
CONTROLLED RANDOMNESS IMPROVES THE PERFORMANCE OF TRANSFORMER MODELS

A PREPRINT

Tobias Deußer^{*1,2}, Cong Zhao^{1,3}, Wolfgang Krämer³, David Leonhard^{1,2},
Christian Bauckhage^{1,2}, and Rafet Sifa^{1,2}

¹University of Bonn, Bonn, Germany

²Fraunhofer IAIS, Sankt Augustin, Germany

³Deutsche Telekom, Bonn, Germany

ABSTRACT

During the pre-training step of natural language models, the main objective is to learn a general representation of the pre-training dataset, usually requiring large amounts of textual data to capture the complexity and diversity of natural language. Contrasting this, in most cases, the size of the data available to solve the specific downstream task is often dwarfed by the aforementioned pre-training dataset, especially in domains where data is scarce. We introduce controlled randomness, i.e. noise, into the training process to improve fine-tuning language models and explore the performance of targeted noise in addition to the parameters of these models. We find that adding such noise can improve the performance in our two downstream tasks of *joint named entity recognition and relation extraction* and *text summarization*.

Keywords Natural Language Processing · Regularization · Transformer · Machine Learning

1 Introduction

The emergence of pre-trained transformer models brought a massive breakthrough in the field of natural language processing. During pre-training, such transformer models can learn generic language representations with strong generalization capabilities by applying a self-supervised learning approach and leveraging large text corpora. These pre-trained language models can be fine-tuned in various downstream tasks without needing to train from scratch compared to traditional training methods, significantly reducing training costs while achieving excellent performance. Models like BERT Devlin et al. (2019), ELECTRA Clark et al. (2020), or T5 Raffel et al. (2020) have achieved remarkable results on several language processing tasks and the most recent developments of even larger language models, made prominent by GPT-3 Brown et al. (2020) and GPT-4 OpenAI (2023) but not limited to these two¹, improved on these even further. These models have enabled researchers and developers to exploit existing computational linguistic knowledge more conveniently, which in turn has dramatically accelerated the development of natural language processing research and applications.

One of the keys to the success of these models is the ability to adapt to data not encountered during the pre-training phase, i.e. when a downstream task is tackled and the model fine-tuned, as shown in virtually all such applications of transformer models (e.g. Chalkidis et al. (2019), Beltagy et al. (2020), Wang et al. (2021), Aghajanyan et al. (2021), Hillebrand et al. (2022), Ye et al. (2022), Deußer et al. (2023) Ramamurthy et al. (2023), Deußer et al. (2023)). To avoid overfitting and instability during this process, one can apply various regularization techniques and data augmentation, which both might help stabilize fine-tuning and improve performance.

*tdeusser@uni-bonn.de, ORCID-ID: 0000-0003-4685-0847

¹See for example OPT Zhang et al. (2022), Bloom Scao et al. (2023), LLaMA Touvron et al. (2023), or Falcon Almazrouei et al. (2023).

In this work, we build upon the concept introduced in Wu et al. (2022), namely NoisyTune, which describes the process of adding noise to all the parameters of language models in order to regularize it.

We introduce such noise to more parts of the model to examine how this influences the performance of two downstream tasks, specifically *joint named entity recognition and relation extraction* and *text summarization*. We observe that we can enhance the F_1 score of the *joint named entity recognition and relation extraction* model on the KPI-EDGAR dataset Deußer et al. (2022) by **2.977** and the ROUGE–Average² score of the *text summarization* model by **1.387** on the BillSum dataset Kornilova and Eidelman (2019).

Our contribution is thus investigating how noise added to various parts can improve the performance of the considered downstream task. Therein, we find that the model’s performance is significantly improved by adding controlled noise to certain components. We further theorize that this novel approach can help improve the generalisation of large language models to small datasets or low-resource languages.

In the following, we first review related work with a focus on regularization techniques for natural language processing as well as previous studies on *joint named entity recognition and relation extraction* and *text summarization*. Section 3 describes our methodology, i.e., how and where we introduce noise to the model and the general model architecture of both downstream tasks. Thereafter, in Section 4, we outline our dataset, present our experiments, and discuss the results. Section 5 then adds concluding remarks and an outlook into conceivable future work.

