

Rethinking the Role of Scale for In-Context Learning: An Interpretability-based Case Study at 66 Billion Scale

Hritik Bansal^{1*} Karthik Gopalakrishnan^{2†} Saket Dingliwal² Sravan Bodapati²
Katrin Kirchhoff² Dan Roth²

¹University of California, Los Angeles ²AWS AI Labs

hbansal@cs.ucla.edu

{karthgop, skdin, sravanb, katrinki, drot}@amazon.com

Abstract

Language models have been shown to perform better with an increase in scale on a wide variety of tasks via the in-context learning paradigm. In this paper, we investigate the hypothesis that the ability of a large language model to in-context learn-perform a task is not uniformly spread across all of its underlying components. Using a 66 billion parameter language model (OPT-66B) across a diverse set of 14 downstream tasks, we find this is indeed the case: $\sim 70\%$ of the attention heads and $\sim 20\%$ of the feed forward networks can be removed with minimal decline in task performance. We find substantial overlap in the set of attention heads (un)important for in-context learning across tasks and number of in-context examples. We also address our hypothesis through a task-agnostic lens, finding that a small set of attention heads in OPT-66B score highly on their ability to perform primitive *induction* operations associated with in-context learning, namely, prefix matching and copying. These *induction* heads overlap with task-specific important heads, reinforcing arguments by Olsson et al. (2022) regarding induction head generality to more sophisticated behaviors associated with in-context learning. Overall, our study provides several insights that indicate large language models may be under-trained for in-context learning and opens up questions on how to pre-train language models to more effectively perform in-context learning.

1 Introduction

In recent years, large language models (LLMs) (Brown et al., 2020; Rae et al., 2021; Lieber et al., 2021; Black et al., 2022; Zhang et al., 2022; Chowdhery et al., 2022; Hoffmann et al., 2022; Smith et al., 2022) based on the Transformer architecture (Vaswani et al., 2017) pre-trained using self-

supervision on web-scale textual corpora have revolutionized the field of natural language processing (NLP). At larger scales, these models demonstrate remarkable *emergent* (Wei et al., 2022) prowess in performing a wide variety of tasks without any form of fine-tuning, via the zero/few-shot in-context learning paradigm (Brown et al., 2020).

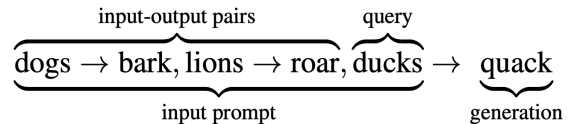


Figure 1: A sample input prompt for in-context learning and the model output.

How in-context learning works has been an open question since its advent and recent studies (Xie et al., 2021; Garg et al., 2022; Olsson et al., 2022; Min et al., 2022b) have begun scratching the surface toward better understanding the paradigm. In this paper, we empirically address the following key question:

Are all LLM components really needed to perform in-context learning?

The first way we address the aforementioned question is through the lens of task-specific importance scores and structured pruning (Li et al., 2016; Molchanov et al., 2016; Anwar et al., 2017) of components underlying modern LLMs, which are primarily stacks composed of multiple high-dimensional self-attention blocks that form *multi-headed attention* and densely activated *feed forward networks* (FFNs). We pick the Open Pre-trained Transformer (OPT) (Zhang et al., 2022) model with 66B parameters for our analyses, which yield several surprising observations. We find that important attention heads are primarily clustered in the intermediate layers and important FFNs are primarily in later layers of the model (§4). We find that the ability to perform zero/few-shot in-context

*Work done as an intern at AWS AI Labs.

†Corresponding author.

Code: github.com/amazon-science/llm-interpret

learning on a variety of 14 NLP datasets/tasks stays nearly intact when up to 70% ($\sim 15.7\text{B}$ parameters in OPT-66B) of the attention heads are removed (§5.1). The attention heads that are (un)important for in-context learning also seem to overlap across tasks (§6.1) and shots (§6.2), and pruning attention heads based on a "universal" importance order computed using all 14 datasets generalizes to varying degrees on out-of-distribution datasets (§6.1.2). These observations indicate that a common task-agnostic subset of the attention heads are responsible for in-context learning. We also find that only up to 20% of the FFNs ($\sim 8.5\text{B}$ parameters) can be removed with minimal decline in zero/few-shot in-context learning performance (§5.2), indicating the importance of FFNs toward in-context learning.

The second way we address the aforementioned question is by quantifying the capacity of all attention heads in OPT-66B to perform a subset of task-agnostic primitive operations associated with in-context learning, namely, *prefix matching* and *copying*: explicitly searching for a prior occurrence of the current token in-context and copying over its suffix. Elhage et al. (2021) and Olsson et al. (2022) developed a mathematical framework to reverse-engineer a Transformer and also find such heads, termed *induction heads*, and explored the hypothesis that such heads drive in-context learning with model sizes up to 13B parameters in a mostly task-agnostic fashion. Using this framework, we compute task-agnostic scores for prefix matching and copying for each attention head and find that a small set of heads in OPT-66B have non-trivial scores for both primitives (§6.3). Qualitative inspection and quantitative analyses show that these heads overlap (to varying degrees) with the ones identified earlier to be important for in-context learning via our set of 14 NLP datasets/tasks, indicating that induction heads are capable of more sophisticated behaviors associated with in-context learning such as latent concept matching but are not the only heads with such capabilities (§6.3.1).

Overall, our study provides several insights about in-context learning at massive scale using both task-specific and task-agnostic settings. In a world of ever increasing language model sizes, we believe these insights serve as a strong foundation for researchers and practitioners in language modeling to build and leverage compact language models that can also demonstrate emergent abilities.

2 Background & Methods

In this section, we establish notation and methods with the Open Pre-trained Transformer (OPT) (Zhang et al., 2022) model used for our study, provide background on in-context learning and the mathematical formulation of induction heads by Olsson et al. (2022) that we build on, and describe our adaptation of oracle and gradient-based importance score formulations for in-context learning.

2.1 Open Pre-trained Transformer (OPT)

OPT is a suite of language models of varying sizes aimed at serving as open replicas of GPT-3. The largest openly accessible model from this suite is OPT-66B with 66 billion parameters.

Architecture: Consider a tokenized input sentence to OPT, $\mathbf{X} \in \mathbb{R}^{N \times d_e}$, where N is the number of tokens in the sentence and d_e is the embedding dimension. The input is processed by multiple decoder layers consisting of *multi-headed attention* (MHA) blocks, *layer norm* (LN) and *feed forward networks* (FFN), followed by a linear layer to produce logits over the vocabulary. The decoder layers can be formally expressed as follows:

$$\mathbf{t}^{(\ell+1)} = \mathbf{z}^\ell + \text{MHA}^\ell(\text{LN}^\ell(\mathbf{z}^\ell)) \quad (1)$$

$$\mathbf{z}^{(\ell+1)} = \mathbf{t}^{(\ell+1)} + \text{FFN}^\ell(\mathbf{t}^{(\ell+1)}) \quad (2)$$

where $\mathbf{z}^1 = \mathbf{X}$, and (1) & (2) are the residual connections corresponding to the MHA and FFN in layer $\ell \geq 1$ respectively. OPT-66B was pre-trained with a maximum sequence length of 2048 and embedding dimension $d_e = 9216$.

MHA: In an MHA block, H attention heads are applied in parallel to the input and their outputs are concatenated. In OPT-66B, there are $H = 72$ attention heads of dimension $d_h = 128$ in every layer ℓ . An individual attention head h in layer ℓ consists of three learnable matrices, $\mathbf{W}_k^h, \mathbf{W}_q^h, \mathbf{W}_v^h \in \mathbb{R}^{d_e \times d_h}$, all unique to the head, such that it applies self-attention $A^h(\cdot)$ on the input, where $d_h = d_e/H$. Formally, for input \mathbf{M} in layer ℓ :

$$\text{MHA}^\ell(\mathbf{M}) = [A^1(\mathbf{M}); \dots; A^H(\mathbf{M})] \mathbf{W}_o^\ell \quad (3)$$

$$A^h(\mathbf{M}) = s^h(\mathbf{M}) \mathbf{M} \mathbf{W}_v^h \quad (4)$$

$$s^h(\mathbf{M}) = \sigma\left(\frac{\mathbf{M} \mathbf{W}_q^h (\mathbf{W}_k^h)^T \mathbf{M}^T}{\sqrt{d_h}}\right) \quad (5)$$

where σ is the softmax function and $\mathbf{W}_o^\ell \in \mathbb{R}^{d_e \times d_e}$ is a learnable output matrix unique to the MHA

block in layer ℓ . To ensure OPT is auto-regressive, the output of $s^h(\cdot)$ is masked to prevent the dependence of the hidden state of the token i , $z_i^\ell \in \mathbb{R}^{d_e}$, on future tokens in indices $\{i + 1, \dots, N\}$.

To remove a head h in layer ℓ in practice, we set $A^h(\mathbf{M})$ to be the zero matrix in Equation (3). This implies that $\mathbf{W}_k^h, \mathbf{W}_q^h, \mathbf{W}_v^h$ can be entirely removed, and the corresponding d_h rows in \mathbf{W}_o^ℓ can also be removed. In total, there are 4608 attention heads across 64 layers in OPT-66B that constitute 21.7B of the total 66B parameters.

FFN: Each layer ℓ consists of a feed forward network (FFN) parameterized by a high-dimensional projection matrix, $\mathbf{W}_1^\ell \in \mathbb{R}^{d_e \times d}$ followed by a low-dimensional projection matrix, $\mathbf{W}_2^\ell \in \mathbb{R}^{d \times d_e}$ where $d = 36864$ for OPT-66B. Formally, for input \mathbf{M} in layer ℓ :

$$\text{FFN}^\ell(\mathbf{M}) = \text{ReLU}(\text{LN}^\ell(\mathbf{M})\mathbf{W}_1^\ell)\mathbf{W}_2^\ell \quad (6)$$

where ReLU is the rectified linear unit activation function and LN is the layer norm.

To remove an FFN in layer ℓ in practice, we set $\text{FFN}^\ell(\mathbf{M})$ to be the zero matrix in Equation (6). This implies $\mathbf{W}_1^\ell, \mathbf{W}_2^\ell$ and the layer norm $\text{LN}^\ell(\cdot)$ for the FFN can be entirely removed. In total, FFNs constitute 43.4B parameters in OPT-66B.

