# PICKING THE UNDERUSED HEADS:
# A NETWORK PRUNING PERSPECTIVE OF ATTENTION HEAD SELECTION
# FOR FUSING DIALOGUE COREFERENCE INFORMATION

*Zhengyuan Liu, Nancy F. Chen*

Institute for Infocomm Research, A*STAR, Singapore

## ABSTRACT

The Transformer-based models with the multi-head self-attention mechanism are widely used in natural language processing, and provide state-of-the-art results. While the pre-trained language backbones are shown to implicitly capture certain linguistic knowledge, explicitly incorporating structure-aware features can bring about further improvement on the downstream tasks. However, such enhancement often requires additional neural components and increases training parameter size. In this work, we investigate the attention head selection and manipulation strategy for feature injection from a network pruning perspective, and conduct a case study on dialogue summarization. We first rank attention heads in a Transformer-based summarizer with layer-wise importance. We then select the underused heads through extensive analysis, and inject structure-aware features by manipulating the selected heads. Experimental results show that the importance-based head selection is effective for feature injection, and dialogue summarization can be improved by incorporating coreference information via head manipulation.

***Index Terms***— Dialogue Summarization, Transformers, Attention Mechanism

## 1. INTRODUCTION

Recently, the Transformer-based models have shown state-of-the-art performance across a variety of Natural Language Processing (NLP) tasks, including, but not limited to, machine translation and reading comprehension [1]. One key component of the Transformer architecture [2] is the layer stacking of multi-head attention that allows the model to capture both local and global information to build feature-rich contextualized representations. In large-scale pre-trained language backbones, it is shown that attention heads at different layers play different roles, and potentially capture grammatical features such as part-of-speech [3] and structural information like syntactic dependency [4]. However, without directly training on corpora that provide explicit and specific linguistic annotation such as coreference and discourse information, model performance remains subpar for language generation tasks that require high-level semantic reasoning [1]. Thus, incorporating such features in a more explicit way raises emerg-

ing research interest [5, 6], including adding attention constrain [7] and adopting separate graph neural components [5].

When fine-tuned on downstream tasks, previous studies show that Transformer-based models are over-parameterized, and can be compressed via structured pruning or knowledge distillation. For instance, previous work showed that a few attention heads do the "heavy lifting" whereas others contribute very little or nothing at all [8]. In practice, in a well-trained multi-layer Transformer, a large percentage of attention heads can be masked at the inference stage without significantly affecting performance, and some layers can even be reduced to only one head [9]. Inspired by this observation, we rethink the strategy of incorporating structure-aware information in a network pruning perspective: redundant attention heads can be replaced with featured weights, and it is much more computationally efficient than introducing additional neural components. In this paper, we conduct a case study on abstractive dialogue summarization, where high-level semantic features are necessary for achieving optimal performance. We investigate the following two research questions:

- Are some attention heads less important or redundant in a well-trained dialogue summarizer?

- Can we manipulate the underused heads with coreference information to improve the summarization model?

Experiments are conducted on a benchmark dialogue summarization corpus SAMSum [10]. We first take a quantitative observation on the importance of attention heads by scoring and ranking them with a gradient-based algorithm, and conducting structured head pruning at the fine-tuning and inference stage. We empirically find that masking a set of lowest-ranking heads does not affect the model performance on the downstream task, regardless of different training settings. We then evaluate two head manipulation methods to incorporate the coreference information into the neural dialogue summarizer, and experimental results show that the model can obtain improved performance, and the manipulated heads are effectively utilized with higher importance.

