

MECHANISTIC INSIGHTS: CIRCUIT TRANSFORMATIONS ACROSS INPUT AND FINE-TUNING LANDSCAPES

Anonymous authors

Paper under double-blind review

ABSTRACT

Mechanistic interpretability seeks to uncover the internal mechanisms of Large Language Models (LLMs) by identifying circuits—subgraphs in the model’s computational graph that correspond to specific behaviors—while ensuring sparsity and maintaining task performance. Although automated methods have made massive circuit discovery feasible, determining the functionalities of circuit components still requires manual effort, limiting scalability and efficiency. To address this, we propose a novel framework that accelerates circuit discovery and analysis. Building on methods like edge pruning, our framework introduces circuit selection, comparison, attention grouping, and logit clustering to investigate the intended functionalities of circuit components. By focusing on what components aim to achieve, rather than their direct causal effects, this framework streamlines the process of understanding interpretability, reduces manual labor, and scales the analysis of model behaviors across various tasks. Inspired by observing circuit variations when models are fine-tuned or prompts are tweaked (while maintaining the same task type), we apply our framework to explore these variations across four PEFT methods and full fine-tuning on two well-known tasks. Our results suggest that while fine-tuning generally preserves the structure of the mechanism for solving tasks, individual circuit components may not retain their original intended functionalities.

1 INTRODUCTION

The rapid advancement of large language models (LLMs) has established them as powerful tools across a wide range of tasks for (Vaswani, 2017; Devlin, 2018; Achiam et al., 2023) natural language and beyond. However, their vast scale—often involving billions of parameters—and black-box nature pose significant challenges in understanding their decision-making processes, leading to unpredictable risks in critical domains where accuracy, fairness, and transparency are essential. The field of interpretability aims to address these challenges by developing methods to uncover the internal reasoning of these models. Mechanistic interpretability, a subfield within this area, seeks to reverse-engineer LLMs into human-understandable algorithms (Conmy et al., 2023; Bereska & Gavves, 2024). A key objective is to identify **circuits**, subgraphs within the model’s computational graph that correspond to specific behaviors, while maintaining a high level of sparsity and preserving the model’s performance on the given task.

Recent progress in mechanistic interpretability has shed light on the inner workings of language models through meticulous human inspections, uncovering key circuits responsible for specific tasks (Wang et al., 2022; Hanna et al., 2024; Merullo et al., 2023). To accelerate this process and reduce reliance on costly manual efforts, automated tools have been developed to systematically identify circuits driving certain behaviors (Conmy et al., 2023; Bhaskar et al., 2024; Syed et al., 2023). However, these methods are either time-consuming (Conmy et al., 2023) or rely on approximations that prioritize speed over reliability (Syed et al., 2023). As a result, further investigations into the relationships between circuits (Tigges et al., 2024; Prakash et al., 2024; Jain et al., 2024) are limited in either scale or precision in finding circuits. The recent introduction of techniques like Edge Pruning (Bhaskar et al., 2024) has made it possible to automatically identify more accurate and faithful circuits across larger datasets while keeping computational costs manageable.

054 Despite advancements in circuit discovery, significant challenges remain in understanding the func-
 055 tionality of the circuit components. Existing studies often assign specific functions to individual
 056 nodes within circuits through manual inspection. These nodes and their roles are typically identified
 057 during the discovery process by observing changes in the model’s output (e.g., variations in logit
 058 values) when specific node activations are perturbed. This approach generally requires extensive
 059 path interventions to determine the direct effects of nodes on the final output and careful design of
 060 such intervention. The process demands considerable manual effort, leading to a limited scope (e.g.
 061 only focus on the change in target output logits) and placing a heavy burden on researchers, ulti-
 062 mately slowing down the progress of Mechanistic Interpretability. Moreover, as circuit discovery
 063 methods become automated, the identification of circuits and their functionalities no longer hap-
 064 pen in parallel. While progress has been made in circuit identification, a growing gap persists in
 065 understanding the broader implications —*after circuits are identified, what coming next?* With
 066 automated methodologies generating vast amounts of circuits, investigating their internal func-
 067 tionalities has become increasingly laborious and challenging.

068 This challenge is evident in our observations from the IOI task. Modifying objects within the task
 069 results in significant changes to the discovered circuits, as shown in Fig. 1. Intuitively, circuits
 070 responsible for the same reasoning across similar tasks should remain consistent. However, the pre-
 071 trained model struggles to maintain this consistency when performing identical tasks across various
 072 types of objects (e.g., changing “Mary” to “Dog”). Moreover, applying different supervised fine-
 073 tuning methods introduces further variations in the circuits. However, investigating these changes in
 074 circuits and their functionalities becomes increasingly demanding and labor-intensive.

075 Inspired by these challenges and building on these findings, this work explores how circuits vary in
 076 both structure and functionality across different ablated prompts, and how model fine-tuning meth-
 077 ods affect these circuits. To address this, we propose a framework for circuit analysis that integrates
 078 **efficient discovery** with **automated interpretability**. To broaden the scope of functionality ex-
 079 ploration and accelerate the process, we make a trade-off by foregoing functionality conclusions
 080 based on causal effects. Instead, we integrated the attention pattern and logit lens with clustering
 081 techniques to describe the **intended functionalities** of the circuits components. By ‘intended func-
 082 tionality’, we refer to what the circuit components are attempting to do in the context of a given
 083 task, rather than what they ultimately contribute to the final output. Our framework extends the
 084 edge pruning algorithm (Bhaskar et al., 2024) for circuit discovery by incorporating stages for **cir-
 085 cuit selection**, **circuit comparison**, **attention grouping**, and **logit clustering**, as outlined in Fig. 2.
 086 We further summarize our findings and contributions as following: (1). We introduce a novel frame-
 087 work that accelerates the process of estimating the functionality of circuit components, reducing the
 088 reliance on manual efforts. (2). We demonstrate that the functionality of circuit components does not
 089 necessarily remain consistent across ablated prompts. (3). We show that while fine-tuning methods
 090 maintain the overall functional structure, the specific components performing those functions may
 091 change.

092 2 EXPERIMENT SETUP AND PIPELINE

093 In this section, we first outline the experimental setup, providing an overview of the prompt settings
 094 for the tasks and fine-tuning methodologies involved in our work. Next, we describe the process
 095 for identifying and selecting circuits, along with the metrics used for circuit validations. We then
 096 introduce the methods for evaluating the intended functionality of the circuit components. Lastly,
 097 we discuss how comparative circuit analysis is conducted, following the pipeline in Fig. 2.

099 2.1 TASKS AND MODELS

100 We explore circuit variations within the context of two specific tasks: Indirect Object Identifica-
 101 tion (IOI) and Greater Than (GT). These tasks were initially examined by Wang et al. and Hanna
 102 et al., respectively. We focused on IOI and GT due to their well-established, human-inspected cir-
 103 cuits, making them particularly suitable for implementing our framework and conducting an in-depth
 104 analysis of circuit variations through different fine-tuning methods and ablated prompts. Addition-
 105 ally, recent work by Merullo et al. offers promising insights through the study of these fundamental
 106 tasks. Building on top of these tasks, Merullo et al. has studied the transferrability of these circuits
 107 beyond the base syntactic structure of these tasks. To gain a more granular understanding of the

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

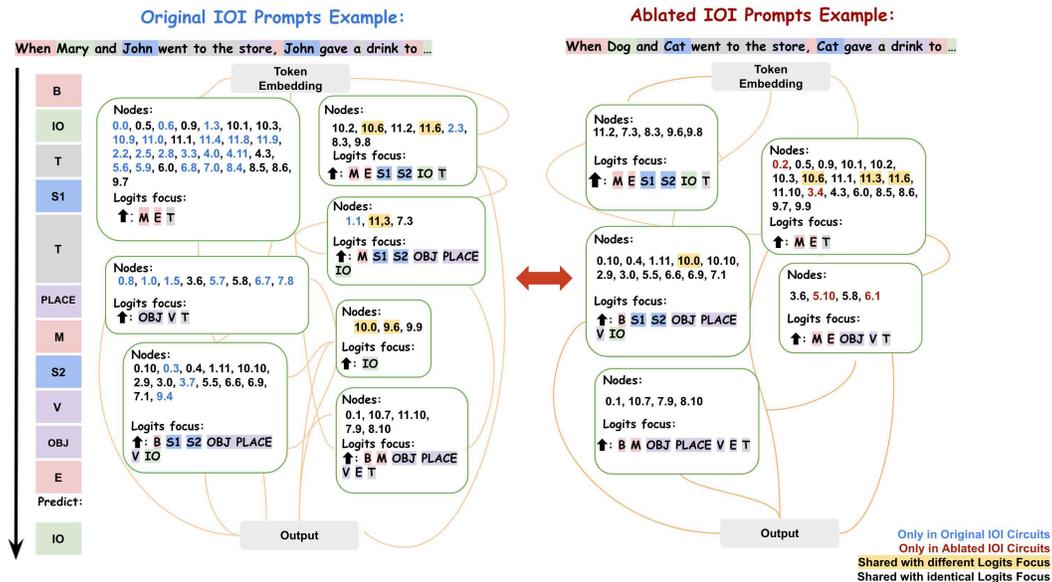


Figure 1: We demonstrate how circuits vary between IOI prompts and ablated prompts. When human names are replaced with animals, the newly discovered circuits become almost a subset of the IOI circuits. Nodes are clustered based on their logit outputs. We focus on logit outputs regarding to the prompt components. For each cluster, the major up-weights on logits for the sentence components are listed. Upon further analysis, we found that even overlapping nodes can perform different functions, colored in yellow.

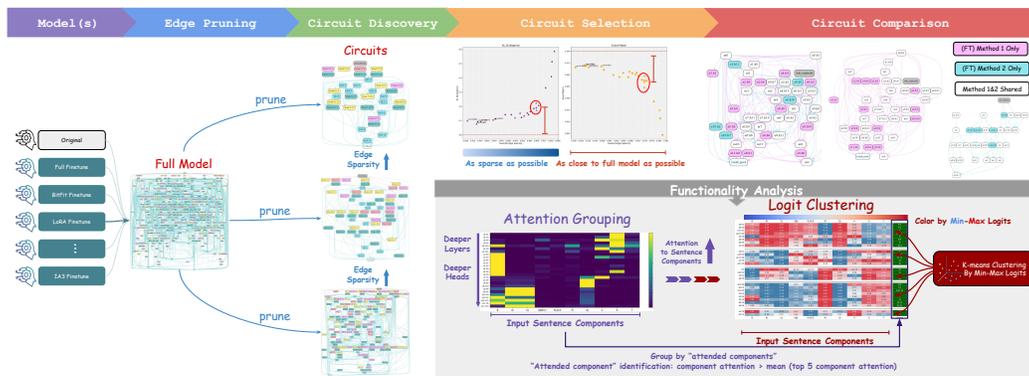


Figure 2: Pipeline of experiments: (1) Fine-tuned models on certain tasks; control model being the original GPT-2 model (2) Apply Edge-pruning to identify a series of candidate circuits (3) Select ideal circuits that balance sparsity and performance recovery, evaluated using KL Divergence and Exact Match metrics. (4) Compare circuits across models fine-tuned with different methods. (5) Conduct functionality analysis by: (5a). Grouping the heads based on their primary attended input sentence components - what the head is focusing on. (5b). Computing logits for each input sentence component per head - which components the head up-weights or down-weights for generation. (5c). Clustering heads by their min-maxed logits across input sentence components - which heads share or differ in functionality.

model’s mechanisms, we focus on GPT-2 small, as its resulted circuits are more manageable and human-readable. More details are discussed as follows:

Indirect Object Identification (IOI): IOI (Wang et al., 2022) is the task of predicting the name of the indirect object in a sentence. The task follows the format of “When $\{A\}$ and $\{B\}$ went to $\{PLACE\}$, $\{B\}$ bought a $\{OBJECT\}$ to $\rightarrow \{A\}$ ”. The datasets are generated by filling in A, B, PLACE, and OBJECT with a list of nouns. The model then completes the sentence by filling in the predicted indirect object, ideally $\{A\}$. Additionally, the prompts include variations in both ABBA and BABA formats. In the original IOI settings, A and B are filled by human names such as “John” and “Mary” as shown in Fig. 1. Following the setting by Edge Pruning (Bhaskar et al., 2024), we adopt the prompts on a variant with 30 templates from Hugging Face. Similarly, we randomly select 200 examples each for the train and validations; and 36,084 instances for the test sets.

Ablated IOI: In addition to the IOI dataset, we introduce an ablated IOI dataset to investigate whether the functionality and structure within the circuits are preserved across different types of objects, inspired by the observations in Fig. 1. We retain the template structure from the original IOI tasks; however, instead of populating A and B with human names, we use capitalized names of animals, cities, and colors. An example of the ablated prompt can be found in Fig. 1. To avoid tokenization issues that could affect model performance, we only include names that map to a single token. All other settings, such as train, test, and validation splits, remain the same as in the original IOI tasks.

Greater Than (GT): GT tasks (Hanna et al., 2024) follows the format of “*The war lasted from the year 1743 to 17 $\rightarrow xy$.*” The tasks requires the model to place a higher probability on the continuations 44, 45, ..., 99, compared to 00, 01, ..., 42. We adopted the version of dataset proposed in Edge Pruning (Bhaskar et al., 2024) which has 150 examples in the train and validation, and 12,240 instances in the test. This task is generally considered simpler than IOI, as it involves fewer logical steps to complete.

PEFT Methods: We primarily focus on Parameter-Efficient Fine-Tuning (PEFT) methods that preserve the original model structure. Specifically, we fine-tuned the GPT-2 small model with supervision, using Bitfit (Zaken et al., 2021), LoRA (Hu et al., 2021), IA3 (Liu et al., 2022), and AdaLoRA (Zhang et al., 2023), on the IOI and GT tasks. A fully fine-tuned model was used as a control for comparison. To ensure a fair comparison, all fine-tuned models are controlled to achieve nearly identical performances on the same test set. For the IOI task, the goal was to minimize the prediction loss for the indirect object. For the GT task, the goal was to maximize the gap between the cumulative probabilities of correct and incorrect years.

