
Image-based Natural Language Understanding Using 2D Convolutional Neural Networks

Erinç Merdivan^{*†} Anastasios Vafeiadis^{*§} Dimitrios Kalatzis^{*§} Sten Hanke[†]
Johannes Kropf[†] Konstantinos Votis[§] Dimitrios Giakoumis[§]
Dimitrios Tzovaras[§] Liming Chen[¶] Raouf Hamzaoui[¶] Matthieu Geist^{||}

Abstract

We propose a new approach to natural language understanding in which we consider the input text as an image and apply 2D Convolutional Neural Networks to learn the local and global semantics of the sentences from the variations of the visual patterns of words. Our approach demonstrates that it is possible to get semantically meaningful features from images with text without using optical character recognition and sequential processing pipelines, techniques that traditional Natural Language Understanding algorithms require. To validate our approach, we present results for two applications: text classification and dialog modeling. Using a 2D Convolutional Neural Network, we were able to outperform the state-of-art accuracy results of non-Latin alphabet-based text classification and achieved promising results for eight text classification datasets. Furthermore, our approach outperformed the memory networks when using out of vocabulary entities from task 4 of the bAbI dialog dataset.

1 Introduction

Recent advances in natural language processing make heavy use of neural network models. Solutions for tasks such as semantic tagging [8], text classification [26] and sentiment analysis [9] rely on either Recurrent Neural Network (RNN) or Convolutional Neural Network (CNN) variants. In the latter case, the vast majority of the proposed models are based on character-level CNNs applied on one-hot vectors of text or 1D CNNs [18]. Although the results are promising, having either surpassed or equaled the previous state of the art, there are a few issues regarding the proposed models, which are all related to the fundamental inductive bias underlying these models' architectural design. Whether working at the word- or character-level, language processing with most neural network models almost always translates to sequential processing of a string of abstract discrete symbols.

CNNs based on 1D or character convolutions constitute the vast majority of CNN models used in language processing. These networks are fast if the dictionary size is small. However, for some languages, the one-hot encoding vector dimension for input sequences can be very large (e.g., over 3000 for Chinese characters). Furthermore, and specifically for RNN variants, training for long input sequences is difficult due to the well-known problem of vanishing gradients. While architectures like Long Short-Term Memory (LSTM) [17] and Gated Recurrent Units (GRU) [6] were specifically designed to tackle this problem, stable training on long sequences remains an elusive goal, with recent

^{*}The authors contributed equally to this work. Corresponding author erinc.merdivan@ait.ac.at

[†]Department of Health & Environment - Austrian Institute of Technology, GmbH

[‡]CentraleSupélec, Université de Lorraine, CNRS, LORIA, F-57000 Metz, France

[§]Information Technologies Institute (ITI) - Center for Research & Technology Hellas (CERTH)

[¶]De Montfort University (DMU)

^{||}Université de Lorraine, CNRS, LIEC, F-57000 Metz, France (now at Google Brain)

works devising yet more ways to improve performance in recurrent models [27, 30, 33]. Moreover, many state of the art recurrent models rely on the attention mechanism to improve performance [1, 23, 31], which places an additional computational burden on the overall method.

To tackle the above problems, we use CNNs to process the entire text at once as an image. In other words, we convert our textual datasets into images of the relevant documents and apply our model on raw pixel values. This allows us to sensibly apply 2D convolutional layers on text, taking advantage of advances in neural network models designed for and targeting computer vision problems. Doing so, allows us to bypass the issues stated earlier relating to the use of 1D character-level CNNs and RNNs, since now the processing of documents relies on parallel extraction of visual features of many lines (depending on filter size) of text. As far as the vanishing/exploding gradient is concerned, for large CNN architectures, we can easily take advantage of recent architectural advances [15, 16, 7], which specifically aim to improve its effects. In terms of linguistics, our approach is based on the distributional hypothesis [14], where our model produces compositional hierarchies of document semantics by way of its hierarchical architecture. Beyond providing an alternative computational method to deal with the problems described above, our approach is also motivated by findings in neuroscience, cognitive science and the medical sciences where the link between visual perception and recognition of words and semantic processing of language has long been established [22, 28]. Our approach is robust to textual anomalies, such as spelling mistakes, unconventional use of punctuation (e.g., multiple exclamation marks), etc. which factors in during feature extraction. As a result, not only is the need of laborious text preprocessing removed, but the derived models are able to understand the semantic significance of the occurrence of such phenomena (e.g., multiple exclamation marks to denote emphasis), which proves to be especially helpful in tasks such as text classification and/or sentiment analysis. Moreover, our approach can work with any alphabet (latin and non-latin), text font, misspellings and punctuation. Furthermore, it can be extended to handwriting, background colors and table formatted text naturally. It also removes the need of pre-processing real-world documents (and thus the need for optical character recognition, spell check, stemming, and character encoding correction).