## 2. TRANSFORMER-BASED MODELS

The Transformer [2] utilizes self-attention instead of recurrent or convolutional neural components. It can be in the

form of encoder-only and encoder-decoder architectures, and Transformer-based sequence-to-sequence models are popular in various language generation tasks such as machine translation and summarization [11]. The encoder consists of stacked Transformer layers, and in each of them, there are two sub-components: a multi-head self-attention layer and a position-wise feed-forward layer. Between these two sub-components, residual connection and layer normalization are added. The $u$-th encoding layer is formulated as:

$$\widetilde{h}^u = \text{LayerNorm}(h^{u-1} + \text{MultiHeadAttn}(h^{u-1})) \quad (1)$$

$$h^u = \text{LayerNorm}(\widetilde{h}^u + \text{FFN}(\widetilde{h}^u)) \quad (2)$$

where $h^u$ is the input to $u$-th layer. MultiHeadAttn, FNN, and LayerNorm are multi-head attention, feed-forward, and layer normalization, respectively. The decoder consists of stacked Transformer layers as well. In addition to the two sub-components in encoder, the decoder performs another multi-head attention over the previous decoding hidden states, and over all encoded representations (i.e., cross-attention). Generally, the decoder generates tokens in an auto-regressive manner from left to right.

One sophisticated design of the Transformer for a strong representation learning capability is the multi-head attention mechanism. More specifically, instead of performing a single attention calculation on the input tuple (*i.e.,* key, value, and query) in a $d$-dimension, multi-head attention projects them into $N_h$ different sub-spaces (each sub-space is expected to capture different features [2, 4]). After calculating attention for every head, where produces a $d/N_h$-dimensional output, we aggregate and project the vectors, and obtain the final contextualized representation. The multi-head attention of one layer is formulated as:

$$\text{Attention}(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d/N_h}})V \quad (3)$$

$$head_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (4)$$

$$\text{MultiHeadAttn}(Q, K, V) = \text{Concat}(head_1, ..., head_{N_h}) \quad (5)$$

## 3. DIALOGUE SUMMARIZATION

Abstractive dialogue summarization has raised much research interest in recent years [12, 13]. Unlike documents, conversations are interactions among multiple speakers, they are less structured and are interspersed with more informal linguistic usage [14, 15], making dialogue summarization more challenging. In common human-to-human conversations, the useful information (which usually focuses on some dialogue topics) is exchanged back and forth across multiple speakers (*i.e.*, interlocutors) and dialogue turns. Aside from speakers referring to themselves and each other, they also mention third-party persons, concepts, and objects, resulting in ubiquitous
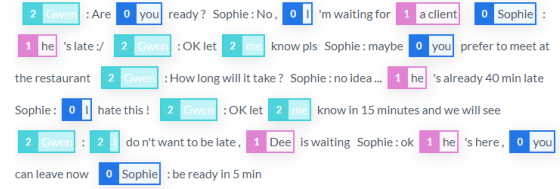


**Fig. 1**. A dialogue example processed by coreference resolution. Colored spans are items in coreference clusters with ID numbers. Each cluster is in one color for better readability.

| Training Set (14,732 Samples) | Mean/Std. of Dialogue Turns | 11.7 (6.45) |
| | Mean/Std. of Dialogue Length | 124.5 (94.2) |
| | Mean/Std. of Summary Length | 23.44 (12.72) |
| Validation Set (818 Samples) | Mean/Std. of Dialogue Turns | 10.83 (6.37) |
| | Mean/Std. of Dialogue Length | 121.6 (94.6) |
| | Mean/Std. of Summary Length | 23.42 (12.71) |
| Test Set (819 Samples) | Mean/Std. of Dialogue Turns | 11.25 (6.35) |
| | Mean/Std. of Dialogue Length | 126.7 (95.7) |
| | Mean/Std. of Summary Length | 23.12 (12.20) |

**Table 1**. Data statistics of the SAMSum corpus.

coreferential expressions [6]. Moreover, the implicit referring such as anaphora or cataphora makes coreference chains more elusive to track. For instance, as the dialogue shown in Figure 1, two speakers exchange information among interactive turns, where the pronoun *"he"* is used multiple times, referring to the word *"client"*. Without sufficient modeling of the referring information, a summarizer cannot link mentions with their antecedents, and produces outputs with factual inconsistency [16]. Therefore, enhancing the model with coreference information is beneficial for dialogue summarization to more appropriately context comprehension, and precisely track the interactive flow throughout a conversation. In this work, we conduct a case study on abstractive dialogue summarization, and assess our proposed method of improving context understanding by incorporating conference features.

