yet, we can still benefit from the tried and tested methods of the early information retrieval. Considering the small collection of documents at hand, first we will investigate if we can handle the task using approaches that were available to the researchers when the size of corpora that were available to them were small as well [75]. However, we will get leverage from a state of the art tool from the modern computer science; *Google Translate*.

4.1 Monolingual Retrieval

First of all, we created a corpora suitable for the task by translating the target language’s definitions to English using Google Translate. We used the Google Cloud API¹ in order to automate the process.

Table 4.1 presents some definitions and their automated translations to English. While the translations can inform the user about the general sense of the word, sentence

¹<https://cloud.google.com/translate>

Original Definition	Translated Definition
bila pri starih Grkih in Rimljanih vojna ladja s tremi vrstami vesel	with the ancient Greeks and Romans, a three-wheeled war ship was happy
znanstvena veja matematike, ki se ukvarja s prostorskimi lastnostmi teles in njihovimi medsebojnimi odnosi	the scientific branch of mathematics, which deals with the spatial properties of the bodies and their interrelations
circumstanțe, stări de fapte, lucruri luate în calcul într-o discuție sau dezbatere	circumstances, facts, things taken into account in a discussion or discussion
Fază a lunii în care este complet iluminată	Phase of the month in which it is fully illuminated
Persoană care se ocupă cu pescuitul și uneori cu conservarea peștelui pescuit	A person who deals with fishing and sometimes with the conservation of fish fishing
E bëj diçka që të shkojë e të përputhet me një qëllim të caktuar, i bëj ndryshimet e ndreqjet e nevojshme, që t'u përgjigjet kushteve e rrëthanave të caktuara.	I do something to go and match a certain purpose, make the necessary adjustments and adjustments, to respond to certain conditions and circumstances.
pohoj, a paraqit dikujt një gjë për ta parë, për t'u njojur me të a për të gjykuar për të; zbuloj diçka që të shihet; tregoj	I assure you, did someone present a thing to see, to know him or to judge him? I find something to be seen; show

TABLE 4.1: Example definitions and translations

structure is flawed at times. We can also see word duplication when synonyms collapse into single words. These shortcomings are also reported by Groves & Mundt [77].

With the English Princeton WordNet definitions and 6 target wordnet definitions at hand, we can handle the task as *monolingual document retrieval*. We will use the vector space model that was mentioned in Chapter 2 to have real valued vectors that represent the definitions.

To begin with, the definitions that were extracted from English Princeton WordNet and the target wordnet have been collocated to form a corpora. This process includes one target wordnet at a time and experiments were simply repeated for each wordnet. The vocabulary of this corpora is shared between the wordnet pairs as they have been joined in a single language.

4.1.1 Term Document Matrix

With the definition collection and the vocabulary at hand, we created a *term-document* matrix. In this instance of the term-document matrix, rows represent the documents and columns denote individual vocabulary entries. Such a matrix C is m by n while m is the number of documents and n is the size of the vocabulary $|V|$. Our definitions are just short documents so we use the terms *definition* and *document* interchangeably with the added benefit of using d to denote either of them.

For any element of C , $c_{i,j}$ is the number of times j occurs in the document i . This matrix is sparse since documents use a fraction of the available vocabulary V . Our case with definitions is no different with an average of 10.62 words per definition with vocabulary sizes well over 10,000s.

4.1.2 Term Weighting

Elements of a term-document matrix is scaled in order to assign weights to elements according to their discriminative power [78]. Stopwords like “the”, “and” have virtually no significance and terms that are common in the corpora’s domain will have no discriminating power either, as explained by Manning *et al.* [78]. We have opted out of stopword removal considering our short definitions.

What is a good scaling method for term weighing? $tf\text{-}idf$ is composed of two scores; term frequency tf and inverse document frequency idf . idf was initially

suggested by Jones [27] and is highlighted by Manning *et al.* [78]. Term frequency $tf_{w,d}$ is the number of times term w occurs in a document d . Inverse document frequency idf_w is calculated using the number of documents that contain the term w which is also known as the document frequency of w , df_w . We have used the smooth variant of idf as follows;

$$idf_w = \log \frac{N + 1}{df_w + 1} \quad (4.1)$$

Where N is the number of documents in the corpus. $\log(\cdot)$ is used to dampen the effect of the idf .

Manning *et al.* explains the intuition behind the $tf\text{-}idf$ in his *Introduction to Information Retrieval* as follows;

[...], $tf\text{-}idf_{t,d}$ assigns to term t a weight in document d that is

1. highest when t occurs many times within a small number of documents (thus lending high discriminating power to those documents);
2. lower when the term occurs fewer times in a document, or occurs in many documents (thus offering a less pronounced relevance signal);
3. lowest when the term occurs in virtually all documents.

By weighing our term-document matrix C , for any term $x_{i,j}$

$$x_{i,j} = tf_{td} \cdot idf_t \quad (4.2)$$

Such that term $x_{i,j}$ shows the importance of term t with relation to its general significance throughout the corpus.

4.1.3 Similarity Measure

Using the term-document matrix, similarity between two documents can be calculated by assigning a scoring scheme that uses rows as document vectors. *Cosine similarity* is such a measure that can be used to find the cosine of the angle between two vectors and scales to multidimensional case trivially. Moreover, since the majority of the dimensions between the vectors are zero, cosine similarity ignores these dimensions and implicitly prioritizes the non-zero dimensions.

For the row vector \vec{d}_t and \vec{d}_p , cosine similarity between definitions t, p is

$$\cos(\theta) = \text{sim}(d_t, d_p) = \frac{\vec{d}_t \cdot \vec{d}_p}{|\vec{d}_t| |\vec{d}_p|} \quad (4.3)$$

4.1.4 Retrieval

With definition vectors and a similarity measure at hand, we split the *tf-idf* weighted term-document matrix C' as English Princeton WordNet definitions and target wordnet definitions. The target wordnet definitions were used as pseudo-queries and English Princeton WordNet definitions were used as the document collection to retrieve definitions from.