Our approach is based on the hypothesis that more semantic information can be extracted from features derived from the visual processing of text than by processing strings of abstract discrete symbols. We test this hypothesis on Natural Language Processing (NLP) tasks and show that a solid understanding of text semantics leads to better model performance. Our contributions are summarized as follows:

- a proof of concept that text classification can be achieved over an image of the text;
- a proof of concept that basic dialogue modeling (restaurant booking), in an information retrieval setting, can be completed using only image processing methods;

The remainder of the paper is organized as follows: Section 2 positions our approach compared to related work, Section 3 introduces the proposed method, Section 4 presents the experimental results and Section 5 draws the conclusions.

2 Related Work

The use of convolutional neural networks for natural language processing has attracted increasing attention in recent years. For sentence classification, Kim [21] used a simple CNN architecture consisting of one convolutional layer with multiple filters of different sizes, followed by max-pooling. The feature maps produced are then fed to a softmax layer for classification. Despite its simplicity, this architecture exhibited good performance. Sentence modeling was further explored by Blunsom et al. [3] who used an extended application, which they call Dynamic Convolutional Neural Network (DCNN) to deal with various input lengths and short- and long-term linguistic dependencies. Wang et al. [32] perform clustering in an embedding space to derive semantic features which they then feed to a CNN with a convolutional layer, followed by k-max pooling and a softmax layer for classification.

Character-level (as opposed to word- or sentence-level) feature extraction was investigated by Zhang et al. [34] who used a standard deep convolutional architecture for text classification. Dos Santos and Gatti [13] carried out sentiment analysis on sentences taken from text corpora, using a CNN architecture which derives input representations that are hierarchically built from the character to the sentence level. Johnson and Zhang [19] used a CNN for text categorization. Their method does not

rely on pre-trained word embeddings, but rather computes convolutions directly on high-dimensional text data represented by one-hot vectors. An architectural variation was also proposed for adapting a bag-of-words model in the convolutional layers. Johnson and Zhang [20] used CNNs for sentiment and topic classification in a semi-supervised framework, where they retained the representations derived by a CNN over text regions, and which they then integrated into the supervised CNN classifier. Ruder et al. [25] employed a novel architecture combining character- and word-level channels to determine an unseen text’s author among a large number of potential candidates, a task they called large-scale authorship attribution. Bjerva et al. [2] introduced a semantic tagging method, which combines (a) stacked neural network models, consisting of a vanilla CNN or a ResNet [16] in the lower level for character-/word-level feature extraction and a bidirectional Gated Recurrent Unit (GRU) in the higher level, with (b) a residual bypass function which preserves the saliency of lower-level features that could be potentially lost in the processing chain of intermediate layers.

While all the aforementioned works used CNNs for NLP tasks, all essentially used text data as input, either pre-trained word embeddings or simply one-hot vector representations.

3 Method

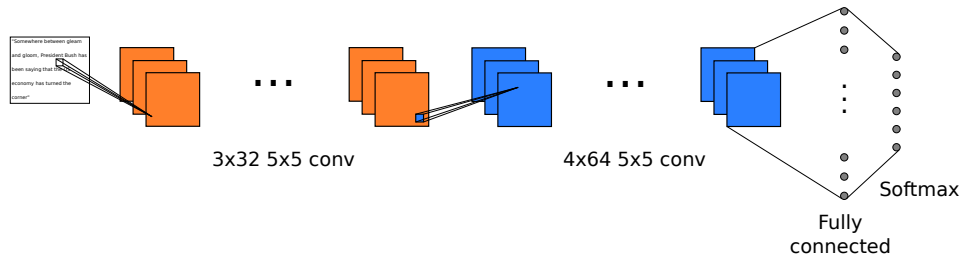


Figure 1: Proposed model: 3 convolutional layers consisting of 32 5x5 filters each, are followed by 4 convolutional layers consisting of 64 5x5 filters each. A linear fully connected layer and a classification output layer complete the model.