## 4. EXPERIMENTAL SETTING

**Corpus** In our setting, we conduct experiments on the SAMSum [10], a benchmark dialogue summarization dataset consisting of 16,369 daily conversations with human-written summaries. Dataset statistics are listed in Table 1.

**Evaluation Metrics** We quantitatively evaluated model outputs with the standard metric ROUGE [17], and reported ROUGE-1, ROUGE-2, and ROUGE-L. All reported results use the same evaluation criteria following previous works [10, 13] for the benchmarked comparison.

**Model Configuration** The baselines and proposed models were implemented in PyTorch and Hugging Face Transformers. AdamW optimizer was used, and the initial learning rate was set at 1e-5. Beam search size was 5. We trained each model for 15 epochs and selected the best checkpoints on the validation set with ROUGE-2 score. All experiments were run on a single Tesla A100 GPU with 40GB memory.
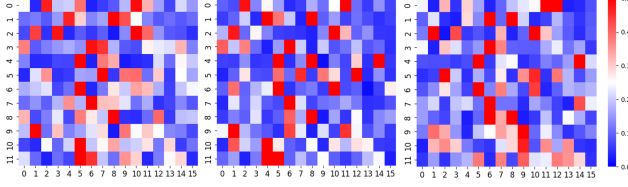
**Fig. 2**. Head importance heatmaps of the *BART-large* model trained with three different training configurations.
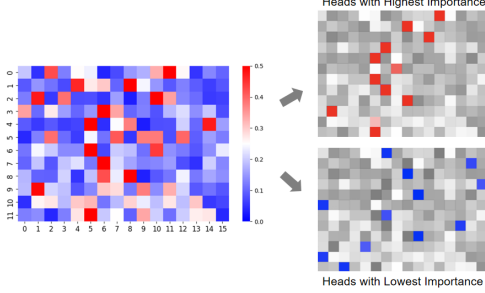


**Fig. 3**. Importance heatmap after averaging operation, and heads with highest/lowest scores are selected in each layer.

## 5. HEAD IMPORTANCE FOR SPECIFIC DOWNSTREAM TASKS

To assess the importance of attention heads in a Transformer-based model, we rank the heads from a network pruning perspective. Here the structured pruning is adopted, which is based on the hypothesis that there is redundancy in the attention heads [8, 18]. To obtain importance scores, we follow [9, 18] and calculate the expected sensitivity (gradient) of the attention heads to the mask variable $\xi^{(i,u)}$ (i.e., $\{0, 1\}$):

$$\text{MultiHeadAttn}(Q, K, V) = \text{Concat}(\xi_1 head_1, .., \xi_{N_h} head_{N_h}) \tag{6}$$

$$S^{(i,u)} = E_{x \sim X} \left| \frac{\partial \mathcal{L}(x)}{\partial \xi^{(i,u)}} \right| \tag{7}$$

where $X$ is the data distribution and $\mathcal{L}(x)$ the loss on sample $x$. Intuitively, if $S^{(i,u)}$ has a high value then changing $\xi^{(i,u)}$ is liable to have a large effect on the model (denotes the contribution score for attention head $i$ at layer $u$). After calculating the contribution scores, we rank the attention heads of each Transformer encoder layer after layer normalization, and obtain those with the highest/lowest scores.

In our preliminary fine-tuning experiments, we observed that there are some ranking variance. As shown in Figure 2, the head importance heatmaps of a *BART-large* model [11] upon different training configurations of random seeds and learning rates are not exactly the same. We then use averaging to exclude heads with high deviation. As shown in Figure 3, some heads show consistently higher/lower layerwise importance scores; It indicates that the head importance has some correlation with the downstream task and training corpus.