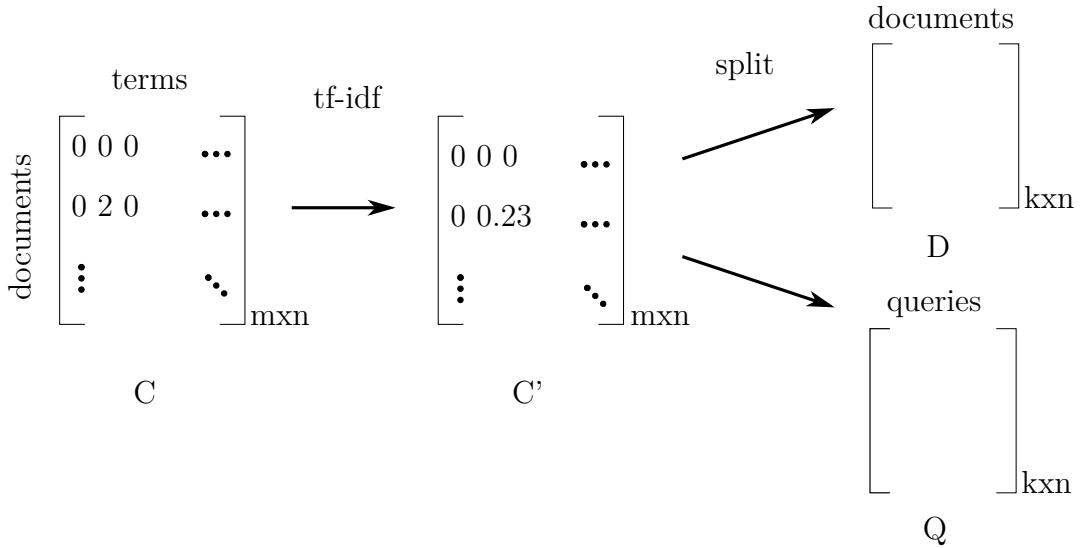


FIGURE 4.1: Term document matrix to queries and documents

In order to retrieve and rank definitions given the query definition, cosine similarity between query and each corpus definition is calculated. The similarity scores that range between 0 and 1 are sorted in descending order with the definition index attached. The retrieved documents are ranked according to their similarity score. Since the definitions are aligned across wordnets, a correctly retrieved definition will have the same index as the query definition.

4.1.5 Evaluation

Information retrieval literature usually uses *precision* and *recall* as the evaluation metric with the assumption that the user wants to receive as many relevant documents as possible [79]. However, we only have one correct definition we are interested in and cannot access relevance any further for any other retrieved definitions anyway. As a result, we used a tweaked version of *Mean reciprocal rank* (MRR) as introduced by Voorhees [80] in “The TREC-8 Question Answering Track Report”. In the report, mean reciprocal rank was used in order to evaluate a question answering track where a singular correct answer was sought among other answers and is rewarded according to the returned rank. In other words, mean reciprocal rank evaluates a system’s retrieval score with respect to the rank of a single correct answer.

$$\text{MRR} = \frac{1}{|Q|} \cdot \sum_q^{|Q|} \frac{1}{\text{rank}_q} \quad (4.4)$$

Where rank_q is the ranking of the *correct* definition for the query q .

4.1.6 Results

4.2 Cross Lingual Document Retrieval

4.2.1 Word Movers Distance

Following the popularization of the word2vec [38], Kusner *et al.* [81] introduced *Word Mover’s Distance*. By adapting the optimal transport theorem into using distances between word embeddings and documents as probability distributions that are composed of words, he suggested the following;

For an embedding matrix $X \in R^{n \times d}$ d is the dimension the embedding was learned at and n is the size of the vocabulary. In this configuration, a row $\vec{x}_i \in R^d$ of the matrix X is the embedding for the i^{th} word w_i .

Given a separate term-document matrix $C \in R^{m \times n}$, a row $\vec{d}_j \in R^n$ is the row vector for the j^{th} document d_j . In the original paper “From Word Embeddings to Document

English Princeton WordNet Definition	Translated Definition
ancient Greek or Roman galley or warship having three tiers of oars on each side	with the ancient Greeks and Romans, a three-wheeled war ship was happy
the pure mathematics of points and lines and curves and surfaces	the scientific branch of mathematics, which deals with the spatial properties of the bodies and their interrelations
everything stated or assumed in a given discussion	circumstances, facts, things taken into account in a discussion or discussion
the time when the Moon is fully illuminated	Phase of the month in which it is fully illuminated
someone whose occupation is catching fish	A person who deals with fishing and sometimes with the conservation of fish fishing
make fit for, or change to suit a new purpose	I do something to go and match a certain purpose, make the necessary adjustments and adjustments, to respond to certain conditions and circumstances.
admit (to a wrongdoing)	I assure you, did someone present a thing to see, to know him or to judge him? I find something to be seen; show

TABLE 4.2: English Princeton WordNet definitions and the target wordnet definitions we want to match

Language	MRR	Top 10	Top 10 %	Top 1	Top 1 %
bg	0.087	674	0.337	403	0.202
el	0.122	1006	0.503	709	0.355
it	0.016	476	0.238	250	0.125
ro	0.051	988	0.494	728	0.364
sl	0.073	584	0.292	317	0.159
sq	0.056	979	0.490	767	0.384

TABLE 4.3: Experiment results for monolingual retrieval

Distances”, the elements of the term-document matrix C are given as normalized bag-of-words representations. Given the i^{th} word in a document;

$$d_i = \frac{tf}{\sum_{j=1}^n tf_j} \quad (4.5)$$

Where $tf_{i,j}$ is the term frequency or the number of times term i appears in document j .