In our approach, we treat text understanding as a problem which concerns the learning of context-dependent semantic conventions of language use in a given corpus of text. We treat this complex problem as an image processing problem, where the model processes an image with the text body (Figure 2), learning both the local (word- and sentence-level) and the global semantics of the corpus. In this way, the domain or context dependent meaning of sentences is implicitly contained in the variations of the visual patterns given by the distribution of words in sentences. As such, the problem reduces to a single, where the model needs to observe as many variations of in-domain text as possible to be able to generalize adequately. This process is similar to the analytical method of learning to read [5], where the global meaning of a body of text is acquired first and learning of the text’s meaning moves to hierarchically lower linguistic units. In our case, this translates to understanding the structure and context of the whole corpus first, then the sentences, and finally the words that constitute these sentences.

3.1 Models

For the tasks of (English and Chinese) text classification we used a vanilla CNN and Xception architecture [7] to compare if better vision deep networks can increase performance. We call this network architecture Xception in paper.

The vanilla CNN consists of 7 convolutional layers, a fully connected layer and an output layer containing as many units as classes (e.g., for a classification problem with 4 classes, the output layer would contain 4 units). All filters in the convolutional layers are 5x5 with stride 2. The first 3 layers use 32 filters, while the rest use 64 filters. The fully connected layer consists of 128 units. All units in all layers use the rectifier function, apart from the output layer, which uses a softmax output. Figure 1 shows the architecture of our model.

For the task of dialog modeling we used version 4 of the recently proposed deep Inception network (Inception-V4) [29]. Our choice was motivated by the fact that the vanilla CNN model was too

simple to effectively model the dialog structure, as well as its pragmatics (i.e., the use of language in discourse within the context of a given domain), a problem which Inception-V4 seems to have tackled, at least to a certain extent. We selected the Inception-V4 against the Xception because we wanted to experiment with different advanced architectures for similar tasks.

3.2 Data Augmentation

Data augmentation has been shown to be essential for training robust models [10, 11]. For image recognition, augmentation is applied using simple transformations such as shifting the width and the height of images by a certain percentage, scaling, or randomly extracting sub-windows from a sample of images [12].

For the task of English and Chinese text classification, we used the *ImageDataGenerator* function provided by Keras. The input image was shifted in width and height by 20%, rotated by 15° and flipped horizontally, using a batch size of 50. For the task of dialogue modeling, we applied the same augmentation techniques and random character flipping. Character flip and in particular changing the rating of a restaurant improved the per-response and per-dialog accuracy, especially for difficult sentences, such as booking a 4 star restaurant.

4 Results

To validate our approach, we ran experiments for two separate tasks: text classification and dialog modeling, using a single NVIDIA GTX 1080 Ti GPU.

4.1 Text classification

In this task we trained our model on an array of datasets which contained text related to news (AG’s News and Sogou’s News), structured ontologies on Wikipedia (DBPedia), reviews (Yelp and Amazon) and question answering (Yahoo! answers). Details about the datasets can be found in [34]. For this task, Zhang et al. [34] tested CNNs that use 1D convolutions in the task of text classification, which may more broadly include natural language understanding, as well as sentiment analysis. While the model used in [34] uses text as input vectors, our proposed method uses image data of text. In other words, whereas Zhang et al. [34] use one-hot vector representations of words or word embeddings, we use binarized pixel values of grayscale images of text corpora.

