

Domain-Specific Word Embeddings for Patent Classification

Julian Risch
Hasso Plattner Institute
University of Potsdam
Potsdam, Germany
julian.risch@hpi.de

Ralf Krestel
Hasso Plattner Institute
University of Potsdam
Potsdam, Germany
ralf.krestel@hpi.de

Abstract—Patent offices and other stakeholders in the patent domain need to classify patent applications according to a standardized classification scheme. To examine the novelty of an application it can then be compared to previously granted patents in the same class. Automatic classification would be highly beneficial, because of the large volume of patents and the domain-specific knowledge needed to accomplish this costly manual task. However, a challenge for the automation is patent-specific language use, such as special vocabulary and phrases. To account for this language use, we present domain-specific pre-trained word embeddings for the patent domain. We train our model on a very large dataset of more than 5 million patents and evaluate it at the task of patent classification. To this end, we propose a deep learning approach based on gated recurrent units for automatic patent classification built on the trained word embeddings. Experiments on a standardized evaluation dataset show that our approach increases average precision for patent classification by 17 percent compared to state-of-the-art approaches. In this paper, we further investigate the model’s strengths and weaknesses. An extensive error analysis reveals that the learned embeddings indeed mirror patent-specific language use. The imbalanced training data and underrepresented classes are the most difficult remaining challenge.

Index Terms—Document Classification, Deep Learning, Word Embedding, Patents

I. INTRODUCTION

In 2018, 308,853 U.S. patents have been granted by the U.S. Patent and Trademark Office, which is the second-largest number of grants ever¹. All granted U.S. patents since 1976 are publicly available as full text². These large text collections represent an extensive amount of human knowledge in an almost unstructured form. This makes mining information from them challenging and automatic classification and retrieval a hard problem.

Not only the number of documents but also the patent-specific vocabulary makes the tasks more difficult. Because of the underlying legal purpose of patent documents, they follow a specific writing style. Patent applications need to define the scope of an invention and need to delimit it from others whilst covering as much variation as possible. As a consequence, patent descriptions use vague language. For example, a patent calls an invention “electronic still camera” and “electronic

imaging apparatus”, whereas such a device is called “digital camera” in colloquial speech (Fig. 1). A patent’s claims are a controversial subject, because a patent grants rights and also limits the rights of others. Patents grant a monopoly for a limited time in exchange for the disclosure of the invention so that others can license it.

Unstructured text sections, such as abstracts, descriptions, and claims, make up the largest part of a patent. The claims section is essential for defining the scope of an invention. It describes the extent of the monopoly rights granted by the patent. Court decisions of the past precisely define the meaning of “patent speak”. An example are the slight differences of “consist of” and “comprise”³: “consist of” implies an exhaustive enumeration, whereas “comprise” commences an enumeration that is not necessarily exhaustive. Classifying patents is challenging because of patent-specific language use — even for domain experts.

Another difficulty are special technical terms and long lists of synonyms, such as light-sensitive, photosensitive, and photoreceptive. Depending on the context these synonyms actually might or might not have slightly different meaning. Patent applicants might come up with a new term to describe their invention to underline its novelty.

The International Patent Classification (IPC) is a hierarchical classification system for patents. It has been periodically revised and adapted to the upcoming of new fields of invention. The system considers 4 levels of hierarchy: sections, classes, subclasses, and group. For example, the U.S. patent no. 4131919 with the IPC code H04N 1/21 is in group H04N 1/21, which is in the subclass H04N, the class H04, and section H. The subparts of this code correspond to the section “electricity”, class “electric communication technique”, subclass “pictorial communication, e.g. television”, and group “Intermediate information storage”. An excerpt of this patent is depicted in Fig. 1 with the deprecated IPC code H04N 005/79.

This complicated classification system is applied at several different steps in the patenting process. On the one hand, patent applicants need to search for prior art, if they file a patent. They need to retrieve patents about similar inventions

¹<https://www.ificclaims.com/rankings-trends-2018.htm>

²<https://bulkdata.uspto.gov/>

³https://www.epo.org/law-practice/legal-texts/html/guidelines/e/f_iv_4_21.htm

United States Patent [19]		[11]	4,131,919
Lloyd et al.		[45]	Dec. 26, 1978
[54]	ELECTRONIC STILL CAMERA	[57]	ABSTRACT
[75]	Inventors: Gareth A. Lloyd; Steven J. Sasson, both of Rochester, N.Y.		Electronic imaging apparatus, preferably an electronic still camera, employs an inexpensive information-recording medium such as audio-grade magnetic tape for "capturing" scene images. The camera includes a charge coupled device comprised of an array of photo-sensitive elements which form a charge pattern corresponding to an optical image projected onto the elements during an exposure interval. A charge transfer circuit converts the charge pattern into a high frequency pulsed electrical signal immediately following the exposure interval to remove the charge from the device in a short period of time to maintain unwanted "dark current" at a low level. Each pulse represents the image-forming light projected onto a particular photo-sensitive element. A high speed analog-to-digital converter converts these pulses to multi-bit digital words in real time. A digital buffer memory temporarily stores these words, then retransmits them at a rate that is compatible for recording on the audio-grade tape. The image can be displayed on a conventional television receiver by reading the recorded words from the tape and converting them to a format compatible with the signal-receiving circuitry of the television.
[73]	Assignee: Eastman Kodak Company, Rochester, N.Y.		
[21]	Appl. No.: 798,956		
[22]	Filed: May 20, 1977		
[51]	Int. Cl. ²	H04N 5/79	
[52]	U.S. Cl.	360/9; 360/35; 358/127; 358/134; 358/213	
[58]	Field of Search	360/9, 10, 8, 35, 33; 179/2 TV; 358/127, 134, 213, 85, 133, 78	
[56]	References Cited		
	U.S. PATENT DOCUMENTS		
	3,858,232 12/1974 Boyle	357/24	
	3,911,467 10/1975 Levine	358/213	
	3,962,725 6/1976 Lemke	360/37	
	4,016,361 4/1977 Pandey	360/9	
	4,057,830 11/1977 Adcock	360/35	
	Primary Examiner—Bernard Konick Assistant Examiner—Alan Faber Attorney, Agent, or Firm—D. P. Monteith		
			8 Claims, 4 Drawing Figures

Fig. 1: Patent documents follow a standardized structure and consist of several fields, such as title, abstract, and claims, but also references. This example is an excerpt of U.S. patent no. 4131919.

although they might use different words for description. On the other hand, patent examiners in a patent office need to check a patent application for its “inventive step or non-obviousness” and its “novelty”. A patent examiner specialized in the field of the invention needs to be matched to the patent application. Finally, patent courts and patent attorneys deal with the infringement and validity of granted patents. All three scenarios involve an information retrieval task, where patents similar to a given patent need to be found. Based on their similarity, similar patents mutually limit their scopes.