2.2 CIRCUIT DISCOVERY AND SELECTION

We adopted edge-pruning (Bhaskar et al., 2024) to automate the circuits discovery step for each finetuned models along with the pretrained version across the above mentioned tasks. To evaluate and ensure the faithfulness and preciseness of the identified circuits, we mainly relied on the measure of KL-Divergence and Exact Match, adopted by Bhaskar et al. and Conmy et al..

Circuit Discovery Algorithm: We implemented the edge pruning algorithm for automated circuit discovery (Bhaskar et al., 2024). This method addresses circuit discovery through gradient-based pruning on the edges of the model’s computational graph over hyperparameter *edge-sparsity* (*es*). Similar to path-patching (Wang et al., 2022) and the other automated methodologies (Conmy et al., 2023; Syed et al., 2023), it involves causal interventions on edges by substituting the target edges with counterfactual activations from corrupted examples. This process further generates sparse circuits by masking out edges that has no causal effect on the target tasks.

In this work, we first reproduced the circuits for the IOI and GT tasks. Next, we retrieved the circuits from the models fine-tuned with supervision on the two tasks. Finally, we retrieved the circuits from the models on the ablated IOI dataset using all the described models.

Circuit Evaluation and Selection: A circuit is considered accurate and faithful when its output closely aligns with that of the full model, even at a high level of graph sparsity, as highlighted in previous works (Hanna et al., 2024; Conmy et al., 2023; Bhaskar et al., 2024). Specifically, we utilized metrics, such as Exact Match for IOI, Kendall’s τ for GT, and KL Divergence for both, from the edge pruning approach to assess how closely the circuit’s output aligns with the full model. Addi-

tionally, we measured the differences between target and distractor outputs, such as logit difference for IOI and probability difference for GT, both traditionally used for these tasks.

As shown in the circuit selection stage of Fig. 2, we observed an exponential relationship between edge sparsity and KL divergence, consistent across all models and tasks. More results for the metrics are provided in the Appendix B. For this work, we selected a group of ‘best’ circuits for each model by identifying those at the ‘knee’ of the KL-divergence curve, as highlighted in red in the circuit selection stage of Fig. 2. Using this selection criterion, we chose the circuits with the highest possible sparsity, just before the KL divergence sharply increases. Among the circuits in this selected group, we considered them to be equivalently ‘best’. Therefore, we randomly selected one circuit from the group for further comparison and analysis. A sample circuit from this group represents the best tradeoff between performance and sparsity.

2.3 CIRCUIT FUNCTIONALITY DISCOVERY

We explore the intended functionalities of circuit components using a combination of attention heat maps and the logit lens. Unlike previous works, we conceptualize the mechanism within the model as either reducing the logits of incorrect choices or enhancing the logits of correct ones. In doing so, we broaden the focus beyond just the target outputs, providing a more comprehensive view of the model’s inner mechanism. For IOI-related tasks, we categorize the prompts into ten groups: “B, IO, S1, PLACE, M, S2, V, OBJ, E, and T,” following the structure of the IOI prompts. “B” represents the beginning of the sentence (e.g., “When” in Fig. 3), “M” refers to the middle of the sentence (e.g., a comma), and “T” denotes the end of the sentence (e.g., “to”). “V” represents the verb in the second half of the sentence, such as “gave,” while “OBJ” refers to the object before “E.” “S1” and “S2” refer to the distractor such as “John”, while “IO” represents the indirect object the model needs to predict, such as “Mary,”. We assign the rest of the templates into “T”. For GT tasks, we divide the prompts into six parts: “B, N, V, S, E, and T.” “B” represents the beginning of the sentence, “N” refers to the noun that the task focuses on (e.g., “war”), “V” indicates the verb such as “lasts” (suggesting the predicted numbers should be larger), “S” represents the start of the year, and “E” refers to the end of the sentence. We assign the rest of the templates into “T”. It is important to note that the intended functionalities do not directly reflect the ultimate contributions of the circuit components to the final output. However, understanding these intended functionalities still provides valuable insights and allows us to estimate their direct effects. As shown in Fig. 3, we found that the previously identified Name Mover Head formed its own group of intended functionality, exhibiting shared similarities in attention patterns. In fact, most of the previously identified nodes with the same functionalities cluster well together when analyzed using the logit lens.

Attention Grouping is a technique used to analyze and categorize transformer heads based on their attention patterns. In transformers, each head focuses on different components of the input, reflecting what it “intends” to process. By examining the attention values, we can infer which parts of the sentence each head is attending to and determine its specific intended functionality. Since not all attended tokens carry equal importance, the method focus on the most relevant components. Heads that focus on similar elements, such as an indirect object (IO), are grouped together, suggesting they contribute similarly to the model’s decision-making process.

Attention grouping helps uncover the various strategies a model uses to achieve the same outcome. For example, when the indirect object (IO) is emphasized, the model can do this by either increasing attention to the IO or by down-weighting attention to other components. Without grouping heads by their attention patterns, these diverse approaches could be overlooked. By clustering heads that focus on similar components, attention grouping allows for a deeper understanding of how the model processes information and how different heads contribute to the final prediction. This is especially useful in tasks like IOI, where multiple heads may contribute to the same result through distinct ways, offering valuable insights into the model’s internal workings.

The attention grouping process involves computing the attention of each head for all tokens in a sentence and mapping them to sentence components (e.g., B, IO). The average attention for each component is calculated, and the top k components with the highest attention values are selected. If a component’s attention exceeds the mean across these top components, the head is considered to be attending to that component. Heads are grouped if they attend to the same components. An example can be found in Fig. 3.

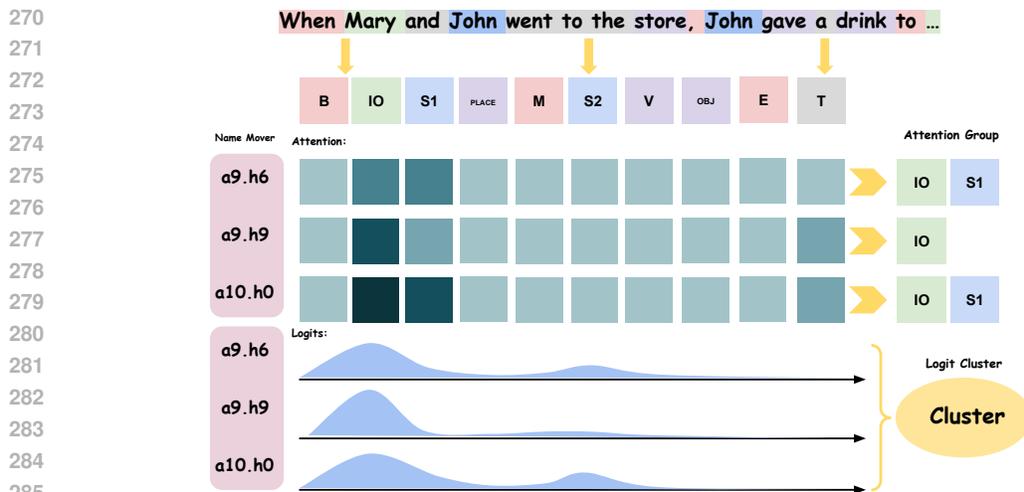


Figure 3: An example of attention grouping and logit clustering from IOI circuit. We assign the IOI prompts into several categories following their general pattern. We found that the previously identified Name Remover Head form its own cluster where the logits on IO get highly up-weighted over the other components. These heads’ attention focus on IO and S1/S2 as well, indicating their intentional functionality of attending to and up-weighting IO for correct prediction.

For the IOI task, with 9 sentence components, $k = 5$ is chosen based on observations that heads primarily focus on 4 or less components. For tasks with a different number of sentence components, the value of k can be adjusted accordingly.

Logit Clustering is an essential technique used to group transformer heads based on their logit patterns, providing insights into how different heads intended to contribute to the model’s decision-making process. One key method used in logit clustering is the *logit lens*, which allows researchers to probe intermediate representations in neural networks (NNs) and transformers. The logit lens works by examining the logits produced at each layer, representing the model’s confidence about which token it might output at that specific stage. This technique reveals how early layers begin forming predictions and how these predictions evolve and refine as more layers process the input.

The *logit patterns* generated by different heads give a sense of how each head contributes to the model’s predictions by up-weighting or down-weighting specific components. For example, in tasks such as the IOI task, where the indirect object (IO) is the correct target and the first subject (S1) is a distractor, a head that increases the logit on IO while decreasing the logit on S1 plays a crucial role in making the correct prediction. Similarly, other heads may contribute by focusing on irrelevant tokens such as punctuation or template words, which might help the model understand sentence structure or task objectives. Comparing the logit patterns across different heads and layers, thus, helps identify where significant changes in token predictions occur and how intermediate layers potentially contribute to the final outcome.

By clustering heads based on their logit patterns, we can group heads that are performing similar functions. This *logit clustering* allows for the identification of heads that collaborate in the model’s decision-making process, providing a clearer view of the functional roles played by different heads. The use of logit clustering is particularly useful when logit patterns vary across layers or even repeat in heads across multiple layers, as we have observed. Applying clustering algorithms to these patterns automates the grouping process, facilitating the analysis of head functions. For example, clustering can reveal whether heads within certain layers, such as deeper layers, are grouped together or whether logit patterns are distributed across all layers.

In this work, we used K-means algorithm for logits clustering. While other clustering methods or human calibrations may provide more precise results, K-means has proven to deliver satisfying results in grouping logit patterns, reducing the amount of manual effort required in functional analysis. This method allows for a scalable approach to understanding how different transformer heads influence the model’s output and decisions.

Attention Grouping + Logit Clustering: Both attention grouping and logit clustering are needed because each method has limitations when used in isolation. Attention grouping only indicates which components the heads focus on but doesn't reveal whether the heads are attempting to up-weight or down-weight those components. On the other hand, logit clustering identifies which components the heads are likely up-weighting or down-weighting but cannot clarify whether this effect is intentional or a byproduct from another function. For example, an up-weighted logit on IO might either reflect direct attention and up-weighting on this component, or it could be the result of down-weighting other components.

By integrating these two methods, we can better understand the intentional functionality of the heads. Attention grouping helps reveal the intention behind the model's focus, while logit clustering estimates the functionality — the effect the heads are having on specific components. When observing attention groups within the same logit cluster, it's easier to discern how nodes with similar functions differ in their intentions. Similarly, identical attention groups may belong to different logit clusters, indicating varied roles in the model's decision-making process. This combined framework simplifies the investigation of intentional functionalities, helping to clarify how different components intended to contribute to the overall behavior of the model.

3 EXPERIMENTAL RESULTS

In this section, we provide a detailed analysis of how the structure and intended functionality of IOI circuits differ from those retrieved from ablated prompts, where we focus on the animal objects since the original GPT2-small is incapable of performing IOI tasks when changing IO and S to cities and colors. We then explore how various finetuning methods enhance performance on both the IOI and GT tasks. Lastly, we investigate how finetuning on the original IOI task improves the model's capability of performing the IOI task when applied to ablated prompts. Evaluations on model performance can be found in Appendix. Sec. B. For the following results, we focus on attention head only by collapsing the QKV nodes into the corresponding attention heads.

Table 1: The model's circuit differences on the GT and IOI tasks are compared in terms of the number of nodes. For clarity, we have consolidated the QKV (Query, Key, Value) nodes of the IOI task into a single attention head. The results are averaged over the selected group of the "best" circuits for each model under each task. A 95% confidence interval is calculated to demonstrate statistically significant variations in circuit sizes. Blocks highlighted in green indicate a statistically significant larger number of nodes. Our analysis shows that finetuning methods generally reduce the circuit sizes for the GT task while introducing more components to the IOI task.