We have studied the *tf-idf* weighted case for using term-document matrices for document retrieval tasks. Kusner *et al.* gives a brilliant example on a case where word overlap scoring nature of term-document matrices fall short;

Recall the earlier example of two similar, but word-different sentences in one document: “Obama speaks to the media in Illinois” and in another: “The President greets the press in Chicago”. After stop-word removal, the two corresponding nBOW vectors d and d' have no common non-zero dimensions and therefore have close to maximum simplex distance, although their true distance is small.

In order to use the innate similarity between words provided by word embeddings, Kusner *et al.* first introduces a flow matrix $T \in R^{n \times n}$.

4.3 Balikas *et al.* Implementation

4.3.1 Optimal Transport

4.3.2 Sinkhorn

Chapter 5

Supervised Validation

The approaches we have presented so far work fully unsupervised. Given two collections of definitions, we have tried to retrieve the definition(s) that represented the same meaning. However, given the moderately sized data in our hands we have accepted as “golden”, we studied the feasibility of training an encoder [83], where the objective is to learn whether the pair of definitions entail the same sense across languages.

Recurrent Neural Network (RNN) architectures improve upon the prototypical neural network model by introducing a *memory* for the connections in the network [84]. By updating the hidden unit over time using the output of the previous time-step, the model can *remember* features of the input signal for later inputs. One particular architecture of RNNs proposed by Hochreiter & Schmidhuber [85] is *long short-term memory* (LSTM). LSTM models have been successful on language modelling tasks [83], handwriting recognition [86, 87] and machine translation with a focus on rare words [88] to name a few. Highlight of these results are that LSTM has an advantage on tasks that require contextual information to persist over long periods of time [89]. Furthermore, LSTMs do not require fixed input vectors, which is a necessity for us since our definition pairs do not have to be the same size.

5.0.1 Vanishing Gradient Problem

LSTM is borne out of the need to address the *vanishing gradient problem* [85, 90]. On the original publication by Hochreiter & Schmidhuber [85], a crucial shortcoming of RNNs have been identified as their slow rate of training which may not converge in the end at all. Independently, Bengio *et al.* [90] suggested that the problem stems from the choice between the conservation of the previous inputs versus resisting against the

noise they accumulate. Figure 5.1, adapted from Graves [89] illustrates the problem using shades as the influence of input over neural network units.

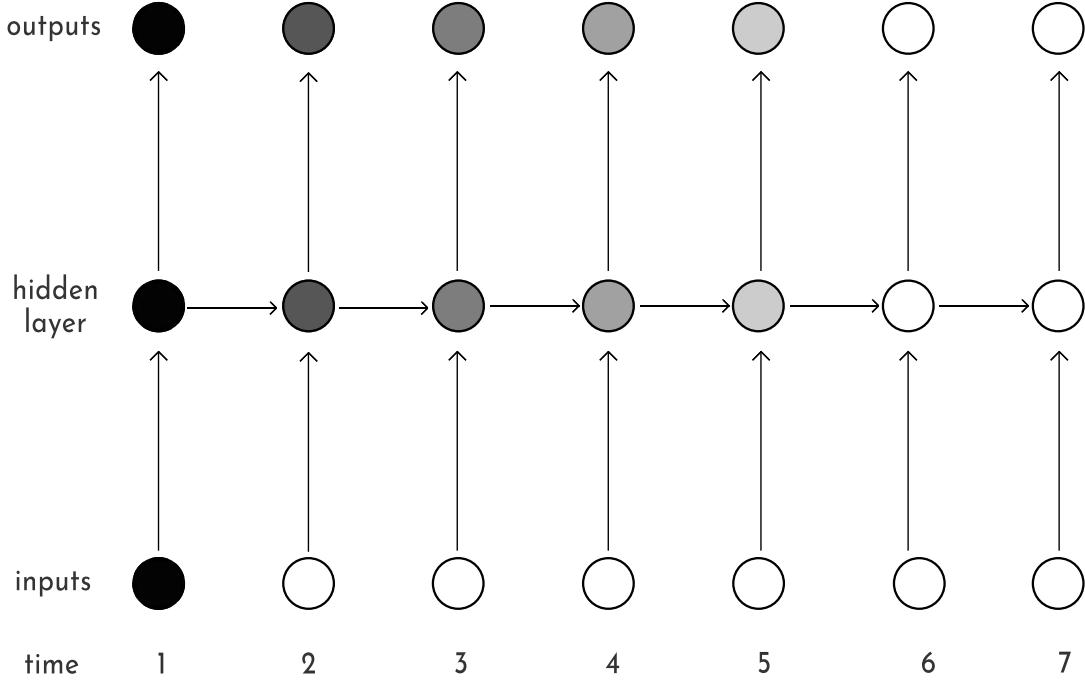


FIGURE 5.1: Graphical representation of vanishing gradient problem where the shades of the nodes represent the influence of the input signal, adapted from Figure 4.1 of Graves [89]

As the input signal traverses the units of an RNN, it either diminishes or blows up [89]. A long line of solutions were proposed such as decomposing structures so that only unexpected ones can be relevant. LSTM is the solution highlighted by Graves [89] as a recurrent neural network model that can work over temporally distant input signals while preserving their influence or diminishing their noise. The centrepiece idea is to use a *constant error carousel*, special cells that enforce a constant error flow. In the original paper, this complex unit is named *memory cell* [85]. Using an *input gate*, the cell is updated if the current input is relevant and using an *output gate*, the unit will not update other cells if the current input is not relevant. A simplified overview of the suggested model is presented in Figure 5.2.