The IPC systematically classifies patents into topical subclasses. Thereby the retrieval of similar patents can be performed by looking up patents in the same subclass. However, manually classifying patents into such subclasses is costly in terms of working power and needs domain-specific knowledge due to the complexity of the IPC. The goal of automated patent classification is to save these costs and associate a given patent document with its correct subclasses automatically. Smith summarizes the applications of automated patent classification as (1) matching patent applications with a patent examiner who is a domain expert for the field of invention, (2) classification of external documents so that they can easily be retrieved during the patent examination process, and reclassification of older patents labeled with outdated classification schemes [1]. In practice, patents can be associated with multiple subclasses. Therefore, patent classification is not a multi-class but a multi-label classification task. In fact, our example patent in Fig. 1 is associated with 2 IPC subclasses. In total, the IPC knows 637 subclasses.

In this paper, we propose to improve automatic patent classification by leveraging recent deep learning techniques. In particular, we train fastText word embeddings on a large dataset of more than 5 million patents. We use these embeddings together with bi-directional Gated Recurrent Units (GRUs) to

classify patents. Experiments show that our approach is superior to state-of-the-art approaches in terms of three evaluation measures. For example, we increase micro-average precision at predicting a patent’s subclass by 17 percent. Further, we find that domain-specific word embeddings trained on patent documents outperform standard word embeddings trained on Wikipedia pages by 9 percent when combined with a GRU-based neural network.

Our contributions from previous work are the computation of word embeddings on the second largest corpus ever used for training and providing these word embeddings for download⁴. Further, we propose a deep neural network architecture based on bi-directional Gated Recurrent Units (GRUs) for patent classification. In this paper, we further investigate the model’s strengths and weaknesses. An extensive error analysis reveals that the learned embeddings indeed mirror patent-specific language use. The imbalanced training data and underrepresented classes are the most difficult remaining challenge.

Section II summarizes related work in the field of automatic patent classification and gives an overview of different word embedding approaches. The three datasets used in this paper are described in Section III and Section IV describes our approach to capture semantics in patent language by domain-specific word embeddings and automatically classify patents. We evaluate our approach with three experiments in Section V and conduct an error analysis in Section VI. We conclude in Section VII.

II. RELATED WORK

Fall et al. established a collection of around 75,000 excerpts of English-language patent applications as a de-facto standard dataset for the evaluation of automatic patent classification [2]. The dataset is called WIPO-alpha⁵ and is provided by the World Intellectual Property Office (WIPO). Fall et al. further propose three evaluation measures that are tailored to the patent classification task, where a patent is typically associated with a main subclass, but also with several incidental subclasses. We apply the three measures in our experiments and describe them in detail in Section V. In general, the micro-precision of assigning the correct class to a given patent is evaluated.

Table I gives an overview of related work in the field of patent classification.

Seneviratne et al. propose to generate signatures from patents instead of using the full vocabulary as features [3]. They evaluate their patent classification approach on IPC class level (114 classes) and subclass level (451 subclasses) on the WIPO-alpha dataset. While they improve classification performance in comparison to Fall et al., they optimize also the time required to index and search a patent collection. Other results on the WIPO-alpha dataset have been published by Nguyen (macro-f1: 0.452, micro-f1: 0.755) [4], Rousu et al. (micro-f1: 0.767) [5], and Qiu et al. (macro-f1: 0.418) [6].

⁴<https://hpi.de/naumann/projects/repeatability/text%2Dmining.html>

⁵WIPO-en-alpha dataset, World Intellectual Property Office, Geneva, Switzerland, 2002

TABLE I: Related Work Approaches use the Datasets WIPO-alpha, USPTO-2M, and CLEF-IP for Evaluation.

Approach	WIPO	USPTO	CLEF
Fall et al. [2]	x	-	-
Seneviratne et al. [3]	x	-	-
Nguyen [4]	x	-	-
Rousu et al. [5]	x	-	-
Qiu et al. [6]	x	-	-
Derieux et al. [7]	-	-	x
Verberne and Dhondt [8]	-	-	x
Li et al. [9]	-	x	x
Risch and Krestel [10]	x	x	-

Several researchers conducted their experiments on other datasets, which makes a direct comparison with their results impossible. An ensemble of different classifiers slightly improves micro-F1 score on a refined version of the WIPO-alpha dataset according to Mathiassen and Ortiz-Arroyo [11]. They report a micro-f1 of 0.867. Instead of IPC, Tran and Kavuluru use the Cooperative Patent Classification (CPC) system, which replaces the earlier the U.S. Patent Classification (USPC) system [12]. They report a micro-f1 of 0.700 on a dataset of patents with 654 subclasses. Dhondt et al. report a micro-f1 0.751 and a micro-precision of 0.800 on a subset of 532,264 English abstracts from the so called CLEF-IP 2010 corpus [13].

The CLEF-IP 2010 corpus from the Conference and Labs of the Evaluation Forum’s track for retrieval experiments in the intellectual property domain (CLEF-IP) considered two tasks: (1) recommending patents as prior art for another patent and (2) patent classification according to the International Patent Classification system (IPC). As its predecessor, the CLEF-IP track of 2011 [14] provided datasets and tasks for a large number of publications concerning retrieval in the intellectual property domain.