2 Related Work

In this section, we discuss various other studies conducted on the effect of regularization techniques on machine learning models. Afterwards, we shortly introduce the most relevant advances in both of our tasks, *joint named entity recognition and relation extraction* and *text summarization*.

2.1 Regularization

Overfitting and thus regularization of machine learning models have been thoroughly studied since the emergence of the field. Early examples of regularization include the work by Hoerl and Kennard (1970), Hanson and Pratt (1988), McCloskey and Cohen (1989), Breiman (1995), Girosi et al. (1995), and Prechelt (2002). Using noise as a regularization technique for training machine learning models was first considered in Bishop (1995).

Hinton et al. (2012) introduced the popular dropout method, which randomly omits a certain part of the feature detectors on each training case. Building on the idea of weight decay Hanson and Pratt (1988), Loshchilov and Hutter (2019) investigated the effects of decoupled weight decay regularization on the training of deep neural networks.

Applying such techniques to natural language processing has gained more importance with the ever-increasing size of the number of parameters of language models. Merity et al. (2018) considered the challenge of word-level language modelling and investigated strategies for regularizing Long Short Term Memory-based Hochreiter and Schmidhuber (1997) models. Lee et al. (2020) proposed to add dropout to randomly mix pre-trained parameters into the downstream model to reduce forgetting in BERT fine-tuning. Furthermore, Dodge et al. (2020) proposed an early stopping method to filter out poor-performing random seeds. On the topic of machine translation, Ott et al. (2018) analyzed uncertainty in machine translation and proposed tools to assess model calibration. Jiang et al. (2020) used smoothness-inducing regularization, which tries to effectively manage the complexity of the model to encourage models to be smooth within neighbourhoods of all the inputs. It was theorized by Sun et al. (2019) that smaller learning rates during fine-tuning help the model retain prior knowledge while still adapting to the new task. In Howard and Ruder (2018), the authors proposed a gradual unfreezing strategy in which layers of the pre-trained language model are unfrozen one at a time during fine-tuning to cope with catastrophic forgetting. Pan et al. (2023) introduced an extra class-aware initialization stage before fine-tuning and concluded that, in this way, self-supervised models should be easier to train to discriminate between different classes. Finally, Wu et al. (2022) described how to regularize a language model by adding noise to its weight parameters, which we will use as a further baseline for our approach. By investigating the effect of noise injection into the individual parts of the model explicitly and measuring the performance on two discrete downstream tasks, our work succeeds Wu et al. (2022) significantly by adding specificity to the analysis.

2.2 Joint Named Entity Recognition and Relation Extraction

When it comes to the first of the evaluated downstream natural language processing tasks, there is introductory work by Miwa and Sasaki (2014), Li and Ji (2014), and Gupta et al. (2016) that was further confirmed by Kamar et al. (2022),

²See subsection 4.1 for how we define this metric.

demonstrating the advantage of joining the subtasks of named entity recognition (NER) and relation extraction (RE) together. Studies, evaluating further concepts to improve the performance of joint NER and RE (also called JNERE), were carried out by Zheng et al. (2017), who introduced a novel tagging scheme, Bekoulis et al. (2018) and Geng et al. (2021), who treated the task as a multi-head selection problem, Fu et al. (2019), who used graph convolutional networks (GCNs), Giorgi et al. (2019) and Xue et al. (2019), who leveraged the pre-trained model BERT, Yu et al. (2019), who improved performances by introducing a new task decomposition strategy, and Shang et al. (2022), who focused on the problems of cascading errors and redundant information in previous models by creating their own, OneRel, which treats joint extraction as a fine-grained triple classification problem.

There is also insight on the application of JNERE regarding work by Bhatia et al. (2019), Chen et al. (2020), and Jabbari et al. (2020) for the medical, legal, and financial language domain respectively.

The work by Deußer et al. (2022), which was already mentioned in the introduction, also lays groundwork for the examinations undertaken in this study. However, to the best of our knowledge, the effect of deliberately incorporating noise into the training process of JNERE models specifically has not yet been systematically researched.