2.2 In-Context Learning & Induction Heads

With increasingly larger language models being trained in recent years, a new paradigm of learning termed *in-context learning* (Brown et al., 2020) has become popular. In this paradigm, language models perform tasks by being prompted to generate output text conditioned on a few (or zero) in-context training examples that form solved "input-output" pairs for the task along with a query input. Figure 1 illustrates the paradigm for the task of identifying the sound that an animal makes. In some cases, tasks can also be accompanied by task descriptions/templates to help prime the language model better, e.g., zero-shot translating from English to German using the prompt:

English phrase: I like dogs. German phrase:

While these examples involve learning and relying on latent concepts during inference, few-shot in-context learning can additionally involve explicit primitive interactions between the in-context examples. For example, with the prompt:

English phrase: I like dogs. German phrase: ich mag Hunde.

English phrase: I like ducks. German phrase:

the model may rely on prior in-context translations of the tokens *I* and *like* when performing the task for the query input. Olsson et al. (2022) developed a mathematical framework toward better understanding such mechanics, starting off with a task-agnostic formulation of in-context learning as the ability of a model to better predict tokens later in the context than the tokens earlier (Kaplan et al., 2020). They define a set of task-agnostic primitive operations that reflect the kind of interactions we refer to in the above example, namely, *prefix matching* and *copying*. These operations are defined in a simplistic fashion on a repeated sequence of randomly generated tokens: explicitly searching for a prior occurrence of the current token in-context and copying over its suffix. The heads that are capable of performing these operations are termed *induction heads*. Figure 2 depicts these operations for a repeated sequence of tokens. While these operations are intertwined in practice, the capacity of attention heads to *independently* perform them is computed with the scoring algorithms described in detail in Appendix A.8.

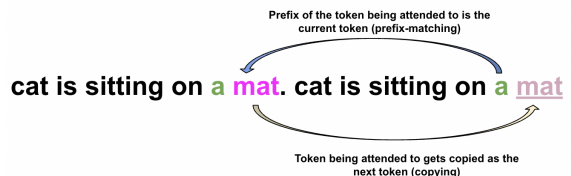


Figure 2: Prefix matching and copying depicted at a given time-step for a repeated sequence of tokens.

2.3 Importance Scores

Consider a model \mathcal{M} and a dataset $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$, where $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$ and $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_L\}$ such that \mathbf{x}_i represents a *prompt* with few (or zero) in-context training examples along with a query input and \mathbf{y}_i represents the corresponding target output sequence. We define and compute importance scores for model components using such datasets to quantify their relative contributions to the model's ability to perform in-context learning.

2.3.1 Oracle

Let $\mathcal{P}_{\mathcal{M}}(\mathcal{D})$ denote a dataset/task-specific performance metric, e.g., accuracy. Given dataset \mathcal{D} , the oracle importance score of a component \mathcal{C} in \mathcal{M} is

computed as follows:

$$IS_{\mathcal{C}}(\mathcal{D}) = \mathcal{P}_{\mathcal{M}}(\mathcal{D}) - \mathcal{P}_{\mathcal{M}_{\setminus \mathcal{C}}}(\mathcal{D}) \quad (7)$$

where $\mathcal{M}_{\setminus \mathcal{C}}$ denotes the resultant model when \mathcal{C} is pruned from \mathcal{M} . Clearly, if pruning a component leads to poor model performance on the task, it must be important for the task. Similarly, if there is no difference or an improvement in performance upon pruning a component, it must be unimportant. Computing oracle importance scores for K model components requires us to perform $\mathcal{O}(K)$ evaluations for each dataset \mathcal{D} .

2.3.2 Gradient-based

Given dataset \mathcal{D} , the gradient-based importance score (Molchanov et al., 2016; Michel et al., 2019) of an attention head h captures the expected sensitivity of the model to h and is computed as follows:

$$IS_h(\mathcal{D}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y})} \left| A^h([\mathbf{x}; \mathbf{y}])^T \frac{\partial \mathcal{L}(\mathbf{y}|\mathbf{x})}{\partial A^h([\mathbf{x}; \mathbf{y}])} \right| \quad (8)$$

where $;$ is the concatenation operator, $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$ such that \mathbf{x} is a sequence of T_x tokens $x_{1:T_x}$, \mathbf{y} is a sequence of T_y tokens $y_{1:T_y}$, A^h is the output of head h defined in (4) and the loss term in (8) is computed using the auto-regressive decomposition of the log-likelihood:

$$\mathcal{L}(\mathbf{y}|\mathbf{x}) = -\frac{1}{T_y} \sum_{j=1}^{j=T_y} \log(p(y_j|\mathbf{x}, y_{1:j-1})) \quad (9)$$

These importance scores can be efficiently computed for all heads by simply performing a single forward and backward pass over the model with \mathcal{D} .

We also define the aggregated importance score of an attention head on a set of datasets $\mathbb{S} = \{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ as follows:

$$IS_h(\mathbb{S}) = \mathbb{E}_{\mathcal{D} \sim \mathbb{S}} [IS_h(\mathcal{D})] \quad (10)$$

3 Experimental Setup

We conducted our experiments on OPT-66B, which was the largest publicly available dense decoder-only language model at the time of our experiments. We efficiently compute gradient-based importance scores for the 4608 attention heads and oracle importance scores for the 64 feed forward networks (FFNs) in OPT-66B. We experiment with a variety of 14 NLP datasets/tasks. For consistency in the evaluation metric, we report

accuracy on all tasks. Our choice of datasets and metric is in line with Zhang et al. (2022). The datasets include ARC Easy and Challenge (Clark et al., 2018) and OpenBookQA (Mihaylov et al., 2018) for advanced question-answering, HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020) and Winogrande (Sakaguchi et al., 2021) for various forms of commonsense reasoning, and the following datasets from the standard SuperGLUE benchmark (Wang et al., 2019): BoolQ, CB, COPA, MultiRC, ReCoRD, RTE, WiC, and WSC. For a subset of experiments involving evaluation of out-of-distribution generalization, we also use 2 additional datasets: MathQA (Amini et al., 2019) and LAMBADA (Paperno et al., 2016). We use a modified version of the *lm-evaluation-harness* framework (Gao et al., 2021) for our experiments. The default framework samples in-context examples at random, which we use without modification.

4 Importance Scores for OPT-66B

4.1 Attention Heads

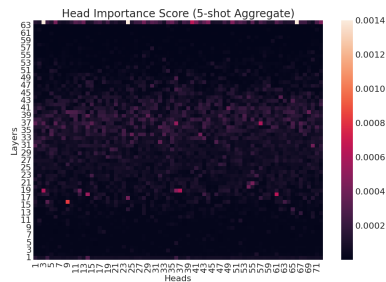


Figure 3: Attention head aggregate importance score heatmap for 5-shot in-context learning with OPT-66B.

Figure 3 depicts a heatmap of the head importance scores averaged across all tasks (as described in §2.3.2) in the 5-shot setting. Task-averaged heatmaps for the 0-shot and 1-shot settings and all task-specific heatmaps are provided in Appendix A.1. We observe that the important attention heads are primarily clustered in the intermediate layers of OPT-66B in both the task-averaged and task-specific cases. We also observe overlap in the important heads across the different zero/few-shot settings, confirmed in follow-up analysis in §6.2.

4.2 Feed Forward Networks

We compute oracle importance scores (both task-specific and averaged across tasks) for each FFN as described in §2.3.1 in the zero/few-shot settings.

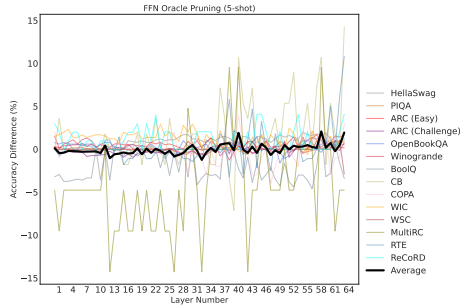


Figure 4: Feed forward network (FFN) oracle importance scores for 5-shot in-context learning with OPT-66B computed by knocking off FFNs independently, i.e., the curves are discrete and **not** cumulative.

Figure 4 depicts the scores in the 5-shot setting. The 0/1-shot scores are depicted in Appendix A.2.

We observe in the 0/1-shot settings that the removal of any FFN in the early (1-30) layers of OPT-66B either gives comparable or better performance for a vast majority of tasks. In the 5-shot setting however, both the early and later layers seem to have important FFNs for most tasks. We also generally observe high variance in FFN importance scores in later layers. We particularly note high variance for WSC and MultiRC, observing that removal of some individual FFNs can lead to absolute accuracy improvements/degradation of up to 20%! We leave further investigation into the cause for this variance for future work.

5 Iterative Pruning

We now assess to what extent we can remove *multiple* attention heads and/or FFNs with minimal decline in task performance. For each task in each (0/1/5-shot) in-context learning setting, we create separate rankings of attention heads and FFNs in OPT-66B by separately sorting them in ascending order by their importance scores (§4.1 and §4.2). We then remove unimportant attention heads or FFNs in an iterative fashion using these rankings, 10% at a time, and re-evaluate task performance after each removal.¹

5.1 Removing Attention Heads

Figure 5 depicts the resulting task-specific and task-averaged accuracy trends in the 5-shot setting. Corresponding 0/1-shot trends are depicted in Appendix A.3. We observe that the average accuracy across tasks does not change much up until

¹We do not remove attention heads one at a time and re-evaluate given the number of heads and evaluation cost.

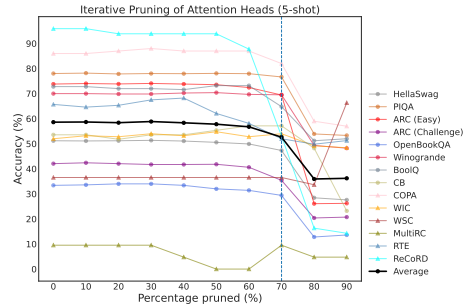


Figure 5: Effect on 5-shot accuracy when removing attention heads in OPT-66B in an iterative manner based on task-specific and shot-specific importance scores.

~70% of the attention heads are removed. A fine-grained look at the individual tasks also mostly shows similar trends, with accuracy staying fairly intact until a large proportion of the heads are removed. Some oddities include tasks such as WSC and CB, wherein we see that the 0-shot accuracy actually increases after removal of 70% of the heads.