Model A	Model B	GT task			IOI task		
		A only	B only	Shared	A only	B only	Shared
Original	IA3	17.2 ± 1.22	11.2 ± 2.10	38.6 ± 2.56	16.2 ± 1.83	20.8 ± 1.99	74.2 ± 1.34
Original	AdaLoRA	20.0 ± 1.69	13.8 ± 1.72	35.8 ± 2.01	13.0 ± 2.13	18.2 ± 1.09	77.4 ± 1.93
Original	Bitfit	25.0 ± 2.43	11.0 ± 1.82	30.8 ± 2.98	21.2 ± 1.58	20.2 ± 0.94	69.2 ± 1.65
Original	Full	19.2 ± 1.72	8.2 ± 1.78	36.6 ± 2.86	18.0 ± 1.84	22.0 ± 1.07	72.4 ± 1.37
Original	LoRA	20.2 ± 1.74	14.0 ± 2.69	35.6 ± 1.37	14.8 ± 2.33	22.0 ± 1.96	75.6 ± 1.48
IA3	AdaLoRA	10.4 ± 1.63	10.2 ± 1.28	39.4 ± 2.36	17.8 ± 1.22	18.4 ± 1.23	77.2 ± 1.01
IA3	Bitfit	17.8 ± 2.87	9.8 ± 1.74	32.0 ± 3.05	22.8 ± 1.91	15.4 ± 1.13	72.2 ± 1.40
IA3	Full	13.0 ± 2.46	8.0 ± 2.57	36.8 ± 1.95	19.2 ± 0.94	18.6 ± 1.29	75.8 ± 0.58
IA3	LoRA	10.8 ± 2.37	10.6 ± 2.07	39.0 ± 1.57	17.0 ± 1.07	19.6 ± 1.43	78.0 ± 1.04
LoRA	AdaLoRA	9.6 ± 1.81	9.6 ± 1.58	40.0 ± 1.90	17.6 ± 1.35	15.6 ± 1.99	80.0 ± 1.33
LoRA	Bitfit	18.0 ± 2.16	10.2 ± 2.99	31.6 ± 1.85	22.4 ± 2.68	12.4 ± 1.09	75.2 ± 1.01
LoRA	Full	14.0 ± 2.35	9.2 ± 2.85	35.6 ± 1.72	19.2 ± 1.60	16.0 ± 1.07	78.4 ± 0.90
AdaLoRA	Bitfit	17.4 ± 1.81	9.6 ± 1.97	32.2 ± 2.87	22.0 ± 0.78	14.0 ± 1.21	73.6 ± 0.66
AdaLoRA	Full	14.4 ± 1.56	9.6 ± 2.03	35.2 ± 2.55	20.2 ± 1.22	19.0 ± 1.21	75.4 ± 0.53
Bitfit	Full	11.2 ± 2.39	14.2 ± 1.95	30.6 ± 2.89	14.6 ± 1.72	21.4 ± 1.85	73.0 ± 1.07

3.1 IOI CIRCUITS AND ABLATED PROMPTS

As shown in Fig. 1, substituting human names with other types of objects, such as animals, leads to significantly smaller circuits. Additionally, as illustrated in Fig. 4, we observed three distinct

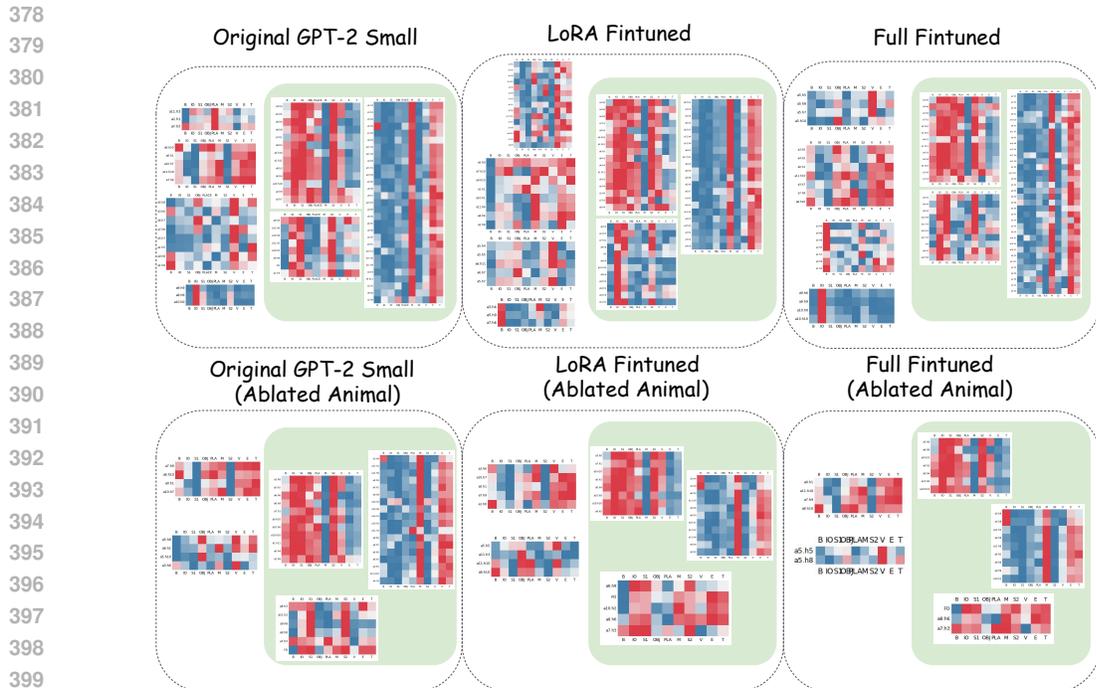


Figure 4: Logit clusters across fine-tuning methods and ablated prompts. The clusters are mostly uniform across fine-tuning methods for either prompt setting, indicating that fine-tuning does not significantly modify circuit’s main functionalities; however, variations can be observed, both in terms of the general sizes of each cluster and individual head’s intentional functionalities. Comparing IOI and ablated prompts, the ablation circuits consist of fewer functional clusters than IOI, and the number of heads for the shared clusters generally reduces. This is likely due to the complexity of processing human names in a semantically meaningful sentence which requires more functional nodes.

clustering groups that consistently appear across all circuits in both IOI and IOI-ablated tasks. One of the clusters places a strong emphasis on the M, E, and T parts of the sentence, in terms of logit, while also attending to the entire sentence. We assume that these nodes intend to complete the sentence with the simplest solutions, such as adding a comma or repeating the final word. Another cluster, in contrast, emphasizes all parts of the sentence except for M, E, and T, while still attending to the entire sentence. This cluster appears to function in opposition to the sentence completion cluster. This cluster also aligns with where all previously identified induction heads (Wang et al., 2022) are clustered. The final preserved cluster focuses heavily on the IO and S tokens, both in terms of logits and attention. This cluster generally down-weights OBJECT and PLACE, components that are irrelevant to correct prediction, which makes it distinguished from the induction heads cluster who usually up-weights them.

In terms of the general cluster analysis, the cluster of induction heads are mostly preserved, with only three additional nodes in IOI. The sentence completion cluster has significantly more heads in IOI, contributing to the difference in circuit sparsity. The cluster that focuses on IO and S is mostly preserved. S-inhibition and the negative name mover (downweighting IO and S, upweighting others) are also preserved.

Interestingly, there is no cluster found in the animals-ablated circuit that purely concentrates on IO in terms of logits and attention. It is worth noting that most of the attention patterns from ablated prompts do not involve any exclusive focus on IO or S tokens in terms of logits and attentions. Instead, they tend to focus on IO and S along with the beginning words. Furthermore, the name mover heads, which are well-clustered in the IOI circuits shown in Fig. 3, now appear to be more dispersed. For instance, node 10.0 still upweights IO over S1, but the effect is less pronounced. It also shows a strong upweighting of PLACE. Node 9.6 shifts to upweight S over IO, while node 9.9 joins the

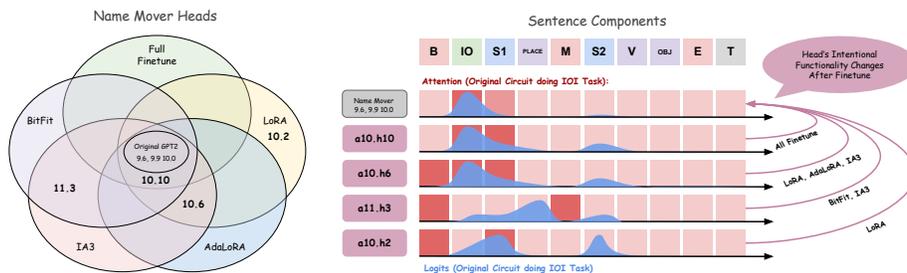
432 sentence completion cluster, upweighting M, E, and T but downweighting the other elements. These
 433 findings overall suggest that the model struggles to isolate IO and S from the prompt. Most of the
 434 time, it relies on the beginning of the sentence to identify the IO and S pair. Although the behavior
 435 of indirect object identification persists, the model’s limited ability to accurately identify IO and
 436 S hinders its performance compared to the original IOI task. These findings are consistent across
 437 finetuning methods and ablation settings.

438
 439 **3.2 PEFT METHODS COMPARISON**
 440

441 By fine-tuning GPT-2 small on IOI and GT tasks, we observed that the circuits retrieved from GT
 442 tasks are generally statistically significantly smaller than those from the original GT tasks, as shown
 443 in Table 1. An analysis of logits and attention patterns on GT revealed that the fine-tuned model
 444 tends to focus more on the starting year, while the MLP in deeper layers increasingly up-weights
 445 target number continuations. These findings suggest that the reduction in circuit size for GT tasks is
 446 likely due to the simplicity of the task —simply selecting numbers greater than the starting year, as
 447 the task only has one direction (Greater Than)— and a more focused approach in solving it. More
 448 results on GT tasks can be found in Appendix. Sec. D. In contrast, for the more complex IOI task,
 449 fine-tuning methods generally result in a statistically significantly larger circuit to address the task’s
 450 increased complexity.

451 For the IOI task, the overall structure of functionality for circuits is largely preserved, aligning with
 452 findings from Prakash et al.. Specifically, many of the heads identified by Wang et al. as serving
 453 certain functions continue to perform similar roles in fine-tuned cases. For example, Duplicate
 454 Token Heads, Previous Token Heads, some S-inhibition Heads, and certain Backup Name Mover
 455 Heads remain consistent across all fine-tuning methods, maintaining their functionalities despite the
 456 fine-tuning adjustments. This stability in key components ensures that model behavior is preserved,
 457 allowing fine-tuned heads to focus more effectively on specific task elements, thereby enhancing
 458 overall task performance. Furthermore, this consistency not only highlights the critical contribution
 459 of these heads to solving the IOI task but also validates our framework, which integrates attention
 460 grouping and logit clustering to explore the intended functionalities of LLM attention heads.

461 While many individual heads’ intentional functionalities are generally preserved, fine-tuning often
 462 causes shifts in their behavior. Different fine-tuning methods may or may not achieve the same effect
 463 on this change in behavior. For instance, in the two Induction Heads identified by Wang et al., a5.h5
 464 and a6.h9, we observe that while fine-tuning preserves the intentional functionality of a6.h9, a5.h5
 465 undergoes significant modifications with different fine-tuning methods. In the original circuit, a5.h5
 466 attends to the beginning of the sentence, up-weighting B, IO, S, OBJ, and down-weighting M, E,
 467 and T. However, after fine-tuning, Full Finetune, BitFit, and LoRA shift it to primarily up-weight V,
 468 whereas AdaLoRA and IA3 reverse the effect, down-weighting B, IO, S, OBJ, and PLACE, while
 469 up-weighting M, V, E, and T. Additional example of how fine-tuning methods align or differ in
 470 modifying the Name Mover Heads on IOI circuits is shown in Fig. 5.



481 Figure 5: A comparison of how different fine-tuning methods modify the cluster of Name Mover
 482 Heads claimed by prior works. The Name Mover Heads significantly attend to and up-weight IO
 483 over other components. All fine-tuning methods increase the number of these heads to enhance
 484 prediction accuracy, shifting the functionalities of some particular heads in the original GPT-2 circuit
 485 for better performance. While some fine-tuning methods share similarities in the additional Name
 Mover Heads, variations are common.

486 Therefore, the differences between circuits may help explain how fine-tuning enhances the perfor-
487 mance of LLMs on certain tasks. From our observations, fine-tuning increases the number of heads
488 that up-weight logits on IO compared to the original GPT-2 circuit. Specifically, the improvement
489 is primarily due to more “focused” up-weighting on IO relative to other components. As shown in
490 Fig. 5, all circuits from fine-tuned models enforce certain heads from the original circuit to func-
491 tion similarly to Name Mover Heads—those that primarily attend to and up-weight IO—leading to
492 improved task performance. Notably, all fine-tuning methods modify the functionality of a10.h10,
493 likely due to its high similarity to Name Mover Heads. However, PEFT requires a larger number
494 of functionality-shifted heads to achieve comparable performance, which is unsurprising since full
495 fine-tuning updates all weights and biases, thoroughly converting the functionalities of individual
496 heads. In contrast, PEFT only updates a small subset of parameters, necessitating more heads with
497 similar intentional functionalities to match the performance of fully fine-tuned circuits.

498 In conclusion, both PEFT methods and full fine-tuning improve performance not only by increasing
499 the presence of highly specialized heads, such as Name Mover and Induction Heads, but also by sim-
500 plifying and refining circuits through the pruning of irrelevant or less important components. This
501 dual effect—enhancing the specificity and accuracy of critical heads while reducing unnecessary
502 complexity—highlights the vital role of fine-tuning in optimizing circuit performance for specific
503 tasks. These adjustments lead to more efficient and interpretable circuits that maintain high levels
504 of task-specific accuracy.

505 4 CONCLUSIONS AND LIMITATIONS

506 Overall, our framework significantly reduces the reliance on manual efforts, enhancing the efficiency
507 of discovering the intended functionality of circuit components. Through our analysis of how the
508 circuits for the IOI and GT tasks vary across different fine-tuning methods and ablated prompts, we
509 found that while the overall structure of intended functionality is preserved, the specific components
510 responsible for these functions may change. This finding suggests that pre-identified circuits and
511 functionalities are subject to variations depending on the ablated prompts and fine-tuning methods
512 used. We hope that our findings and framework will inspire further exploration of circuit utilization
513 and its interpretability in LLMs.

514 As discussed in the paper, rather than estimating the direct effect of individual nodes on target logit
515 values, we focus on their intended functionalities. While exploring intended functionalities can pro-
516 vide a reasonable approximation of direct effects—since pre-identified heads with similar functions
517 are generally well-clustered—it is important to recognize that they are not equivalent. Drawing an
518 analogy to “correlation is not causation,” we emphasize that intended functionality does not nec-
519 essarily reflect the “true” functionality of circuit components. Furthermore, intended functionality
520 may depend on the performance of previous states, meaning it can become inconsistent if those ear-
521 lier states are perturbed. However, this challenge could also arise in causal exploration methods like
522 Patch patching.

