

Deeper Insights Without Updates: The Power of In-Context Learning Over Fine-Tuning

Anonymous ACL submission

Abstract

Fine-tuning and in-context learning (ICL) are two prevalent methods in imbuing large language models with task-specific knowledge. It is commonly believed that fine-tuning can surpass ICL given sufficient training samples as it allows the model to adjust its internal parameters based on the data. However, this paper presents a counterintuitive finding: For tasks with implicit patterns, ICL captures these patterns significantly better than fine-tuning. We developed several datasets featuring implicit patterns, such as sequences determining answers through parity or identifying reducible terms in calculations. We then evaluated the models' understanding of these patterns under both fine-tuning and ICL across models ranging from 0.5B to 7B parameters. The results indicate that models employing ICL can quickly grasp deep patterns and significantly improve accuracy. In contrast, fine-tuning, despite utilizing thousands of times more training samples than ICL, achieved only limited improvements. We also proposed circuit shift theory from a mechanistic interpretability's view to explain why ICL wins.

1 Introduction

Adapting pre-trained models to specific tasks or domains is commonly achieved through fine-tuning (Hu et al., 2023; Peters et al., 2019) or in-context learning (Gan and Mori, 2023). Fine-tuning, a well-established method, involves further training a pre-trained model on a smaller, domain-specific dataset, directly updating the model's parameters to retain improvements across various contexts and scenarios. In contrast, in-context learning (ICL) enhances task performance by incorporating task-specific examples into prompts, guiding the model in task completion without altering its parameters during training.

There has been much debate about the pros and cons of fine-tuning and in-context learning. Fine-

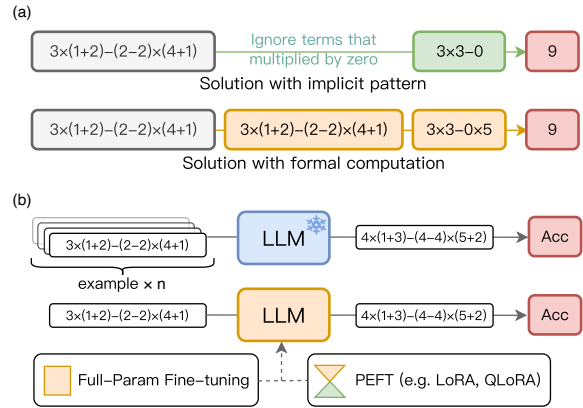


Figure 1: (a) A simple example of an implicit pattern detection task. The given problem (arithmetic expression calculation task in this figure) can be solved in either a formal way, e.g., directly calculating, or by exploiting the detected implicit pattern as a shortcut. (b) Illustration of implicit pattern detection for in-context learning and fine-tuning. For ICL, several examples with answers are given in context, and a further new question is used to test accuracy. For fine-tuning, LLM learns from single examples using parameter update methods like full-parameter fine-tuning or PEFT methods.

tuning is praised for its ability to bring permanent memorization to models (Hu et al., 2023), and it can perform well even with a small amount of training data (Liu et al., 2022). However, critics argue that fine-tuning demands substantial computational resources (Hu et al., 2021) and can encounter issues such as catastrophic forgetting (Zhai et al., 2023). This conserves computational resources but necessitates longer prompts and incurs higher inference costs.

How about ICL? It is favored for its training-free nature (Dong et al., 2022), allowing prompts to be easily changed for adaptation to other domains without re-training (Min et al., 2022). Other works (Bhattamishra et al., 2023) showed that ICL can help the model uniquely identify a discrete function sample-efficiently. (Reddy, 2023) showed

that ICL is driven by the abrupt emergence of an induction head, which subsequently competes with in-weights learning. (Shen et al., 2024) observes that ICL and gradient descent modify the output distribution of language models differently. Despite these advantages, ICL is limited by context length restrictions and incurs higher costs during each inference stage due to the longer prompts required.