5.2 Removing FFNs

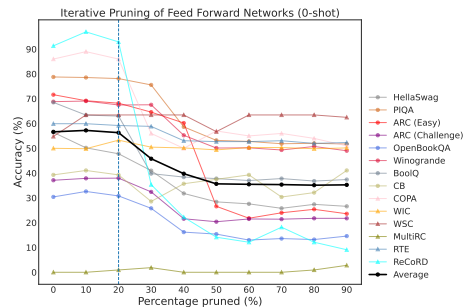


Figure 6: Effect on 0-shot accuracy when removing feed forward networks (FFNs) in OPT-66B in an iterative manner based on task-specific and shot-specific importance scores.

Figure 6 depicts the resulting task-specific and task-averaged accuracy trends in the 0-shot setting. Corresponding 1/5-shot trends are depicted in Appendix A.4. We observe that in the 0-shot setting, the average accuracy across tasks does not change up until ~20% of the FFNs are removed. For some tasks such as PIQA, Winogrande and RTE, the accuracy does not change even if 30% of the FFNs (~13B of the 66B parameters) are removed. We also observe that the inflection point after which we observe a sharp decline in accuracy changes to 10% for the few-shot settings. Overall, these observations indicate that FFNs play a critical role toward in-context learning.

5.3 Combined Removal of Heads & FFNs

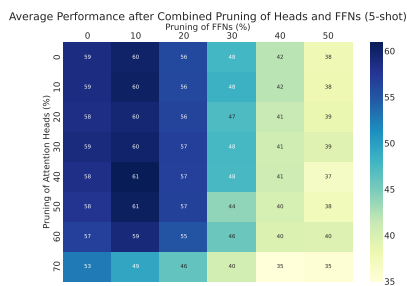


Figure 7: Effect on average 5-shot accuracy when removing **both** attention heads and feed forward networks (FFNs) in OPT-66B in an iterative manner based on shot-specific task-aggregate importance scores.

We now investigate whether the inflection points to in-context learning performance when removing either attention heads or FFNs in an iterative fashion still hold when removing them in *tandem*. Figure 7 depicts the average 5-shot accuracy of all tasks on joint iterative removal of attention heads and FFNs. Corresponding 0/1-shot trends are depicted in Appendix A.5. We observe that the removal of 70% of the attention heads ($\sim 15.7B$ parameters) and 20% of the FFNs ($\sim 8.5B$ parameters) leads to a mere 5% absolute drop in the average 0-shot accuracy. In the 1-shot setting, the drop in accuracy is 6% on removing 70% of the attention heads and 10% of the FFNs. In the 5-shot setting, the drop in accuracy is 4% on removing 60% of the attention heads and 20% of the FFNs. Overall, these new inflection points have deviated by at most 10% absolute, which may be attributed to the interplay between heads and FFNs.

6 Detailed Analysis of Attention Heads

In this section, we perform a detailed analysis of the attention heads in OPT-66B, given that in-context learning is auto-regressive in nature and attention heads explicitly encode cross-token interactions.

6.1 Cross-Task Analysis

Michel et al. (2019) found preliminary empirical evidence of the existence of "universally" important attention heads in trained task-specific Transformer and BERT models via evaluating on out-of-domain test sets for machine translation and natural language inference respectively. With similar motivation, we study if the (un)important attention heads identified in various in-context learning settings for OPT-66B are shared across tasks.

6.1.1 Spearman's Rank Correlation

We assess overlap in (un)important attention heads across tasks by sorting task-specific head importance scores to get head importance rankings and computing the Spearman's rank correlation coefficient (SRCC) between the rankings for every pair of tasks in the zero-shot and few-shot settings. We also sort the task-aggregate head importance scores to get the aggregate ranking and compute the SRCC against the ranking for every constituent task. All correlations are depicted in Figure 8 for the 5-shot setting and Appendix A.6 for the 0/1-shot settings.

In both zero and few-shot settings, we observe statistically significant ($p < 0.01$) positive correlations in the head importance rankings for every pair of tasks, as well as between every task's ranking and the aggregate ranking. This indicates that the set of (un)important attention heads are clustered together across tasks. We also observe seemingly lower magnitude SRCC values between every task and ReCoRD, a long reading comprehension task which requires commonsense reasoning, indicating the amount of head overlap is proportionally lower.

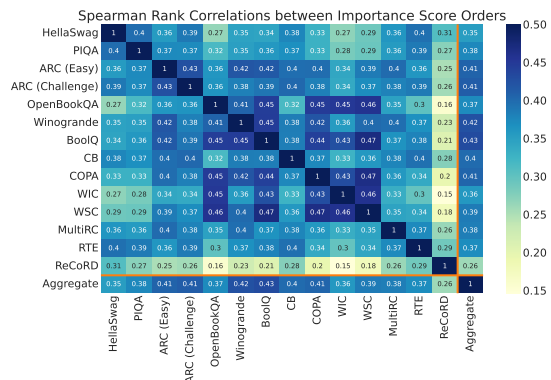


Figure 8: Spearman's rank correlation coefficients between the attention head importance rankings for various tasks in the 5-shot setting. All p -values < 0.01 .

6.1.2 Generalization Trends

To understand how well head importance rankings generalize across tasks, we study accuracy trends for tasks when pruning using various head importance rankings. We study two sets of tasks.

The first set of tasks we study were used to compute the aggregate ranking: COPA, Winogrande and ReCoRD. For each of these 3 tasks, we consider the impact of pruning based on the self-ranking, aggregate ranking and the rankings from the tasks which share the *highest* and *lowest* SRCC with them. Figures 9a, 9b and 9c de-

pict the accuracy trends for these 3 tasks in the 5-shot setting. Corresponding trends in the 0/1-shot settings are in Appendix A.7. In the 0-shot setting, we observe that the accuracy on all 3 tasks when pruning using the rankings described is almost unaffected up to the 50% mark. We then observe a sharp decline in accuracy on COPA and Winogrande when the model is pruned to the 70% mark using the ranking identified via ReCoRD, the task with the lowest SRCC (0.13) with both COPA and Winogrande. This indicates that even if the rankings vary between ReCoRD and COPA/Winogrande (as reflected in the low magnitude of the SRCC score), the set of attention heads important for 0-shot learning with ReCoRD are important for COPA/Winogrande too. To further verify this, we calculated and found 71% and 76% overlap between the top 30% important attention heads for ReCoRD-COPA and ReCoRD-Winogrande respectively. Comparing the zero-shot setting against the few-shot settings, we note that the decline/divergence in accuracy beyond the 50% pruning mark using the ReCoRD ranking is less sharp for COPA and Winogrande in the 1-shot setting and fades away in the 5-shot setting, indicating a convergence of important heads across tasks.

The second set of tasks we study are unseen, i.e., not used to compute the aggregate ranking: MathQA and LAMBADA. For these tasks, we analyze accuracy trends when pruning using the self-ranking and aggregate ranking. Figures 9d and 9e depict their accuracy trends in the 5-shot setting. Corresponding trends in the 0/1-shot settings are in Appendix A.7. As expected, we observe that the self-ranking accuracy curves are somewhat higher than the aggregate ranking accuracy curves in general across both tasks. For MathQA, we also observe that the absolute difference in accuracy for both cases is within 1-2%. These indicate that the aggregate rankings generalize well to MathQA but not as much to LAMBADA.

6.2 Cross-Shot Analysis

To see if the attention heads identified to be (un)important for a task are shared across the different zero and few-shot settings, we compute Spearman’s rank correlation coefficient (SRCC) between the cross-shot head importance rankings for each task and compute the mean and variance across all 14 tasks. We observe that the mean SRCC is higher for rankings *within* the few-shot setting (0.41 for

1-shot vs. 5-shot) than for rankings *across* the zero and few-shot settings (0.39 for 0-shot vs. 1-shot and 0.37 for 0-shot vs. 5-shot), with low variance (0.001) and p -value < 0.01 . This matches the intuition that a similar set of heads must be important within the different few-shot settings than across the zero-shot and any of the few-shot settings. However, we also see that the SRCC magnitudes for the latter are not very far off. In totality, these indicate non-trivial overlap in the (un)important attention heads for tasks across shots.

6.3 Induction Heads in OPT-66B

We look for induction heads in OPT-66B by quantifying the capacity of all attention heads to perform prefix matching and copying using random input sequences in a task-agnostic fashion, following the definition and algorithms by Olsson et al. (2022) discussed in §2.2 and Appendix A.8.

Figures 10a and 10b depict the prefix matching and copying score heatmaps respectively for OPT-66B. We observe that a small subset of attention heads in OPT-66B have high prefix matching scores, located in the upper layers (31+) of the model. On the other hand, there are a relatively larger number of attention heads with high copying scores, although the vast majority of these are also located in the upper layers (41+). When seen in conjunction, these observations indicate that there is a sparse set of attention heads that are capable of performing both primitive operations and thus can be deemed plausible induction heads.

6.3.1 Are Induction Heads Important?

We now study whether induction heads (which encode the basic in-context learning primitives of explicit prefix matching and copying) overlap with attention heads identified to be important (and consequently capable of sophisticated and latent behaviors associated with in-context learning) for our chosen downstream tasks.

