

# Multi-Task Deep Learning for Legal Document Translation, Summarization and Multi-Label Classification

Ahmed ELNAGGAR<sup>1</sup>, Christoph GEBENDORFER<sup>a</sup>, Ingo GLASER<sup>a</sup>, and Florian MATTHES<sup>a</sup>

<sup>a</sup>*Software Engineering for Business Information Systems, Technische Universität München, Germany*

**Abstract.** The digitalization of the legal domain has been ongoing for a couple of years. In that process, the application of different machine learning (ML) techniques is crucial. Tasks such as the classification of legal documents or contract clauses as well as the translation of those are highly relevant. On the other side, digitized documents are barely accessible in this field, particularly in Germany. Today, deep learning (DL) is one of the hot topics with many publications and various applications. Sometimes it provides results outperforming the human level. Hence this technique may be feasible for the legal domain as well. However, DL requires thousands of samples to provide decent results. A potential solution to this problem is multi-task DL to enable transfer learning. This approach may be able to overcome the data scarcity problem in the legal domain, specifically for the German language. We applied the state of the art multi-task model on three tasks: translation, summarization, and multi-label classification. The experiments were conducted on legal document corpora utilizing several task combinations as well as various model parameters. The goal was to find the optimal configuration for the tasks at hand within the legal domain. The multi-task DL approach outperformed the state of the art results in all three tasks. This opens a new direction to integrate DL technology more efficiently in the legal domain.

**Keywords.** multi-task deep learning, translation, summarization, multi-label classification

## 1. Introduction

On the past few years, deep learning yielded to great results in many fields, including computer vision, natural language processing (NLP), speech and robotics. In many areas, it was able to out-perform humans including, image classification [1], health [2] and reading comprehension [3]. The availability of large amount of annotated data and fast computing power are the two main reasons behind this big hype. In the legal domain, legal professionals are doing a lot of tasks related to natural language processing daily, which could be replaced by ML algorithm, but that didn't happen yet deeply because the

---

<sup>1</sup>Corresponding Author: Ahmed Elnaggar, Software Engineering for Business Information Systems, Boltzmannstr. 3, 85748 Garching bei München, Germany; E-mail: ahmed.elnaggar@tum.de.

scarcity of annotated data. Despite that fact that there are exceptionally large text base in the legal domain, it was not preprocessed and structured in a format that could be used for ML technology. The use of ML in the legal domain started to take the attentions of the legal profession and some work has already been done, like translating legal documents [4] or classifying verdicts of the French Supreme Court [5]. However, a lot of possible use cases are not exploited yet.

Generating annotated datasets is generally a costly process. It is even more difficult in the legal domain because we can't easily crowd source it. For example "image net" the biggest image classification dataset and "SQUAD" the biggest reading comprehension dataset, were created through Amazon Mechanical Turk, by sourcing people without very specific knowledge. In the legal domain, we need people with very specific knowledge and education to annotate these unstructured data. Which is hard to crowd source it and even more costly. This leads to very circuitual problem on the legal domain:

- NLP is highly required for the legal domain, but the annotated datasets barely exist at all.

One way to overcome this problem is by using multi-task deep learning [6]. In this approach, we train multiple tasks using only one model to provide better results of these problems through transfer learning, especially, tasks that suffers from data scarcity. Therefore, in our work, we needed to achieve two goals:

1. Investigate the effect of transfer learning in the legal problems.
2. Find a big legal text dataset that could be used for transfer learning for any other legal task.

Furthermore, we want to answer three questions regarding the usage of the multi-task deep learning in the legal domain:

1. Is transfer learning through multi-tasking benefits tasks in the legal domain?
2. What are the results of training multiple problems jointly versus separately?
3. Can the multi-task approach outperform the state of the art in the legal domain?

## 2. Related Work

The usage of deep learning has not used intensively in the legal domain. Furthermore, according to our knowledge, the multi-task deep learning was not deeply investigated by researchers and has not been applied in the legal domain. However, we will try to cover the most related research to our work.