523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545
546 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
547 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
548 2022.
- 549
550 Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety—a review. *arXiv*
preprint arXiv:2404.14082, 2024.
- 551
552 Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. Finding transformer circuits
553 with edge pruning. *arXiv preprint arXiv:2406.16778*, 2024.
- 554
555 Lawrence Chan, Adrià Garriga-Alonso, Nicholas Goldwosky-Dill, Ryan Greenblatt, Jenny
556 Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal scrub-
557 bing, a method for rigorously testing interpretability hypotheses. *AI Alignment For-*
558 *um*, 2022. [https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/
causal-scrubbing-a-method-for-rigorously-testing](https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/causal-scrubbing-a-method-for-rigorously-testing).
- 559
560 Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-
561 Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural*
Information Processing Systems, 36:16318–16352, 2023.
- 562
563 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding.
564 *arXiv preprint arXiv:1810.04805*, 2018.
- 565
566 Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang.
567 Black-box prompt learning for pre-trained language models. *arXiv preprint arXiv:2201.08531*,
568 2022.
- 569
570 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin
571 Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained
language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- 572
573 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
574 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for
575 transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
- 576
577 Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order op-
578 timization methods for in-context learning: A study with linear models. *arXiv preprint*
arXiv:2310.17086, 2023.
- 579
580 Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model
581 behavior with path patching. *arXiv preprint arXiv:2304.05969*, 2023.
- 582
583 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?: Inter-
584 preting mathematical abilities in a pre-trained language model. *Advances in Neural Information*
Processing Systems, 36, 2024.
- 585
586 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
587 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
588 *arXiv:2106.09685*, 2021.
- 589
590 Samyak Jain, Ekdeep Singh Lubana, Kemal Oksuz, Tom Joy, Philip HS Torr, Amartya Sanyal, and
591 Puneet K Dokania. What makes and breaks safety fine-tuning? mechanistic study. *arXiv preprint*
arXiv:2407.10264, 2024.
- 592
593 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt
tuning. *arXiv preprint arXiv:2104.08691*, 2021.

- 594 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*
595 *preprint arXiv:2101.00190*, 2021.
596
- 597 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
598 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
599 learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- 600 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin
601 Bossan. Pefit: State-of-the-art parameter-efficient fine-tuning methods. [https://github.](https://github.com/huggingface/peft)
602 [com/huggingface/peft](https://github.com/huggingface/peft), 2022.
603
- 604 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual asso-
605 ciations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022a.
- 606 Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing
607 memory in a transformer. *arXiv preprint arXiv:2210.07229*, 2022b.
608
- 609 Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Circuit component reuse across tasks in trans-
610 former language models. *arXiv preprint arXiv:2310.08744*, 2023.
- 611 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
612 for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
613
- 614 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
615 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction
616 heads. *arXiv preprint arXiv:2209.11895*, 2022.
- 617 Nikhil Prakash, Tamar Rott Shaham, Tal Haklay, Yonatan Belinkov, and David Bau. Fine-tuning
618 enhances existing mechanisms: A case study on entity tracking. *arXiv preprint arXiv:2402.14811*,
619 2024.
- 620 Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou,
621 Banghua Zhu, Lianmin Zheng, Kurt Keutzer, et al. S-lora: Serving thousands of concurrent lora
622 adapters. *arXiv preprint arXiv:2311.03285*, 2023.
623
- 624 Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks.
625 *Advances in Neural Information Processing Systems*, 34:24193–24205, 2021.
626
- 627 Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit
628 discovery. *arXiv preprint arXiv:2310.10348*, 2023.
- 629 Curt Tigges, Michael Hanna, Qinan Yu, and Stella Biderman. Llm circuit analyses are consistent
630 across training and scale. *arXiv preprint arXiv:2407.10827*, 2024.
- 631 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
632
- 633 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-
634 pretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint*
635 *arXiv:2211.00593*, 2022.
- 636 Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning
637 for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.
638
- 639 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He,
640 Yu Cheng, Weizhu Chen, and Tuo Zhao. Adalora: Adaptive budget allocation for parameter-
641 efficient fine-tuning. *arXiv preprint arXiv:2303.10512*, 2023.
642
643
644
645
646
647

648	APPENDIX TABLE OF CONTENTS	
649		
650		
651	A Model Performance	14
652		
653	B Background and Related Works	14
654		
655		
656	C IOI circuits analysis	14
657		
658	C.1 Circuit Comparison	14
659	C.2 Logit clustering and attention Grouping	23
660		
661		
662	D GT circuit analysis	24
663		
664	D.1 Circuit comparison	24
665	D.2 Logit clustering and attention Grouping	24
666		
667		
668	E Reproducibility	24
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		

A MODEL PERFORMANCE

B BACKGROUND AND RELATED WORKS

Mechanistic Interpretability (MI) (Elhage et al., 2021; Olsson et al., 2022) aims to discover interpretable components of an otherwise blackbox neural network. Such analyses are typically performed via a series of input perturbations (also known as *patching*) to ablate the effect of individual model components to its predictive behavior (Chan et al., 2022; Meng et al., 2022a;b; Goldowsky-Dill et al., 2023). MI has been successfully applied to natural language tasks that feature a controlled output space on small-scale pre-trained LMs (Wang et al., 2022; Hanna et al., 2024); it has also been applied to study in-context learning and algorithmic behaviors on stylized transformers (Akyürek et al., 2022; Fu et al., 2023; Nanda et al., 2023). While early works put equal focus on circuit discovery and interpretability, recent ones have emphasized scalable circuit discovery (Conmy et al., 2023; Syed et al., 2023; Bhaskar et al., 2024) over extracting human-understandable algorithms, which requires extensive manual effort to examine the computational behavior of model components. Our work integrates recent best practices in circuit discovery to *automate* interpretability.

Parameter-Efficient Fine-tuning (Mangrulkar et al., 2022) aims to improve model performance on downstream tasks by training only a small portion of parameter relative to the full model. Referring to Ding et al. (2023) for a more detailed survey, most popular approaches exploits the low-rank structure of projection matrices (Hu et al., 2021; Zhang et al., 2023) or introduce a fixed set of scaling and/or bias parameters (Zaken et al., 2021; Liu et al., 2022). Other representative approaches include prompt tuning (Lester et al., 2021; Li & Liang, 2021; Diao et al., 2022) and dynamically identifying tuning parameters via influence functions (Sung et al., 2021). In this work, we primarily investigate LoRA (Hu et al., 2021), AdaLora (Zhang et al., 2023), BitFit (Zaken et al., 2021), and IA3 (Liu et al., 2022). These methods are broadly applied in many production settings, since they are more scalable with commercial hardware and can be served with thousands of replicas simultaneously (Sheng et al., 2023).

Bhaskar et al. has studied the change effects of circuits induced by fine-tuning on an entity tracking task, and found that fine-tuning *enhances* existing mechanisms for billion-scale LMs. Our work extends this study with a more diverse set of tasks and PEFT methods; and more importantly, we have identified that fine-tuning *modifies* existing mechanisms for small LMs.

This section presents the evaluation results of all models across various tasks and ablated prompts. Although the model performs well on the IOI ablated tasks after fine-tuning, GPT-2 small struggles with indirect object identification when the main objects are replaced with colors and cities. Specifically, the model’s ability to handle these ablations remains limited, indicating that fine-tuning has not fully generalized the model to different types of input modifications. Moreover, even after fine-tuning, the performance on the IOI task ablated with cities is still suboptimal, suggesting that the model’s understanding of abstract entities such as locations remains insufficient. These findings highlight the need for more targeted interventions or further fine-tuning strategies to improve model robustness across diverse ablations.

C IOI CIRCUITS ANALYSIS

This section documents the detailed qualitative results for IOI analysis. Sec. C.1 lists the visualization of comparative circuit analysis. The overlapped nodes are in white colors while the unique nodes to each circuits are shown in either pink or blue. Sec. C.2 lists the results of logit clustering and attention group visualizations.

C.1 CIRCUIT COMPARISON

Comparative circuit analysis on the original IOI tasks can be found in Fig. 6, Fig. 7, Fig. 8, Fig. 9, Fig. 10, and Fig. 11. The analysis with ablated prompts with animals can be found in Fig. 12, Fig. 13, Fig. 14, and Fig. 17. The analysis with ablated prompts with cities can be found in Fig. 18 and Fig. 19. Analysis on prompts ablated with colors can be found in Fig. 21.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

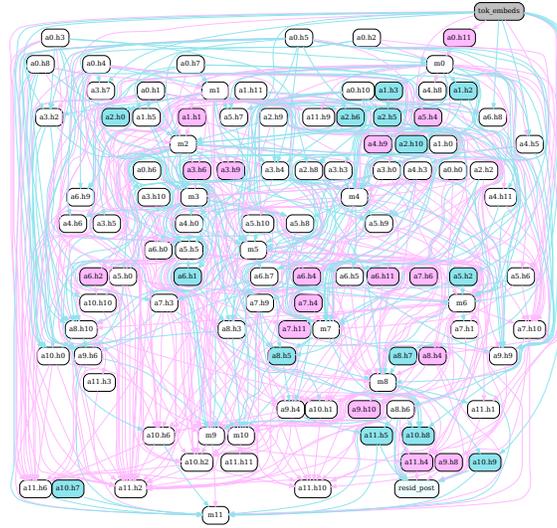


Figure 6: Circuit difference: LoRA vs Full. Red/pink indicates LoRA-specific circuits, blue indicates Full-specific circuits, and white represents shared circuits.

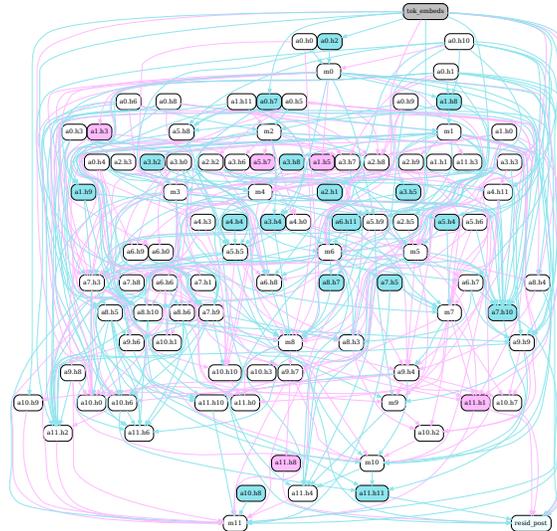


Figure 7: Circuit difference: Original vs AdaLoRA. Red/pink indicates Original-specific circuits, blue indicates AdaLoRA-specific circuits, and white represents shared circuits.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

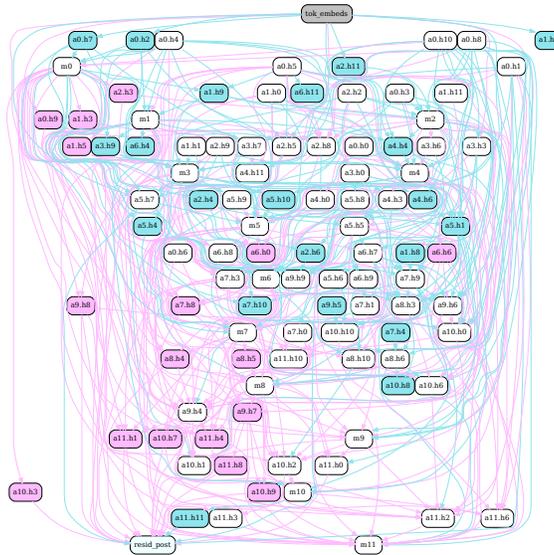


Figure 8: Circuit difference: Original vs BitFit. Red/pink indicates Original-specific circuits, blue indicates BitFit-specific circuits, and white represents shared circuits.

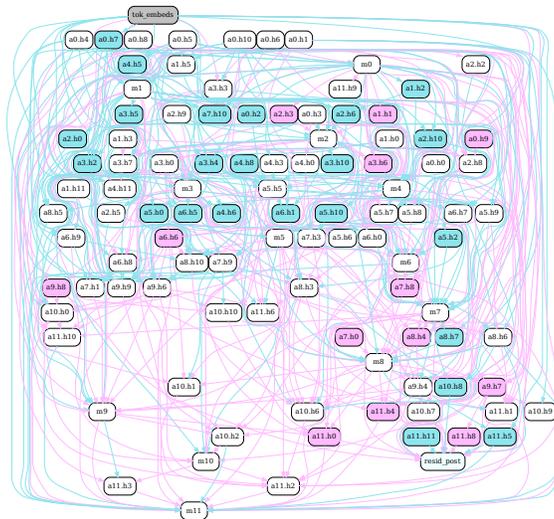


Figure 9: Circuit difference: Original vs Full. Red/pink indicates Original-specific circuits, blue indicates Full-specific circuits, and white represents shared circuits.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

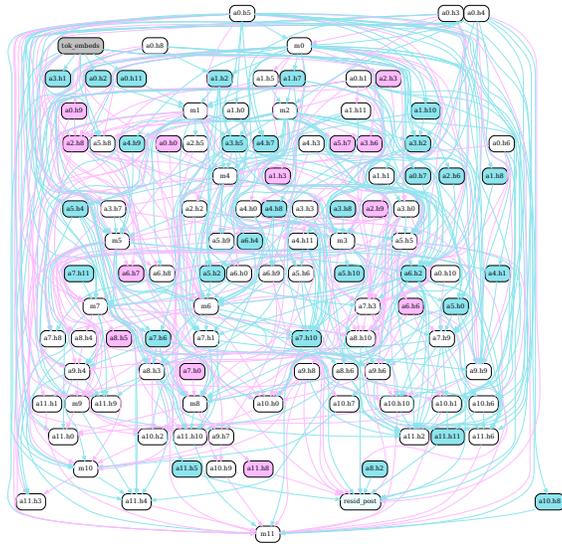


Figure 10: Circuit difference: Original vs IA3. Red/pink indicates Original-specific circuits, blue indicates IA3-specific circuits, and white represents shared circuits.