A qualitative comparison of the heatmaps in Figure 10 against the heatmaps referenced in §4.1 indicates that induction heads do overlap with task-aggregated important attention heads. To better facilitate this comparison, we first formalize the total capacity of a model to perform prefix matching (or copying) to be the sum of the respective scores for individual attention heads in the model. We then investigate how much of this capacity is retained when attention heads are pruned in the order of least important heads first. Figure 11 de-

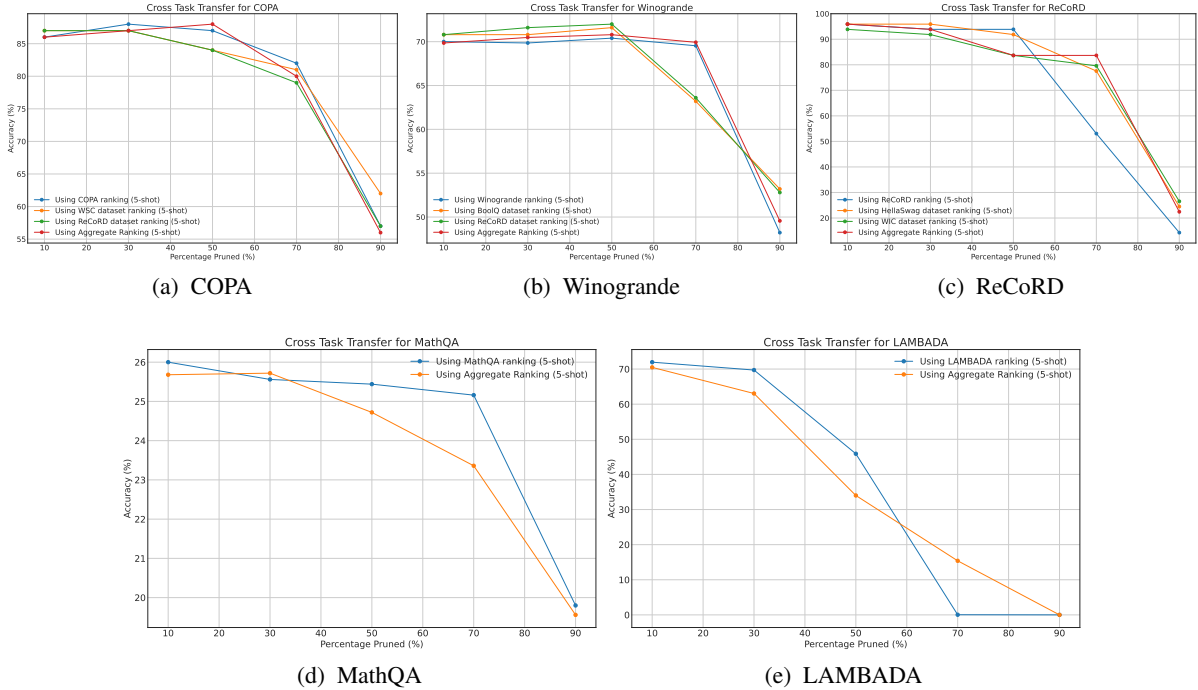


Figure 9: Cross-task transfer of attention head importance rankings as measured by impact of pruning on accuracy in the 5-shot setting.

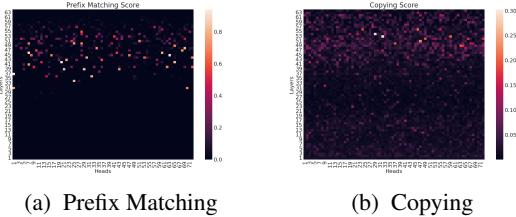


Figure 10: Attention head prefix matching and copying score heatmaps for OPT-66B.

picts this comparison. We observe that much of the total prefix matching score is retained when 20% of the least important heads are removed, with the slope of decline becoming sharp only after the 40% pruning mark. This indicates that unimportant heads also have low prefix matching scores. We also observe that the prefix matching scores are generally higher for heads important for few-shot in-context learning than for heads important for zero-shot learning. On the other hand, we observe across the zero-shot and few-shot settings that the total copying score retained on pruning attention heads rapidly and consistently declines, indicating that even unimportant heads have a non-trivial capacity to perform copying. When seen in conjunction, these observations indicate that induction

heads in OPT-66B are capable of sophisticated behaviors associated with in-context learning popular downstream NLP tasks and reinforce the induction head generality arguments [Olsson et al. \(2022\)](#) make in the context of smaller models with stylized and synthetic tasks. We also provide per-task plots in [Appendix A.9](#) which showcase that some tasks rely on induction heads more than other tasks.

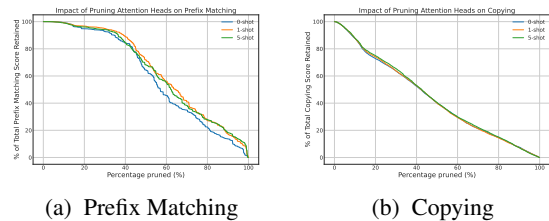


Figure 11: Total inductive (prefix matching / copying) capacity retained as a function of percentage of attention heads pruned, where heads are pruned based on task-aggregate importance score rankings in the order of least important first.

7 Related Work

There has been an interest in effectively *leveraging* the in-context learning paradigm ([Zhao et al., 2021](#); [Holtzman et al., 2021](#); [Min et al., 2022a](#); [Liu et al., 2022](#); [Lu et al., 2022](#); [Rubin et al., 2022](#);

Mishra et al., 2022) ever since its introduction by Brown et al. (2020), but there have been relatively fewer studies toward better *understanding* the paradigm itself. Xie et al. (2021) cast in-context learning as implicit Bayesian inference where the language model implicitly infers a shared concept among in-context examples when making a prediction. Min et al. (2022b) study the role of the in-context examples themselves, finding that the ground-truth labels are not needed in the examples and that the more important drivers are provision of the label space, the distribution of the input text and the overall format of the sequence. Garg et al. (2022) showcase that Transformer models trained from scratch can in-context learn the class of linear functions with performance comparable to the optimal least squares estimator even under distribution shifts. Razeghi et al. (2022) showcase that in-context learning performance is correlated strongly with term frequencies in the pre-training corpora used. Olsson et al. (2022) consider an alternate framing of in-context learning as the ability of a language model to better predict tokens later in the context than tokens earlier and hypothesize the existence of induction heads that are responsible for in-context learning. Chan et al. (2022) show that Transformers exhibit striking differences in generalizing from in-context vs. in-weights information.

Several works have also focused on analyzing and interpreting how attention works. Vig and Belinkov (2019) performed a study on GPT-2, finding that attention targets different parts of speech at different layer depths and aligns with dependency relations most strongly in the middle layers. Tenney et al. (2019) showcase that BERT encodes the classical NLP pipeline in an interpretable way across layers. There are works relying on different formulations for head importance, such as layer-wise relevance propagation (Voita et al., 2019), gradient-based importance and oracle knock-off importance (Michel et al., 2019), with small task-specific trained models and report the existence of specialized heads. Given the recent trend of increasing model scale (Lieber et al., 2021; Chowdhery et al., 2022; Smith et al., 2022; Rae et al., 2021) toward tuning-free general-purpose language models that exhibit emergent in-context learning abilities, we draw and build on prior work to understand just how much scale is really needed and/or used for in-context learning *downstream*, an aspect somewhat eclipsed by the focus on the pre-training loss

curve in scaling laws (Hoffmann et al., 2022). It is also worth noting that some of our empirical observations rely on a simple greedy approach to training-free pruning since our focus was not to optimally prune a language model with respect to performing in-context learning. Li et al. (2021) show the greedy approach is sub-optimal and produces under-estimates and Halabi et al. (2022) account for the need to re-compute importance scores after removal of each attention head or FFN by formulating pruning as weakly sub-modular maximization.

8 Conclusion & Future Work

In this paper, we studied the efficacy of attention heads and feed forward networks (FFNs) in a large language model (OPT-66B) in performing in-context learning in both task-specific and task-agnostic settings. We observed that while in-context learning may have emerged via self-supervised pre-training at scale, only a core nucleus of attention heads and FFNs seem to be important for in-context learning across a wide variety of downstream tasks. We observed that a small set of attention heads have the capacity to perform task-agnostic primitive induction operations associated with in-context learning, namely, prefix matching and copying. We also saw that these induction heads overlap with task-specific important attention heads, indicating that induction heads are capable of more sophisticated forms of in-context learning and reinforcing arguments (Olsson et al., 2022) about their generality. Overall, our in-context learning-centric observations complement recent work (Hoffmann et al., 2022) in indicating that large language models may be under-trained and motivate several interesting directions for future work. While induction heads are formed naturally during self-supervised pre-training in its current form, we believe it may be possible to increase the number and strength of induction heads formed by defining auxiliary pre-training objectives for primitives like prefix matching and copying. More generally, it may also be prudent to investigate and improve (pre-)training regimes to increase the number of important model components to in-context learn-perform a wide variety of downstream tasks. Multi-task instruction-tuning likely belongs to this category and it would be interesting to replicate our study with now increasingly accessible instruction-tuned model variants (such as OPT’s instruction meta-learned variant OPT-IML).

9 Limitations

Our work is a comprehensive empirical study of a popular large language model’s capacity to perform in-context learning, relying on both task-specific (via a wide variety of challenging and practically relevant downstream tasks) and task-agnostic (via looking for induction heads) analyses and connecting the two via correlation/overlap investigations. We do not claim a causal link, i.e., we do not claim that an attention head that acquires the capacity to be an induction head will become capable of more sophisticated in-context learning associated with our downstream tasks. Making this claim will require a more deeper investigation that is outside the scope of this paper. We also do not fully understand *why* most attention heads seem to be unimportant for in-context learning and *why* there is an overlap in (un)important attention heads across tasks and shots, which warrant further investigation. Other more obvious limitations to our work include our use of only up to 5 in-context examples, random selection of in-context examples for a query input and our choice of all monolingual downstream tasks.

10 Impact Statement

The findings in our work have significant implications for the design, development and deployment of large language models, known to have a very high carbon footprint as well as training and inference costs. Having identified that a core nucleus of model parameters seem to be important for in-context learning, it may be possible to reduce these models’ carbon footprint and mitigate these costs. Our findings provide architectural transparency and may also be helpful in identifying targeted improvements for downstream tasks as well as for more broader facets such as bias and fairness.

References

Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*.

Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. 2017. Structured pruning of deep convolutional neural networks. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 13(3):1–18.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the*

AAAI conference on artificial intelligence, volume 34, pages 7432–7439.

- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv:2204.06745*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Stephanie CY Chan, Ishita Dasgupta, Junkyung Kim, Dharshan Kumaran, Andrew K Lampinen, and Felix Hill. 2022. Transformers generalize differently from information stored in context vs in weights. *arXiv preprint arXiv:2210.05675*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. A framework for few-shot language model evaluation.
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. 2022. What can transformers learn in-context? a case study of simple function classes. *arXiv preprint arXiv:2208.01066*.
- Marwa El Halabi, Suraj Srinivas, and Simon Lacoste-Julien. 2022. Data-efficient structured pruning via submodular optimization. *arXiv preprint arXiv:2203.04940*.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks,

- Johannes Welbl, Aidan Clark, et al. 2022. [Training compute-optimal large language models](#). *arXiv preprint arXiv:2203.15556*.
- Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. [Surface form competition: Why the highest probability answer isn't always right](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7038–7051, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. 2016. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*.
- Jiaoda Li, Ryan Cotterell, and Mrinmaya Sachan. 2021. [Differentiable subset pruning of transformer heads](#). *Transactions of the Association for Computational Linguistics*, 9:1442–1459.
- Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. 2021. [Jurassic-1: Technical details and evaluation](#). *White Paper. AI21 Labs*.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. [Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. [Are sixteen heads really better than one?](#) *Advances in neural information processing systems*, 32.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. [Can a suit of armor conduct electricity? a new dataset for open book question answering](#). *arXiv preprint arXiv:1809.02789*.
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022a. [Noisy channel language model prompting for few-shot text classification](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5316–5330, Dublin, Ireland. Association for Computational Linguistics.
- Sewon Min, Xinxin Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022b. [Rethinking the role of demonstrations: What makes in-context learning work?](#) *arXiv preprint arXiv:2202.12837*.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2022. [Reframing instructional prompts to GPTk's language](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 589–612, Dublin, Ireland. Association for Computational Linguistics.
- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. 2016. [Pruning convolutional neural networks for resource efficient inference](#). *arXiv preprint arXiv:1611.06440*.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2022. [In-context learning and induction heads](#). *Transformer Circuits Thread*.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The lambda dataset: Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. [Scaling language models: Methods, analysis & insights from training gopher](#). *arXiv preprint arXiv:2112.11446*.
- Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. 2022. Impact of pretraining term frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. [Learning to retrieve prompts for in-context learning](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. [Winogrande: An adversarial winograd schema challenge at scale](#). *Communications of the ACM*, 64(9):99–106.
- Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhunoye, George Zerveas, Vijay Korthikanti, et al. 2022. [Using deep-speed and megatron to train megatron-turing nl](#)

530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*.

Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. [BERT rediscovers the classical NLP pipeline](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). *Advances in neural information processing systems*, 30.

Jesse Vig and Yonatan Belinkov. 2019. Analyzing the structure of attention in a transformer language model. *arXiv preprint arXiv:1906.04284*.

Elena Voita, David Talbot, Fedor Moiseev, Rico Senrich, and Ivan Titov. 2019. [Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned](#). *arXiv preprint arXiv:1905.09418*.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). *Advances in neural information processing systems*, 32.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.

Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [Hellaswag: Can a machine really finish your sentence?](#) *arXiv preprint arXiv:1905.07830*.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. [Opt: Open pre-trained transformer language models](#). *arXiv preprint arXiv:2205.01068*.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR.

A Appendix

A.1 Head Importance Scores

Figure 12 depicts the attention head aggregate importance score heatmaps in the 0-shot and 1-shot

settings. Figures 14, 15 and 16 depict the attention head importance scores for each task in the 0-shot, 1-shot and 5-shot settings respectively.

A.2 FFN Importance Scores

Figure 13 depicts the task-specific and task-averaged importance scores for feed forward networks in the 0-shot and 1-shot settings.

A.3 Removing Attention Heads

Figure 17 depicts the task-specific and task-averaged accuracy trends on iterative removal of attention heads in the order of least important first in the 0-shot and 1-shot settings.

A.4 Removing FFNs

Figure 18 depicts the task-specific and task-averaged accuracy trends on iterative removal of feed forward networks in the order of least important first in the 1-shot and 5-shot settings.

A.5 Combined Removal of Heads & FFNs

Figure 19 depicts the average accuracy of all tasks on joint iterative removal of attention heads and feed forward networks in the order of least important first in the 0-shot and 1-shot settings.

A.6 Cross-Task Analysis: Spearman’s Rank Correlation

Figure 20 depicts the Spearman’s rank correlation coefficients (SRCC) between the attention head importance rankings for every pair of tasks in the 0-shot and 1-shot settings. It also depicts the SRCC between the aggregate ranking and the ranking for each constituent task.

A.7 Cross-Task Analysis: Generalization Trends

Figures 21 and 22 depict the cross-task head importance ranking generalization plots in the 0-shot and 1-shot settings.

A.8 Details of Prefix Matching and Copying Scores

Algorithms 1 and 2 contain pseudo-code to compute prefix matching and copying scores respectively for each attention head in OPT-66B. We follow the approach described by Olsson et al. (2022), but instead of computing scores using 10 sequences with fixed length of 25, we compute these scores using 100 sequences with varying lengths to account for OPT-66B’s large maximum sequence length.

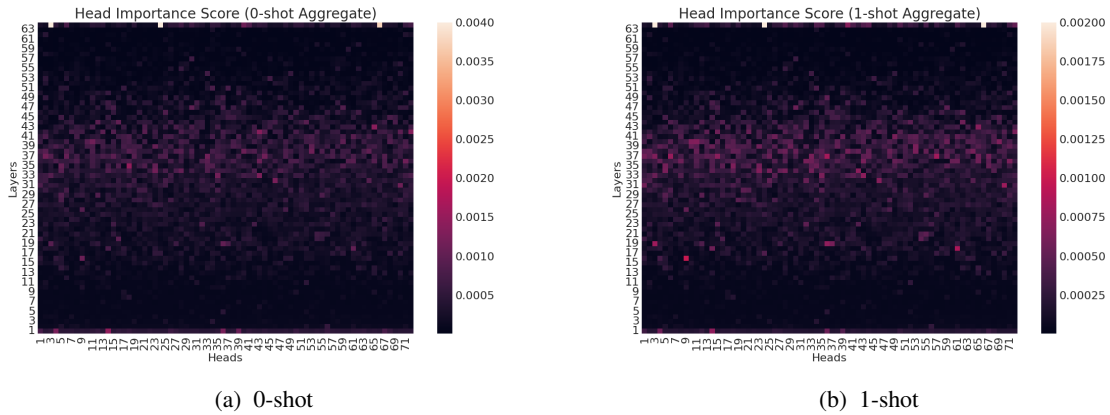


Figure 12: Attention head aggregate importance score heatmaps for 0/1-shot in-context learning with OPT-66B.

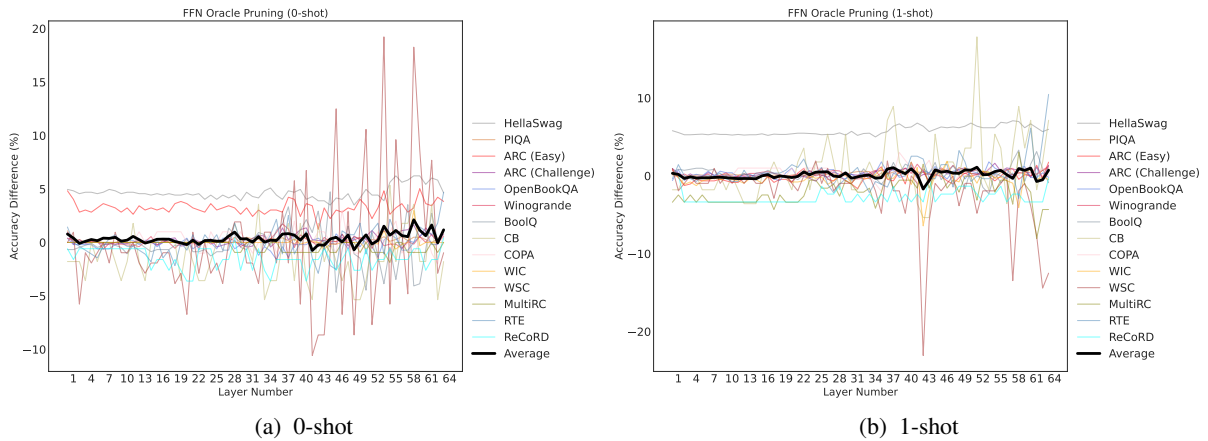


Figure 13: Feed forward network (FFN) oracle importance scores for 0/1-shot in-context learning with OPT-66B. Each FFN is knocked off independently to compute these scores, i.e., the curves are discrete and **not** cumulative.

As in [Olsson et al. \(2022\)](#), we exclude a small fraction of the most and least common tokens from the model’s vocabulary and randomly sample tokens for these sequences to strip out the effects of pre-training corpora memorization from our scores and inductive behavior analyses.

For prefix matching, the high-level approach is the following: take a random sequence, repeat it 4 times, perform a forward pass and then for each head, compute the attention pattern and take the average of all attention pattern entries attending from a given token back to tokens that succeeded the same token in earlier repeats.

For copying, the high-level approach is the following: take a random sequence, directly feed the sequence through each head and compute the contribution of the head to the output logits, and then measure how much the head increased the logit of the maximally attended to token over increasing

the logits of other attendable tokens at each time-step. Unlike [Olsson et al. \(2022\)](#), we do not scale the raw scores to be in the range of -1 to 1.

A.9 Importance of Induction Heads to Each Task

Figures 23 and 24 showcase the importance of induction heads to each task via measuring the percentage of the total prefix matching and copying capacities retained as a function of percentage of attention heads pruned, where heads are pruned based on each task’s head importance ranking for each in-context learning setting (zero-shot, one-shot and five-shot) in the order of least important first. A small initial slope of decline implies that unimportant heads also have low prefix matching or copying scores while a steep initial slope of decline implies unimportant heads also have high prefix matching or copying scores. We observe differ-

ences in the slopes of decline across different tasks, with tasks like HellaSwag and ReCoRD (which have high accuracies in Figure 5) having smaller initial slopes than a task like OpenBookQA (which has relatively lower accuracy in Figure 5). When seen in conjunction, these plots not only point to the generality of induction heads to more sophisticated behaviors associated with in-context learning but also indicate that some tasks rely on induction heads more than others.