To evaluate the impact of pruning heads on the downstream task at the inference stage, we prune the heads of a well-trained model on the summarization corpus, and com-

| Model | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---|---|---|---|
| Baseline (*BART-large*) | 53.14 | 28.58 | 49.69 |
| *Pruning Heads at Inference Stage* | | | |
| Highest-Ranking Heads | 51.72 [↓2.7%] | 27.10 [↓5.2%] | 48.04 [↓3.4%] |
| Lowest-Ranking Heads | 52.89 [↓0.5%] | 27.88 [↓2.5%] | 49.46 [↓0.5%] |
| *Pruning Heads at Training Stage* | | | |
| Highest-Ranking Heads | 52.70 [↓0.9%] | 28.05 [↓1.9%] | 49.14 [↓1.2%] |
| Lowest-Ranking Heads | 53.16 [↑0.1%] | 28.59 [↑0.1%] | 49.73 [↑0.1%] |

**Table 2**. ROUGE F1 scores on attention head pruning at training and inference stage. Relative changes are in brackets.

pare results on the test set. As shown in Table 2, masking the highest-ranking heads of all Transformer layers leads to a 3.4% relative decrease on ROUGE-L, while masking lowest-ranking heads only brings a 0.5% drop. To evaluate the impact of pruning heads at the training stage, we mask the attention heads based on their importance during the fine-tuning process. As shown in Table 2, the model can achieve a comparable result when we mask the highest-ranking heads of all layers. In contrast, the evaluation performance upon masking lowest-ranking (underused) heads even becomes slightly higher. We postulate that the model turns to exploit the rest heads. Therefore, at both the training and inference stages, the head importance is effective to indicate its contribution to the task, and some heads provide limited contribution and can be pruned before fine-tuning.

## 6. HEAD MANIPULATION FOR LINGUISTIC FEATURE INJECTION

### 6.1. Constructing Structure-Aware Matrix

Given dialogue content after coreference resolution, to build the chain of a coreference cluster, we add links between each item and their mentions. Following previous works [6], to better retain local information, we connect all adjacent items in one cluster. More specifically, given a cluster $C_i$ of $m$ items $\{E_1^i, E_2^i..., E_m^i\}$, we add a bi-directional link of each $E$ to its precedent. To construct a structure-aware matrix upon coreference chains for enhancing the Transformer encoder, here we investigate two methods:

**Full-link Matrix** Given a dialogue input $x$ of $n$ tokens (a subword tokenization is utilized), a structure-aware coreference matrix $A_x$ is initialized in a $n^2$ dimension. Iterating each coreference cluster $C$, the first token $t_i$ of each item (*e.g.*, word and text span) is connected with the first token $t_j$ of its antecedent in the same cluster with a bi-directional edge (*i.e.*, $A_x[i][j] = 1$ and $A_x[j][i] = 1$). Then the weights are re-scaled by averaging on the cluster size ($m$ items).

**Adjacent-link Matrix** When the size of one cluster is big, its averaged weights in a full-link matrix will be very small and cause gradient vanishing. Therefore, following the feature aggregation from neighbors in graph neural networks (GNNs), we construct the adjacent-link matrix by only connecting each item with its nearest neighbors.
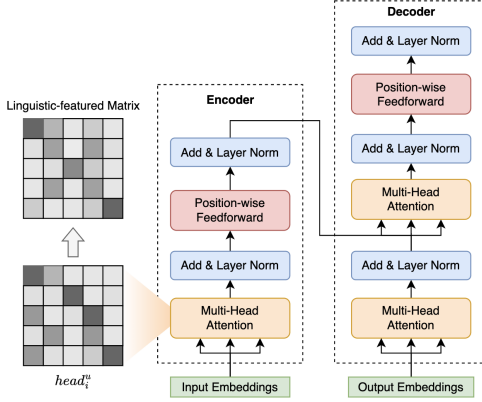
**Fig. 4**. Architecture overview of the coreference-aware Transformer with attention head manipulation.

| Model | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---|---|---|---|
| Baseline (*BART-large*) | 53.14 | 28.58 | 49.69 |
| *MV-BART-Large* | 53.42 | 27.98 | 49.97 |
| *LM-Annotator ($\mathcal{D}_{All}$)* | 53.70 | 28.79 | 50.81 |
| Probing-based Head Selection | | | |
| Full-link Matrix | 53.80 | 28.58 | 50.25 |
| Adjacent-link Matrix | 53.58 | 28.83 | 50.12 |
| Importance-based Head Selection | | | |
| Full-link Matrix | 53.68 | 28.71 | 50.03 |
| Adjacent-link Matrix | 53.98 | 29.15 | 50.73 |