Multiple arrows denote input from current time frame and recurrent connections. \odot symbol denotes multiplication. By weighing the input and output gates between 0 and 1, the impact of the current input can be adjusted. Overall, input gate controls how

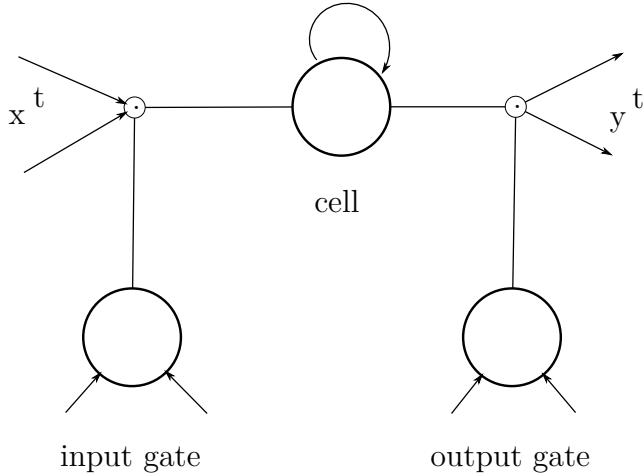


FIGURE 5.2: Simplified long short-term memory cell architecture

much cell will learn and output gate controls how much the cell will propagate. Figure 5.3 adapted from Graves [89] illustrates how the LSTMs operate against vanishing gradient.

Two gates on top of the cell structure model got extended with a third *forget gate* by Gers *et al.* [91]. The aim was to handle input sequences that are not segmented in a predictable manner. The proposed forget gate is implemented to *reset* the cell. When a cell state got irrelevant due to a change in problem domain, forget gate gradually resets the cell state instead of erroneous activations from the input gate.

Another extension came in the form of *peephole connections* by Gers *et al.* [92]. By allowing internal gates to inspect the cell state, they have shown improvements on non-linear tasks.

The finalized model with input, forget and output gates as well as the internal peephole connections was debuted in Graves & Schmidhuber [93]. In their expansive study comparing 8 LSTM variants over 15 years of CPU time, Greff *et al.* [94] named this model the “vanilla LSTM”.

Recently, Mueller & Thyagarajan [95] proposed a *siamese* model using two LSTMs.

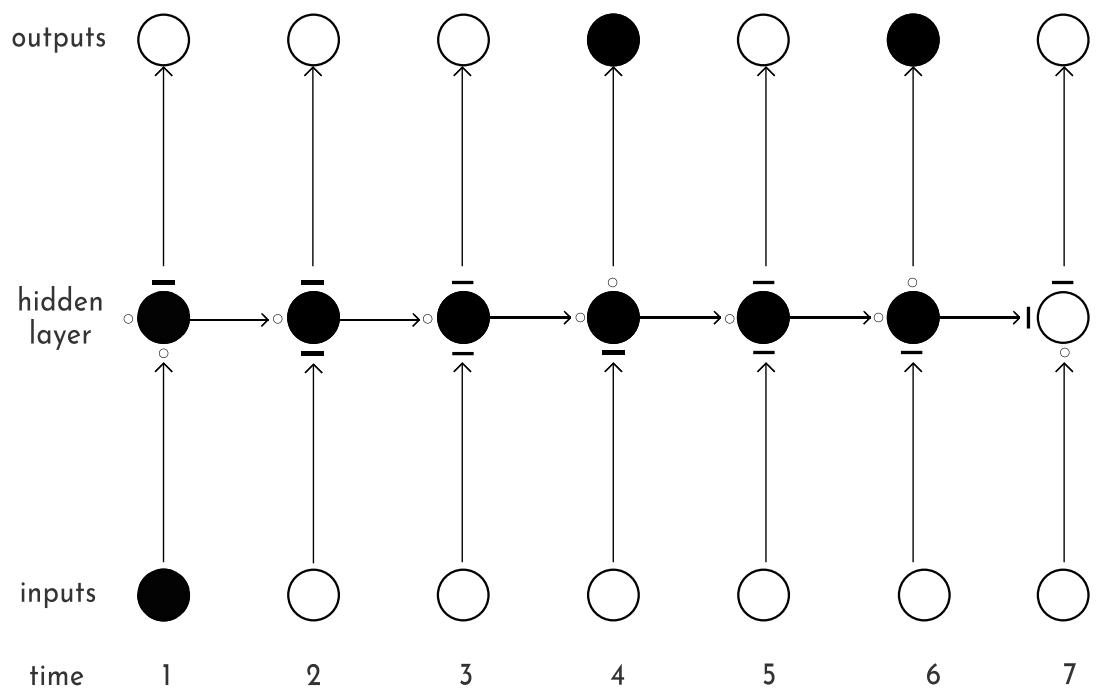


FIGURE 5.3: Preserving the input signal through blocking (-) or allowing (O) the input signal, adapted from Figure 4.4 of Graves [89]

Chapter 6

Experiments and Evaluation

We experimented with in domain *fasttext* embeddings. We have trained Romanian and Bulgarian embeddings and ran the WMD and Sinkhorn experiments. The performance dropped to the quarter of pre-trained embeddings so we have not repeated the experiments for other language corpora.

We have evaluated our embeddings using the standard bilingual lexicon extraction as a general measure for their performance. Although Ruder *et al.* [49] and Glavas *et al.* [96] both say you should evaluate them using downstream tasks.

Language	FastText 1M	FastText 500k	numberbatch
bg	33.61	35.17	51.97
el	37.37	39.58	30.35
it	58.20	59.28	50.37
ro	37.33	38.71	64.17
sl	21.42	22.91	74.74
sq	24.46	25.36	58.63

TABLE 6.1: Accuracy Scores of the Vectors Aligned Using VecMap

Language	FastText 1M	FastText 500k	numberbatch
bg	96.43	93.36	17.53
el	94.44	90.28	12.15
it	97.93	95.97	41.08
ro	97.06	94.91	16.4
sl	94.67	90.73	9.23
sq	83.59	80.92	9.51

TABLE 6.2: Coverage Scores for the Vectors Aligned Using VecMap

Chapter 7

Conclusion

Italian is a core language in numberbatch, the performance skyrockets!