Table 1: Results of Latin and non-Latin alphabet-based text classification in terms of held-out accuracy. *Worst-Best Performance* reports the results of the worst and best performing baselines from Table 4 of Zhang et al. [34] and Conneau et al. [10]. Results reported for *TI-CNN* were obtained in 10 epochs

Dataset	Worst-best Performance (%)	TI-CNN (%)	Xception (%)	Number of Classes
AG’s News	83.1-92.3	80.0	91.8	4
Sogou News (Pinyin)	89.2-97.2	90.2	94.6	5
Sogou News (Chinese)	93.1-94.5	-	98.0	5
DBPedia	91.4-98.7	91.7	94.5	14
Yelp Review Polarity	87.3-95.7	90.3	92.8	2
Yelp Review Full	52.6-64.8	55.1	55.7	5
Yahoo! Answers	61.6-73.4	57.6	73.0	10
Amazon Review Full ⁷	44.1-63	50.2	57.9	5
Amazon Review Polarity ¹	81.6-95.7	88.6	94.0	2

Table 1 shows our method’s held-out accuracy in the task of Latin and non-Latin (Sogou News in Chinese) alphabet-based text classification for each of the datasets. All baselines are derived

⁷Amazon datasets were large and we did not have enough computational resources to achieve comparable results to SoA with the Xception

Table 2: *TI-CNN* sentiment prediction for human-generated input text. The model was trained on Amazon Review Polarity dataset

Sample No.	Text Sample	Positivity Score
1	'this product is mediocre'	0.60
2	'this product is excelent'	0.91
3	'this product is excellent'	0.96
4	this product is excellent!!!!	0.98
5	'I love this product it is great'	0.99
6	'I like this product it is ok'	0.78
7	'I don't know'	0.56
8	'as;kdna;sdn nokorgmnsd kasdn;laknsdnaf'	0.51
9	'I recommended this product to every one in the beginning, but it turned out to be horrible later'	0.64
10	'I recommended this product to every one, in the beginning it was working great, I was in love it'	0.96

from Table 4 of Zhang et al. [34] and Conneau et al. [10]. We denote the vanilla CNN used by *TI-CNN* (which stands for Text-to-Image Convolutional Neural Networks). The column *Worst-Best Performance* shows the worst and best held-out accuracy achieved by the baseline models. Since this is a proof-of-concept of a visual approach to natural language understanding, to further contextualize it, we have also included the equivalent of random classification for each dataset. Our approach achieved comparable results to most of the best performing baselines. For the Yelp Review Full, the results were below the state-of-the-art, since both the *TI-CNN* and Xception were not fine-tuned for each dataset.

Table 2 shows human generated text (not included in the training set) the model used for testing. For these examples, the table shows model predictions after the model was trained on the *Amazon Review Polarity* data set [24] containing reviews of products in various product categories. The dataset is used for binary (positive/negative) sentiment classification of text and the metric (*positivity score*) is the class probability of being the positive class. The text samples were meant to illustrate our model's merits, compared to pure NLP methods. In detail, the model is able to discriminate between words expressing different degrees of the same sentiment (e.g. samples 1,6 compared to samples 2-5). Sample 2 (compared to samples 3-4) illustrates our method's robustness to anomalies like spelling mistakes. In a traditional NLP setting the misspelled word would have a different representation from the respective correctly-spelled word. Unless the model was trained on data that contained many (and many variations) of these anomalies, or engineered by a human, it would not necessarily correlate the misspelled word with the sentiment it expressed. In our model the misspellings are handled naturally by the network, since the similarity level of misspelled words is high. We note that while this can be alleviated by preprocessing procedures or character-level models, these require more pre-processing or human intervention than our method, which focuses on visual patterns of words.

As discussed before, the model builds these visual representations in a bottom-up fashion, creating a semantic hierarchy which is derived from language use within the context of the corpus domain. Sample 4 shows another interesting characteristic of our model which is the understanding of punctuation even if used informally. The exclamation marks used in sample 4 generated the highest prediction score for positive sentiment among all variations of the same phrase (samples 2-4), indicating that the model is capable of understanding the emphasis. Samples 5 and 6 have a similar structure but the different choice of words to describe positive sentiment affects prediction score. This also exhibits the model's capacity to build meaningful hierarchical representations, as it has learned to discriminate between the small nuances (e.g. choice of words) encountered in (visually and semantically) similar textual structures (sentences). Interestingly, an input which expresses a "neutral" sentiment, such as sample 7, has an analogous prediction score (0.56) that is closer to random guessing in a model that was trained in binary sentiment prediction, which is reasonable behavior. The model is also robust to nonsensical text such as sample 8. The prediction score is, again, close to random guessing.