**Translation:** A. Vaswani [7] proposed the transformer which represents the current state of the art in general translation, with a BLEU [8] score of 41.8. P. Koehn [4] built 462 machine translation systems for all language pairs of the Acquis Communautaire corpus, which is the body of common rights and obligations which have been adopted by all European Union (EU) Member States.

**Summarization:** AM. Rush [9] initiated work on abstractive summarization with neural networks and induced researchers to continue with sequence-to-sequence models. Additional variants were proposed after that for both extractive and abstractive summarization [10]. C. Grover [11] build the HOLJ corpus for extractive summa-

rization of British judgments. B. Hachey [12] used machine learning for extractive summarization using a corpus of judgments of the UK House of Lords.

**Classification:** [13] Multi-label classification of legal document of the JRC-Acquis using the EuroVoc thesaurus [14,15] is one of the difficult tasks because it has more than 6000 labels and low number of samples per label. R. Steinberger [13] achieved a respectable accuracy of 47.3% on German and 48% on English documents of the JRC-Acquis involving the EuroVoc thesaurus.

**Multi-Task:** R. Collobert [16] build a unified multi-task architecture for various NLP tasks such as SRL, NER, POS, chunking and language modeling. They demonstrated that learning tasks simultaneously can improve performance, and they achieved state-of-the-art performance in SRL by training the SRL task jointly with language model. X. Liu [17] successfully develop a multi-task DNN to combine tasks as disparate as classification and web page ranking. The experimental results demonstrate that the model consistently outperforms strong baselines. P. Liu [18] proposed three RNN based architectures to model text sequence with multi-task learning. They focused their work on four different text classification tasks about movies reviews. H. Zhang [19] proposed a multi-task learning architecture for text classification with four types of recurrent neural layers. Their model outperforms the single task models for various datasets for products and movies reviews. L. Kaiser [20] took the next step of multi-task learning by combining tasks from different modalities including image classification, image caption generation, text translation, text parsing and speech recognition. They showed that adding these tasks together never hurts performance and in most cases improves it on all tasks. They also showed that tasks with less data benefit largely from joint training with other tasks, while performance on large tasks degrades only slightly if at all.

### 3. Legal Corpora

The three datasets that were used are the proceedings of the European parliament (Europarl) [21], digital corpus of the European parliament (DCEP) [22] and Joint Research Centre - Acquis Communautaire (JRC-Acquis) [23].

The Europarl corpus provides the proceedings of the European Parliament between the years 1996 and 2011 for 20 languages. Usually, the documents cover the discussions of political topics. Frequently, sentences contain first-person narrative text expressing political opinions and positions. The DCEP covers different areas including press releases, session protocols, reports of the parliamentary committees and written questions for 23 languages. The JRC-Acquis is a collection of legislative documents, retrieved from the European Union (EU) law, stating EU laws and policies for 22 languages, which have to be implemented by each member state.

Only seven major languages were selected for training as proof of concept including English, German, French, Italian, Spanish, Czech and Swedish. Furthermore, the three datasets were preprocessed from their original format to Moses format [24], which ease the integration of any machine learning platform or library.

**Translation Dataset:** The three datasets were used for the translation is Europarl<sup>2</sup>, DCEP<sup>3</sup> and JRC-Acquis<sup>4</sup>. Including only the previous mentioned 7 languages. The final combined translation dataset contains 4 to 8 million samples sentences per language pair. It is consider a good source for transfer learning, because Summarization and multi-labeling datasets are only 0.5% and 0.3% of it is size.

**Summarization Dataset:** The JRC-Acquis<sup>5</sup> dataset was used for summarization where each document contains a title element holding a short description of the document body. This summary usually varies between one to three sentences representing the semantic core of each document. The dataset contains between 18k to 22k samples per each language.

**Multi-Labeling Dataset:** The JRC-Acquis<sup>6</sup> dataset was used for Multi-Labeling, where each document is assigned to EuroVoc thesaurus annotations. These EuroVoc thesaurus has a hierarchical structure with over than 6000 classes, for example: agriculture, food, health, information technology, law or politics. Furthermore, Each document is usually assigned to various classes that ranges between one and seven classes. The dataset contains between 11k to 14k samples per each language.