Bibliography

1. *A Practical Guide to Lexicography* (ed van Sterkenburg, P.) *Terminology and Lexicography Research and Practice* 6. OCLC: 249659375 (Benjamins, Amsterdam, 2003). 459 pp. ISBN: 978-90-272-2329-6 978-90-272-2330-2 978-1-58811-380-1 978-1-58811-381-8.
2. Uzun, E. N. Modern Dilbilim Bulguları Işığında Türkçe Sözlüğe Bir Bakış. *Çukurova Üniversitesi Türkoloji Araştırmaları Merkezi* (2005).
3. İbrahim USTA, H. Türkçe Sözlük Hazırlamada Yöntem Sorunları. *Ankara Üniversitesi Dil ve Tarih-Coğrafya Fakültesi Dergisi*, 223–242 (Jan. 1, 2006).
4. Uzun, L. 1945’TEN BU YANA TÜRKÇE SÖZLÜKLER. *KEBİKEÇ İnsan Bilimleri İçin Kaynak Araştırmaları Dergisi*, 53–57. ISSN: 1300-2864 (1999).
5. Kendall, J. *The Forgotten Founding Father: Noah Webster’s Obsession and the Creation of an American Culture* 1st Edition. 368 pp. ISBN: 0-399-15699-2 (G.P. Putnam’s Sons, Apr. 14, 2011).
6. Fellbaum, C. *WordNet : An Electronic Lexical Database* ISBN: 978-0-262-27255-1 (MIT Press, 1998).
7. Miller, G. A. Nouns in WordNet: A Lexical Inheritance System. *Int J Lexicography* 3, 245–264. ISSN: 0950-3846. <https://academic.oup.com/ijl/article/3/4/245/923281> (2019) (Dec. 1, 1990).
8. Winston, M. E., Chaffin, R. & Herrmann, D. A Taxonomy of Part-Whole Relations. *Cognitive Science* 11, 417–444. ISSN: 1551-6709. https://onlinelibrary.wiley.com/doi/abs/10.1207/s15516709cog1104_2 (2019) (1987).
9. Sagot, B. & Fišer, D. Building a Free French Wordnet from Multilingual Resources (May 31, 2008).
10. Fellbaum, C. in *EuroWordNet: A Multilingual Database with Lexical Semantic Networks* (ed Vossen, P.) 137–148 (Springer Netherlands, Dordrecht, 1998). ISBN: 978-94-017-1491-4. https://doi.org/10.1007/978-94-017-1491-4_6 (2019).
11. Bond, F. & Paik, K. A Survey of WordNets and Their Licenses (Jan. 1, 2012).

12. Ruci, E. *On the Current State of Albanet and Related Applications* (Technical report, University of Vlora.(<http://fjalnet.com> ..., 2008).
13. Simov, K. I. & Osenova, P. *Constructing of an Ontology-Based Lexicon for Bulgarian.* in *LREC* (Citeseer, 2010).
14. Stamou, S., Nenadic, G. & Christodoulakis, D. *Exploring Balkanet Shared Ontology for Multilingual Conceptual Indexing.* in *LREC* (2004).
15. Pianta, E., Bentivogli, L. & Girardi, C. MultiWordNet: Developing an Aligned Multilingual Database (Jan. 1, 2002).
16. Fišer, D., Novak, J. & Erjavec, T. *sloWNet 3.0: Development, Extension and Cleaning in Proceedings of 6th International Global Wordnet Conference (GWC 2012)* (2012), 113–117.
17. Tufiş, D., Ion, R., Bozianu, L., Ceauşu, A. & řtefanescu, D. *Romanian Wordnet: Current State, New Applications and Prospects* in *Proceedings of 4th Global WordNet Conference, GWC* (2008), 441–452.
18. Somers, J. *You're Probably Using the Wrong Dictionary* <http://jsomers.net/blog/dictionary> (2019).
19. Almeida, F. & Xexéo, G. Word Embeddings: A Survey. arXiv: [1901.09069](https://arxiv.org/abs/1901.09069) [cs, stat]. [http://arxiv.org/abs/1901.09069](https://arxiv.org/abs/1901.09069) (2019) (Jan. 25, 2019).
20. Collobert, R. & Weston, J. *A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning* in. Proceedings of the 25th International Conference on Machine Learning (ACM, May 7, 2008), 160–167. ISBN: 978-1-60558-205-4. [http://dl.acm.org/citation.cfm?id=1390156.1390177](https://dl.acm.org/citation.cfm?id=1390156.1390177) (2019).
21. Harris, Z. S. Distributional Structure. *WORD* **10**, 146–162. ISSN: 0043-7956. <https://doi.org/10.1080/00437956.1954.11659520> (2019) (Aug. 1, 1954).
22. Firth, J. R. A Synopsis of Linguistic Theory 1930–1955. *Studies in linguistic analysis Special volume of the Philological Society*, 11 (1957).
23. Osgood, C. E., Suci, G. J. & Tannenbaum, P. H. *The Measurement of Meaning* 358 pp. ISBN: 978-0-252-74539-3 (University of Illinois Press, 1957).
24. Lund, K. & Burgess, C. Producing High-Dimensional Semantic Spaces from Lexical Co-Occurrence. *Behavior Research Methods, Instruments, & Computers* **28**, 203–208. ISSN: 1532-5970. <https://doi.org/10.3758/BF03204766> (2019) (June 1, 1996).