D’hondt et al. find that bigrams are important phrasal features to capture multi-word terms, which are frequent in patents [15]. However, their work considers only the class level (120) classes and not the more diverse and thus more difficult subclass level. Verberne and D’Hondt investigate the usefulness of different text sections of patents, such as title, abstract, claims, and description in context of CLEF-IP 2010 patent retrieval and re-ranking tasks [16]. Abstract and description achieve best precision and recall at the retrieval task and significantly outperform title and claims [8], [16]. With regard to the usefulness of metadata, such as applicants, inventors, and address, they conclude that it does not improve classification [8]. This result contradicts Beney, who finds that applicant and address improves classification [17]. They argue that names and addresses identify companies, which work in restricted domains. Derieux et al., find that results are language specific, classification on English patents is at least 10% better than on German and French patents [7]. Guyot et al. aim at building a single patent classifier for German, French, and English patents [18]. They avoid sophisticated preprocessing steps to be as language-independent as possible.

The observation that language has a strong influence on the classification motivates further investigation of patent-specific language use. In this paper, we consider to model this language use with patent-specific word embeddings. To this end, we summarize work in the field of word embeddings.

The upcoming of word embeddings or, more general speaking, dense vector representations to capture the semantic meaning of words influences many natural language processing tasks. With Word2Vec, Mikolov et al. propose an efficient way to train word embeddings [19]. As a consequence, they are able to train embeddings on large datasets with billions of words. A similar approach, termed global vectors (GloVe), trains word embeddings on global word-word co-occurrence counts rather than on context windows of limited size [20]. A disadvantage of both Word2Vec and GloVe is the inherent out-of-vocabulary problem: a word that occurs only in the test data but not in the training data has no vector representation in the word embedding space. To overcome this problem, Bojanowski et al. introduce another context-window-based approach, which they call fastText [21]. fastText word embeddings incorporate information about character n-grams as subparts of a word. As a consequence, they overcome the out-of-vocabulary problem of other word embedding approaches by falling back to embeddings of character n-grams if a word is unknown.

Recently, deep learning approaches for patent classification have been proposed. Xia et al. outline a general deep learning approach for patent classification based on sparse auto-encoders and deep belief networks [22]. However, their proposal is limited to a theoretical approach and lacks practical experiments. Grawe et al. automatically classify patents based on word embeddings and long-short term memory units (LSTMs) in a neural network [23]. Their approach is similar to ours but has several limitations: (1) it considers only 50 different classes, (2) it achieves only 63% accuracy, and (3) as opposed to our approach it suffers from out-of-vocabulary problems, which is inherent to the applied Word2Vec model. Li et al. proposed a convolutional neural network for patent classification [9]. The first 100 words of each patent’s title and abstract are represented with 200-dimensional word embeddings. Convolutions with filter size 3, 4, and 5 are used in combination with a max-pooling layer. Their approach achieves only slightly better precision and recall as a random forest baseline. The applied word embedding technique suffers from out-of-vocabulary problems, which we overcome with our approach. However, the comparison of convolutional neural networks and recurrent neural networks at the task of patent classification remains an open task for future research.

Instead of content-based approaches, which consider only a patent’s text sections, Li et al. propose a citation-based approach [24]. They exploit co-citation relations among patents. Further, they leverage the fact that patents reference other patents in the same field to explain the novelty of their ideas. These references are not limited to patents, but also include scientific papers. Cross-collection topic models can be used to recommend references across these different document collec-

tions [25]. While these approaches can help to retrieve similar patents and can therefore be of help in the patenting process, we solely focus on content-based classification of patents in this paper. Similar to the IPC system in the patent domain is the Medical Subject Headings (MeSH) ontology in the medical domain. Eisinger et al. compare automatic document classification for the two classification schemes [26]. They leverage class co-occurrence frequencies to enrich labeled classes and propose a guided search as an application. In contrast to most other work, Chen et al. do not classify on the level of subclasses but on the level of subgroups [27]. Subgroups are the lowest level in the patent classification hierarchy. While there are 648 subclasses, there are approximately 72,000 subgroups. Benzineb and Guyot describe current challenges of automated patent classification and also its historical background [28]. Two examples are classification consistency issues arising because of rapidly growing number of human examiners or the so-called “horizontal nature” of patents, which refers to multiple correct labels per patent.

In our previous work, we proposed domain-specific word embeddings for the patent domain and showed that these embeddings outperform common word embeddings trained on Wikipedia [10]. In this work, we conduct more experiments with regard to two different aspects. First, we investigate the process to obtain domain-specific word embeddings for the patent domain. With a new experiment we investigate on the influence of the number of training samples on the model’s performance. Furthermore, we conduct an error analysis, where we investigate at which level of the classification hierarchy mis-classification happens.

III. DATASET

In this paper, we consider three different datasets of patent documents. Tab. II gives an overview of the datasets, their number of documents, and their number of tokens. The first dataset is the WIPO-alpha dataset established by Fall et al., which is a de-facto standard for the evaluation of automated patent classification and has been widely used [2]–[6]. The dataset contains more than 75,000 patents with title, abstract, claims, and full description. Further, each patent is associated with a main subclass and incidental subclasses.

The second dataset is much larger and contains 5.4 million patents granted by the United States Patent and Trademark Office (USPTO). The USPTO keeps records of all U.S. patent activity since 1790. On their website⁶, they provide free bulk downloads of full text patent publications from 1976 to 2016. We use this full dataset and refer to it as USPTO-5M.