#### 4. Multi-Task Legal System

The algorithm we used for multi-task learning is MultiModel algorithm. The algorithm was proposed by Google [20] to create a single generalized deep learning model which is capable of solving tasks across multiple areas (natural language processing, computer vision and speech recognition). This single model was originally trained concurrently on general tasks including image classification, image captioning generation, language translation, English parsing task and speech recognition. However, in our work we used the algorithm for language translation, summarization and document classification specifically in the legal domain.

##### 4.1. MultiModel Architecture

The multi-model uses the sequence to sequence approach based on convolutional neural network. The model consists of four building blocks (Modality Nets, Encoder, Decoder and mixer) as shown on figure 1 and briefly presented [20] bellow.

**Modality Nets:** The model uses four different modality nets (language, image, audio, categorical). This allows it to accept and produce different inputs and outputs types. However, it uses only the language, image and audio for inputs and language and categorical for outputs. Furthermore, it produces a unified representation for all of the tasks. In our case, only the language modality was used.

**Encoder:** The encoder takes the unified output of the modality nets and process it with six custom built convolutional blocks with one mixture-of-expert layer in between.

---

<sup>2</sup><https://mediatum.ub.tum.de/1446650>

<sup>3</sup><https://mediatum.ub.tum.de/1446648>

<sup>4</sup><https://mediatum.ub.tum.de/1446655>

<sup>5</sup><https://mediatum.ub.tum.de/1446654>

<sup>6</sup><https://mediatum.ub.tum.de/1446653>

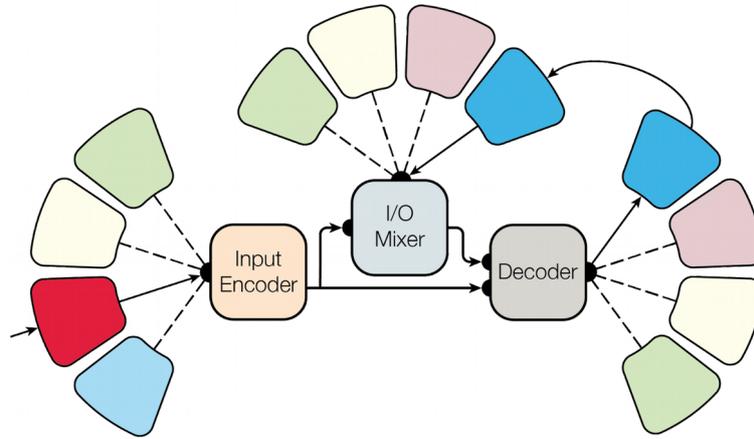


Figure 1. Google Multi-Model Building Blocks

**Decoder:** The decoder takes the output from the encoder and produce the final output using modality nets. It consists of four convolutional attention blocks with one mixture of expert layer in the middle. Furthermore, at each training step a token is passed to it, which allows it learn different representations for different tasks.

**I/O Mixer:** The mixer takes the output of the encoder and the previous output from the decoder. This allows it to learn long term dependencies. It consists of two convolutional blocks and one attention block.

## 5. Experimental Settings

### 5.1. Training Details

Generally, every model was trained until it converged, and sometimes we used early stop to prevent over-fitting. For the multimodel, we have used two configurations the base (MM-B) configuration as was described in the paper and light version (MM-L) configuration. The light version has fewer parameters and was used to test the effect of number of parameters on the result. The transformer, multi-model base and multi-model light were trained with batch size 2048, 2048 and 1048, while the hidden size of each layer was 128, 512 and 512, and the filter size was 1024, 2084 and 2048. In case of our Multi-Task model at each training step we trained the model for the same batch size of each problem sequentially. To speed up the process, we trained the algorithms on four machines. The transformer model on a machine with 4x Tesla K80, the multimodel base version on two machines the first was NVIDIA DGX-1 with 8x Tesla V100 and the second with 5x Titan XP, and the multimodel light version on a machine with 4x Titan 1080Ti.