Samples 9 and 10 are meant to illustrate the model's ability to deal with relatively long sequences of words, which express different (even contradictory) sentiments. Whereas in sample 10 positive sentiment is constantly expressed throughout the sequence, in sample 9 the sentiment switches from



Figure 2: Top: Sogou News dataset with Chinese characters. Bottom: Sogou News dataset with pinyin

positive to negative mid-sequence. While the model erroneously predicts positive sentiment (the prediction score is greater than 0.5), the prediction score is significantly lower for sample 9, compared to sample 10, which indicates that the model has factored in the sentiment switch in its understanding of the text’s semantics.

Finally, we applied the Xception architecture to the Sogou News dataset, using the original Chinese characters (Figure 2). We did not run experiments with the *TI-CNN* for the Sogou News (Chinese), since we only picked the best performing network. Huang and Wang [18] used 1D CNNs for text classification with Chinese characters and showed that the accuracy recognition was higher than the traditional conversion to the pinyin romanization system. We extended this work by using the Xception architecture in the 2D image to achieve almost the same result (Table 1). This proves that regardless of how many Chinese words we fit in a 300x300 or a 200x200 image, our approach outperformed the NLP sequential CNNs. Furthermore, the performance improved when using the Chinese characters instead of the pinyin.

4.2 Dialog modeling

For the dialog modeling task, we tested our Inception-V4-based agent in task 4 (Figure 3) of the bAbI dialog dataset from Bordes et al. [4]. The bAbI dialog dataset consists of 1000 training, 1000 validation and 1000 test dialogs in the restaurant booking domain. Each dialog is divided in four different tasks. Here we focus on task 4, where the dialog agent should be able to read entries about restaurants from a relevant knowledge base and provide the user the information they requested, such as restaurant address or telephone. We note that restaurant telephone numbers and addresses have been delexicalized and replaced by tokens representing this information. We chose to focus on this task to demonstrate the increased effectiveness of visual processing of dialog as opposed to purely linguistic processing, due to the high number of different lexical tokens. In our approach the agent needs to correlate the visual pattern of a knowledge base entry to the relevant request. While in principle this should be easy to achieve using artificial delexicalized tokens, as in this benchmark task, it would be far more difficult to do so in the real world, with non-standard sequences of words (such as restaurant names, addresses, telephone numbers, etc). However, given the results of the text classification tasks, we hypothesize that given enough data, our visual approach can create semantic models that encapsulate such correlations.

5 Conclusion

We presented a proof-of-concept for natural language understanding, relying only on visual features of text. For non-dialog text, images of text as input to the CNN models can build hierarchical semantic representations which let them detect various subtleties in language use, irrespective of the language of the input data. For dialog text, we showed that CNN models learn both the structure of discourse and the implied dialog pragmatics implicitly contained within the training data. Crucially, our approach does not require any preprocessing of natural language data, traditionally found in NLP applications, such as tokenization, optical character recognition, stemming, or spell checking. Our method can work using different computer fonts, background colors and can be expanded to human handwriting. Unlike traditional and language agnostic models, our method can perform NLP tasks on real-world documents that include tables, bold, underlined and colored text, where traditional NLP methods, as well as language agnostic models (1D CNN) fail.

Our work is a first step towards expanding the methods for natural language understanding, exploiting recent advances in image recognition and computer vision. Initial results of this approach are promising for a wide range of NLP tasks, such as text classification, sentiment analysis, dialog modeling and natural language understanding. Future work will study the effect of pre-trained models on non-dialog corpora with regard to modeling performance, incorporation within generative frameworks for text generation for tasks like text summarization or performance of more complex applications (such as recurrent CNNs). For the task of text classification, the recognition accuracy of the Xception will increase, since the reported results are achieved without any fine-tuning for the specific datasets. As computer vision deep learning models continue to improve, we expect our results for the NLP task to follow suit (Table 1).