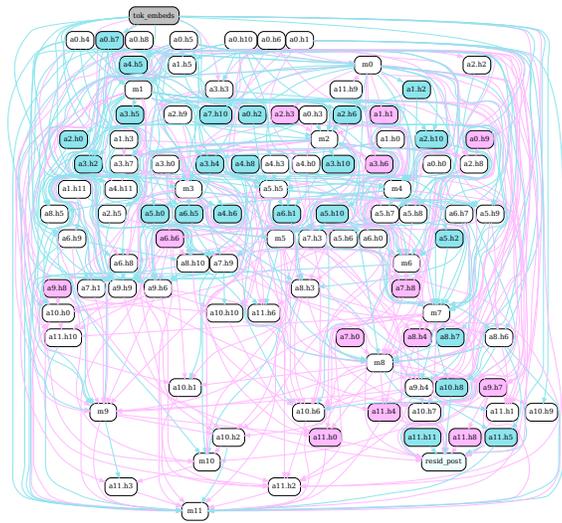


Figure 11: Circuit difference: Original vs LoRA. Red/pink indicates Original-specific circuits, blue indicates LoRA-specific circuits, and white represents shared circuits.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

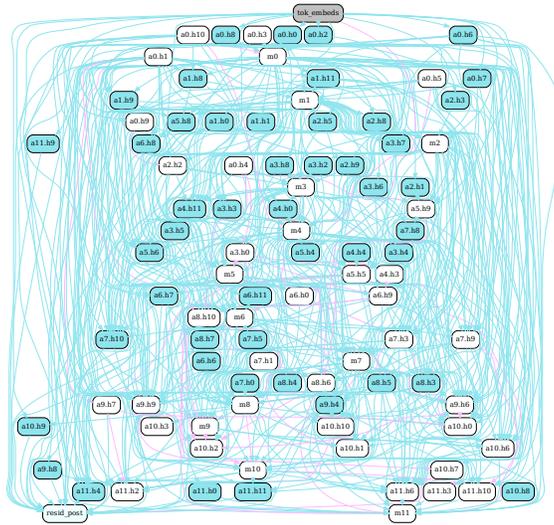


Figure 12: AdaLoRA circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

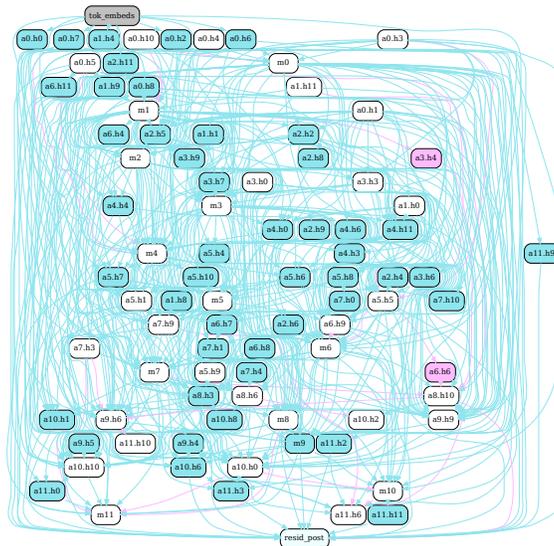


Figure 13: BitFit circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

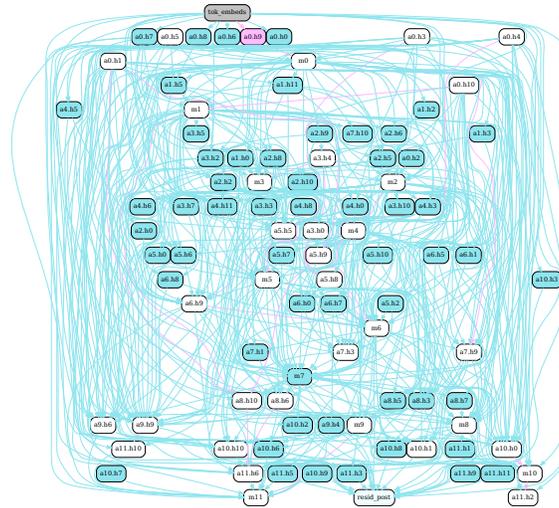


Figure 14: Full circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

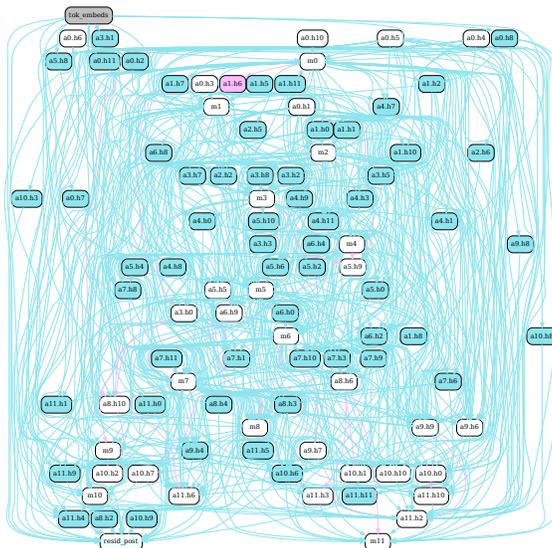


Figure 15: IA3 circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047

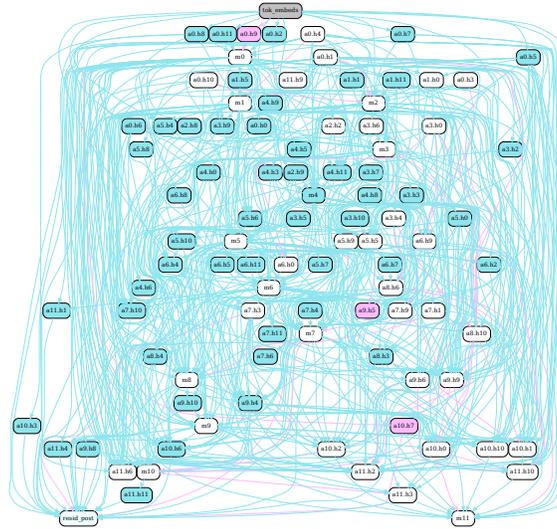


Figure 16: LoRA circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074

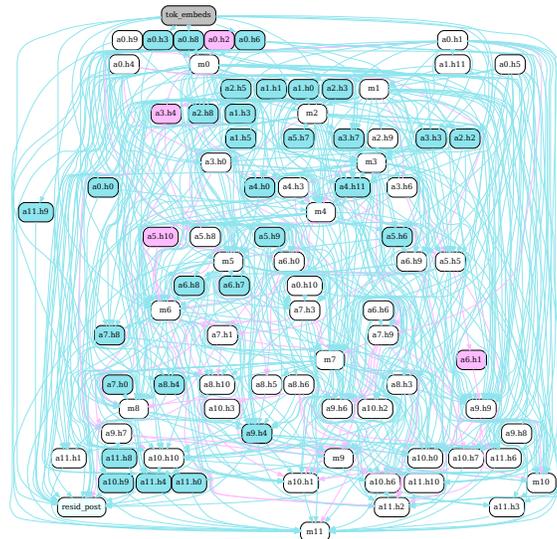


Figure 17: Original circuit difference between IOI and IOI-animals, where blue corresponds to IOI only and red corresponds to IOI-animals only circuit components. White refers to shared circuit components.

1075
1076
1077
1078
1079

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100

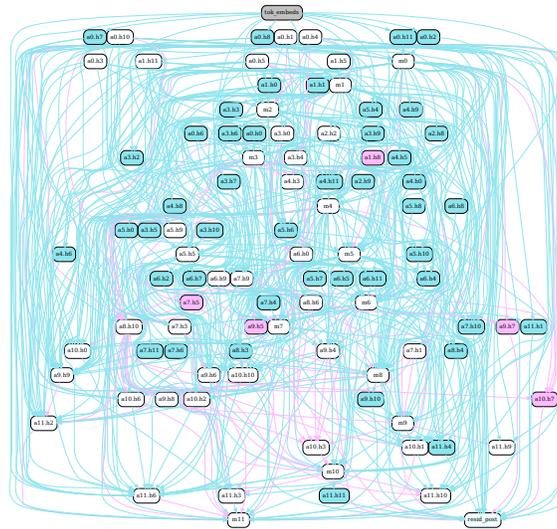


Figure 18: LoRA circuit difference between IOI and IOI-cities, where blue corresponds to IOI only and red corresponds to IOI-cities only circuit components. White refers to shared circuit components.

1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128

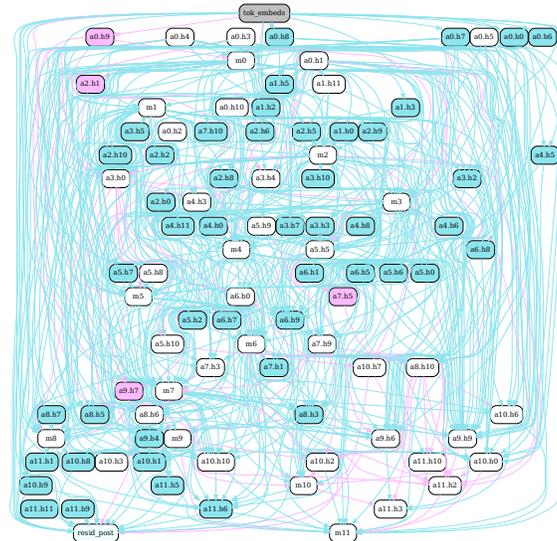


Figure 19: Full finetune circuit difference between IOI and IOI-cities, where blue corresponds to IOI only and red corresponds to IOI-cities only circuit components. White refers to shared circuit components.

1129
1130
1131
1132
1133

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

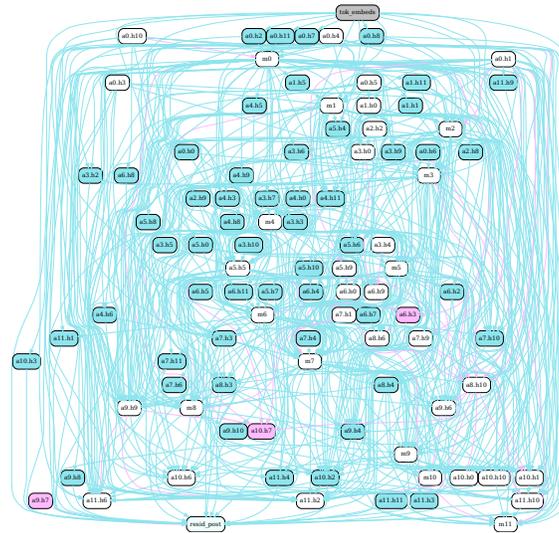


Figure 20: LoRA circuit difference between IOI and IOI-colors, where blue corresponds to IOI only and red corresponds to IOI-colors only circuit components. White refers to shared circuit components.

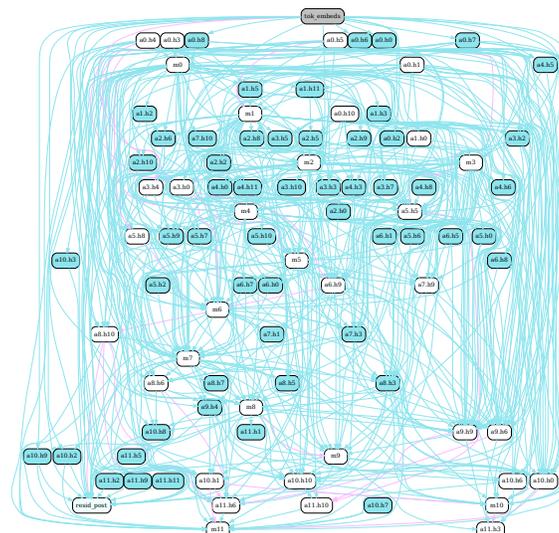


Figure 21: Full finetune circuit difference between IOI and IOI-colors, where blue corresponds to IOI only and red corresponds to IOI-colors only circuit components. White refers to shared circuit components.

Table 2: IOI circuit performances, rounded to the nearest hundredth. The original GPT-2 model, due to poor accuracy on colors and cities ablated settings, is not satisfiable for meaningful circuit investigations on these two tasks. Among fine-tuning methods, cities-ablated prompts pose the most difficulty for the circuit to preserve full model performance.