Different combinations of the jointly tasks has been tested. For translation, we choose two combinations, the pool combination (jt-pool-5) consists of the five available German translation pairs "de-en, de-es, de-fr, de-it, de-sv", and the chain combination

(jt-chain-7) which consists of a chain of language "cs-de, de-en, en-es, es-fr, fr-it, it-sv". For summarization, we had one combination (js-7) which joint all the summarization languages. For multi-labeling, we had one combination (jl-7) which joint all the multi-labeling languages. Finally, we had a last combination which combines different tasks with the same language (ja-3). It combines the translation, summarization and multi-labeling tasks of the same source language together. All of these combinations of the Multi-Model was compared with the result of the state of the art models, which is the transformer for general translation and summarization, and JEX for JRC-Acquis multi-label classification. Finally, due to the time and the number of pages constrains we only report the result of the German language.

## 5.2. Metrics

We report our results with common task-dependent metrics. In the follow sections we cover each task metrics.

### 5.2.1. Translation

The BLEU [8] score was used to evaluate the translation results. It measure the quality of the translation based in the n-grams overlaps between the predicted translation and the target translation.

$$BLEU = \min\left(1, \frac{\text{hypothesis length}}{\text{reference length}}\right) \left(\prod_{i=1}^4 \text{precision}_i\right)^{\frac{1}{4}} \quad (1)$$

### 5.2.2. Summarization

The standard metric for evaluating the summarization is ROUGE [25] score, which we used in for the summary evaluation. We only evaluated the results based on 1-gram, 2-grams and the longest n-gram. They simply called ROUGE-1, ROUGE-2 and ROUGE-L.

$$ROUGE_N = \frac{\sum_{S \in \text{reference.summaries}} \sum_{\text{gram}_n \in S} \text{count\_match}(\text{gram}_n)}{\sum_{S \in \text{reference.summaries}} \sum_{\text{gram}_n \in S} \text{count}(\text{gram}_n)} \quad (2)$$

### 5.2.3. Multi-Label Classification

For multi-label classification, we report precision, recall and F1 score.

$$\text{Precision} = \frac{\#TruePositive}{\#TruePositive + \#FalsePositive} \quad (3)$$

$$\text{Recall} = \frac{\#TruePositive}{\#TruePositive + \#FalseNegative} \quad (4)$$

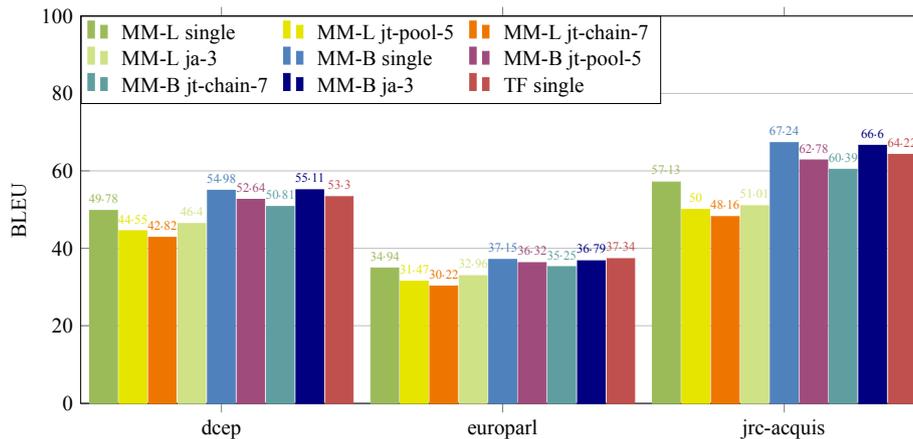
$$F_1 = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

## 6. Result and Discussion

Figure 2 shows the translation results. Generally, all models had better result on both dcep and jrc-acquis datasets than the europarl. This might be because these two datasets contain a lot of cross references, sentence fragments and enumerations compared to the europarl.