Essentially, the primary distinction between fine-tuning and ICL lies in parameter updating; all fine-tuning methods modify the model’s parameters. It might seem, therefore, that ICL’s impact is less profound. However, our research reveals a counterintuitive finding: **for datasets with implicit patterns, ICL is more adept at uncovering these latent patterns than fine-tuning.**

To investigate this phenomenon, we designed datasets containing implicit patterns across various domains, including two mathematical tasks: expression calculation and boolean function, one textual task: relation reasoning, and one code reading task. These domains share a common trait: the presence of implicit patterns that can simplify problem-solving. We evaluated LLMs’ capability to recognize such patterns with these datasets. Our findings include: (1) Both fine-tuning and ICL could detect and utilize implicit patterns, resulting in increased test accuracy. (2) ICL performed much better than fine-tuning in implicit pattern detection, *e.g.*, ICL-based models enjoyed higher test accuracy. (3) ICL also showed strong performance in robustness tests and OOD data tests. Our experiments demonstrate that the ability of LLMs to leverage implicit patterns significantly enhances their problem-solving capabilities, providing a clear advantage for tasks involving complex data structures.

Understanding the operational principles of LLMs is crucial for their safety and ethical implications and can further promote improvements. Therefore, we delved deeper into the mechanisms behind this phenomenon. From a mechanistic interpretability perspective (Reddy, 2023), we proposed the **Circuit Shift** theory. Circuits are certain groups of attention heads and MLP layers (Conmy et al., 2023). A shift in circuits typically represents the model adopting a different method in problem-solving. Our findings indicated that ICL resulted in a larger-scale circuit shift compared to fine-tuning, which means that with ICL, the model changed its problem-solving method more significantly for implicit pattern detection and utilization. We also

provided a visualized heatmap of circuits for detailed observation. In summary, our contributions are threefold:

Implicit Pattern Detection dataset. We defined and illustrated the implicit pattern detection task, then developed a dataset across mathematics (expression calculation, boolean function), textual reasoning (relation test) and code (output guessing).

Ability Comparison. We presented a counterintuitive finding: LLMs with in-context learning detected implicit patterns much better than fine-tuned ones. We extensively tested this capability on models ranging from 0.5B to 7B parameters.

Mechanism explanation. We analyzed the principles behind the implicit finding mechanism. And we proposed circuit shift theory to explain why ICL finds implicit patterns better than fine-tuning.

2 Background

Transformer. Transformer (Vaswani et al., 2017) is the cornerstone architecture for LLMs nowadays, with its breathtaking ability in parallel training and sota performance. One Transformer model f_{trf} usually consists multiple of Transformer layers f_{layer} and an embedding layer f_{emb} . For an input sequence (typically IDs after tokenization) $\mathbf{X}_0 \in R^{n \times 1}$ with length n , it first passes through an embedding layer f_{emb} with hidden state size d , then passes all the Transformer layers, and finally gets an output $\mathbf{O}_l \in R^{n \times d}$ with l layers: $\mathbf{O}_l = f_{\text{trf}}(\mathbf{X}_0) = \left(\bigcirc_{i=1}^l f_{\text{layer}}^{(i)} \right) (\mathbf{X}_0)$, where for each layer f_{layer} , it usually contains an Attention block and an MLP block:

$$\mathbf{O}_i^{\text{att}} = \mathbf{X}_i + \text{Attn}(\text{Norm}(\mathbf{X}_i)), \quad (1)$$

$$\mathbf{O}_i = \mathbf{O}_i^{\text{att}} + \text{MLP}(\text{Norm}(\mathbf{O}_i^{\text{att}})). \quad (2)$$

Here, $\mathbf{X}_i^{\text{att}}$ is the output of the attention block, and $\mathbf{X}_i^{\text{mlp}}$ is the output of the MLP block for layer i , with residual connections preventing it from gradient disappearance and normalization (typically pre-norm) for stabilizing the training process.