2.3 Text Summarization

As for the task of text summarization, there are several studies conducted on the same dataset as our experiments, that is BillSum, introduced by Kornilova and Eidelman (2019). Most of these make use of or compare their proposed models to models based on BERT, such as An et al. (2021), who used a model that incorporates additional knowledge into the summarization task, Abdel-Salam and Rafea (2022), who compared various BERT variants, Liang et al. (2022), who proposed the Coarse-to- Fine Facet-Aware Ranking (C2F-FAR) framework for unsupervised long document summarization, and Jain et al. (2022), who leveraged the Kullback-Leibler based summarization. Different but general, model-comparative approaches were, for example, undertaken by Rehman et al. (2022) and Mahmoud and Hafez (2022).

BERT was, although on other datasets, also used for text summarization by Zhang et al. (2019), who used the model for the generation of the abstractive summarization output, and Koniaris et al. (2023), who worked on greek legal texts. Other legal domain automatic text summarization techniques were studied by, among, but not limited to, Liu and Lapata (2019), who proposed a general framework for extractive and abstractive models, and Sheik and Nirmala (2021), who focused on improvements through preparation mechanisms and throughout various baselines.

The effects of deliberately introducing noise to the specific task of text summarization were studied by Yousefi-Azar and Hamey (2017), who presented extractive summarization methods based on term-frequency. Liu et al. (2020) used noise to better model uncertainty during training with a student and a teacher model interacting with each other. Both reported improvements using noise at the respective steps in the model training processes. Our work then takes a more fundamental and systematic approach, thoroughly studying the different effects of noise on the different steps of the training process.

3 Methodology

In this section, we first describe how we add noise to the language models to regularize them. Afterwards, we briefly touch upon the models we use for our two tasks, *joint named entity recognition and relation extraction* and *text summarization*.

3.1 Regularization by adding noise

The parameter matrices of a language model are denoted as $[\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_N]$, where N is the number of parameter matrix types. As shown in Wu et al. (2022), the perturbed version $\tilde{\mathbf{W}}_i$ of the parameter matrix \mathbf{W}_i is defined as

$$\tilde{\mathbf{W}}_i = \mathbf{W}_i + \mathbf{U}\left(-\frac{\lambda}{2}, \frac{\lambda}{2}\right) \cdot \sigma(\mathbf{W}_i), \quad (1)$$

where $\mathbf{U}(a, b)$ represents uniformly distributed noise ranging from a to b , λ is the fine-tunable hyperparameter that controls the noise intensity, and $\sigma(c)$ denotes the standard deviation of c . Clearly, parameters with a larger variance are subject to stronger noise due to the multiplication of the noise term with the standard deviation $\sigma(c)$.

In addition to simply adding such noise to all parameters, as seen in Wu et al. (2022), we investigate how the performance of the downstream task is affected when we only partially inject the noise term into certain parts of the model.

More precisely, we add noise to either the bias term, the weight term, or both with different intensities. Furthermore, during the relation extraction downstream task, we add noise to the residual connection and the layer normalization step. We also divide the BERT Devlin et al. (2019) encoder into two separate *layer zones* to add noise in different intensities, as Tenney et al. (2019) theorized that BERT solves various language understanding tasks at different layer depths. On the other hand, during the text summarization task, we insert noise separately into the encoder and decoder parts of the model.

Therefore, we adjust Equation 1 by only perturbing a certain part of the parameter matrix \mathbf{W}_i :

$$\tilde{\mathbf{W}}_i^{\text{loc}} = \mathbf{W}_i^{\text{loc}} + \text{U} \left(-\frac{\lambda}{2}, \frac{\lambda}{2} \right) \cdot \sigma(\mathbf{W}_i^{\text{loc}}), \quad (2)$$

where $\mathbf{W}_i^{\text{loc}}$ is the localized, i.e. restricted to certain parts of the model, parameter matrix. Then, $\tilde{\mathbf{W}}_i^{\text{loc}}$ is the perturbed localized parameter matrix.

3.2 Joint named entity recognition and relation extraction

The *joint named entity recognition and relation extraction* task is defined as extracting entities from a text segment, mostly sentences, and linking them together afterwards. Given this sentence from Deußler et al. (2022),

“In 2021 and 2020 the **total net revenue** was \$**100** million and \$**80** million, respectively.”

one should extract and find the relations

total net revenue – **100**, **total net revenue** – **80**.