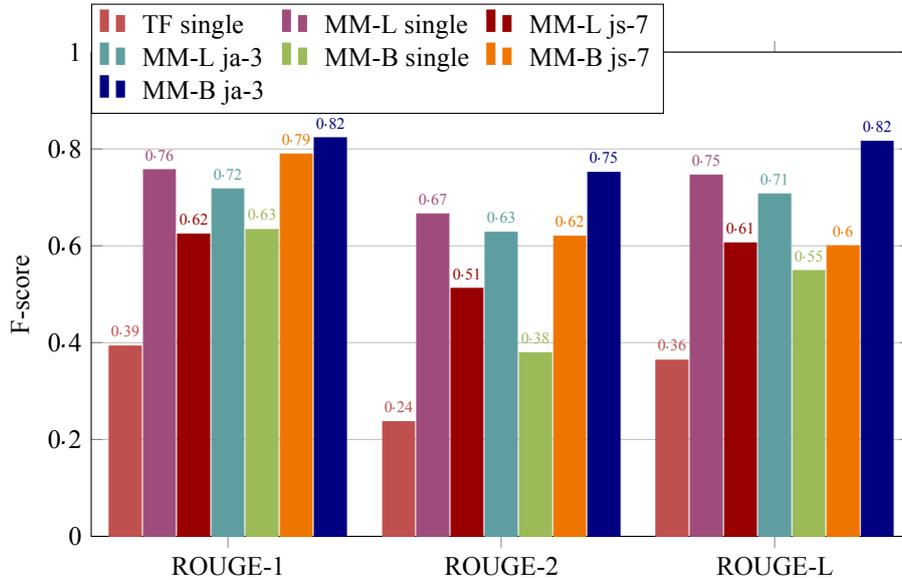
The Multi-model light version (MM-L single, MM-L jt-pool-5, MM-L jt-chain-7 and MML- ja-3) falls behind both the transformer model (TF-B single) and the multi-model base version on the three datasets. This because the number of parameters of the light version is almost half of the number of parameters of the other models. The light version usually produces shorter sentences, however, after manually examine them, we found that semantic meaning remains largely untouched. Another observation that multi-model light single which was trained on a single task outperform the same model but with joint tasks. This is because the limited capacity of the model didn't allow it to learn multi-tasks jointly. By increasing the number of tasks the BLEU score decrease.

The Multi-model base version (MM-B single) outperformed the transformer model for both dcep and jrc-acquis datasets with BLEU score 54.98 and 67.24 compared to 53.3 and 64.22. However, in the case of europarl dataset, the BLEU score was a little bit less, 37.15 compared to 37.34. When the model was trained jointly with other translation languages (MM-B jt-pool-5 and MM-B ja-chain-7) the BLEU score falls behind the transformer. However, the model (MM-B ja-3) which was trained with different tasks (summarization and classification) with the same input language (Germany), outperformed all other models in the dcep dataset with BLEU score 55.11. It was better than the transformer, but less than the multi-model which was trained on a single task (MM-B single) for the acquis with BLEU score 66.6. For the europarl dataset it slightly falls behind both the transformer and the multi-model base single task.



**Figure 2.** German-to-English translation BLEU score performance for all single-task & multi-task translation combinations trained on the MultiModel Light (MM-L), MultiModel Base (MM-B) and Transformer (TF)

Figure 3 shows the summarization results. The transformer model falls behind the multi-model for both the light and base versions. It had almost 50% of the ROUGE of the multi-model. The multi-model light versions had almost always better ROUGE scores than the multi-model base when it trained on either single German summarization or multi-language summarization. The reason is that the dataset of summarization was