Algorithm 1 Prefix Matching Scores for Attention Heads

Arguments: Model \mathcal{M}

```
model  $\leftarrow$  Pretrained( $\mathcal{M}$ )
layers, heads  $\leftarrow$  model.num_layers, model.num_heads_per_layer
ranked_vocab_list  $\leftarrow$  model.tokenizer.vocab  $\triangleright$  term-frequency based vocabulary list of model
exclude_vocab_size  $\leftarrow$   $0.04 \times \text{len}(\textit{ranked\_vocab\_list})$   $\triangleright$  remove 4% most & least common tokens
ranked_vocab_list  $\leftarrow$  ranked_vocab_list[exclude_vocab_size : -exclude_vocab_size]
prefix_matching  $\leftarrow$  []
for seed in {1  $\cdots$  100} do
    L  $\leftarrow$   $2 \times \textit{seed} + 23$   $\triangleright$  ensure  $4L \in [100, 892]$ 
    X  $\leftarrow$  random.choice(ranked_vocab_list, size = L, seed = seed, replace = False)  $\triangleright$  L length
    random sequence with all unique tokens
    X  $\leftarrow$  repeat(X, 4)  $\triangleright$  Repeat it four times
    Y  $\leftarrow$  model.forward(X)  $\triangleright$  Forward pass the repeated sequence
    score  $\leftarrow$  zeros(layers, heads)  $\triangleright$  Zero matrix of shape layers  $\times$  heads
    for layer in layers do
        for head in heads do
            att  $\leftarrow$  model.get_attention(layer, head)  $\triangleright$  Shape:  $4L \times 4L$ 
            for token in { $L + 1 \cdots 4L$ } do  $\triangleright$  Repetition starts from token  $L + 1$ 
                att_token  $\leftarrow$  att[token]  $\triangleright$  Shape:  $4L$ 
                for every prev_token == token do  $\triangleright$  Look at the previous repetitions of the token
                    prefix_score = att_token[prev_token + 1]  $\triangleright$  Attention given to token whose
                    prefix is current token
                    score[layer][head]  $\leftarrow$  score[layer][head] + prefix_score
                end for
            end for
            score[layer][head]  $\leftarrow$  score[layer][head]/ $3L$   $\triangleright$  Normalizing by length of for loop
        end for
    end for
    prefix_matching.append(score)  $\triangleright$  Prefix matching scores via one randomly generated example
end for
prefix_matching  $\leftarrow$  average(prefix_matching)  $\triangleright$  Attention head-wise average over all examples
return prefix_matching
```

Algorithm 2 Copying Scores for Attention Heads

Arguments: Model \mathcal{M}

Definitions: Dimension per Head D , Vocabulary Size V

```
model ← Pretrained( $\mathcal{M}$ )
layers, heads ← model.num_layers, model.num_heads_per_layer
ranked_vocab_list ← model.tokenizer.vocab ▷ term-frequency based vocabulary list of model
exclude_vocab_size ←  $0.04 \times \text{len}(\text{ranked\_vocab\_list})$  ▷ remove 4% most & least common tokens
ranked_vocab_list ← ranked_vocab_list[exclude_vocab_size : -exclude_vocab_size]
copying_score ← []
for seed in {1 ··· 100} do
    L ←  $4 \times (2 \times \text{seed} + 23)$  ▷  $L \in [100, 892]$ 
    X ← random.choice(ranked_vocab_list, size = L, seed = seed, replace = False) ▷ L length
    random sequence with all unique tokens
    score ← zeros(layers, heads) ▷ Zero matrix of shape layers × heads
    for layer in layers do
        for head in heads do
            attn_layer_head ← model.get_attention_head(layer, head)
            out ← attn_layer_head(X) ▷ Shape:  $L \times D$ 
            attention ← model.get_attention(layer, head) ▷ Shape:  $L \times L$ 
            logits ← model.hidden_to_vocab(out) ▷ Shape:  $L \times V$ 
            logits ← softmax(logits, dim = 1)
            for token in {1 ··· L} do
                max_ind ← argmax(attention[token]) ▷ Index of the token being max attended to
                attendable_input ← X[1 : token] ▷ Attendable input tokens
                attendable_logits ← logits[token][attendable_input] ▷ Logits of attendable tokens
                mean_of_logits ← average(attendable_logits)
                raised_logits ← attendable_logits - mean_of_logits
                relu_raised_logits ← ReLU(raised_logits) ▷ Computing raise in logit values
                relu_raised_logit_max_ind ← relu_raised_logits[X[max_ind]]
                temp_score ← relu_raised_logit_max_ind / sum(relu_raised_logits)
                score[layer][head] ← score[layer][head] + temp_score
            end for
            score[layer][head] ← score[layer][head] / L ▷ Normalizing by length of for loop
        end for
    end for
    copying_score.append(score) ▷ Copying scores via one randomly generated example
end for
copying_score ← average(copying_score) ▷ Attention head-wise average over all examples
return copying_score
```

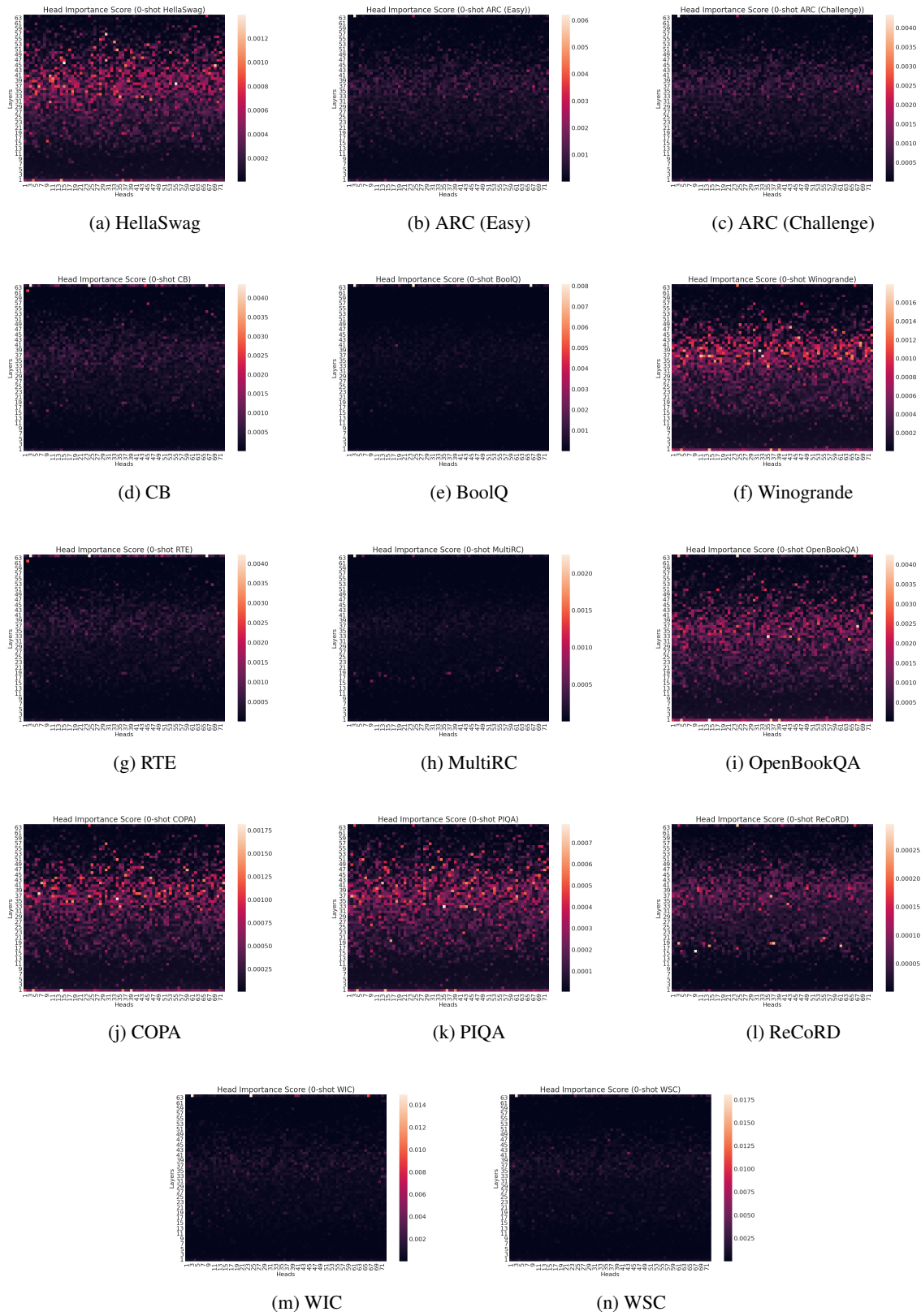


Figure 14: Attention head importance score heatmaps for 0-shot in-context learning with OPT-66B for each task.