Each patent contains bibliographic data, such as title, inventor, owner, filing date, and granting date. Furthermore, author information, patent type classification, claims, abstract, links to other patents or papers, and a detailed description of the invention are provided. For our experiments, we focus on textual data and leave out figures and their captions. In

⁶<https://www.uspto.gov/learning%2Dand%2Dresources/electronic%2Dbulk%2Ddata%2Dproducts>

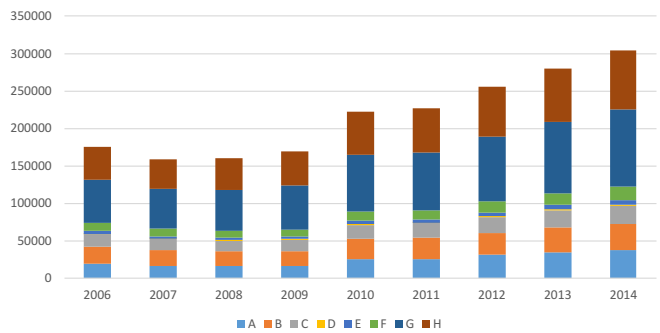


Fig. 2: Number of granted patents in the USPTO-2M dataset per year and per section of the IPC system.

TABLE II: A Comparison of the Three Patent Datasets

Dataset	# Documents	# Tokens
WIPO-alpha	75,250	561 million
USPTO-2M	2 million	235 million
USPTO-5M	5 million	38 billion

comparison to WIPO-alpha, USPTO-5M is 70 times larger in terms of number of documents and also number of tokens.

The third dataset is called USPTO-2M, contains 2 million patents, and goes back to Li et al. [9]. It is publicly available online⁷ in a pre-processed JSON format so that other researchers can use it easily. The dataset is split into a training set with 1.95 million documents and a test set with the remaining 50,000 documents. Further, it is limited to titles, abstracts, document identifiers, and subclasses. In total, there are 637 subclasses. In contrast to WIPO-alpha, USPTO-2M does not distinguish between main subclass and incidental subclasses. Figure 2 shows the number of granted patents in the USPTO-2M dataset per year and per section of the IPC system. For all classes this number is almost constant until 2009 and linearly increases year by year starting from 2010.

IV. DEEP LEARNING FOR PATENT CLASSIFICATION

Our goal is to automatically classify patents into their assigned subclasses. The large amount of available patents and their full text plus the recent success of deep learning for natural language processing motivate to investigate deep learning for patent classification. To this end, we propose to use word embeddings to capture the semantics of the specific language that is used in patents. Further, we propose a neural network architecture to automatically classify patents based on the inferred word embeddings.

A. Domain-Specific Word Embeddings

Word embeddings are a basic ingredient for a variety of tasks in natural language processing. They represent words as dense vectors in a vector space. Pre-trained on a large number of tokens, relations of these representations in a vector space can mirror semantic relations of words [19].

⁷<http://mleg.cse.sc.edu/DeepPatent/>

We propose to train fastText word embeddings based on the method by Bojanowski et al. [21] with 100, 200, and 300 dimensions. We transform all characters to lowercase and discard all words that occur less than ten times. The used context window size is 5.

We train the embeddings on our dataset USPTO-5M, which contains 38 billion tokens and publish the resulting word embeddings online⁸. To the best of our knowledge, this is the second largest number of tokens ever used to train word embeddings. It contains more than twice the number of tokens of the English Wikipedia (16 billion) and is only exceeded by the Common Crawl dataset, which consists of 600 billion tokens. We assume that the embeddings are helpful not only for patent classification but also for other tasks in the patent domain and hope that other researchers can build on our results.

B. Neural Network Architecture

Given a patent document, our goal is to infer its main subclass and also potential incidental subclasses. We investigate how domain-specific word embeddings can help to solve this classification problem. Therefore, we extract a patent document’s title and abstract and consider only the sequence of the first 300 words. We choose this limitation to be comparable to related work in our evaluation [2], [3]. Longer sequences linearly increase runtime and memory need.

Fig. 3 visualizes the network architecture. For each word in the input sequence, we calculate its word embedding based on our pre-trained, domain-specific fastText model. This sequence of word embeddings is processed by a spatial dropout, which randomly masks 10% of the input words to make the neural network more robust. The remaining 90% of the sequence serve as input to the next layer in the neural network. In particular, we propose a deep neural network architecture based on gated recurrent units (GRUs). We decided to use GRUs instead of long short-term memory units (LSTMs), because they have a smaller number of trainable parameters and are thus less likely to overfit on the training data. GRUs and LSTMs are superior to simple recurrent units, because they leverage gates to overcome the vanishing gradient problem. We use bi-directional GRUs so that the input sequence is processed in two directions: correct order and reverse order of the words. The outputs of these two processing steps are averaged and followed by a dropout of 10%, again to make the network more robust.

In Section I, we pointed out that patent classification is not a multi-class but a multi-label classification task. A typical final layer of our neural network would therefore be a dense layer with as many units as subclasses and a sigmoid activation. Instead, we use a dense layer with as many units as subclasses and a softmax activation. Thereby, we train the model for the multi-class classification task only and aim to predict a patent’s main subclass. For training the neural network, the softmax activation together with a categorical loss function

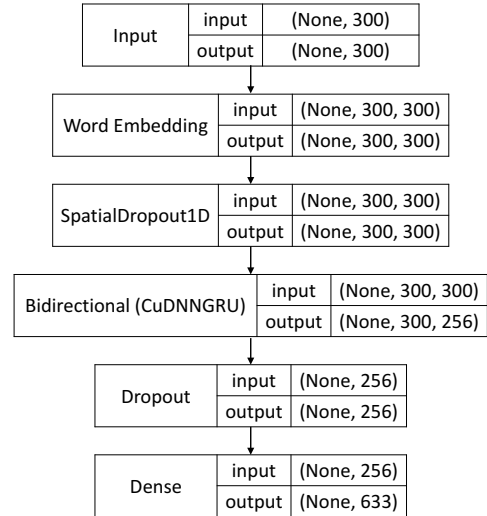


Fig. 3: The neural network uses pre-trained word embeddings, spatial dropout, GRUs, dropout, and a dense layer with softmax activation.

considers only a single subclass as correct. If our model predicts any other subclass, such as any incidental subclasses, the prediction is considered wrong during training.

However, during testing, we consider the probabilities output by the softmax activation for all subclasses. We consider the top three subclasses with the highest probabilities as our final prediction. Although the neural network is trained to predict only the main subclass and not the incidental subclasses, our experiments in Section V show that the model achieves competitive results for both tasks.