Fine-tuning. Fine-tuning is a process where a pre-trained LLM is further trained on a specific task or dataset to improve its performance for that particular application. Suppose there exists a pre-trained Transformer model f_{trf} with learnable parameters θ_{pre} . The goal of fine-tuning is to adjust these parameters to minimize a task-specific loss function

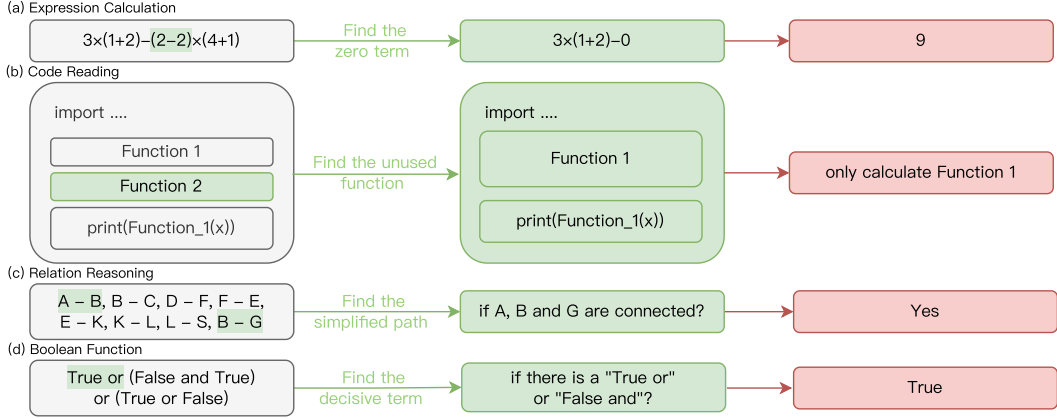


Figure 2: Examples of implicit pattern detection for four reasoning tasks. The implicit pattern, once detected, can reward the model with reduced computation to arrive at the answer.

$\mathcal{L}_{\text{task}}$ on a new dataset $\mathcal{D}_{\text{task}}$. During fine-tuning, the parameters θ_{fine} of the model are updated using gradient descent or one of its variants. The update rule for the parameters at each iteration t can be expressed as:

$$\theta_{\text{fine}}^{(t+1)} = \theta_{\text{fine}}^{(t)} - \eta \nabla_{\theta} \mathcal{L}_{\text{task}}(f_{\text{trf}}(\mathbf{X}_t; \theta_{\text{fine}}^{(t)}), \mathbf{Y}_t), \quad (3)$$

where η is the learning rate, \mathbf{X}_t represents the input data in iteration t , \mathbf{Y}_t represents the target labels in iteration t , and $\nabla_{\theta_{\text{fine}}} \mathcal{L}_{\text{task}}$ denotes the gradient of the loss function with respect to the model parameters. Fine-tuning typically requires substantial computational resources. For instance, full-parameter fine-tuning of LLaMA-3 with 8 billion parameters and an 8K context using the Adam optimizer and gradient checkpointing demands a minimum of 152 GB of VRAM (Rasley et al., 2020), which equates to at least two A100 80 GB GPUs with parallel training. While parameter-efficient fine-tuning (PEFT) is less resource-intensive compared to full-parameter fine-tuning, it still requires 16 GB of VRAM (QLoRA with a 1K context (Detmers et al., 2024)), necessitating at least one RTX 3090 GPU. Additionally, some studies have shown that PEFT can result in a noticeable drop in the model’s performance (Pu et al., 2023; Zou et al., 2023).

In-Context Learning In-Context Learning (ICL) in LLMs is an emergent capability where the model uses the provided context to perform tasks. Given a special task F and a series of prompt inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$, ICL happens when these inputs and their answers $\mathbf{y}_1 = F(\mathbf{x}_1)$ are given in multi-shot, *i.e.*, $(\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{y}_n, \mathbf{x}_{n+1})$. In this scenario, the goal for LLM to do ICL is to learn the task F and accurately predict \mathbf{y}_{n+1} . This phenomenon

allows the model to adaptively handle a variety of tasks, such as translation, question-answering, and more, simply through appropriate prompt engineering. ICL happens in inference-stage without explicit re-training, thus resulting in more friendly requirements for GPUs (Yin et al., 2024; Hong et al., 2023). Even LLaMA-3 70B could run on a single 3090 GPU with PowerInfer (Song et al., 2023).