Metric	Ablation	AdaLoRA	BitFit	Full	IA3	LoRA	Original
Acc	Names	0.95 ± 0.01	0.99 ± 0.00	0.98 ± 0.00	0.97 ± 0.00	0.94 ± 0.00	0.73 ± 0.01
	Animals	0.96 ± 0.02	0.99 ± 0.00	0.99 ± 0.00	0.96 ± 0.01	0.97 ± 0.00	0.52 ± 0.01
	Colors	0.95 ± 0.02	0.99 ± 0.00	0.99 ± 0.00	0.96 ± 0.01	0.96 ± 0.00	–
	Cities	0.74 ± 0.02	0.97 ± 0.01	0.93 ± 0.01	0.76 ± 0.02	0.79 ± 0.00	–
LD	Names	5.95 ± 0.16	11.62 ± 0.15	11.63 ± 0.12	6.15 ± 0.10	5.83 ± 0.12	3.24 ± 0.18
	Animals	4.61 ± 0.41	12.32 ± 0.43	11.03 ± 0.34	4.35 ± 0.30	5.17 ± 0.27	0.86 ± 0.07
	Colors	4.23 ± 0.42	11.50 ± 0.39	10.41 ± 0.36	4.39 ± 0.31	4.91 ± 0.20	–
	Cities	2.57 ± 0.20	9.60 ± 0.38	8.82 ± 0.34	2.19 ± 0.15	3.37 ± 0.08	–
KL	Names	0.21 ± 0.01	0.05 ± 0.00	0.06 ± 0.01	0.16 ± 0.01	0.20 ± 0.01	0.30 ± 0.00
	Animals	0.09 ± 0.02	0.06 ± 0.01	0.04 ± 0.01	0.08 ± 0.01	0.06 ± 0.01	0.28 ± 0.01
	Colors	0.10 ± 0.03	0.05 ± 0.01	0.04 ± 0.00	0.09 ± 0.02	0.07 ± 0.00	–
	Cities	0.38 ± 0.04	0.11 ± 0.02	0.21 ± 0.03	0.39 ± 0.03	0.37 ± 0.01	–
EM	Names	0.95 ± 0.00	0.99 ± 0.00	0.98 ± 0.00	0.96 ± 0.00	0.94 ± 0.00	0.76 ± 0.01
	Animals	0.96 ± 0.02	0.99 ± 0.00	0.99 ± 0.00	0.96 ± 0.01	0.97 ± 0.00	0.66 ± 0.01
	Colors	0.95 ± 0.02	0.99 ± 0.00	0.99 ± 0.00	0.96 ± 0.01	0.97 ± 0.00	–
	Cities	0.77 ± 0.02	0.97 ± 0.01	0.93 ± 0.01	0.77 ± 0.02	0.82 ± 0.00	–

Table 3: GT task performances

Finetune Method	ES	PD	PD(10)	KT	KL
Original	0.99 ± 0.00	0.72 ± 0.00	0.33 ± 0.01	0.78 ± 0.01	0.23 ± 0.02
AdaLoRA	0.99 ± 0.00	0.96 ± 0.00	0.62 ± 0.01	0.84 ± 0.01	0.38 ± 0.02
BitFit	0.99 ± 0.00	1.00 ± 0.00	0.02 ± 0.00	0.88 ± 0.00	0.20 ± 0.01
Full	0.99 ± 0.00	0.98 ± 0.00	0.53 ± 0.03	0.82 ± 0.01	0.74 ± 0.05
IA3	0.99 ± 0.00	0.92 ± 0.00	0.55 ± 0.01	0.83 ± 0.00	0.33 ± 0.01
LoRA	0.99 ± 0.00	0.92 ± 0.00	0.58 ± 0.01	0.83 ± 0.01	0.30 ± 0.02

C.2 LOGIT CLUSTERING AND ATTENTION GROUPING

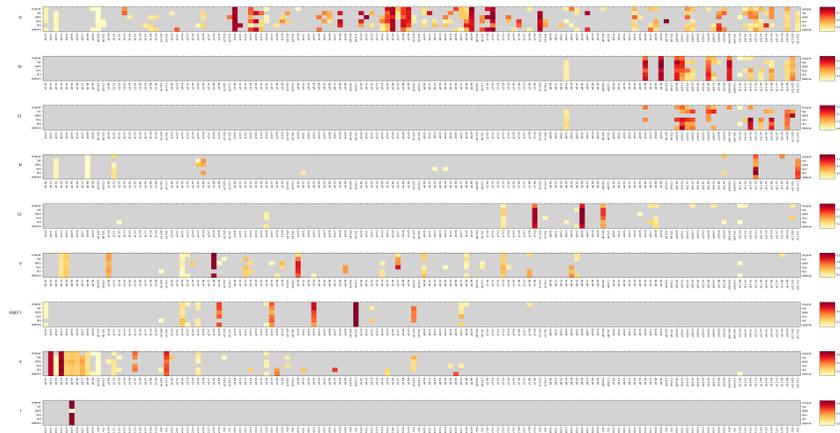


Figure 22: Attention grouping results for all components circuit. Grey cells indicate pruned nodes of minimal attention.

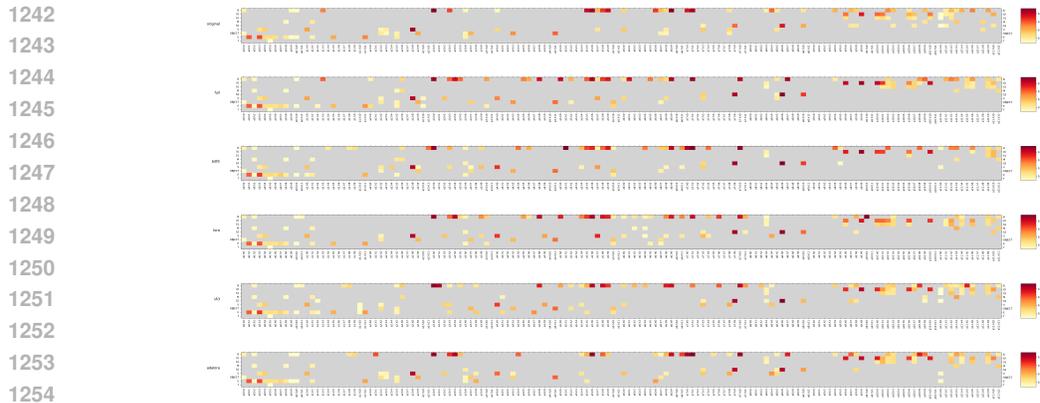


Figure 23: Attention grouping results for all finetuning methods circuit. Grey cells indicate pruned nodes of minimal attention.

D GT CIRCUIT ANALYSIS

D.1 CIRCUIT COMPARISON

D.2 LOGIT CLUSTERING AND ATTENTION GROUPING

E REPRODUCIBILITY

The implementation of circuit discovery mainly depends on the code from Edge Pruning. More details can be found in their github. The code to perform logit lens and attention grouping will be released upon acceptance.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

BitFit Logit Lens Heatmap

	B	IO	S1	OBJ	PLACE	M	S2	V	E	T	
FO	-85.51	-40.48	-51.64	-59.80	-60.72	-56.71	-51.64	58.74	-55.73	-55.66	FO
a10 h2	-2.02	1.04	1.10	-1.75	-1.19	3.44	1.10	0.61	2.61	2.24	B
a10 h6	-1.34	5.22	1.47	-0.61	-2.91	2.56	1.47	0.28	2.90	2.35	IO-S1
a11 h2	-0.85	2.19	3.92	-0.35	-0.91	2.33	3.92	-0.17	1.84	1.27	B
a3 h0	-0.39	0.32	0.33	0.36	0.00	0.16	0.33	0.45	-0.23	0.13	B
a5 h1	-0.50	0.03	-0.08	-0.25	-0.45	0.11	-0.08	-0.46	0.15	-0.01	B
a6 h11	0.09	0.51	0.42	0.49	-0.19	0.39	0.42	0.63	0.64	0.16	B-V
a7 h0	0.03	0.18	0.18	0.04	-0.16	0.04	0.18	-0.03	0.18	0.03	B
a7 h10	-0.06	0.04	0.05	-0.04	-0.28	0.07	0.05	0.16	0.21	0.11	B
a8 h3	1.35	1.96	1.62	-0.80	-1.67	0.17	1.62	0.64	1.54	0.71	IO-S1
a8 h6	-2.78	0.90	1.25	-4.12	-6.58	10.19	1.25	0.30	8.43	7.42	S2
a1 h1	-1.33	-0.79	-0.69	-0.29	0.44	-0.32	-0.69	-0.86	-0.29	-0.57	E-M
a11 h10	-2.05	-3.57	-2.87	-0.83	1.61	-1.64	-2.87	-1.92	-1.87	-1.29	IO-S1
a2 h6	-0.64	-0.68	-0.77	-0.21	-0.00	-0.04	-0.77	-0.64	0.14	-0.34	B-M
a2 h8	-0.17	0.18	0.04	0.54	-0.06	0.45	0.04	-0.19	0.62	0.10	V
a3 h7	0.37	0.32	0.21	0.53	0.46	0.43	0.21	0.64	0.93	0.66	OBJ-V
a5 h10	-1.04	-1.12	-1.27	1.86	-0.82	-0.02	-1.27	0.79	1.02	0.40	B-E-OBJ
a6 h8	0.55	0.29	0.19	0.42	-0.14	0.25	0.19	0.16	1.57	1.01	B-V
a0 h1	-1.50	-2.20	-2.37	-1.38	-0.97	-1.43	-2.37	-0.50	-0.54	-0.24	E
a0 h8	-0.53	-1.06	-1.13	-1.20	-1.03	-1.05	-1.13	-0.76	-0.71	-0.85	E-M
a1 h0	0.54	-0.02	-0.20	0.18	-0.43	0.10	-0.20	0.44	0.18	0.20	E-V
a3 h6	0.38	-0.11	-0.24	-0.01	-0.19	0.23	-0.24	0.31	-0.05	0.05	B-E-S2
a5 h9	0.58	-1.42	-1.47	-1.64	-0.76	0.38	-1.47	-1.89	0.34	-0.15	B
a7 h9	-0.74	-3.22	-6.19	-0.57	0.61	4.38	-6.19	2.61	3.59	3.88	S2
a8 h10	-2.68	-7.13	-11.25	-3.46	-2.24	-4.35	-11.25	-4.30	-3.17	-3.63	S2
a0 h0	-0.59	-0.88	-0.83	-0.60	-0.71	0.71	-0.83	-0.20	0.40	0.35	B-OBJ
a0 h10	-0.43	-0.47	-0.44	-0.14	-0.38	0.24	-0.44	0.15	0.10	0.08	B-E
a0 h2	-0.41	-1.15	-1.16	-0.54	-0.48	1.55	-1.16	0.14	1.36	1.02	B-E-M
a0 h5	-0.80	-1.99	-2.04	-1.42	-0.99	1.22	-2.04	-0.77	0.72	0.99	E-T
a0 h6	-1.73	-1.90	-1.90	-1.47	-1.30	0.78	-1.90	-0.31	0.01	0.03	E
a0 h7	-0.82	-0.58	-0.71	-0.15	-0.19	0.29	-0.71	-0.48	0.53	0.41	E
a1 h8	-0.96	-1.45	-1.32	-1.15	-0.81	1.28	-1.32	-0.43	0.97	0.72	B
a1 h9	-0.65	-1.64	-1.64	-1.12	-0.99	2.49	-1.64	-0.27	2.07	1.72	B
a10 h1	-1.40	2.48	-0.07	-1.54	0.16	7.65	-0.07	2.46	5.16	5.05	IO
a10 h8	-3.46	-4.78	-4.82	-3.42	-2.79	8.14	-4.82	0.11	5.63	5.64	B-V
a11 h0	-51.93	-102.71	-101.47	-65.78	-50.96	201.10	-101.47	0.34	166.50	137.75	B-V
a11 h11	-7.07	-8.77	-8.73	-6.38	-6.11	11.39	-8.73	0.66	9.65	8.08	B-M
a11 h6	-2.94	-1.04	-0.56	-2.97	-2.88	6.90	-0.56	0.84	5.45	4.33	B-IO
a11 h9	-9.95	-14.91	-14.22	-9.44	-7.44	24.97	-14.22	2.84	19.94	16.45	B-IO-S1
a2 h11	-1.52	-1.37	-1.22	-1.20	-1.04	0.60	-1.22	-0.74	0.13	0.21	B
a2 h2	-0.17	-0.09	-0.12	0.09	-0.32	0.71	-0.12	0.32	0.68	0.51	OBJ-V
a2 h5	-0.87	-1.70	-1.61	-0.85	-0.07	0.80	-1.61	-0.81	1.10	0.41	E-M
a3 h3	-0.51	-0.68	-0.71	-0.44	-0.61	0.95	-0.71	0.08	0.40	0.35	B-V
a3 h9	-0.26	-0.62	-0.66	-0.06	-0.36	1.10	-0.66	0.86	0.56	0.64	B-V
a4 h0	-0.87	-0.87	-1.04	0.13	-1.26	1.41	-1.04	0.31	2.02	0.88	V
a4 h11	-0.20	-1.18	-1.34	-0.68	-0.74	3.27	-1.34	0.71	2.85	2.31	OBJ
a4 h3	-0.53	-0.80	-1.09	0.12	-0.31	1.62	-1.09	0.15	2.09	1.28	OBJ
a5 h6	-0.58	-0.73	-0.67	-0.41	-0.50	1.23	-0.67	-0.04	0.84	0.68	B
a10 h0	-0.56	26.75	7.03	2.69	-0.40	0.89	7.03	-0.85	0.11	-0.21	IO
a10 h10	0.53	21.84	5.49	0.63	0.04	8.69	5.49	-1.43	8.38	6.35	IO
a11 h3	-0.38	7.43	5.20	4.21	0.24	-1.10	5.20	0.20	-0.65	-0.80	B-M
a6 h9	-0.16	0.55	0.56	-0.01	-0.33	-0.28	0.56	-0.12	-0.43	-0.27	B
a9 h6	-1.55	21.43	5.85	-0.36	-1.76	3.52	5.85	-0.28	3.33	2.43	IO
a9 h9	0.59	46.30	3.95	-0.31	-0.36	1.58	3.95	-0.70	0.22	0.67	IO
a2 h4	-0.22	-0.19	-0.19	0.15	-0.10	0.41	-0.19	0.43	-0.27	0.11	V
a5 h4	-0.14	-0.04	-0.05	0.23	-0.29	0.17	-0.05	0.35	0.08	0.14	B-V
a5 h5	-0.04	0.13	0.30	0.38	-0.31	0.27	0.30	1.02	0.28	0.07	B
a5 h7	-0.01	0.71	0.62	0.75	-0.46	0.85	0.62	1.92	-0.41	0.12	B-V
a5 h8	0.64	0.85	0.56	0.78	-0.53	0.80	0.56	1.14	0.76	1.04	B
a6 h7	0.27	0.47	0.34	1.89	0.74	-0.48	0.34	0.98	-0.29	0.11	B-OBJ
a7 h4	0.47	0.34	0.40	0.63	0.25	0.06	0.40	1.42	0.55	0.65	B
a0 h3	0.87	2.25	2.25	1.52	1.18	-3.60	2.25	0.42	-2.82	-2.46	E-V
a0 h4	0.11	0.37	0.40	0.11	-0.19	-1.99	0.40	-0.42	-1.71	-1.29	E-V
a1 h11	0.80	2.79	2.76	1.68	1.03	-6.74	2.76	0.25	-5.67	-4.81	E
a1 h4	1.40	2.00	1.77	1.17	0.78	-1.84	1.77	0.54	-0.90	-0.56	B
a2 h9	0.60	0.72	0.48	0.75	0.61	-0.24	0.48	0.85	0.70	0.48	OBJ-V
a4 h4	-0.28	0.34	0.44	0.53	0.27	-1.31	0.44	-0.43	-1.12	-0.69	B
a4 h6	0.26	0.60	0.48	0.49	0.39	-0.60	0.48	-0.50	0.11	0.26	B
a5 h4	0.30	0.78	0.87	0.91	0.34	-0.65	0.87	0.48	0.15	0.10	B
a7 h1	0.32	0.94	0.92	0.85	0.48	-0.54	0.92	0.48	-0.26	-0.12	B
a7 h3	-0.36	-0.24	-0.49	-1.41	0.41	-1.73	-0.49	-3.62	-1.20	-1.19	S2-V
a9 h4	0.61	1.91	2.17	0.16	0.66	-2.12	2.17	1.07	-0.70	-0.63	B
a9 h5	-1.14	-2.29	-2.34	-2.09	-1.99	-12.86	-2.34	-7.00	-10.27	-10.05	B-S2
	B	IO	S1	OBJ	PLACE	M	S2	V	E	T	Group