Training of the neural network until convergence takes 13 epochs with a batch size of 256. With a larger batch size, more subclasses are covered in a particular epoch. The more diverse set of subclasses potentially prevents the model from optimizing for a small subset of all subclasses per epoch only. However, we find no significant difference in classification performance if we train the model with a batch size of 32 until convergence for 5 epochs. We assume that smaller batches, which cover less subclasses, have no negative effect on classification performance at our task, but we did not conduct experiments to further investigate this matter.

V. EXPERIMENTS

For our experiments, we use three evaluation measures as proposed by Fall et al. [2]. Fig. 4 visualizes the three evaluation measures and how they differ in comparing the ranked, top three predicted subclasses with the ground truth main subclass (MC) and incidental subclasses (IC). These measures are tailored to the practical application of patent classification. The measure “top prediction” compares only the top-ranked prediction to the main subclass. The measure “three guesses” compares not only the top-ranked but the three top-ranked ranked predictions to the main subclass. The prediction is successful if one of the top three predictions matches the

⁸<https://hpi.de/naumann/projects/repeatability/text%2Dmining.html>

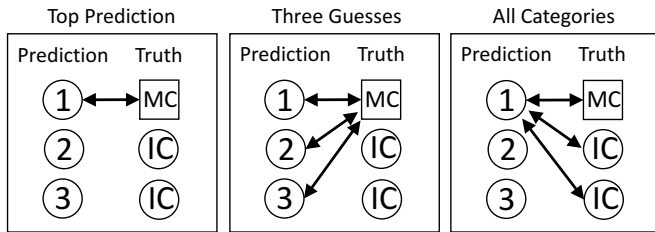


Fig. 4: Three evaluation measures for the task of patent classification. The predicted, ranked subclasses are compared to the ground truth main subclass (MC) and the incidental subclasses (IC). (adapted from Fall et al. [2])

TABLE III: A Comparison of Micro-Average Precision for Different Numbers of Word Embedding Dimensions on the WIPO-alpha dataset.

Evaluation Measure	Word Embedding Dimensions			
	100-patent	200-patent	300-patent	300-wiki
Top-Prediction	45%	48%	49%	42%
Three-Guesses	70%	72%	72%	67%
All-Categories	54%	56%	57%	50%

ground truth main subclass. Both measures, “top prediction” and “three guesses” evaluate only a multi-class classification task. In contrast, the measure “all categories” considers also the incidental subclasses as ground truth information and thus evaluates based on a multi-label ground truth. The measure checks whether the top prediction is included in the set of the main subclass and all incidental subclasses. In theory, this set could contain more than three subclasses. However, in practice the set contains less than two subclasses on average.

We run three experiments to show that our domain-specific word embeddings are beneficial for the task of patent classification. A fourth experiment evaluates how the model’s performance is effected by altering the size of the training data. In the first experiment, we compare the classification performance of four different approaches. Three of them use our pre-trained, patent-specific word embeddings and differ only in the number of word embedding dimensions (either 100, 200, or 300). The fourth approach uses generic 300-dimensional word embeddings trained on Wikipedia pages. All four approaches have the same neural network architecture as described in Section IV-B. We use the WIPO-alpha dataset and apply the three evaluation measures: “top prediction”, “three guesses”, and “all categories”.

Tab. III lists the results of our first experiment. The patent-specific word embeddings, which we trained on 38 billion tokens, outperform word embeddings trained on English Wikipedia pages. This superiority holds if we use 300-dimensional word embeddings for both approaches. However, if we train domain-specific word embeddings with only 100-dimensional vectors, 300-dimensional word embeddings trained on Wikipedia are almost as good as domain-specific word embeddings. The performances of 200- and 300-dimensional domain-specific word embeddings differ only

slightly.

The second experiment compares our best model to state-of-the-art approaches for patent classification to show that our approach achieves competitive results. To this end, we use the same experiment setup as Fall et al., again on the WIPO-alpha dataset and are thereby able to compare with results reported in related work [2], [3]. Tab. IV lists the results of our second experiment. Our best model with domain-specific word embeddings outperforms the best other approach by up to 17 percent (42 percent compared to 49 percent precision).

The third experiment evaluates our approach on a more recent and larger dataset than WIPO-alpha, called USPTO-2M. Unfortunately, this dataset does not distinguish between main subclasses and incidental subclasses. For training our approach, we consider the first listed subclass of each patent as its main subclass. For the majority of patents only one subclass is listed anyways.

The measure “all categories” is not influenced by the fact that the dataset does not explicitly list main subclasses. Both other measures, “top prediction” and “three guesses”, can only be approximated, because we can only guess the ground truth main subclass out of the set of all listed subclasses. Another limitation of the dataset is that it does not contain patents’ full texts but only their abstracts and titles. However, the WIPO-alpha and the USPTO-2M dataset are still quite similar and we assume that the task of patent classification is equally difficult on both datasets. We use the patents of the years 2006 to 2013 as training data and the patents of the year 2014 as test data.

Because of the size of the dataset and memory constraints during training, we can only process the first 30 words of each patent (instead of the first 300 words as in our other experiments). For the same reason, we can use only 100-dimensional and no 300-dimensional word embeddings. Tab. V lists the results of our third experiment. Surprisingly, the classification results are even better on the USPTO-2M dataset with the limited approach than on the WIPO-alpha dataset with our more complex approach. The USPTO-2M dataset contains 25 times more training samples than the WIPO-alpha dataset. We assume that the larger number of training samples is the main reason for the model’s strong performance.

Together, the three experiments show that domain-specific word embeddings together and a GRU-based neural network achieve competitive results at the task of patent classification. In particular, patent-specific word embeddings outperform generic word embeddings trained on Wikipedia pages. However, memory constraints during training limit our approach for the USPTO-2M dataset.