To solve the joint named entity recognition and relation extraction task, we employ the model introduced in Hillebrand et al. (2022) titled KPI-BERT. Said model has three main building blocks: A BERT-based sentence encoder, a named entity recognition decoder, and a relation extraction decoder.

To give more detail, given an input sentence tokenized using WordPiece Schuster and Nakajima (2012), we utilize the pre-trained BERT model to obtain the encoded token embeddings. Subsequently, we employ a pooling function to generate word representations by combining the embeddings of individual subwords with a trainable recurrent neural network (RNN) pooling mechanism introduced by Hillebrand et al. (2022). The RNN pooling mechanism is built on a bidirectional gated recurrent unit (GRU) Cho et al. (2014). Additionally, we employ conditional label masking to sequentially tag entities before classifying their relations.

3.3 Text summarization

The *text summarization* task involves generating a shorter version of a given text while preserving its vital information. It can be approached through extractive or abstractive methods. Extractive methods involve selecting and combining critical sentences from the original text.

In contrast, abstractive methods involve generating new sentences that capture the essence of the original text, which is the approach we choose to solve. We pick the Longformer-Encoder-Decoder Beltagy et al. (2020) setup and pre-trained model to solve the text summarisation task. It is a Longformer variant with encoder and decoder transformer stacks, but the encoder uses Longformer’s efficient local and global attention model instead of the initial fully self-attentive one. Additionally, the decoder uses full self-attention for the encoded token and previously decoded locations.

4 Experiments

Here, we first shed some light on the two datasets for our two tasks, namely *joint named entity recognition and relation extraction* and *text summarization*, and how we evaluated these two. Then, we take a closer look at the results and sum them up. All experiments were conducted on a single Nvidia GeForce RTX 2060 SUPER, and the code is implemented in PyTorch.

We execute the experiments within four phases, starting without noise for the downstream tasks. Second, we add noise to all weights as seen in Wu et al. (2022). These two approaches can be seen as our two baselines. Finally, we arrived at our results by adding noise only to the bias or weight parts. For joint named entity recognition and relation

Noise added to	λ	F ₁ in %
Nothing	0	40.696
All	0.81	42.824
Bias	0.41	43.672
Weights	0.50	43.084
Add&Norm	0.2	43.411
Layer zones	0.9	42.200

Table 1: Results of adding noise to KPI-BERT Hillebrand et al. (2022), evaluated on KPI-EDGAR Deußner et al. (2022). Adding noise to all parameters is equivalent to the approach from Wu et al. (2022). Add&Norm refers to the process of adding noise to residual connections and layer normalization.

extraction, noise is additionally added to the residual connections and layer normalization and the layer zones. For text summarization, noise is added separately to the encoder and decoder.

Furthermore, to not cherry-pick any particularly well-performing results, as warned about in Trosten (2023), we run each configuration five times with a different seed each and compare the average of these runs.

4.1 Data and downstream tasks

The joint named entity recognition and relation extraction, defined in the KPI-EDGAR dataset Deußner et al. (2022), aims to extract information from financial documents, including key performance indicators. It holds a total of 1,355 sentences holding 4,522 entities and 3,841 relations, and is split into a training, validation, and test set, encompassing 969, 146, and 240 sentences each. The named entity recognition part is realized in finding numerical and non-numerical entities, whereas the relation extraction part links the found entities to allow for a meaningful extraction of them. Deußner et al. (2022) also defined an *adjusted* F₁ score capturing when an entity is only partially found, which we will also be using.

For this *adjusted* F₁ score, the *true positives* (tp), *false negatives* (fn), and *false positives* (fp) of a relation r between entity i and j are given by:

$$\text{tp}_r = \frac{1}{2} \left(\frac{o_i}{n_{i,\text{gt}}} + \frac{o_j}{n_{j,\text{gt}}} \right) \quad (3)$$

$$\text{fn}_r = 1 - \text{tp}_r \quad (4)$$

$$\text{fp}_r = \frac{1}{2} \left(\frac{n_{i,\text{pred}} - o_i}{n_{i,\text{pred}}} + \frac{n_{j,\text{pred}} - o_j}{n_{j,\text{pred}}} \right), \quad (5)$$

where

$$o_i := |e_{i,\text{pred}} \cap e_{i,\text{gt}}|, \quad (6)$$

is the overlap/intersection o_i of an entity prediction i and its ground truth and $e_{i,\text{pred}}$ and $e_{i,\text{gt}}$ is the set of all token identifiers for the entity prediction and ground truth, respectively. The operation $|\cdot|$ calculates the size of a given set.