References

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [2] Johannes Bjerva, Barbara Plank, and Johan Bos. Semantic tagging with deep residual networks. *arXiv preprint arXiv:1609.07053*, 2016.
- [3] Phil Blunsom, Edward Grefenstette, and Nal Kalchbrenner. A convolutional neural network for modelling sentences. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, 2014.
- [4] Antoine Bordes and Jason Weston. Learning end-to-end goal-oriented dialog. *arXiv preprint arXiv:1605.07683*, 2016.
- [5] Sylvie Cèbe and Roland Goigoux. *Apprendre à lire à l'école: Tout ce qu'il faut savoir pour accompagner l'enfant*. Retz, 2011.
- [6] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, 2014.
- [7] François Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint*, pages 1610–02357, 2017.
- [8] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [9] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537, 2011.
- [10] Alexis Conneau, Holger Schwenk, Loïc Barrault, and Yann Lecun. Very deep convolutional networks for text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, volume 1, pages 1107–1116, 2017.
- [11] Ann Copestake. Augmented and alternative nlp techniques for augmentative and alternative communication. *Natural Language Processing for Communication Aids*, 1997.

- [12] Ítalo de Pontes Oliveira, Joao Lucas Peixoto Medeiros, and Vinícius Fernandes de Sousa. A data augmentation methodology to improve age estimation using convolutional neural networks. In *Graphics, Patterns and Images (SIBGRAPI), 2016 29th SIBGRAPI Conference on*, pages 88–95. IEEE, 2016.
- [13] Cícero Nogueira Dos Santos and Maira Gatti. Deep convolutional neural networks for sentiment analysis of short texts. In *COLING*, pages 69–78, 2014.
- [14] Zellig S Harris. Distributional structure. *Word*, 10(2-3):146–162, 1954.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630–645. Springer, 2016.
- [17] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [18] Weijie Huang and Jun Wang. Character-level convolutional network for text classification applied to chinese corpus. *arXiv preprint arXiv:1611.04358*, 2016.
- [19] Rie Johnson and Tong Zhang. Effective use of word order for text categorization with convolutional neural networks. *arXiv preprint arXiv:1412.1058*, 2014.
- [20] Rie Johnson and Tong Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In *Advances in neural information processing systems*, pages 919–927, 2015.
- [21] Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- [22] Stefan Koelsch, Elisabeth Kasper, Daniela Sammler, Katrin Schulze, Thomas Gunter, and Angela D Friederici. Music, language and meaning: brain signatures of semantic processing. *Nature neuroscience*, 7(3):302, 2004.
- [23] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.
- [24] Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pages 165–172. ACM, 2013.
- [25] Sebastian Ruder, Parsa Ghaffari, and John G Breslin. Character-level and multi-channel convolutional neural networks for large-scale authorship attribution. *arXiv preprint arXiv:1609.06686*, 2016.
- [26] Fabrizio Sebastiani. Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1):1–47, 2002.
- [27] Minjoon Seo, Sewon Min, Ali Farhadi, and Hannaneh Hajishirzi. Neural speed reading via skim-rnn. *arXiv preprint arXiv:1711.02085*, 2017.
- [28] Panagiotis G Simos, Luis FH Basile, and Andrew C Papanicolaou. Source localization of the n400 response in a sentence-reading paradigm using evoked magnetic fields and magnetic resonance imaging. *Brain research*, 762(1-2):29–39, 1997.
- [29] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, volume 4, page 12, 2017.
- [30] Trieu H Trinh, Andrew M Dai, Thang Luong, and Quoc V Le. Learning longer-term dependencies in rnns with auxiliary losses. *arXiv preprint arXiv:1803.00144*, 2018.
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 6000–6010, 2017.
- [32] Peng Wang, Jiaming Xu, Bo Xu, Chenglin Liu, Heng Zhang, Fangyuan Wang, and Hongwei Hao. Semantic clustering and convolutional neural network for short text categorization. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, volume 2, pages 352–357, 2015.

- [33] Adams Wei Yu, Hongrae Lee, and Quoc V Le. Learning to skim text. *arXiv preprint arXiv:1704.06877*, 2017.
- [34] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In *Advances in neural information processing systems*, pages 649–657, 2015.