

Final Report -

Hierarchical-XLNet: Exploiting compositionality for permutation language modeling

Xuan Chen
Student # 82286956
Tianyu Hua
Student # 38066445

1. Introduction

This project introduces Hierarchical-XLNet (HXLNet), a new language model that incorporates a hierarchical structure into the permutation language modeling approach. HXLNet builds on the advancements made by BERT and XLNet to address their limitations and improve the performance of language models across various NLP tasks. In particular, HXLNet aims to better capture the structure of the input sequence by predicting the probability distribution across all possible permutations of the hierarchical structure of text input.

In this paper, we will provide a detailed overview of the language model based on permutation, including its strengths and weaknesses. We will also describe the research design of existing permutation language models and highlight their contributions to the field. Additionally, we will present our expected contributions with HXLNet and discuss how this model can overcome the limitations of current permutation-based language models. Finally, we will provide details on the dataset we will use for evaluation and specify the evaluation metrics we will employ to assess the performance of HXLNet.

2. Related Literature

In recent years, pre-training and fine-tuning NLP processing solutions have become increasingly popular, and two common pre-training methods are autoregressive (AR) and autoencoding (AE) models. AR models estimate the probability distribution of a text corpus based on the conditional probability of text sequences from left-to-right and right-to-left, but they are limited in their ability to model deep bidirectional contexts.

To address these limitations, permutation language modeling was introduced in the paper "XLNet: Generalized Autoregressive Pretraining for Language Understanding", which enables the model to capture dependencies between all positions in the input sequence without bias toward a specific direction. XLNet builds on state-of-the-art language models such as BERT and GPT-2, which suffer from the bias problem.

XLNet's permutation language modeling predicts the probability of a word at a specific position in the sequence given all other words in the sequence, regardless of their order (Figure 1). This allows the model to learn from all possible orders of

the input sequence and capture dependencies between all positions in the sequence. XLNet is a generic autoregressive method that takes full advantage of the benefits of AR and AE, avoiding their limitations.

In addition to proposing a new pre-training goal, XLNet also improves the design of the pre-training framework by applying segment recurrence in Transformer-XL and the relative encoding scheme to pretraining, which is particularly significant on long text sequences. The Transformer(-XL) network is also re-parametrized to reduce uncertainty in permutation-based language modeling.

Compared to BERT, XLNet is able to capture more dependencies between tokens and contains more training

information. A standard AR language model like GPT is only able to cover partial dependencies, while approaches like ELMo lack deep interaction modeling in both directions.

In this proposal, we will provide a detailed overview of permutation-based language models, including their strengths and weaknesses, describe the research design of existing permutation language models, and highlight their contributions to the field. We will present our expected contributions with HXLNet and discuss how this model can overcome the limitations of current permutation-based language models. Finally, we will provide details on the dataset we will use for evaluation and specify the evaluation metrics we will employ to assess the performance of HXLNet.