Figure 5 shows a 2-dimensional projection of the word embedding space. For reasons of simplification the visualization includes only the 10,000 most frequent words. We apply the t-SNE algorithm [29] for dimensionality reduction from 300 to 2 dimensions. To this end, we use 500 iterations, a learning rate of 10, and a perplexity of 25. This visualization can be explored interactively (search for words and display their closest neighbors) in a web browser. We have prepared a 10,000-word subset of our dataset and made it available

TABLE IV: A Comparison of Micro-Average Precision for State-of-the-Art Approaches [2], [3] and our Neural Network with Wikipedia Word Embeddings (RNN-wiki) and Patent Word Embeddings (RNN-patent) on the WIPO-alpha dataset.

Evaluation Measure	Naive Bayes [2]	k-NN [2]	SVM [2]	SNoW [2]	k-NN [3]	RNN-wiki	RNN-patent
Top-Prediction	33%	39%	41%	36%	42%	45%	49%
Three-Guesses	53%	62%	59%	56%	67%	69%	72%
All-Categories	41%	46%	48%	43%	50%	53%	57%

TABLE V: Micro-Average Precision for our Neural Network with Patent Word Embeddings (RNN-patent) with 100 dimensions (limited to the first 30 words of each Patent) on the USPTO-2M dataset.

Evaluation Measure	RNN-patent
Top-Prediction	53%
Three-Guesses	75%
All-Categories	64%

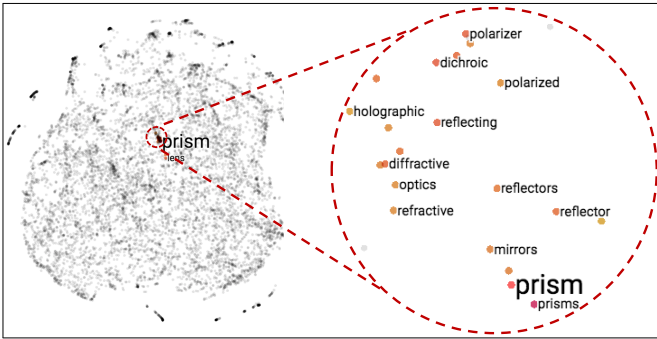


Fig. 5: 2-dimensional projection of the word embedding space with a zoom on the word “prism” and its neighborhood.

online⁹. This subset can be visualized online with the help of tensorflow’s projector¹⁰. The focus of Figure 5 is on the word “prism”. Note that the distance between words in the 2-dimensional projection might slightly differ compared to the distance in the original, 300-dimensional space. Nearest neighbors of the word “prism” in the original space are (from closest to furthest): prisms, dichroic, polarizer, reflecting, lens, reflection, and diffractive. Another example is the word “light-sensitive”, which has the words “photosensitive”, “image-forming”, “photoconductive”, and “photoreceptor” as closest neighbors. The trained word embedding space is correctly taking up the high semantic similarity of these words and represents them with similar embedding vectors.

Our fourth experiment evaluates how the model’s performance is effected by altering the size of the training data. In total there are more than 1.7 million training samples and we speculate that even more data would help performance. To this end we systematically train models with an increasing percentage of the full USPTO-2M training dataset, from 10%, 25%, 50%, and 75% to 100%. Again, we use the patents of the years 2006 to 2013 as training data and sample a subset of each

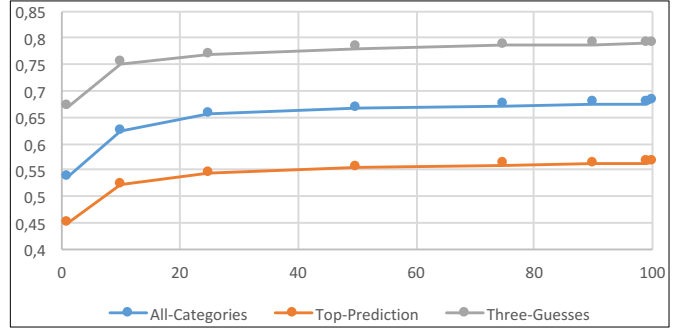


Fig. 6: The micro-average precision at all three evaluated tasks improves with an increasing amount of training data on the USPTO-2M dataset.

year’s patents. Thereby, the smaller samples of the training data cover still the same time span. The test set remains unchanged and contains the patents of the year 2014. Figure 6 shows that more training samples improve the predictive power of the model, however the rate of improvement slows down after training on approximately 25 percent of the training data (400,000 samples). Still the experiment suggests that a slightly better performance is possible with more training data.

VI. ERROR ANALYSIS

We analyze the mis-classifications of our model on different hierarchy levels. Figure 7 visualizes precision scores of the top prediction on the hierarchy level of sections, classes, and subclasses. The performance on the section level is homogeneous and ranges between 0.72 and 0.82, which is visualized by the first column in Figure 7. Although there are many more training samples for section B (10,790 samples) than for section D (1352 samples) or E (2172 samples), the model manages to learn the characteristics of all the different sections with almost equal precision. The second column in Figure 7 reveals that precision on the class level is more heterogeneous. For example, while the model achieves a precision of 0.76 for class A63, the precision for class A62 is only 0.34. One reason is the amount of training samples per class. While there are 565 samples of class A63, there are only 150 samples for class A62. There is also an outlier in section D, where the model achieves a precision of only 0.2 for class D02. This low performance is connected to a small number of training samples for class D02 (29 samples) in contrast to the total number of training samples for section D (1352 samples). The third column in Figure 7 shows the precision for each subclass. Sorting by precision reveals a staircase-shaped sequence of

⁹<https://projector.tensorflow.org/>

¹⁰<https://hpi.de/naumann/projects/repeatability/text-mining.html>

bars in the bar chart. On the subclass D06L the model achieves its worst precision of only 0.06, which we assume is due to the small number of training samples for subclass D06L (20 samples). The classification accuracy on this particular underrepresented subclass (and others, such as A47H, B61B or C21C) could only be improved with more training samples.