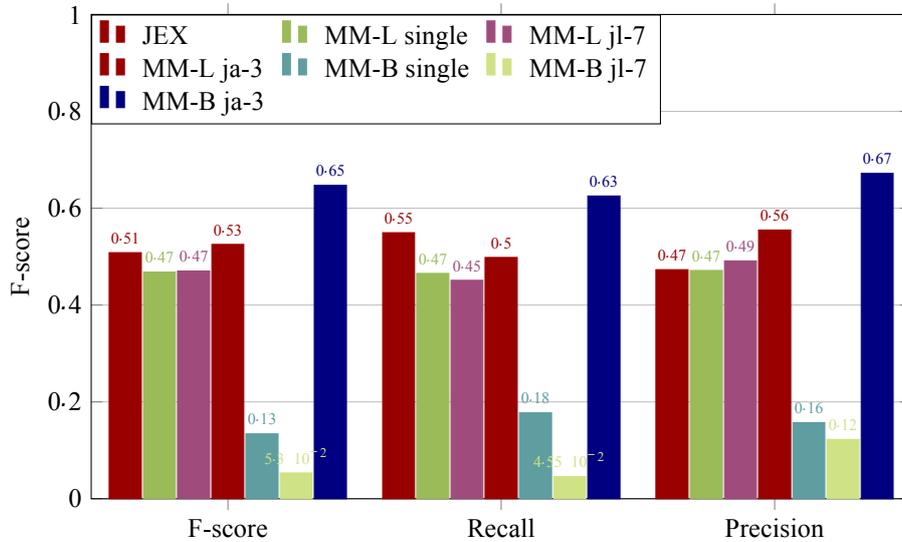
**Figure 3.** German Summarization performance using ROUGE score for all single-task & multi-task translation combinations trained on the MultiModel Light (MM-L), MultiModel Base (MM-B) and Transformer (TF)

small relative to the number of parameters for the base model, which lead to fast overfitting even with using regularization techniques. The best ROUGE scores were obtained from the multi model (MM-B ja-3) which was trained on the three different tasks jointly with the same input language (Germany), with ROUGE-1, ROUGE-2 and ROUGE-L of 0.82, 0.75 and 0.82.

Figure 4 shows the multi-label classification results. The JEX model outperformed the multi-model light (MM-L single, MM-L jl-7 and ML-L ja-3) on both F-score and recall, but it had lower precision. The multi-model base which was trained on single (MM-B single) and all classification languages (MM-B jl-7) falls badly to provide any good classification. The reason behind that is the multi-labeling datasets, which is very small compared to the model capacity, that made the model to over-fit. The best result that outperformed JEX the state of the art model was obtained by combining the three different tasks together with the same language (MM-B ja-3) with F-score, recall and precision of 0.65, 0.63 and 0.67 compared to 0.51, 0.55 and 0.47.

The previous experiments lead to three important points, which answers the three research questions. First, multi-task deep learning outperforms the single task state of the art models, when it is combined with different tasks of the same input language and one of these tasks has a large number of samples. This allows to transfer the knowledge the algorithm learn between these tasks. Second, the greater the number of tasks in a joint task the greater the impact on performance than the relatedness or diversity of the joined tasks. Third, the capacity of the multi-task models must be adopted depending on datasets sizes <sup>7</sup>.

<sup>7</sup>The output of the translation, summarization and classification tasks with the different models can be downloaded from <sup>2</sup>, <sup>3</sup> and <sup>4</sup> in the decodes folder.



**Figure 4.** German multi-label classification performance using F-score, Recall and Precision scores for all single-task & multi-task translation combinations trained on the MultiModel Light (MM-L), MultiModel Base (MM-B) and JEX [26]

## 7. Conclusions & Future Work

We proved that multi-task deep learning can be useful in the legal domain. Of course, the type, the amount of joined tasks and the capacity of the multi-task model are major Influential factors of the result. However, It is an effective approach to solve the data scarcity problem through transfer learning. Using this approach will allow us to outperform the current state of the art results, and allow the usage of the deep learning technology on the legal domain. Our work is a base for further research on the effectiveness and usage of multi-task in the legal domain. However, more experiments are required to test it on other tasks, datasets, languages and training combinations. The provided datasets could be used to test the approach on the rest six languages.

## 8. Acknowledgements

We gratefully acknowledge the support of Leibniz-Rechenzentrum, Microsoft Corporation and NVIDIA Corporation with hardware which were used for this research.

## References

- [1] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. 2015: 1026-1034,” 2017.
- [2] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz, K. Shpanskaya, *et al.*, “Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning,” *arXiv preprint arXiv:1711.05225*, 2017.