For the text summarization task, we consider the BillSum dataset Kornilova and Eidelman (2019), which aggregates U.S. congressional and California state bills and was split into 18,949 training bills and 3,269 test bills. We evaluate our results using the ROUGE metrics proposed in Lin (2004). Furthermore, to make a convenient comparison possible, we average the values of the ROUGE-1 F₁, ROUGE-2 F₁, ROUGE-L F₁ and ROUGE-L-sum and title it *ROUGE-Average*.

4.2 Results

As seen in Table 1 and 2, adding noise, in general, does help the models generalize better. Thus, we can confirm the findings in Wu et al. (2022). Additionally, we are able to improve upon that by exposing only certain parts of the models to noise.

Noise added to	λ	ROUGE–Average in %
Nothing	0	34.148
All	0.3	35.131
Bias	0.4	34.854
Weights	0.1	34.703
Encoder	0.1	34.709
Decoder	0.8	35.534

Table 2: Results of adding noise to the Longformer-Encoder-Decoder Beltagy et al. (2020), evaluated on BillSum Kornilova and Eidelman (2019). Adding noise to all parameters is equivalent to the approach from Wu et al. (2022).

Interestingly, there is no simple “best approach” on where to add noise in a model, as seen in the different best-performing noise locations of Table 1 and 2. Therefore, this is another hyperparameter that has to be fine-tuned. In Wu et al. (2022), they only tested a limited range, i.e. $\lambda \in \{0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3\}$, and they found that $0.1 \leq \lambda \leq 0.15$ were optimal. Contrasting this, in our experiments, the best results are obtained with noise intensities significantly larger than this, as demonstrated in Table 1 and 2.

Still, we found that our approach boosted both performances. In the case of KPI-BERT Hillebrand et al. (2022) on KPI-EDGAR Deußer et al. (2022), we achieve a remarkable increase of **2.977%** F_1 compared to the non-perturbed model and an increase of **0.849%** F_1 compared to the NoisyTune approach Wu et al. (2022). On BillSum Kornilova and Eidelman (2019), our approach improved the performance of the Longformer-Encoder-Decoder Beltagy et al. (2020) by **1.387%** and **0.403%** compared to the non-perturbed model and the NoisyTune baseline, respectively.

5 Conclusion and Future Work

In this paper, we study the effect of controlled randomness on transformer models and how such noise, introduced into various parts of the model as shown in Equation 2, can be seen as a potent regularization tool for ever-increasing language models. We studied two different downstream tasks, namely *joint named entity recognition and relation extraction* and *text summarization*, and found that if certain parts of their respective transformer models are infused with a noise component, we can increase their performance significantly. In the first task, three out of the four scenarios, i.e. adding noise to either the bias term, the weights, or to residual connections and layer normalization, were successful and demonstrated additional enhancements over our already string baseline introduced in Wu et al. (2022). In the second task, we can still improve on this baseline, but only one of our scenarios offers an improvement, namely adding noise the whole decoder and leaving the encoder as it is.

There are two further conclusions from our results. One is that the optimal value of the hyperparameter λ , as defined in Equation (2), depends heavily on the dataset and on the task. The other is that in the paper introducing our baseline NoisyTune Wu et al. (2022), they only tested the range $\lambda \in \{0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3\}$, which should have been extended further to achieve even better results, as seen in Table 1 and 2.

In future work, we plan to tackle even more downstream tasks, like named entity recognition, natural language inference, or sentiment analysis, to see if we can achieve such positive results on these natural language processing applications as well. One could also apply this regularization technique to any pre-trained large language model like Llama Touvron et al. (2023) or Bloom Scao et al. (2023) and experiment on how these can then handle themselves. It might partly alleviate the flaw of requiring large datasets for fine-tuning when training data is sparse.

Another interesting idea would be trying out different noise distributions than the uniformly distributed one seen in Equation (2). One candidate for this could be a Gaussian distribution with fat tails, as this does not have such extreme cut-off points as a uniform distribution.