As examples for a correct classification and for a wrong classification, we use the patent applications WO/2000/035682/A1 “Tabbed divider and pocket construction” and WO/1999/004984/A1 “Index Pocket and Method for Manufacturing the Same” from the WIPO-alpha dataset. Both of these patents are in section B “PERFORMING OPERATIONS; TRANSPORTING”, class B42 “BOOKBINDING; ALBUMS; FILES; SPECIAL PRINTED MATTER”, and subclass B42F “SHEETS TEMPORARILY ATTACHED TOGETHER; FILING APPLIANCES; FILE CARDS; INDEXING”. While our model classifies the patent “Index Pocket and Method for Manufacturing the Same” correctly into this subclass, the other patent, “Tabbed divider and pocket construction” is wrongly classified into subclass B42D “BOOKS; BOOK COVERS; LOOSE LEAVES; PRINTED MATTER CHARACTERISED BY...”.

The classification is solely based on the first 30 words of title and abstract of the patents. Both titles contain indicator words for subclass B42F and indexing in particular: “tabbed divider” and “index pocket”. So the misclassification cannot be explained with the title. However, the abstract of the two patents differs significantly. The correctly classified patent is described: “The present invention provides an index pocket for easy identifying an index means formed not only on the first page pocket but...”. This description is about the purpose of the invention. In contrast to that, the description of the misclassified patent explains not the purpose of the invention but how the invention can be constructed from a plain sheet of paper: “A foldable paper sheet is formed having a sheet body portion, a sheet lower side flap extending out from a lower edge of the sheet...”. We assume that this difference in the abstract of the two patents complicates the classification task and causes the mis-classification in this exemplary case.

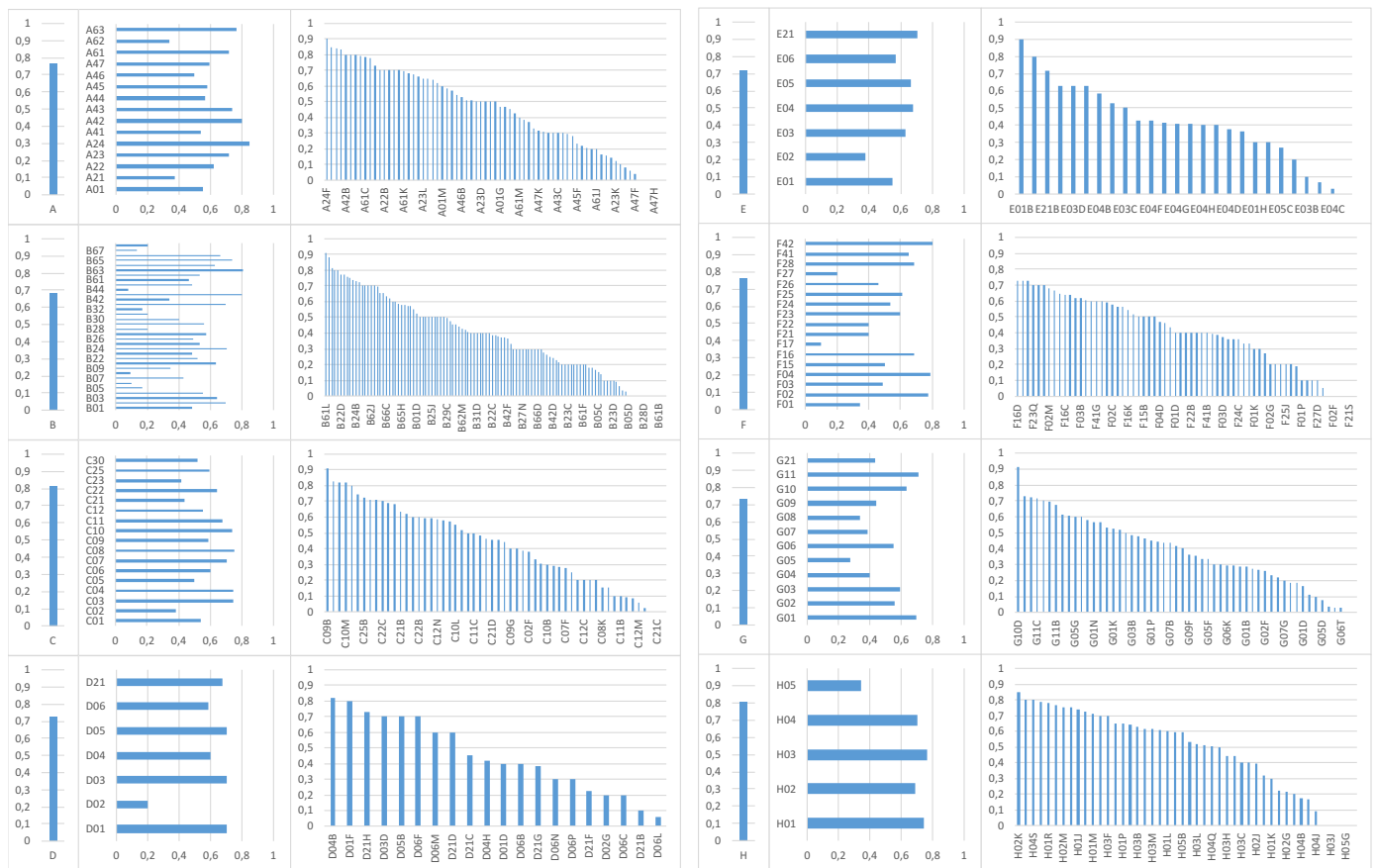
VII. CONCLUSIONS

In this paper, we studied the task of automatic patent classification. We proposed to apply domain-specific fastText word embeddings, which we trained on a large dataset of full texts of more than 5 million patents. Based on these word embeddings that capture the special characteristics of patent speak, we trained a deep neural network with GRUs. Our model is trained with a softmax activation for the task of multi-class classification but is applicable also for multi-label classification. We evaluate our approach with three standard measures in three experiments and improve micro-average precision by 17 percent compared to the state-of-the-art. Further, we find that domain-specific word embeddings, trained specifically on patent documents, outperform generic word embeddings trained in Wikipedia pages. We publish our trained word

embeddings and hope that other researchers can profit from the improved semantic representation of patent language. With an error analysis, we find that mis-classification is often due to low amounts of training data for particular underrepresented subclasses. However, we find that an increasing amount of training data increases overall performance only slightly. The imbalanced training data remains the most difficult challenge. A path for future work is the application of deep learning approaches to other tasks that involve natural language processing in the patent domain, such as classic patent retrieval or reference recommendation. These approaches can surely benefit from pre-trained, domain-specific word embeddings that capture patent speak. Further, an investigation of new neural network architectures tailored to the needs of the patent domain and its hierarchical classification system is promising. The same holds for a comparison of convolutional neural networks and recurrent neural networks at the task of patent classification.