- [3] A. W. Yu, D. Dohan, M.-T. Luong, R. Zhao, K. Chen, M. Norouzi, and Q. V. Le, “Qanet: Combining local convolution with global self-attention for reading comprehension,” *arXiv preprint arXiv:1804.09541*, 2018.
- [4] P. Koehn, A. Birch, and R. Steinberger, “462 machine translation systems for europe,” *Proceedings of MT Summit XII*, pp. 65–72, 2009.
- [5] O.-M. S. ulea, M. Zampieri, S. Malmasi, M. Vela, L. P. Dinu, and J. van Genabith, “Exploring the use of text classification in the legal domain,” 2017.
- [6] S. Ruder, “An overview of multi-task learning in deep neural networks,” *arXiv preprint arXiv:1706.05098*, 2017.
- [7] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- [8] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting on association for computational linguistics*, pp. 311–318, Association for Computational Linguistics, 2002.
- [9] A. M. Rush, S. Chopra, and J. Weston, “A neural attention model for abstractive sentence summarization,” *arXiv preprint arXiv:1509.00685*, 2015.
- [10] L. Liu, Y. Lu, M. Yang, Q. Qu, J. Zhu, and H. Li, “Generative adversarial network for abstractive text summarization,” *arXiv preprint arXiv:1711.09357*, 2017.
- [11] C. Grover, B. Hachey, and I. Hughson, “The holj corpus. supporting summarisation of legal texts,” in *Proceedings of the 5th International Workshop on Linguistically Interpreted Corpora*, 2004.
- [12] B. Hachey and C. Grover, “Automatic legal text summarisation: experiments with summary structuring,” in *Proceedings of the 10th international conference on Artificial intelligence and law*, pp. 75–84, ACM, 2005.
- [13] R. Steinberger, M. Ebrahim, and M. Turchi, “Jrc eurovoc indexer jex-a freely available multi-label categorisation tool,” *arXiv preprint arXiv:1309.5223*, 2013.
- [14] E. L. Mencía and J. Fürnkranz, “Efficient multilabel classification algorithms for large-scale problems in the legal domain,” in *Semantic Processing of Legal Texts*, pp. 192–215, Springer, 2010.
- [15] G. Boella, L. Di Caro, D. Rispoli, and L. Robaldo, “A system for classifying multi-label text into eurovoc,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Law*, pp. 239–240, ACM, 2013.
- [16] R. Collobert and J. Weston, “A unified architecture for natural language processing: Deep neural networks with multitask learning,” in *Proceedings of the 25th international conference on Machine learning*, pp. 160–167, ACM, 2008.
- [17] X. Liu, J. Gao, X. He, L. Deng, K. Duh, and Y.-Y. Wang, “Representation learning using multi-task deep neural networks for semantic classification and information retrieval,” 2015.
- [18] P. Liu, X. Qiu, and X. Huang, “Recurrent neural network for text classification with multi-task learning,” *arXiv preprint arXiv:1605.05101*, 2016.
- [19] H. Zhang, L. Xiao, Y. Wang, and Y. Jin, “A generalized recurrent neural architecture for text classification with multi-task learning,” *arXiv preprint arXiv:1707.02892*, 2017.
- [20] L. Kaiser, A. N. Gomez, N. Shazeer, A. Vaswani, N. Parmar, L. Jones, and J. Uszkoreit, “One model to learn them all,” *arXiv preprint arXiv:1706.05137*, 2017.
- [21] P. Koehn, “Europarl: A parallel corpus for statistical machine translation,” in *MT summit*, vol. 5, pp. 79–86, 2005.
- [22] N. Hajlaoui, D. Kolovratnik, J. Väyrynen, R. Steinberger, and D. Varga, “Dcep-digital corpus of the european parliament,” in *LREC*, pp. 3164–3171, 2014.
- [23] R. Steinberger, B. Pouliquen, A. Widiger, C. Ignat, T. Erjavec, D. Tufis, and D. Varga, “The jrc-acquis: A multilingual aligned parallel corpus with 20+ languages,” *arXiv preprint cs/0609058*, 2006.
- [24] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, *et al.*, “Moses: Open source toolkit for statistical machine translation,” in *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*, pp. 177–180, Association for Computational Linguistics, 2007.
- [25] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pp. 74–81, 2004.
- [26] R. Steinberger, M. Ebrahim, and M. Turchi, “JRC eurovoc indexer JEX - A freely available multi-label categorisation tool,” *CoRR*, vol. abs/1309.5223, 2013.