Figure 25: BitFit Logit Clusters

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

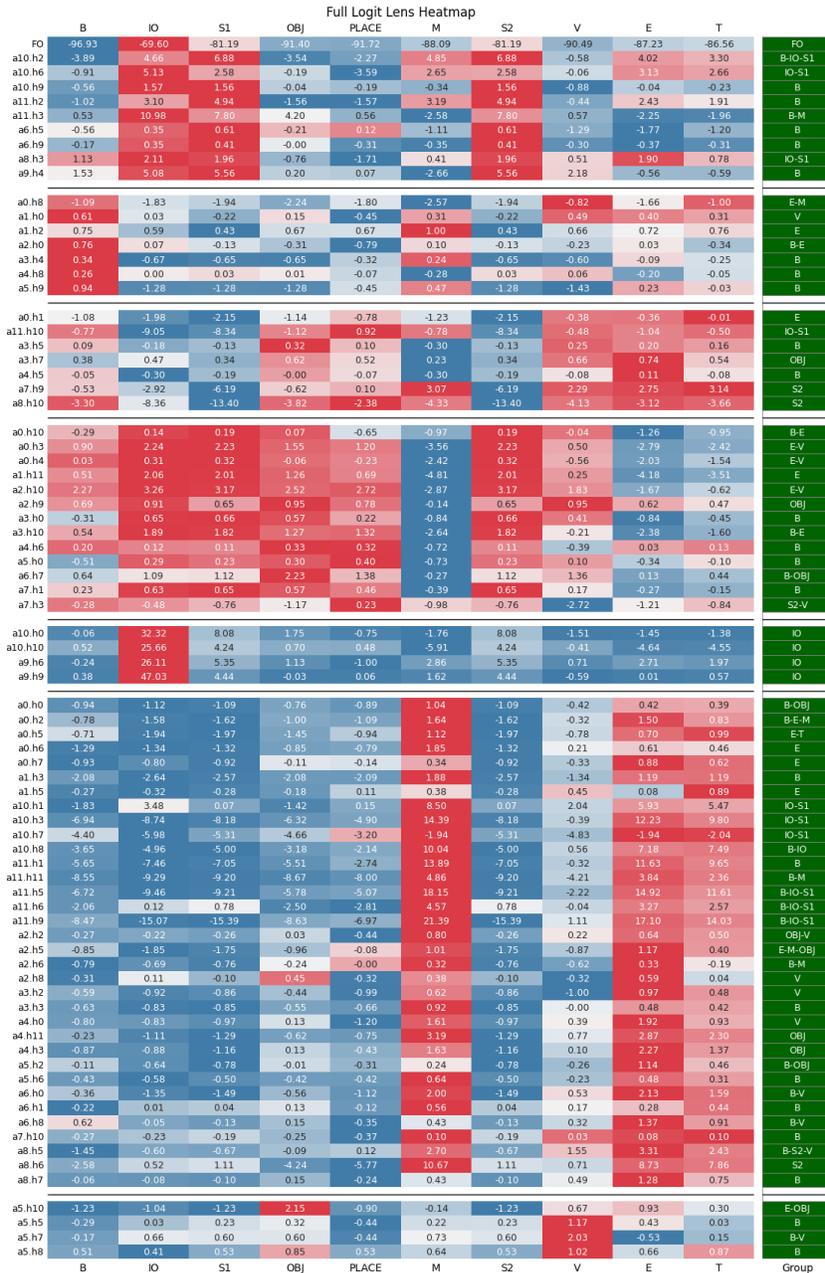


Figure 26: Full Logit Clusters

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

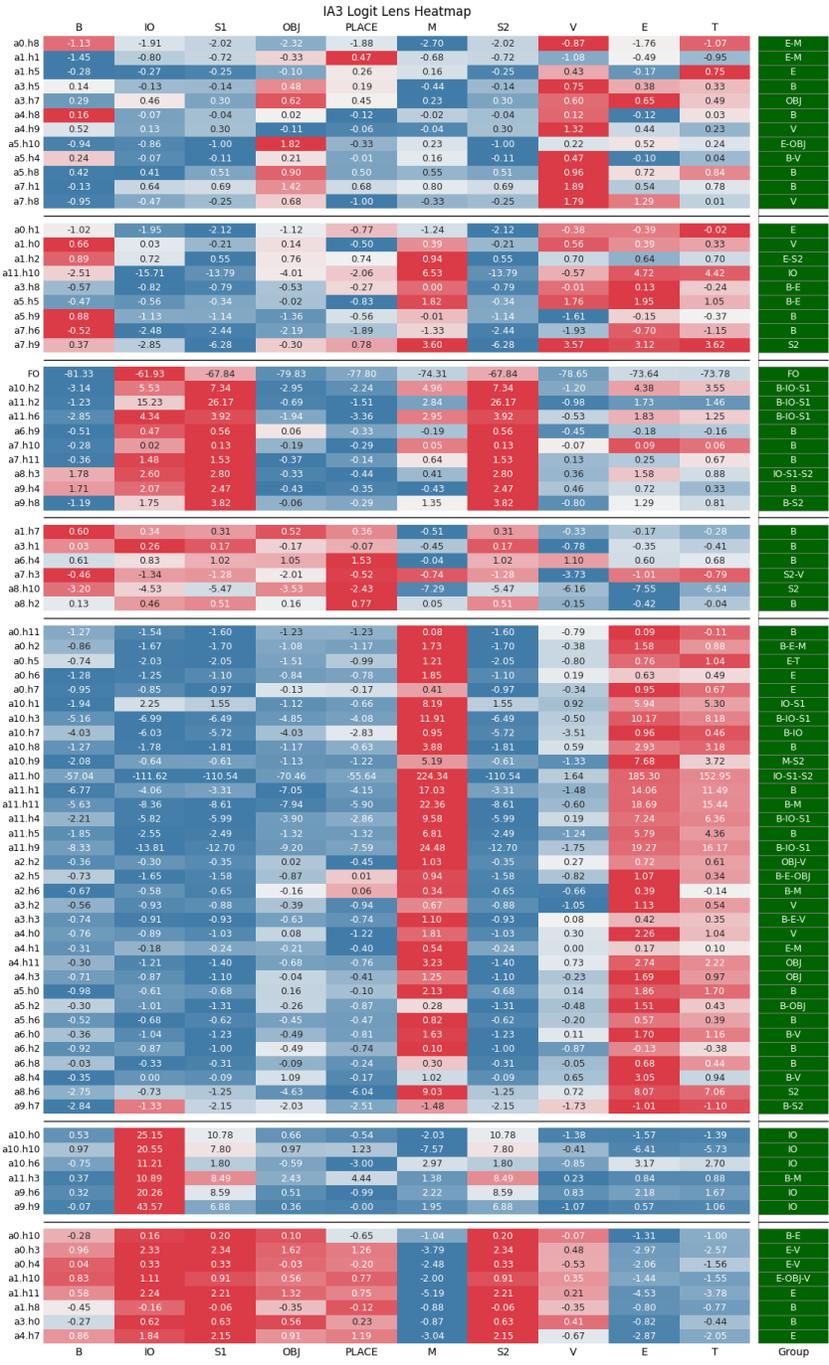


Figure 27: IA3 Logit Clusters

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

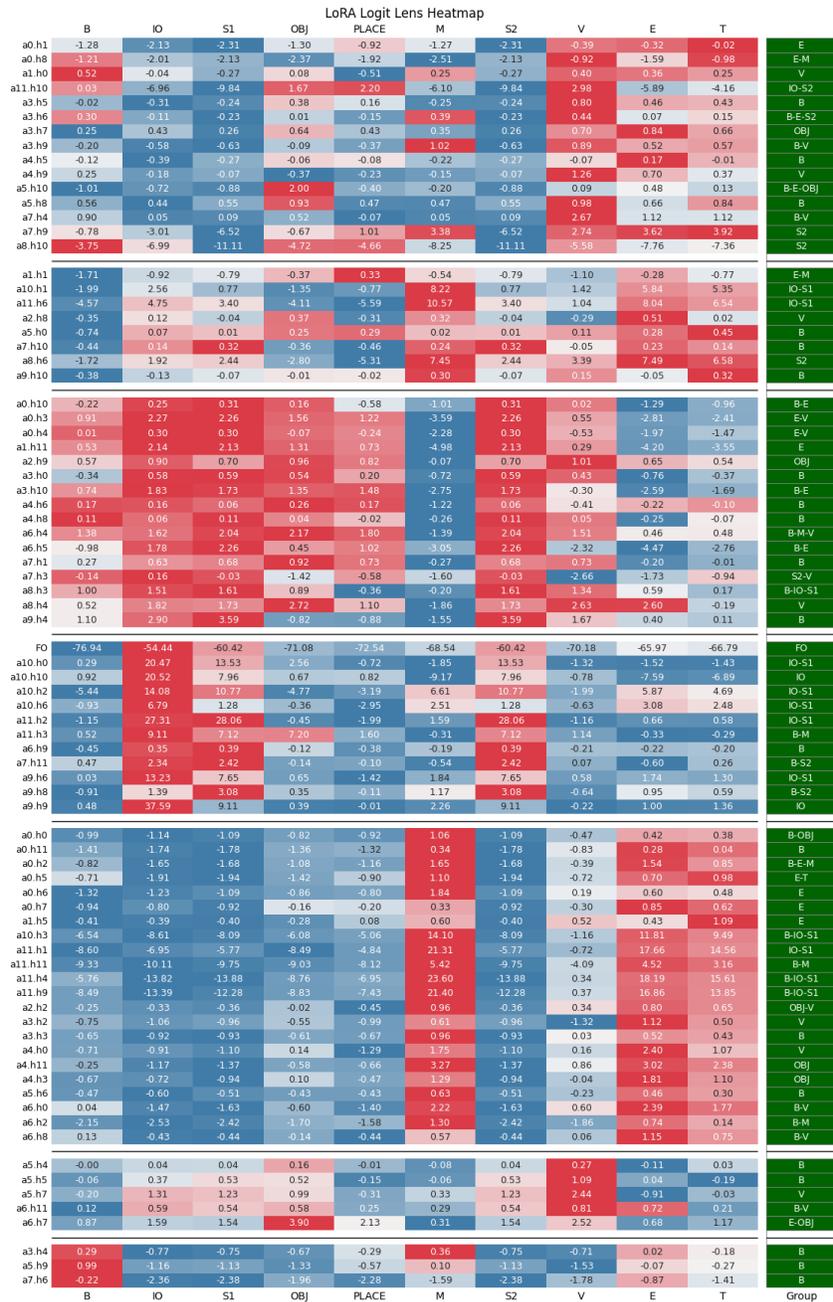


Figure 28: LoRA Logit Clusters

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

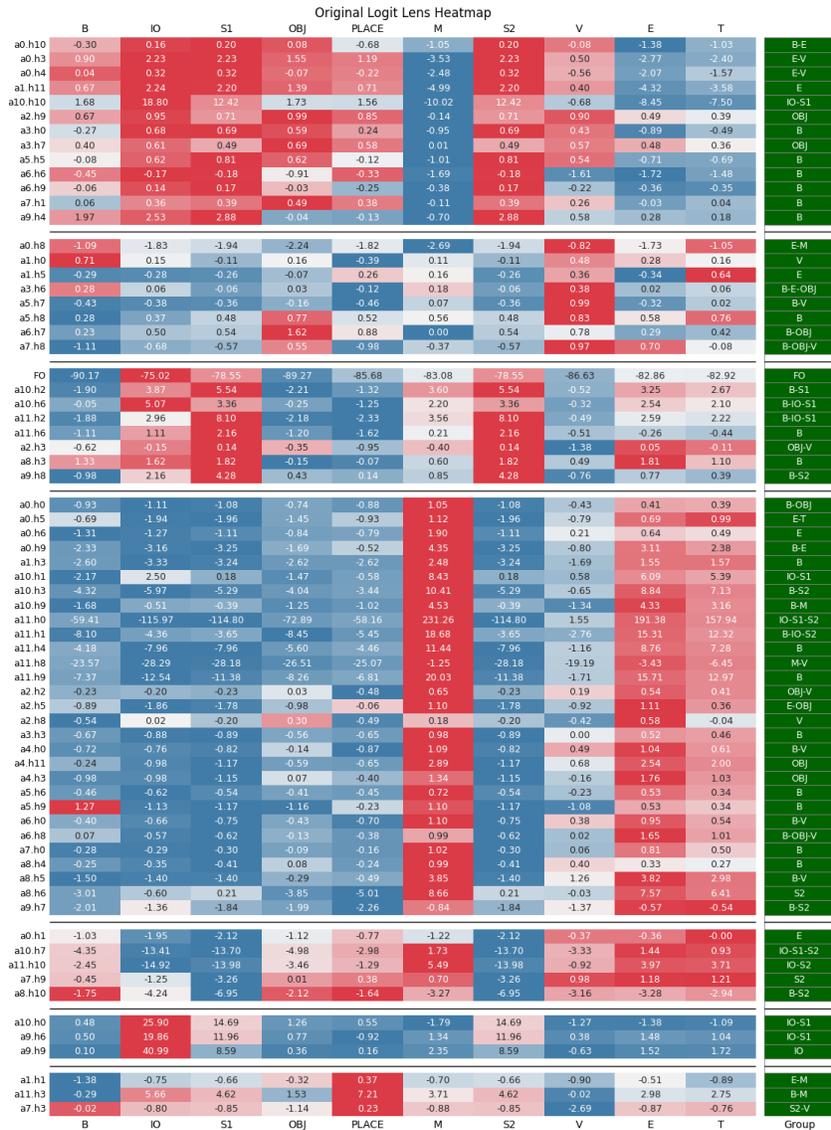


Figure 29: Original Logit Clusters

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

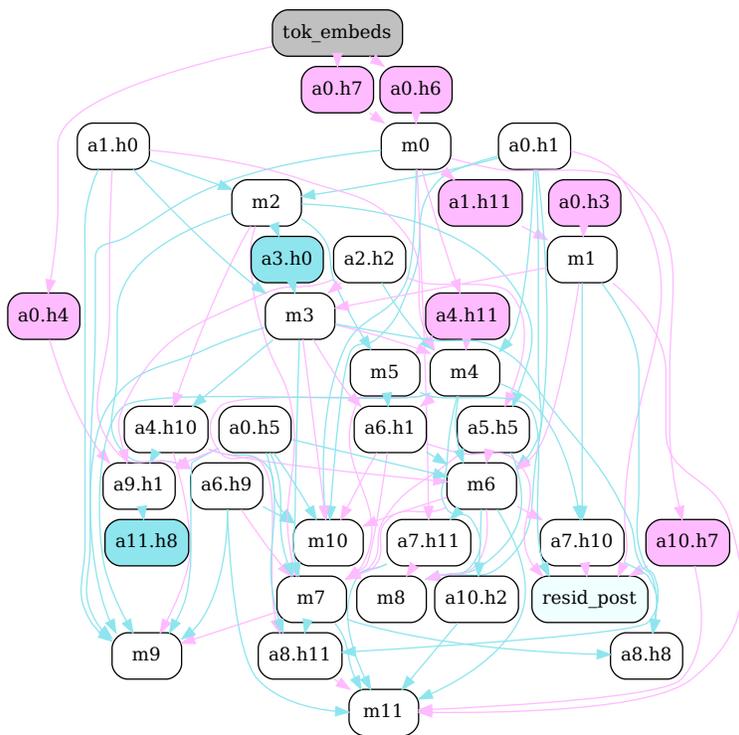


Figure 30: Circuit difference for GT task: LoRA vs Full. Blue indicates LoRA-specific components, red indicates Full-specific components, and white represents shared components.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

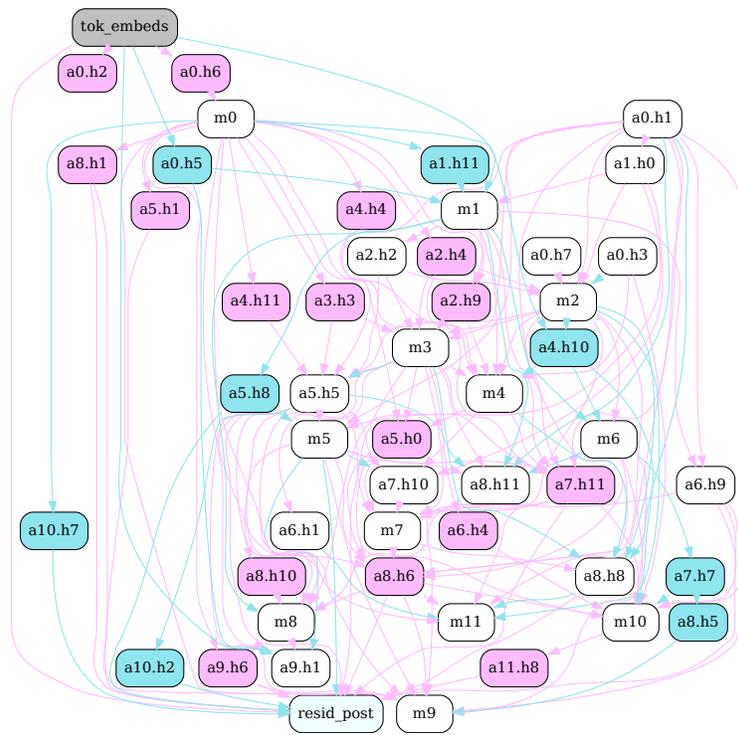


Figure 31: Circuit difference for GT task: Original vs AdaLoRA. Blue indicates Original-specific components, red indicates AdaLoRA-specific components, and white represents shared components.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