REFERENCES

- [1] H. Smith, “Automation of patent classification,” *World Patent Information*, vol. 24, no. 4, pp. 269–271, 2002.
- [2] C. J. Fall, A. Törösvári, K. Benzineb, and G. Karetka, “Automated categorization in the international patent classification,” in *Acm Sigir Forum*, vol. 37, no. 1. ACM, 2003, pp. 10–25.
- [3] D. Seneviratne, S. Geva, G. Zuccon, G. Ferraro, T. Chappell, and M. Meireles, “A signature approach to patent classification,” in *Asia Information Retrieval Symposium*. Springer, 2015, pp. 413–419.
- [4] N. Nguyen, “Improving hierarchical classification with partial labels,” in *ECAI*, 2010, pp. 315–320.
- [5] J. Rousu, C. Saunders, S. Szedmak, and J. Shawe-Taylor, “Kernel-based learning of hierarchical multilabel classification models,” *Journal of Machine Learning Research*, vol. 7, no. Jul, pp. 1601–1626, 2006.
- [6] X. Qiu, X. Huang, Z. Liu, and J. Zhou, “Hierarchical text classification with latent concepts,” in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*. Association for Computational Linguistics, 2011, pp. 598–602.
- [7] F. Derieux, M. Bobeica, D. Pois, and J.-P. Raysz, “Combining semantics and statistics for patent classification,” in *CLEF (Notebook Papers/Labs/Workshop)*, 2010.
- [8] S. Verberne and E. D’hondt, “Patent classification experiments with the linguistic classification system lcs,” in *CLEF (Notebook Papers/Labs/Workshop)*, 2010.
- [9] S. Li, J. Hu, Y. Cui, and J. Hu, “Deepatent: patent classification with convolutional neural networks and word embedding,” *Scientometrics*, vol. 117, no. 2, pp. 721–744, Nov 2018.
- [10] J. Risch and R. Krestel, “Learning patent speak: Investigating domain-specific word embeddings,” in *Proceedings of the Thirteenth International Conference on Digital Information Management (ICDIM)*, September 2018, pp. 1–6.
- [11] H. Mathiassen and D. Ortiz-Arroyo, “Automatic categorization of patent applications using classifier combinations,” in *International Conference on Intelligent Data Engineering and Automated Learning*. Springer, 2006, pp. 1039–1047.
- [12] T. Tran and R. Kavuluru, “Supervised approaches to assign cooperative patent classification (cpc) codes to patents,” in *International Conference on Mining Intelligence and Knowledge Exploration*. Springer, 2017, pp. 22–34.
- [13] E. D’hondt, S. Verberne, C. Koster, and L. Boves, “Text representations for patent classification,” *Computational Linguistics*, vol. 39, no. 3, pp. 755–775, 2013.
- [14] F. Piroi, M. Lupu, A. Hanbury, and V. Zenz, “Clef-ip 2011: Retrieval in the intellectual property domain,” in *CLEF (notebook papers/labs/workshop)*, 2011.



(a) Section A-D

(b) Section E-H

Fig. 7: Precision of the top prediction for section, class, and subclass level.

- [15] E. D'hondt, S. Verberne, N. Weber, C. Koster, and L. Boves, "Using skipgrams and pos-based feature selection for patent classification," *Computational Linguistics in the Netherlands Journal*, vol. 2, pp. 52–70, 2012.
- [16] S. Verberne and E. D'hondt, "Prior art retrieval using the different sections in patent documents," in *CLEF (Notebook Papers/Labs/Workshops)*, 2010.
- [17] J. Beney, "Lci-insa linguistic experiment for clef-ip classification track," in *CLEF (Notebook Papers/Labs/Workshops)*, 2010.
- [18] J. Guyot, K. Benzineb, G. Falquet, and S. Shift, "myclass: A mature tool for patent classification," in *CLEF (Notebook papers/LABS/workshops)*, 2010.
- [19] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [20] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [21] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2017.
- [22] B. Xia, L. Baoan, and X. Lv, "Research on patent document classification based on deep learning," in *Proceedings of the International Conference on Artificial Intelligence and Industrial Engineering (AIIE)*, 2016.
- [23] M. F. Grawe, C. A. Martins, and A. G. Bonfante, "Automated patent classification using word embedding," in *16th IEEE International Conference on Machine Learning and Applications (ICMLA)*, Dec 2017, pp. 408–411.
- [24] X. Li, H. Chen, Z. Zhang, J. Li, and J. F. Nunamaker, "Managing knowledge in light of its evolution process: An empirical study on citation network-based patent classification," *Journal of Management Information Systems*, vol. 26, no. 1, pp. 129–154, 2009.
- [25] J. Risch and R. Krestel, "My approach = your apparatus?" in *Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries*, ser. JCDL '18. New York, NY, USA: ACM, 2018, pp. 283–292.
- [26] D. Eisinger, G. Tsatsaronis, M. Bundschuh, U. Wieneke, and M. Schroeder, "Automated patent categorization and guided patent search using ipc as inspired by mesh and pubmed," in *Journal of biomedical semantics*, vol. 4, no. 1. BioMed Central, 2013, p. S3.
- [27] Y.-L. Chen and Y.-C. Chang, "A three-phase method for patent classification," *Inf. Process. Manage.*, vol. 48, no. 6, pp. 1017–1030, Nov. 2012.
- [28] K. Benzineb and J. Guyot, *Automated Patent Classification*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 239–261.
- [29] L. Van Der Maaten, "Accelerating t-sne using tree-based algorithms," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3221–3245, 2014.