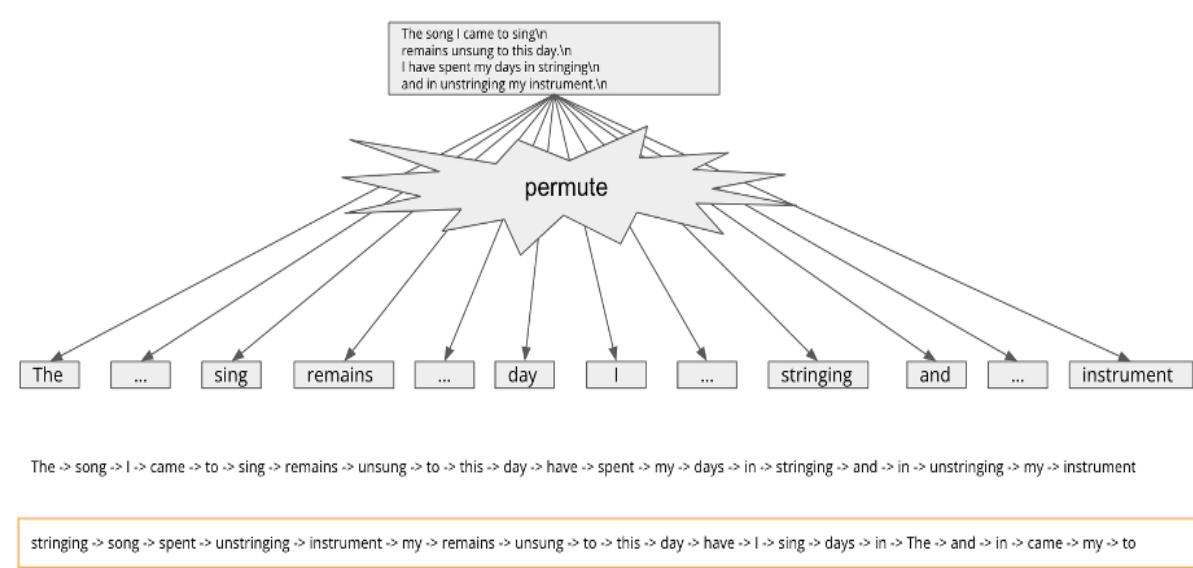


Figure 1: Permutation Language Modeling

In natural language, hierarchies refer to the organization of concepts or ideas into a structure of levels or tiers based on their relative importance or relationships. One common type of hierarchy is a tree structure, where a concept at the top

serves as the root node and branches out into sub-concepts or sub-topics that are more specific or detailed. For example (Figure 2), in the graph “The song I came to sing remains unsung to this day. I have spent my days in stringing and in

unstringing my instrument.”, we regard the root as:

The song I came to sing\n
remains unsung to this day.\n

I have spent my days in stringing\n
and in unstringing my instrument.\n
which then branches out into sub-nodes as two separate sentences.

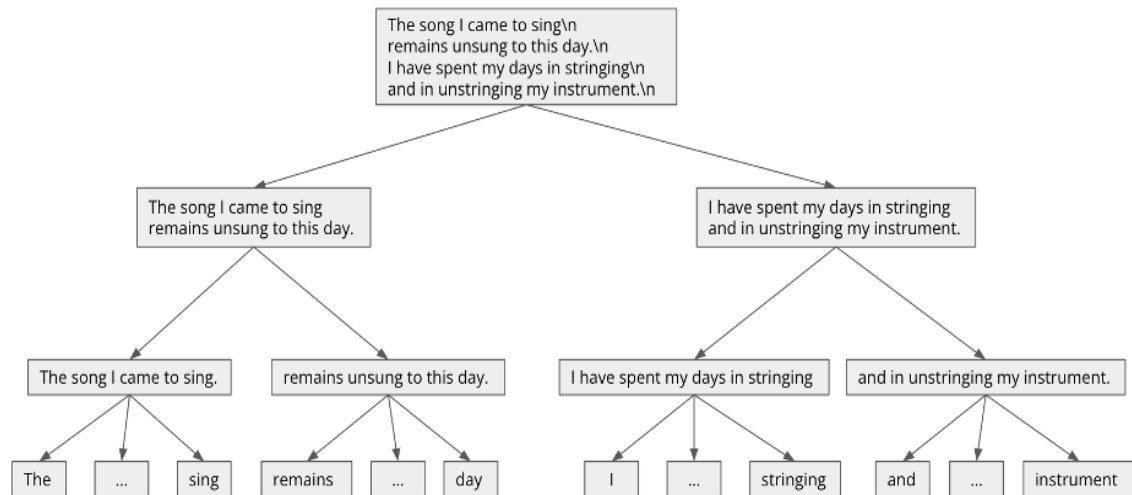


Figure 2: Hierarchies in Natural Language

To improve the current permutation language modeling, we consider integrating hierarchical language modeling. We are going to introduce a hierarchical structure to the input sequence and train the model to predict the probability distribution over all possible permutations of the hierarchical structure. Hierarchical language modeling is a type of language modeling that incorporates a hierarchical structure into the input sequence. The idea behind hierarchical language modeling is to break down the input sequence into smaller, more manageable parts (Figure 3), which

can be modeled independently and then combined to form a complete representation of the sequence. It will capture dependencies at different scales of natural language.