**Table 3**. F1 ROUGE scores on the abstractive dialogue summarization with attention head manipulation.

## 6.2. Attention Head Manipulation

After obtaining structure-aware matrices, we utilize them to enhance the Transformer-based summarizer. Here we directly manipulate attention heads with the featured weights, which is a parameter-free method and more computationally efficient than using additional neural components [5, 19]. It saves 10M parameters (3.1%) and 17% inference time than a GNN-based model. As shown in Figure 4, in $u$-th layer, if one head is lowest-ranking in the importance analysis, we modify it with weights from $A_x$ that present coreference information (*i.e., Importance-based Head Selection*). In our setting, 6 of the 12 *BART* encoding layers were processed for hierarchical modeling, and we empirically found that manipulating heads in higher layers performed better. Additionally, we compare the proposed method to another one named *Probing-based Head Selection* [6], where attention heads that are most similar to the structure-aware matrix $A_x$ (using cosine similarity as measurement) are selected for feature injection (There are no overlapped heads of these two selection strategies).

## 6.3. Results on Dialogue Summarization

Aside from the base model *BART-large* [11], we include two recent state-of-the-art models: *MV-BART-Large* [20] and *LM-Annotator ($D_{All}$)* [13] for extensive comparison. As shown in Table 3, incorporating coreference information helps the

| Model | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---|---|---|---|
| Inference Pruning of Probing-based Heads | | | |
| Full-link Matrix | 52.92 [↓1.7%] | 28.01 [↓2.0%] | 49.17 [↓2.2%] |
| Adjacent-link Matrix | 52.75 [↓1.6%] | 28.10 [↓2.5%] | 49.35 [↓1.6%] |
| Inference Pruning of Importance-based Heads | | | |
| Full-link Matrix | 52.32 [↓2.5%] | 27.41 [↓4.4%] | 48.47 [↓3.1%] |
| Adjacent-link Matrix | 52.44 [↓2.9%] | 27.75 [↓4.9%] | 48.60 [↓4.2%] |

**Table 4**. Ablation study via head pruning at inference stage.

backbone *BART-large*, and makes it comparable to the state-of-the-art models that use additional training data and neural components. In particular, importance-based head selection with adjacent-link matrix performed best with 1.6%, 2.0%, and 2.1% relative F1 score improvement, which is better than the full-link scheme, and the probing-based approach.

## 6.4. Importance Analysis of Manipulated Heads

To qualitatively assess the effectiveness of head manipulation, we conduct an ablation study via head pruning. At the inference stage, we mask the heads that are injected with structure-aware coreference features, and compare it with the unaltered model. As shown in Table 4, pruning the manipulated heads leads to significant performance drop, and the model is affected more by the importance-based than probing-based strategy. Moreover, we compare the importance scores (see Eq. 7) of before and after head manipulation in all encoder layers, and observe that the previously underused heads weigh much higher, demonstrating that the enhanced model effectively utilize the injected features (see Figure 5).
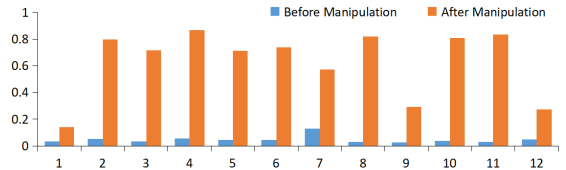


**Fig. 5**. Head importance comparison of before and after the feature injection. Axis X denotes the layer number. Axis Y denotes the normalized importance score.

## 7. CONCLUSION

In this work, we revisited the attention head selection strategy for feature injection from a network pruning perspective. Head importance scoring and ranking of a Transformer-based summarizer showed that some heads are underused after task-specific training. We then manipulated such heads to incorporate structure-aware dialogue coreference features. Experimental results showed that the importance-based head selection is effective for linguistic knowledge injection, and incorporating coreference information is beneficial for dialogue summarization. As a general and computationally efficient approach, this can also be extended to other Transformer-based models and natural language tasks.