25. Salton, G., Wong, A. & Yang, C. S. A Vector Space Model for Automatic Indexing. *Commun. ACM* **18**, 613–620. ISSN: 0001-0782. <http://doi.acm.org/10.1145/361219.361220> (2019) (Nov. 1975).
26. Turney, P. D. & Pantel, P. From Frequency to Meaning: Vector Space Models of Semantics. *Journal of Artificial Intelligence Research* **37**, 141–188. ISSN: 1076-9757. arXiv: [1003.1141](https://arxiv.org/abs/1003.1141). <http://arxiv.org/abs/1003.1141> (2018) (Feb. 27, 2010).
27. Jones, K. S. A Statistical Interpretation of Term Specificity and Its Application in Retrieval. *Journal of Documentation* **28**, 11–21 (1972).
28. Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K. & Harshman, R. Indexing by Latent Semantic Analysis. *Journal of the American society for information science* **41**, 391–407 (1990).
29. Levy, O., Goldberg, Y. & Dagan, I. Improving Distributional Similarity with Lessons Learned from Word Embeddings. *Transactions of the Association for Computational Linguistics* **3**, 211–225 (Dec. 1, 2015).
30. Church, K. W. & Hanks, P. Word Association Norms, Mutual Information, and Lexicography. *Comput. Linguist.* **16**, 22–29. ISSN: 0891-2017. <http://dl.acm.org/citation.cfm?id=89086.89095> (2019) (Mar. 1990).
31. Forsythe, G. E., Malcolm, M. A. & Moler, C. B. *Computer Methods for Mathematical Computations* (Prentice-hall Englewood Cliffs, NJ, 1977).
32. Landauer, T. K. & Dumais, S. T. A Solution to Plato’s Problem: The Latent Semantic Analysis Theory of Acquisition, Induction, and Representation of Knowledge. *Psychological review* **104**, 211 (1997).
33. Schütze, H. *Dimensions of Meaning in Proceedings of the 1992 ACM/IEEE Conference on Supercomputing* Minneapolis, Minnesota, USA (IEEE Computer Society Press, Los Alamitos, CA, USA, 1992), 787–796. ISBN: 978-0-8186-2630-2. <http://dl.acm.org/citation.cfm?id=147877.148132> (2019).
34. Bengio, Y., Ducharme, R., Vincent, P. & Jauvin, C. A Neural Probabilistic Language Model. *Journal of machine learning research* **3**, 1137–1155 (Feb 2003).
35. Xu, W. & Rudnicky, A. *Can Artificial Neural Networks Learn Language Models?* in *Sixth International Conference on Spoken Language Processing* (2000).
36. Turian, J., Ratinov, L.-A. & Bengio, Y. *Word Representations: A Simple and General Method for Semi-Supervised Learning.* in. Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics. **2010** (Jan. 1, 2010), 384–394.

37. Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient Estimation of Word Representations in Vector Space. arXiv: [1301.3781 \[cs\]](https://arxiv.org/abs/1301.3781). <http://arxiv.org/abs/1301.3781> (Jan. 16, 2013).
38. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed Representations of Words and Phrases and Their Compositionality. arXiv: [1310.4546 \[cs, stat\]](https://arxiv.org/abs/1310.4546). <http://arxiv.org/abs/1310.4546> (2019) (Oct. 16, 2013).
39. Mikolov, T., Yih, W.-t. & Zweig, G. *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* (2013), 746–751.
40. Bullinaria, J. A. & Levy, J. P. Extracting Semantic Representations from Word Co-Occurrence Statistics: A Computational Study. *Behavior Research Methods*, 510–526 (2007).
41. Pennington, J., Socher, R. & Manning, C. *Glove: Global Vectors for Word Representation* in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2014), 1532–1543.
42. Baroni, M., Dinu, G. & Kruszewski, G. *Don't Count, Predict! A Systematic Comparison of Context-Counting vs. Context-Predicting Semantic Vectors* in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* Baltimore, Maryland (Association for Computational Linguistics, June 2014), 238–247. <https://www.aclweb.org/anthology/P14-1023>.
43. Mikolov, T., Grave, E., Bojanowski, P., Puhrsch, C. & Joulin, A. *Advances in Pre-Training Distributed Word Representations* in *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)* (2018).
44. Mnih, A. & Kavukcuoglu, K. in *Advances in Neural Information Processing Systems 26* (eds Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z. & Weinberger, K. Q.) 2265–2273 (Curran Associates, Inc., 2013). <http://papers.nips.cc/paper/5165-learning-word-embeddings-efficiently-with-noise-contrastive-estimation.pdf>.
45. Bojanowski, P., Grave, E., Joulin, A. & Mikolov, T. Enriching Word Vectors with Subword Information. <https://arxiv.org/abs/1607.04606> (2018) (July 15, 2016).

46. Grave, E., Bojanowski, P., Gupta, P., Joulin, A. & Mikolov, T. Learning Word Vectors for 157 Languages. arXiv: [1802.06893 \[cs\]](https://arxiv.org/abs/1802.06893). <http://arxiv.org/abs/1802.06893> (2019) (Feb. 19, 2018).
47. Speer, R., Chin, J. & Havasi, C. *ConceptNet 5.5: An Open Multilingual Graph of General Knowledge* in. AAAI Conference on Artificial Intelligence (2017), 4444–4451. <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14972>.
48. Faruqui, M. & Dyer, C. *Improving Vector Space Word Representations Using Multilingual Correlation* in *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics* (Association for Computational Linguistics, Gothenburg, Sweden, 2014), 462–471. <http://aclweb.org/anthology/E14-1049> (2018).
49. Ruder, S., Vulić, I. & Søgaard, A. A Survey Of Cross-Lingual Word Embedding Models. arXiv: [1706.04902 \[cs\]](https://arxiv.org/abs/1706.04902). <http://arxiv.org/abs/1706.04902> (2018) (June 15, 2017).
50. Neale, S. *A Survey on Automatically-Constructed WordNets and Their Evaluation: Lexical and Word Embedding-Based Approaches* in *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)* Miyazaki, Japan (eds Calzolari, N. et al.) (European Language Resources Association (ELRA), May 7–12, 2018). ISBN: 979-10-95546-00-9.
51. Ercan, G. & Cicekli, I. Using Lexical Chains for Keyword Extraction. *Information Processing & Management* **43**, 1705–1714 (2007).
52. Vossen, P. Introduction to EuroWordNet. *Computers and the Humanities* **32**, 73–89. ISSN: 0010-4817. <https://www.jstor.org/stable/30200456> (2019) (1998).
53. Vossen, P. EUROWORDNET: A MULTILINGUAL DATABASE OF AUTONOMOUS AND LANGUAGE-SPECIFIC WORDNETS CONNECTED VIA AN INTER-LINGUALINDEX. *Int J Lexicography* **17**, 161–173. ISSN: 0950-3846. <https://academic.oup.com/ijl/article/17/2/161/969685> (2019) (June 1, 2004).
54. Gonzalo, J., Verdejo, F., Peters, C. & Calzolari, N. in *EuroWordNet: A Multilingual Database with Lexical Semantic Networks* (ed Vossen, P.) 113–135 (Springer Netherlands, Dordrecht, 1998). ISBN: 978-94-017-1491-4. https://doi.org/10.1007/978-94-017-1491-4_5 (2019).
55. Kitamura, K. Cultural Untranslatability. *Translation Journal* **13** (2009).
56. Knight, K. & Luk, S. K. *Building a Large-Scale Knowledge Base for Machine Translation* in *In Proceedings of AAAI* (1994).