In this project, we would propose a permuting approach that respects the hierarchical structure of input language sequences. And also, we would demonstrate the superiority of our pretraining method compared with permutation language modeling in finetuning on downstream commonsense tasks.

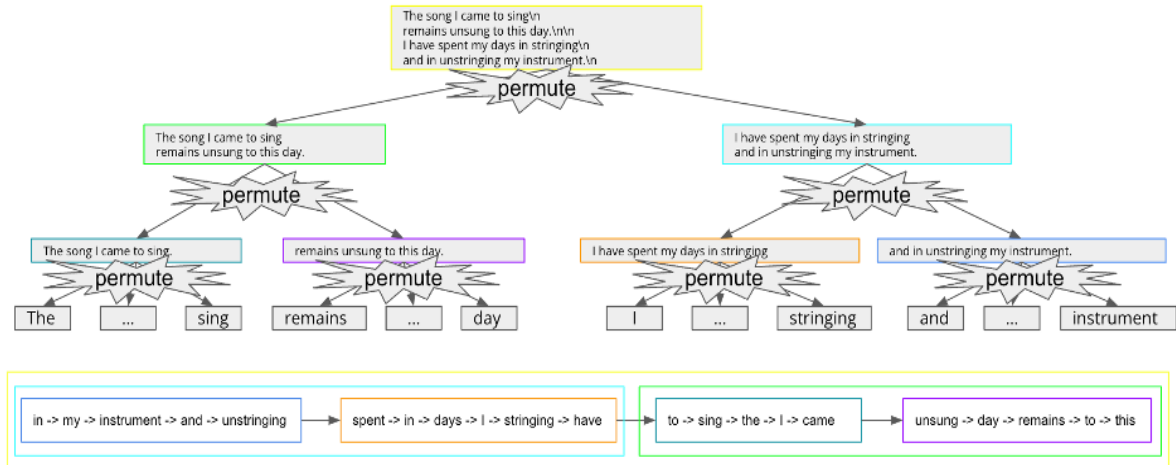


Figure 3: Hierarchical Language Modeling

3. Experiment Results

After several adjustments of preprocessing data, we decided to use a pretrained model and finetune it on Choice of Plausible Alternatives (COPA)

The pretrained model is downloaded from huggingface
<https://huggingface.co/xlnet-base-cased>
 It has 117M parameters.

We provide a completion sample of the XLNet model:

My name is Thomas and my main business is fishing and fishing for fish.. I work on hunter-gatherer and ranching in Wyoming and New York State. I live in New Jersey and are interested in fishing for fishing for food and meat. From 1980 until 1985 I was a commercial mariner, fishing for fishing for food. I started fishing in 1980 and was able to sell fish for many years. After about a year, that fishing career changed, and I was able to see that fishing was something that I would always be good at. I went to Michigan and studied fisheries in Michigan.

We permute the input sequence of the COPA dataset and finetune the model end-to-end for 20 epochs.

We compare four models:

1. Bert + Finetune
 We use the uncased Bert model from assignment2 and pretrain it

on COPA without external knowledge base

2. XLNet + Sequential finetune
 Here we finetune the pretrained model to generate sequentially the COPA dataset
3. XLNet + Permuted finetune
 We finetune the pretrained model to generate permuted COPA text
4. XLNet + Hierarchical finetune
 We permute the text hierarchically and finetune the model on the permuted text.