Furthermore, it would be immensely interesting to see if some low-resource languages can benefit from such a regularization approach. We theorize that in this case, one might see even larger increases in performance, as the “seen language” imbalance of multilingual models, i.e. the distribution of different languages in the training dataset, is usually quite pronounced and heavily skewed towards English, as seen in e.g. Conneau et al. (2020) or Laurençon et al. (2022).

Acknowledgments

This research has been funded by the Federal Ministry of Education and Research of Germany and the state of North-Rhine Westphalia as part of the Lamarr-Institute for Machine Learning and Artificial Intelligence.

References

- Shehab Abdel-Salam and Ahmed Rafea. Performance study on extractive text summarization using bert models. *Information*, 13(2):67, 2022.
- Armen Aghajanyan, Anchit Gupta, Akshat Shrivastava, Xilun Chen, Luke Zettlemoyer, and Sonal Gupta. Muppet: Massive multi-task representations with pre-finetuning. In *Proc. EMNLP*, pages 5799–5811, 2021. doi:10.18653/v1/2021.emnlp-main.468.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. Falcon-40B: an open large language model with state-of-the-art performance, 2023. Placeholder, paper not yet published.
- Chenxin An, Ming Zhong, Zhichao Geng, Jianqiang Yang, and Xipeng Qiu. Retrievalsum: A retrieval enhanced framework for abstractive summarization. *arXiv preprint arXiv:2109.07943*, 2021.
- Giannis Bekoulis, Johannes Deleu, Thomas Demeester, and Chris Develder. Joint entity recognition and relation extraction as a multi-head selection problem. *Expert Systems with Applications*, 114:34–45, 2018.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.
- Parminder Bhatia, Busra Celikkaya, Mohammed Khalilia, and Selvan Senthivel. Comprehend medical: a named entity recognition and relationship extraction web service. In *ICMLA*, pages 1844–1851. IEEE, 2019.
- Chris M Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1):108–116, 1995. doi:10.1162/neco.1995.7.1.108.
- Leo Breiman. Better subset regression using the nonnegative garrote. *Technometrics*, 37(4):373–384, 1995. doi:10.1080/00401706.1995.10484371.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proc. NeurIPS*, volume 33, pages 1877–1901, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.
- Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. Neural legal judgment prediction in english. In *Proc. ACL*, pages 4317–4323, 2019. doi:10.18653/v1/P19-1424.
- Yanguang Chen, Yuanyuan Sun, Zhihao Yang, and Hongfei Lin. Joint entity and relation extraction for legal documents with legal feature enhancement. In *Proc. COLING*, pages 1561–1571, 2020.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *Proc. SSST-8*, pages 103–111, 2014.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. In *Proc. ICLR*, 2020. doi:10.48550/arXiv.2003.10555.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proc. ACL*, pages 8440–8451, 2020. doi:10.18653/v1/2020.acl-main.747.
- Tobias Deußer, Syed Musharraf Ali, Lars Hillebrand, Desiana Nurchalifah, Basil Jacob, Christian Bauckhage, and Rafet Sifa. KPI-EDGAR: A novel dataset and accompanying metric for relation extraction from financial documents. In *Proc. ICMLA*, pages 1654–1659, 2022. doi:10.1109/ICMLA55696.2022.00254.
- Tobias Deußer, Maren Pielka, Lisa Pucknat, Basil Jacob, Tim Dilmaghani, Mahdis Nourimand, Bernd Kliem, Rüdiger Loitz, Christian Bauckhage, and Rafet Sifa. Contradiction detection in financial reports. In *Proc. NLDL*, volume 4, 2023. doi:10.7557/18.6799.
- Tobias Deußer, Lars Hillebrand, Christian Bauckhage, and Rafet Sifa. Informed named entity recognition decoding for generative language models, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proc. NAACL-HLT*, 2019. doi:10.18653/v1/N19-1423.

- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. Fine-tuning pre-trained language models: Weight initializations, data orders, and early stopping, 2020.
- Tsu-Jui Fu, Peng-Hsuan Li, and Wei-Yun Ma. Graphrel: Modeling text as relational graphs for joint entity and relation extraction. In *Proc. ACL*, pages 1409–1418, 2019.
- Zhiqiang Geng, Yanhui Zhang, and Yongming Han. Joint entity and relation extraction model based on rich semantics. *Neurocomputing*, 429:132–140, 2021.
- John Giorgi, Xindi Wang, Nicola Sahar, Won Young Shin, Gary D Bader, and Bo Wang. End-to-end named entity recognition and relation extraction using pre-trained language models. *arXiv preprint arXiv:1912.13415*, 2019.
- Federico Girosi, Michael Jones, and Tomaso Poggio. Regularization theory and neural networks architectures. *Neural computation*, 7(2):219–269, 1995. doi:10.1162/neco.1995.7.2.219.
- Pankaj Gupta, Hinrich Schütze, and Bernt Andrassy. Table filling multi-task recurrent neural network for joint entity and relation extraction. In *Proc. COLING*, pages 2537–2547, 2016.
- Stephen Hanson and Lorien Pratt. Comparing biases for minimal network construction with back-propagation. In D. Touretzky, editor, *Proc. NeurIPS*, volume 1, 1988.
- Lars Hillebrand, Tobias Deußner, Tim Dilmaghani, Bernd Kliem, Rüdiger Loitz, Christian Bauckhage, and Rafet Sifa. KPI-BERT: A joint named entity recognition and relation extraction model for financial reports. In *Proc. ICPR*, pages 606–612, 2022. doi:10.1109/ICPR56361.2022.9956191.
- Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors, 2012.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. doi:10.1162/neco.1997.9.8.1735.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970. doi:10.1080/00401706.1970.10488634.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proc. ACL*, pages 328–339, 2018. doi:10.18653/v1/P18-1031.
- Ali Jabbari, Olivier Sauvage, Hamada Zeine, and Hamza Chergui. A french corpus and annotation schema for named entity recognition and relation extraction of financial news. In *Proc. LREC*, pages 2293–2299, 2020.
- Deepali Jain, Malaya Dutta Borah, and Anupam Biswas. Improving kullback-leibler based legal document summarization using enhanced text representation. In *2022 IEEE SILCON*, pages 1–5. IEEE, 2022.
- Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. SMART: Robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization. In *Proc. ACL*, pages 2177–2190, 2020. doi:10.18653/v1/2020.acl-main.197.
- Mina Esmail Zadeh Nojoo Kamar, Armin Esmaeilzadeh, and Maryam Heidari. A survey on deep learning techniques for joint named entities and relation extraction. In *IEEE AIoT*, pages 218–224. IEEE, 2022.
- Marios Koniaris, Dimitris Galanis, Eugenia Giannini, and Panayiotis Tsanakas. Evaluation of automatic legal text summarization techniques for greek case law. *Information*, 14(4):250, 2023.
- Anastassia Kornilova and Vladimir Eidelman. BillSum: A corpus for automatic summarization of US legislation. In *Proc. New Frontiers in Summarization*, pages 48–56, 2019. doi:10.18653/v1/D19-5406.
- Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. The bigscience roots corpus: A 1.6tb composite multilingual dataset. In *Proc. NeurIPS*, volume 35, pages 31809–31826, 2022.
- Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. Mixout: Effective regularization to finetune large-scale pre-trained language models. In *Proc. ICLR*, 2020. URL <https://openreview.net/forum?id=HkgaETNtDB>.
- Qi Li and Heng Ji. Incremental joint extraction of entity mentions and relations. In *Proc. ACL*, pages 402–412, 2014.
- Xinnian Liang, Jing Li, Shuangzhi Wu, Jiali Zeng, Yufan Jiang, Mu Li, and Zhoujun Li. An efficient coarse-to-fine facet-aware unsupervised summarization framework based on semantic blocks. *arXiv preprint arXiv:2208.08253*, 2022.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, 2004. URL <https://aclanthology.org/W04-1013>.
- Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*, 2019.