57. Boyd-Graber, J., Fellbaum, C., Osherson, D. & Schapire, R. *Adding Dense, Weighted Connections to WordNet* in *Proceedings of the Third International WordNet Conference* (Citeseer, 2006), 29–36.
58. Diab, M. *The Feasibility of Bootstrapping an Arabic Wordnet Leveraging Parallel Corpora and an English Wordnet* in *Proceedings of the Arabic Language Technologies and Resources, NEMLAR, Cairo* (2004).
59. Fiser, D. *Leveraging Parallel Corpora and Existing Wordnets for Automatic Construction of the Slovene Wordnet* in (Aug. 25, 2009), 359–368.
60. Sand, H., Velldal, E. & Øvrelid, L. *Wordnet Extension via Word Embeddings: Experiments on the Norwegian Wordnet* in. Proceedings of the 21st Nordic Conference on Computational Linguistics (May 2017), 298–302. <https://aclweb.org/anthology/papers/W/W17/W17-0242/> (2019).
61. Khodak, M., Risteski, A., Fellbaum, C. & Arora, S. *Automated WordNet Construction Using Word Embeddings* in *Proceedings of the 1st Workshop on Sense, Concept and Entity Representations and Their Applications* (2017), 12–23.
62. Arora, S., Liang, Y. & Ma, T. A Simple but Tough-to-Beat Baseline for Sentence Embeddings. <https://openreview.net/forum?id=SyK00v5xx> (2019) (Nov. 4, 2016).
63. Arora, S., Li, Y., Liang, Y., Ma, T. & Risteski, A. Linear Algebraic Structure of Word Senses, with Applications to Polysemy. *Transactions of the Association for Computational Linguistics* **6**, 483–495. ISSN: 2307-387X. <https://transacl.org/ojs/index.php/tacl/article/view/1346> (2019) (July 20, 2018).
64. Gordeev, D., Rey, A. & Shagarov, D. *Unsupervised Cross-Lingual Matching of Product Classifications* in *Proceedings of the 23rd Conference of Open Innovations Association FRUCT Bologna, Italy* (FRUCT Oy, 2018), 62:459–62:464. <http://dl.acm.org/citation.cfm?id=3299905.3299967>.
65. Lesk, M. *Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone* in *Proceedings of the 5th Annual International Conference on Systems Documentation* (ACM, New York, NY, USA, 1986), 24–26. ISBN: 978-0-89791-224-2. <http://doi.acm.org/10.1145/318723.318728>.
66. Banerjee, S. & Pedersen, T. *An Adapted Lesk Algorithm for Word Sense Disambiguation Using WordNet* in *Computational Linguistics and Intelligent Text Processing* (ed Gelbukh, A.) (Springer Berlin Heidelberg, 2002), 136–145. ISBN: 978-3-540-45715-2.

67. Metzler, D., Dumais, S. & Meek, C. *Similarity Measures for Short Segments of Text* in *Advances in Information Retrieval* (eds Amati, G., Carpineto, C. & Romano, G.) (Springer Berlin Heidelberg, 2007), 16–27. ISBN: 978-3-540-71496-5.
68. Le, Q. V. & Mikolov, T. Distributed Representations of Sentences and Documents. arXiv: [1405.4053 \[cs\]](https://arxiv.org/abs/1405.4053). <http://arxiv.org/abs/1405.4053> (May 16, 2014).
69. Kiros, R. *et al.* Skip-Thought Vectors. arXiv: [1506.06726 \[cs\]](https://arxiv.org/abs/1506.06726). <http://arxiv.org/abs/1506.06726> (2019) (June 22, 2015).
70. Wieting, J., Bansal, M., Gimpel, K. & Livescu, K. Towards Universal Paraphrastic Sentence Embeddings. arXiv: [1511.08198 \[cs\]](https://arxiv.org/abs/1511.08198). <http://arxiv.org/abs/1511.08198> (2019) (Nov. 25, 2015).
71. Zhao, J., Lan, M. & Tian, J. *ECNU: Using Traditional Similarity Measurements and Word Embedding for Semantic Textual Similarity Estimation* in *SemEval@NAACL-HLT* (2015).
72. Šarić, F., Glavaš, G., Karan, M., Šnajder, J. & Bašić, B. D. *TakeLab: Systems for Measuring Semantic Text Similarity* in *Proceedings of the First Joint Conference on Lexical and Computational Semantics - Volume 1: Proceedings of the Main Conference and the Shared Task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation* Montréal, Canada (Association for Computational Linguistics, Stroudsburg, PA, USA, 2012), 441–448. <http://dl.acm.org/citation.cfm?id=2387636.2387708> (2019).
73. Edilson A. Corrêa, J., Marinho, V. & Borges dos Santos, L. NILC-USP at SemEval-2017 Task 4: A Multi-View Ensemble for Twitter Sentiment Analysis (Apr. 7, 2017).
74. Bush, V. As We May Think. *The atlantic monthly* **176**, 101–108 (1945).
75. Singhal, A. Modern Information Retrieval: A Brief Overview. *Bulletin of the Ieee Computer Society Technical Committee on Data Engineering* **24**, 2001 (2001).
76. Luhn, H. P. A Statistical Approach to Mechanized Encoding and Searching of Literary Information. *IBM Journal of research and development* **1**, 309–317 (1957).
77. Groves, M. & Mundt, K. Friend or Foe? Google Translate in Language for Academic Purposes. *English for Specific Purposes* **37**, 112–121. ISSN: 0889-4906. <http://www.sciencedirect.com/science/article/pii/S088949061400060X> (2019) (Jan. 1, 2015).
78. Manning, C. D., Raghavan, P. & Schütze, H. *Introduction to Information Retrieval* Reprinted. OCLC: 549201180. 482 pp. ISBN: 978-0-521-86571-5 (Cambridge Univ. Press, Cambridge, 2009).