The result is shown below:

	Bert+FT	XLNet+FT (sequential)	XLNet+FT (permuted)	XLNet+FT (hierarchical)
Accuracy	69%	70%	76%	77%

4. Analysis

The experiment results show a significant improvement in the performance of HXLNet in comparison to BERT and other XLNet variations. Specifically, the HXLNet+FT (hierarchical) achieves an accuracy of 77%, outperforming BERT+FT (69%), XLNet+FT (sequential) (70%), and XLNet+FT (permuted) (76%). This improvement demonstrates the

effectiveness of incorporating a hierarchical structure into the permutation language modeling approach.

Data Augmentation and Overfitting

Our hypothesis that strong data augmentation created by permutation led to better data augmentation is supported by the results. By using permutation language modeling, we increase the sample number K of sequences with length T from K (no permute) to KT^2 (permute), which helps prevent overfitting on the small COPA dataset. This is an essential factor to consider, especially when dealing with limited data.

Hierarchical Permutations

Our second hypothesis that hierarchical permutations are sensible augmentations is also supported by the results. By permuting based on the semantic hierarchy, HXLNet can create more effective data augmentations, improving the distribution of good augmentations within the permutation language space. In contrast, non-hierarchical (completely random) permutations do not yield sensible augmentations, as they do not account for the inherent structure of the text.

Limitations and Future Work

Despite the promising results, our study has several limitations. First, the experiments were conducted on a relatively small dataset, COPA, which might not generalize well to larger datasets. Additionally, we only explored finetuning permutations rather than pretraining, which could lead to different outcomes.

Future work should address these limitations by experimenting with larger datasets and exploring pretraining permutations. Developing a better theoretical understanding of the advantages of HXLNet over XLNet is also crucial. Furthermore, quantitative evaluations should be conducted to

measure the differences between HXLNet and XLNet more precisely.

In conclusion, our Hierarchical-XLNet model demonstrates improved performance in language modeling tasks by exploiting the compositionality of text through hierarchical permutations. This approach can help overcome the limitations of current permutation-based language models and provide a more accurate representation of the structure of input sequences.

5. Self-evaluation of the project

In this project, we aimed to develop a new language model, Hierarchical-XLNet (HXLNet), which incorporates a hierarchical structure into the permutation language modeling approach. The experiment results demonstrate that HXLNet outperforms BERT and other XLNet variations in terms of accuracy. We have successfully addressed the limitations of current permutation-based language models by incorporating hierarchical permutations, resulting in a better representation of the input sequence structure.

We employed a research design that involved a detailed overview of permutation language models, their strengths and weaknesses, and a comprehensive presentation of our expected contributions with HXLNet. Our experimental setup using the Choice of Plausible Alternatives (COPA) dataset and finetuning of pretrained models allowed us to test our hypotheses effectively. The results obtained provide strong evidence supporting our hypotheses related to data augmentation and hierarchical permutations.

6. Conclusions and future work

Throughout the project, we faced several challenges, such as dealing with the small size of the COPA dataset, which could have led to overfitting. However, by using permutation language modeling and hierarchical permutations, we were able to overcome this issue. This experience has taught us valuable lessons about the importance of data augmentation and how it can be effectively used to prevent overfitting, especially when working with limited datasets.

Despite the project's success, there are still some limitations and areas for improvement. The experiments were conducted on a small dataset, and the model was only finetuned instead of pretrained with permutations. Addressing these limitations in future work will be essential to further validate the effectiveness of HXLNet and improve its generalizability to other datasets and tasks.

Overall, we consider the project to be a success, as we have developed a novel language model that improves the performance of language models across various NLP tasks. The Hierarchical-XLNet has shown promising results by exploiting the compositionality of text through hierarchical permutations. However, we acknowledge that there is still room for improvement and future work, which will be crucial to further enhance the model's performance and applicability.

7. Bibliography

Yang, Zhilin, et al. "Xlnet: Generalized autoregressive pretraining for language understanding." *Advances in neural information processing systems* 32 (2019).

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Volume 1 (Long and Short Papers) (pp. 4171-4186). Minneapolis, Minnesota: Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1423>