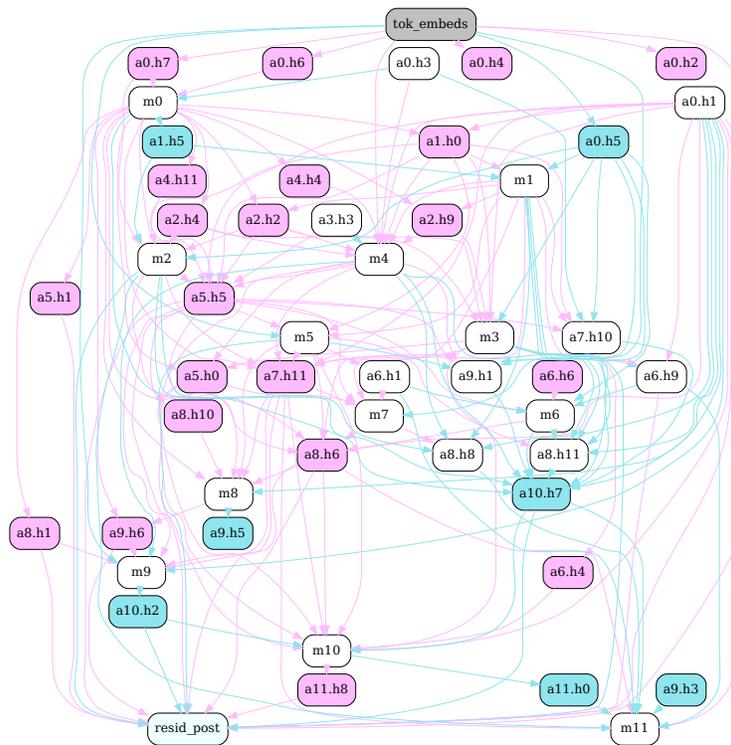


Figure 32: Circuit difference for GT task: Original vs BitFit. Blue indicates Original-specific components, red indicates BitFit-specific components, and white represents shared components.

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

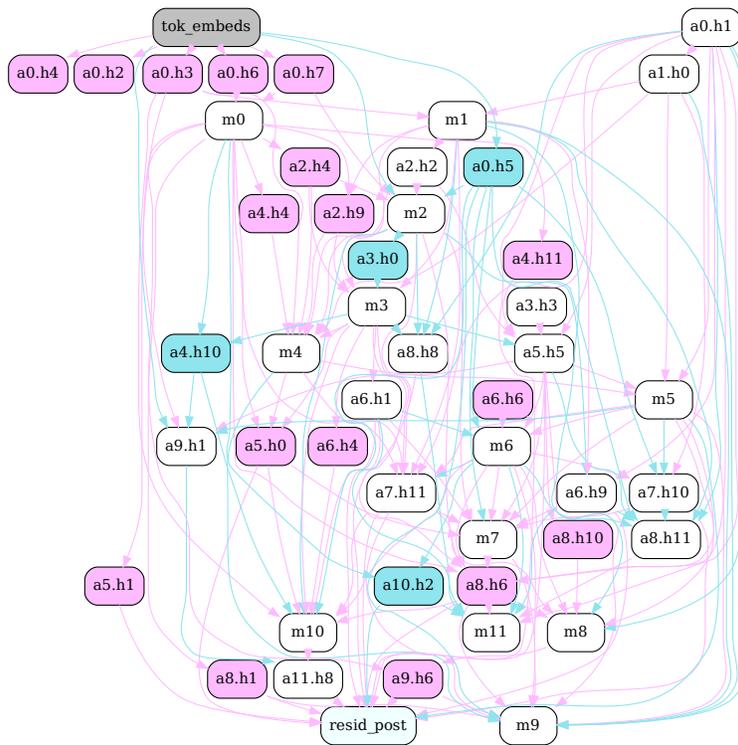


Figure 33: Circuit difference for GT task: Original vs Full. Blue indicates Original-specific components, red indicates Full-specific components, and white represents shared components.

1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889

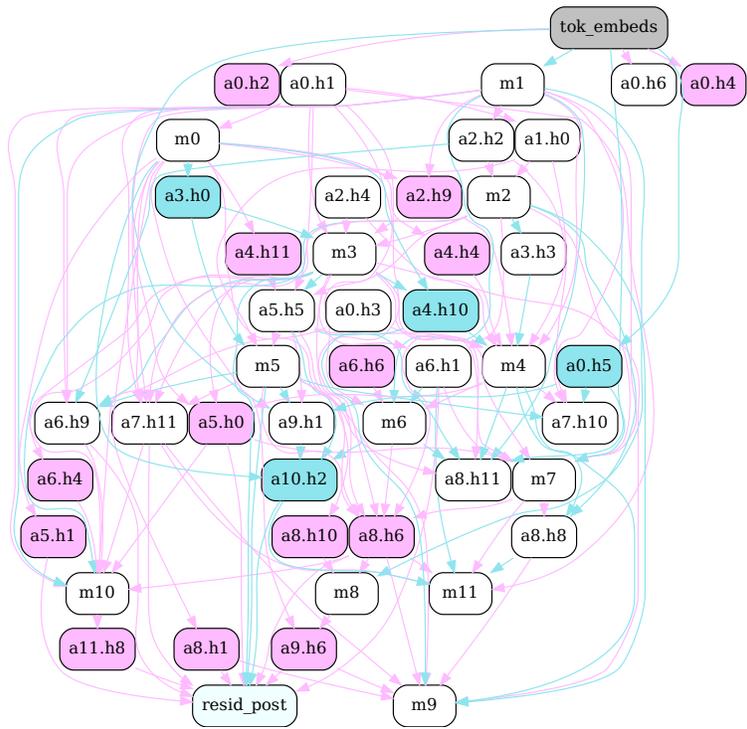


Figure 34: Circuit difference for GT task: Original vs IA3. Blue indicates Original-specific components, red indicates IA3-specific components, and white represents shared components.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

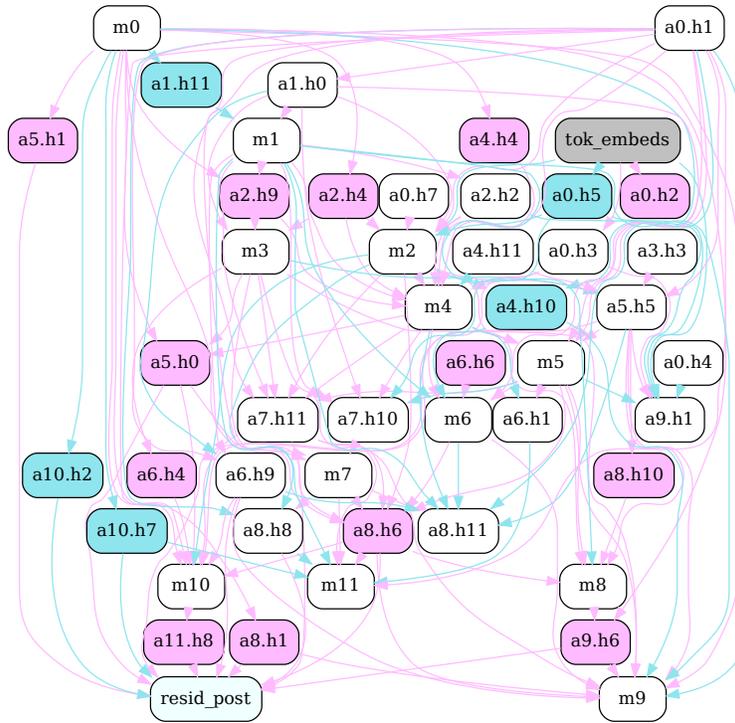


Figure 35: Circuit difference for GT task: Original vs LoRA. Blue indicates Original-specific components, red indicates LoRA-specific components, and white represents shared components.

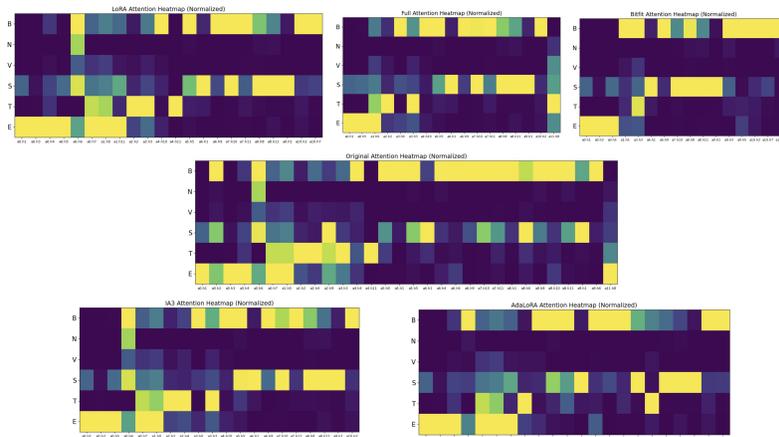


Figure 36: Attention grouping results for all finetuning methods circuit on GT.