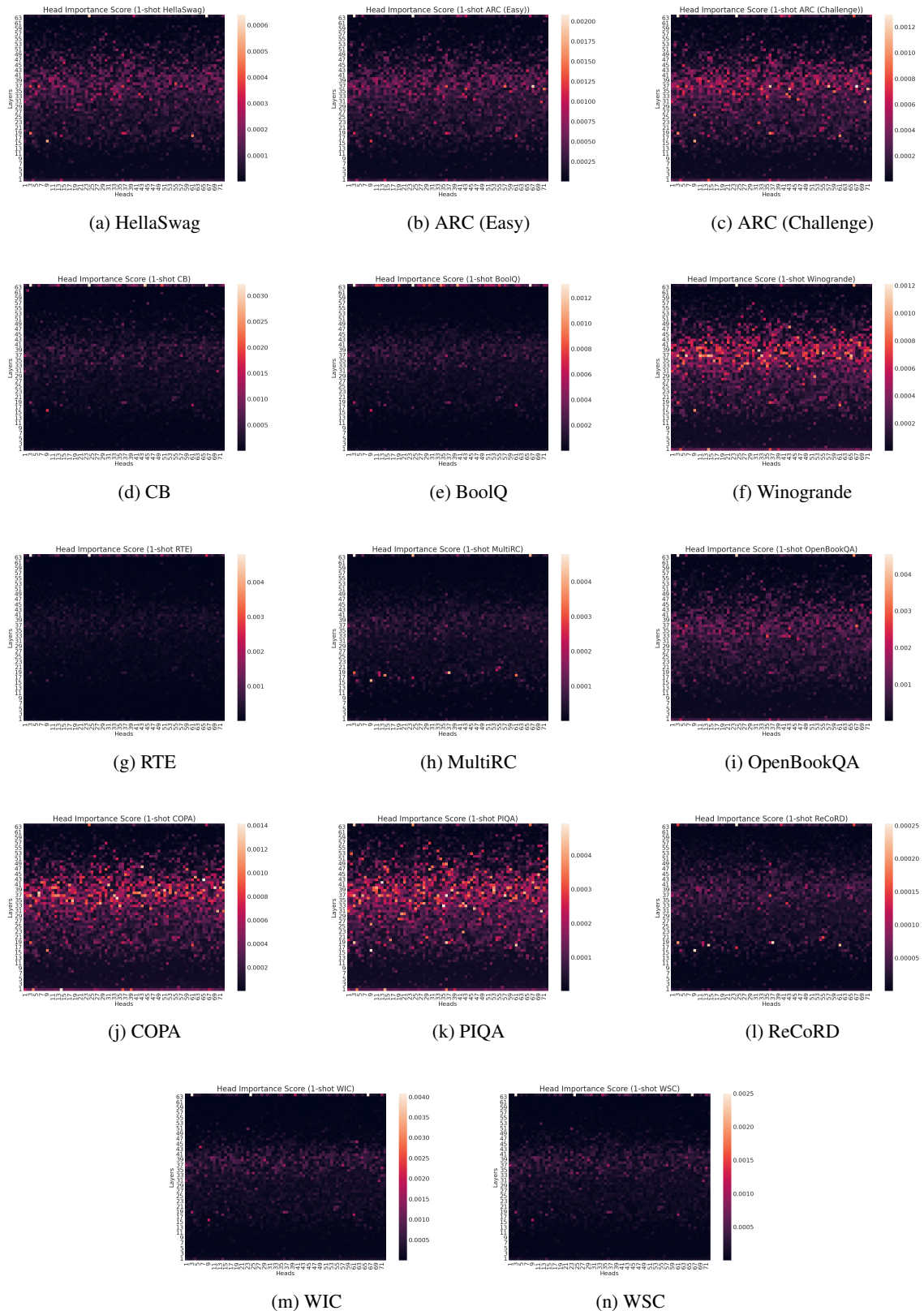


Figure 15: Attention head importance score heatmaps for 1-shot in-context learning with OPT-66B for each task.

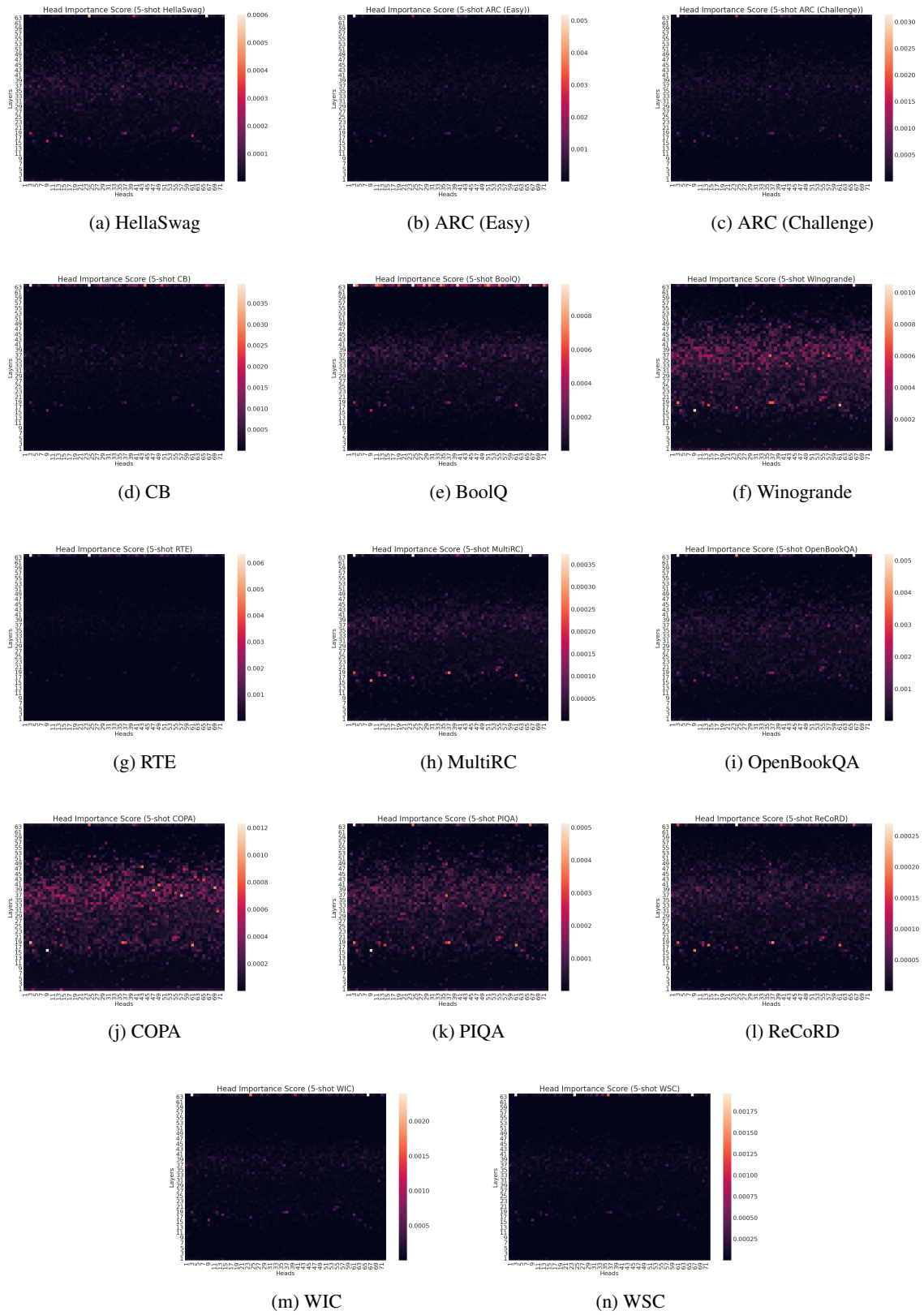
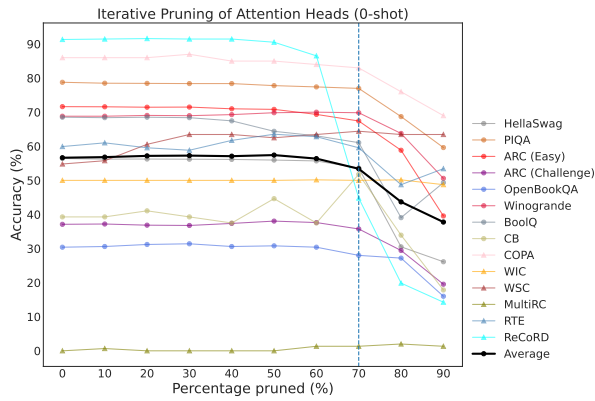
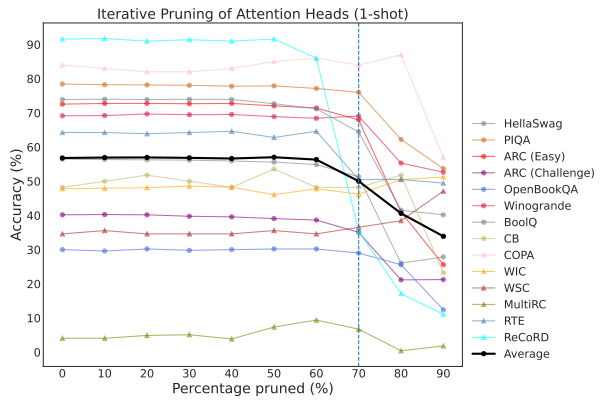


Figure 16: Attention head importance score heatmaps for 5-shot in-context learning with OPT-66B for each task.

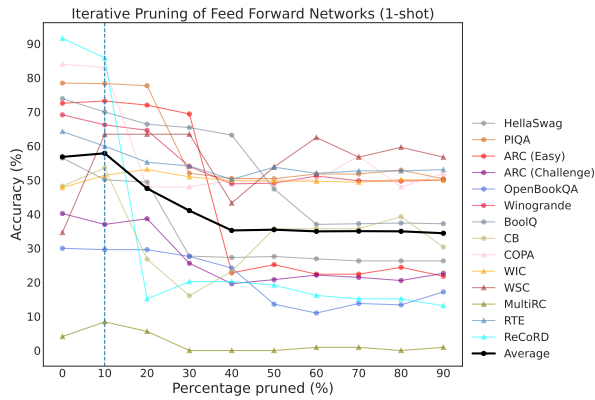


(a) 0-shot

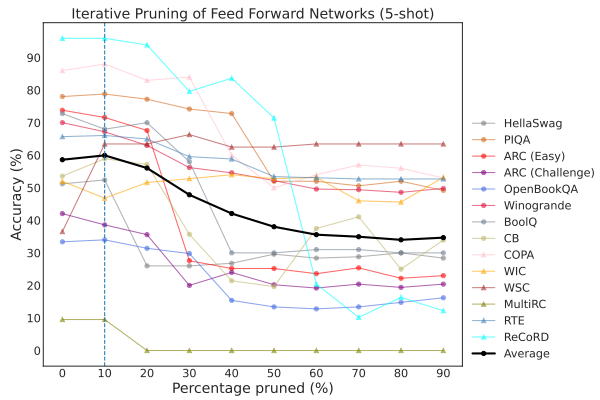


(b) 1-shot

Figure 17: Effect on 0/1-shot accuracy when removing attention heads in OPT-66B in an iterative manner based on task-specific and shot-specific importance scores.

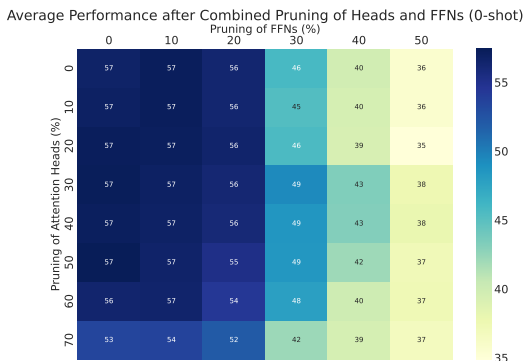


(a) 1-shot

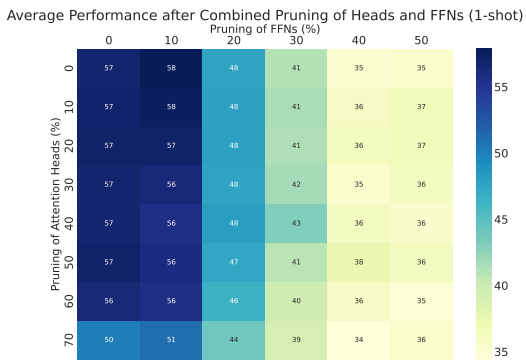


(b) 5-shot

Figure 18: Effect on 1/5-shot accuracy when removing feed forward networks (FFNs) in OPT-66B in an iterative manner based on task-specific and shot-specific importance scores.