79. Salton, G. The State of Retrieval System Evaluation. *Information processing & management* **28**, 441–449 (1992).
80. Voorhees, E. M. *The TREC-8 Question Answering Track Report* in *In Proceedings of TREC-8* (1999), 77–82.
81. Kusner, M. J., Sun, Y., Kolkin, N. I. & Weinberger, K. Q. *From Word Embeddings to Document Distances* in *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37* Lille, France (JMLR.org, 2015), 957–966. <http://dl.acm.org/citation.cfm?id=3045118.3045221> (2019).
82. Balikas, G., Laclau, C., Redko, I. & Amini, M.-R. Cross-Lingual Document Retrieval Using Regularized Wasserstein Distance. arXiv: [1805.04437 \[cs, stat\]](https://arxiv.org/abs/1805.04437). <http://arxiv.org/abs/1805.04437> (2019) (May 11, 2018).
83. Sutskever, I., Vinyals, O. & Le, Q. V. Sequence to Sequence Learning with Neural Networks. arXiv: [1409.3215 \[cs\]](https://arxiv.org/abs/1409.3215). <http://arxiv.org/abs/1409.3215> (2019) (Sept. 10, 2014).
84. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning Representations by Back-Propagating Errors. *Nature* **323**, 533. ISSN: 1476-4687. <https://www.nature.com/articles/323533a0> (2019) (Oct. 1986).
85. Hochreiter, S. & Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **9**, 1735–1780. ISSN: 0899-7667. [http://dx.doi.org/10.1162/neco.1997.9.8.1735](https://dx.doi.org/10.1162/neco.1997.9.8.1735) (2019) (Nov. 1997).
86. Graves, A., Liwicki, M., Bunke, H., Schmidhuber, J. & Fernández, S. in *Advances in Neural Information Processing Systems 20* (eds Platt, J. C., Koller, D., Singer, Y. & Roweis, S. T.) 577–584 (Curran Associates, Inc., 2008). <http://papers.nips.cc/paper/3213-unconstrained-on-line-handwriting-recognition-with-recurrent-neural-networks.pdf>.
87. Graves, A. *et al.* A Novel Connectionist System for Unconstrained Handwriting Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **31**, 855–868. ISSN: 0162-8828 (May 2009).
88. Luong, M.-T., Sutskever, I., Le, Q. V., Vinyals, O. & Zaremba, W. Addressing the Rare Word Problem in Neural Machine Translation. arXiv: [1410.8206 \[cs\]](https://arxiv.org/abs/1410.8206). <http://arxiv.org/abs/1410.8206> (2019) (Oct. 29, 2014).
89. Graves, A. in *Supervised Sequence Labelling with Recurrent Neural Networks* (ed Graves, A.) 37–45 (Springer Berlin Heidelberg, Berlin, Heidelberg, 2012). ISBN: 978-3-642-24797-2. https://doi.org/10.1007/978-3-642-24797-2_4 (2019).

90. Bengio, Y., Simard, P. & Frasconi, P. Learning Long-Term Dependencies with Gradient Descent Is Difficult. *IEEE Transactions on Neural Networks* **5**, 157–166. ISSN: 1045-9227 (Mar. 1994).
91. Gers, F. A., Schmidhuber, J. & Cummins, F. A. Learning to Forget: Continual Prediction with LSTM. *Neural Computation* **12**, 2451–2471 (2000).
92. Gers, F. A., Schraudolph, N. N. & Schmidhuber, J. Learning Precise Timing with Lstm Recurrent Networks. *J. Mach. Learn. Res.* **3**, 115–143. ISSN: 1532-4435. <https://doi.org/10.1162/153244303768966139> (2019) (Mar. 2003).
93. Graves, A. & Schmidhuber, J. Framewise Phoneme Classification with Bidirectional LSTM and Other Neural Network Architectures. *Neural Networks. IJCNN 2005* **18**, 602–610. ISSN: 0893-6080. <http://www.sciencedirect.com/science/article/pii/S0893608005001206> (2019) (July 1, 2005).
94. Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R. & Schmidhuber, J. LSTM: A Search Space Odyssey. *IEEE Trans. Neural Netw. Learning Syst.* **28**, 2222–2232. ISSN: 2162-237X, 2162-2388. arXiv: [1503.04069](https://arxiv.org/abs/1503.04069). <http://arxiv.org/abs/1503.04069> (2019) (Oct. 2017).
95. Mueller, J. & Thyagarajan, A. *Siamese Recurrent Architectures for Learning Sentence Similarity* in *AAAI* (2016).
96. Glavas, G., Litschko, R., Ruder, S. & Vulic, I. How to (Properly) Evaluate Cross-Lingual Word Embeddings: On Strong Baselines, Comparative Analyses, and Some Misconceptions. arXiv: [1902.00508 \[cs\]](https://arxiv.org/abs/1902.00508). <http://arxiv.org/abs/1902.00508> (2019) (Feb. 1, 2019).