- Yang Liu, Sheng Shen, and Mirella Lapata. Noisy self-knowledge distillation for text summarization. *arXiv preprint arXiv:2009.07032*, 2020.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. ICLR*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Hanan A Hosni Mahmoud and Alaaeldin M Hafez. A novel optimized language-independent text summarization technique. *CMC*, 73(3), 2022.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. doi:10.1016/S0079-7421(08)60536-8.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing LSTM language models. In *Proc. ICLR*, 2018. URL <https://openreview.net/forum?id=SyyGPP0TZ>.
- Makoto Miwa and Yutaka Sasaki. Modeling joint entity and relation extraction with table representation. In *Proc. EMNLP*, pages 1858–1869, 2014.
- OpenAI. GPT-4 technical report, 2023.
- Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. Analyzing uncertainty in neural machine translation. In *Proc. ICML*, volume 35, pages 3956–3965. PMLR, 2018. URL <https://proceedings.mlr.press/v80/ott18a.html>.
- Haolin Pan, Yong Guo, Qinyi Deng, Haomin Yang, Jian Chen, and Yiqun Chen. Improving fine-tuning of self-supervised models with contrastive initialization. *Neural Networks*, 159:198–207, 2023. doi:10.1016/j.neunet.2022.12.012.
- Lutz Prechelt. Early stopping-but when? In *Neural Networks: Tricks of the trade*, pages 55–69. Springer, 2002.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1), jan 2020. ISSN 1532-4435. doi:10.48550/arXiv.1910.10683.
- Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization. In *Proc. ICLR*, 2023. URL <https://openreview.net/forum?id=8aHzds2uUyB>.
- Tohida Rehman, Suchandan Das, Debarshi Kumar Sanyal, and Samiran Chattopadhyay. An analysis of abstractive text summarization using pre-trained models. In *Proc. IEM-ICDC*, pages 253–264. Springer, 2022.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. BLOOM: A 176b-parameter open-access multilingual language model, 2023.
- Mike Schuster and Kaisuke Nakajima. Japanese and korean voice search. In *Proc. ICASSP*, 2012.
- Yu-Ming Shang, Heyan Huang, and Xianling Mao. Onerel: Joint entity and relation extraction with one module in one step. In *Proc. AAAI*, pages 11285–11293, 2022.
- Reshma Sheik and S Jaya Nirmala. Deep learning techniques for legal text summarization. In *Proc. UP-CON*, pages 1–5. IEEE, 2021.
- Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. How to fine-tune BERT for text classification? In *Proc. CCL*, pages 194–206. Springer, 2019.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In *Proc. ACL*, pages 4593–4601, 2019. doi:10.18653/v1/P19-1452.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models, 2023.
- Daniel J Trosten. Questionable practices in methodological deep learning research. In *Proc. NLDL*, volume 4, 2023. doi:10.7557/18.6804.
- Xinyu Wang, Yong Jiang, Nguyen Bach, Tao Wang, Zhongqiang Huang, Fei Huang, and Kewei Tu. Automated concatenation of embeddings for structured prediction. In *Proc. ACL-IJCNLP*, pages 2643–2660, 2021. doi:10.18653/v1/2021.acl-long.206.
- Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. NoisyTune: A little noise can help you finetune pretrained language models better. In *Proc. ACL*, pages 680–685, 2022. doi:10.18653/v1/2022.acl-short.76.

- Kui Xue, Yangming Zhou, Zhiyuan Ma, Tong Ruan, Huanhuan Zhang, and Ping He. Fine-tuning bert for joint entity and relation extraction in chinese medical text. In *BIBM*, pages 892–897. IEEE, 2019.
- Deming Ye, Yankai Lin, Peng Li, and Maosong Sun. Packed levitated marker for entity and relation extraction. In *Proc. ACL*, pages 4904–4917, 2022. doi:10.18653/v1/2022.acl-long.337.
- Mahmood Yousefi-Azar and Len Hamey. Text summarization using unsupervised deep learning. *Expert Systems with Applications*, 68:93–105, 2017.
- Bowen Yu, Zhenyu Zhang, Xiaobo Shu, Yubin Wang, Tingwen Liu, Bin Wang, and Sujian Li. Joint extraction of entities and relations based on a novel decomposition strategy. *arXiv preprint arXiv:1909.04273*, 2019.
- Haoyu Zhang, Jianjun Xu, and Ji Wang. Pretraining-based natural language generation for text summarization. *arXiv preprint arXiv:1902.09243*, 2019.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models, 2022.
- Suncong Zheng, Feng Wang, Hongyun Bao, Yuexing Hao, Peng Zhou, and Bo Xu. Joint extraction of entities and relations based on a novel tagging scheme. *arXiv preprint arXiv:1706.05075*, 2017.