(a) 0-shot



(b) 1-shot

Figure 19: Effect on average in-context learning accuracy when removing **both** attention heads and feed forward networks (FFNs) in OPT-66B in an iterative manner based on shot-specific task-aggregate importance scores.

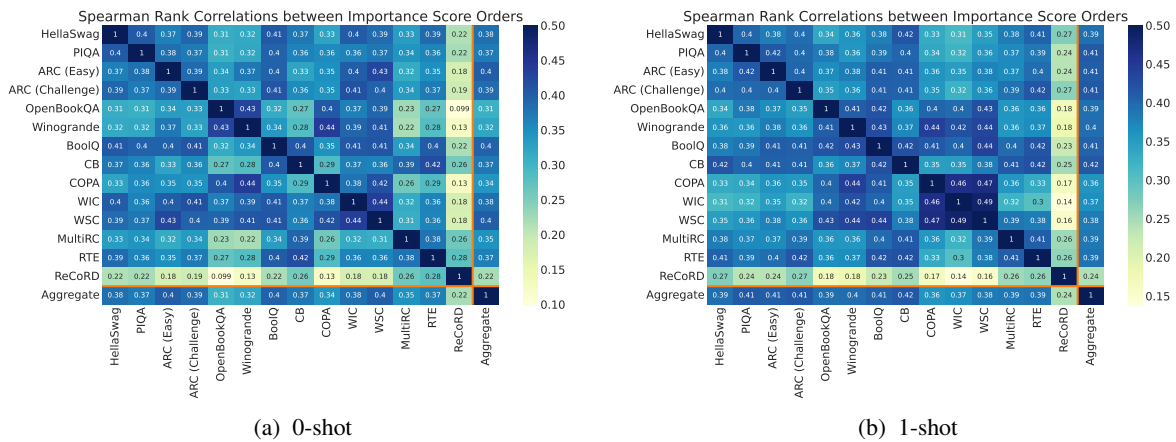


Figure 20: Spearman’s rank correlation coefficients between the attention head importance rankings for various tasks in the 0-shot and 1-shot settings. All p -values < 0.01 .

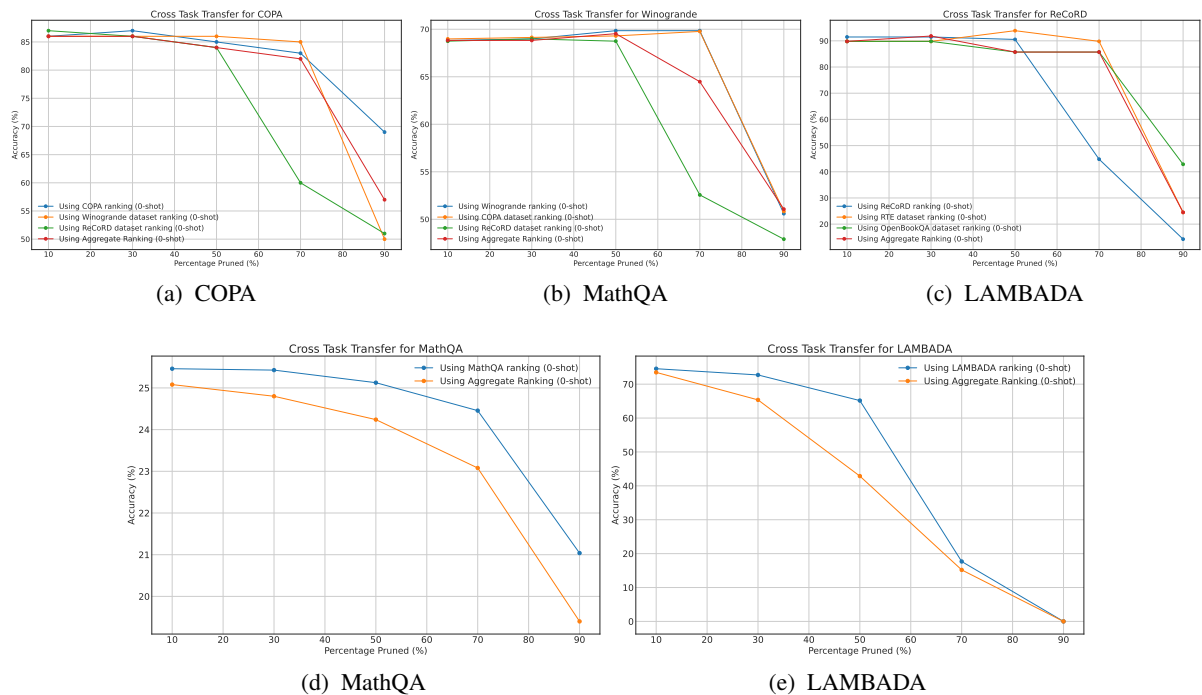


Figure 21: Cross-task transfer of attention head importance rankings as measured by impact of pruning on accuracy in the 0-shot setting.

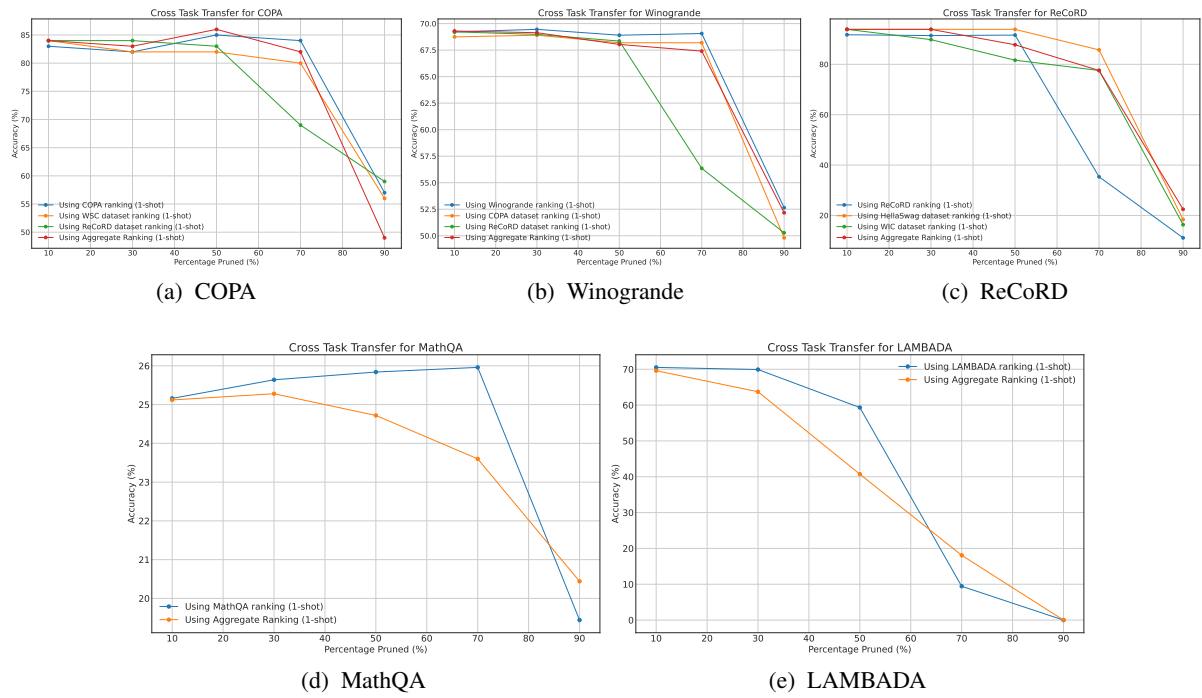


Figure 22: Cross-task transfer of attention head importance rankings as measured by impact of pruning on accuracy in the 1-shot setting.

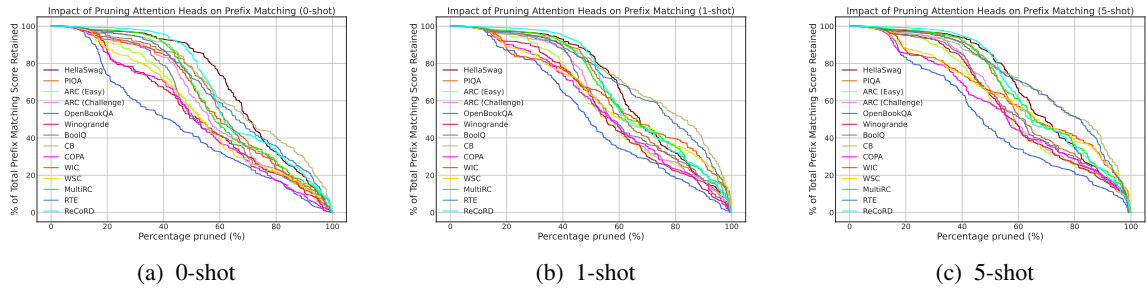


Figure 23: Total prefix matching capacity retained as a function of percentage of attention heads pruned, where heads are pruned based on task-specific importance score rankings in the order of least important first.

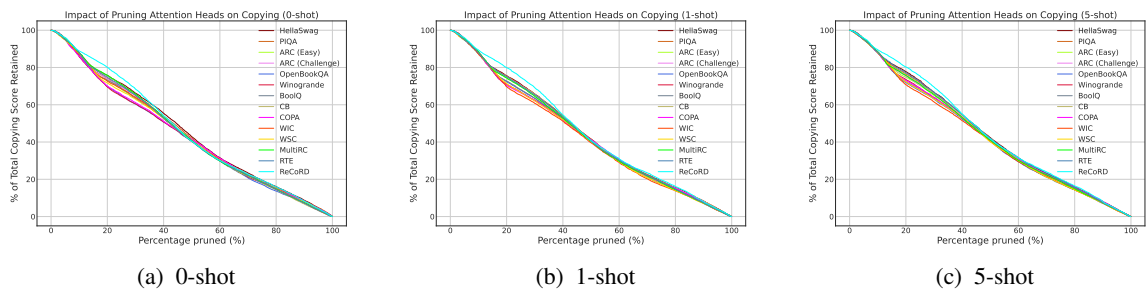


Figure 24: Total copying capacity retained as a function of percentage of attention heads pruned, where heads are pruned based on task-specific importance score rankings in the order of least important first.