## 8. REFERENCES

[1] Pradeep Dasigi, Nelson F. Liu, Ana Marasović, Noah A. Smith, and Matt Gardner, "Quoref: A reading comprehension dataset with questions requiring coreferential reasoning," in *Proceedings of the EMNLP-IJCNLP 2019*, Hong Kong, China, Nov. 2019, pp. 5925–5932, Association for Computational Linguistics.

[2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Proceedings of the NeurIPS*, 2017.

[3] Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick, "What do you learn from context? probing for sentence structure in contextualized word representations," in *Proceedings of the ICLR 2019*, 2019.

[4] John Hewitt and Christopher D. Manning, "A structural probe for finding syntax in word representations," in *Proceedings of the NAACL 2019*, Minneapolis, Minnesota, June 2019, pp. 4129–4138.

[5] Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu, "Discourse-aware neural extractive text summarization," in *Proceedings of the ACL 2020*, Online, July 2020, pp. 5021–5031, Association for Computational Linguistics.

[6] Zhengyuan Liu, Ke Shi, and Nancy Chen, "Coreference-aware dialogue summarization," in *Proceedings of the SIGDIAL 2021*, Singapore and Online, July 2021, pp. 509–519, Association for Computational Linguistics.

[7] Yaushian Wang, Hung-Yi Lee, and Yun-Nung Chen, "Tree transformer: Integrating tree structures into self-attention," in *Proceedings of the EMNLP-IJCNLP 2019*, Hong Kong, China, Nov. 2019, pp. 1061–1070.

[8] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov, "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned," in *Proceedings of the ACL 2019*, Florence, Italy, July 2019, pp. 5797–5808.

[9] Paul Michel, Omer Levy, and Graham Neubig, "Are sixteen heads really better than one?," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[10] Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer, "SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization," in *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, Hong Kong, China, Nov. 2019, pp. 70–79, Association for Computational Linguistics.

[11] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer, "BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proceedings of the ACL 2020*, Online, July 2020, pp. 7871–7880, Association for Computational Linguistics.

[12] Zhengyuan Liu, Angela Ng, Sheldon Lee, Ai Ti Aw, and Nancy F Chen, "Topic-aware pointer-generator networks for summarizing spoken conversations," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 814–821.

[13] Xiachong Feng, Xiaocheng Feng, and Bing Qin, "A survey on dialogue summarization: Recent advances and new frontiers," *arXiv preprint arXiv:2107.03175*, 2021.

[14] HARVEY Sacks, EMANUEL A. SCHEGLOFF, and GAIL JEFFERSON, "A simplest systematics for the organization of turn taking for conversation," in *Studies in the Organization of Conversational Interaction*, JIM SCHENKEIN, Ed., pp. 7 – 55. Academic Press, 1978.

[15] Daniel Jurafsky and James H Martin, "Speech and language processing: An introduction to speech recognition, computational linguistics and natural language processing," *Upper Saddle River, NJ: Prentice Hall*, 2008.

[16] Xiangru Tang, Arjun Nair, Borui Wang, Bingyao Wang, Jai Desai, Aaron Wade, Haoran Li, Asli Celikyilmaz, Yashar Mehdad, and Dragomir Radev, "CONFIT: Toward faithful dialogue summarization with linguistically-informed contrastive fine-tuning," in *Proceedings of the NAACL 2022*, Seattle, United States, July 2022, pp. 5657–5668.

[17] Chin-Yew Lin and Franz Josef Och, "Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics," in *Proceedings of the ACL 2004*, Barcelona, Spain, July 2004.

[18] Sai Prasanna, Anna Rogers, and Anna Rumshisky, "When bert plays the lottery, all tickets are winning," in *Proceedings of the EMNLP 2020*, 2020, pp. 3208–3229.

[19] Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann, "Fixed encoder self-attention patterns in transformer based machine translation," in *Findings of the EMNLP 2020*. 2020, pp. 556–568, Association for Computational Linguistics.

[20] Jiaao Chen and Diyi Yang, "Multi-view sequence-to-sequence models with conversational structure for abstractive dialogue summarization," in *Proceedings of the EMNLP 2020*, Online, Nov. 2020, pp. 4106–4118, Association for Computational Linguistics.