3 Implicit Pattern Detection Test

Through detailed observation and thinking, humans could detect some underlying, non-explicit patterns within the data. This enables us to solve problems more efficiently. Implicit pattern detection refers to the ability of models to recognize underlying, non-explicit patterns within data, enabling them to solve problems more efficiently. This concept is illustrated through tasks such as arithmetic calculations, where the model can bypass complex operations by identifying simplifying patterns. For instance, in mathematical expressions (see Figure 1 and Figure 2), a model might detect that certain terms have negligible impact and can be ignored, leading to quicker computations. We will give a detailed description of our dataset design and experimental settings in the following sections.

3.1 Tasks

To effectively assess the ability of LLMs to identify implicit patterns in data, we have constructed a variety of questions that frequently arise in real-world application scenarios. When the same type of question recurs, we can discover a specific implicit pattern within it to simplify the computational

224	process.	
225	Task 1: Expression Calculation (Imani et al., 2023; Yuan et al., 2023; Yue et al., 2023; He-Yueya et al., 2023)	
226	In the arithmetic calculation	
227	task, the primary focus is on determining whether	
228	certain operations within a given expression can be	
229	disregarded to reduce the complexity of the compu-	
230	tation. The operations considered for these simpli-	
231	fications are limited to addition(+), subtraction(-),	
232	multiplication(\times), and division(/). By exploring	
233	these operations, the model may find that several	
234	terms are multiplied by a continued-to-be-zero	
235	term, and ignoring them could simplify the cal-	
236	culation process and improve the accuracy.	
237		
238	Task 2: Code Reading (Fang et al., 2024)	
239	In the code reading task, LLMs need to analyze and	
240	predict the output of a given piece of code without	
241	executing it, where multiple functions are defined.	
242	Some functions will not influence the final output,	
243	so the key challenge is to determine which func-	
244	tions are essential for producing the output and	
245	which can be disregarded without affecting the re-	
246	sult.	
247	Task 3: Boolean Functions (Zhang et al., 2024)	
248	In the Boolean functions task, the primary objec-	
249	tive is to optimize logical expressions to simplify	
250	their structure without altering the resultant truth	
251	value. The expressions involve logical operators	
252	such as AND (\wedge), OR (\vee), and NOT (\neg). Within	
253	these scenarios, there are specific segments that are	
254	either tautologies, <i>i.e.</i> , always true, or contradic-	
255	tions, <i>i.e.</i> , always false. The model must identify	
256	these segments and bypass their computation.	
257	Task 4: Relation Reasoning (Li et al., 2024)	
258	In the task of relation reasoning, the focus is on de-	
259	termining the relationships between multiple enti-	
260	tities, such as reachability and relative magnitude.	
261	Although the set of relationships involved can be	
262	complex, all queries target fixed entities whose re-	
263	lationships are relatively straightforward. Therefore,	
264	most of the complex relationships can be disre-	
265	garded, simplifying the problem-solving process.	
266		
267	3.2 Settings	
268	Accuracy. Our tasks were constructed such that	
269	implicit patterns can help solve problems more	
270	easily. For example, if an LLM identifies a term	
271	that continues to be zero in arithmetic calculations,	
	it can ignore terms multiplied by it, thereby saving	
	computation. Therefore, we evaluate the model’s	272
	performance with Accuracy .	273
	Misleading Data. LLMs can detect the inner im-	274
	plicit patterns in data and utilize them for simpli-	275
	fying problem-solving. The misleading data is de-	276
	signed to test if LLMs can tackle situations in the	277
	absence of implicit patterns. While implicit pat-	278
	terns are still provided in training or ICL data, mis-	279
	leading data, <i>i.e.</i> , data with no implicit patterns,	280
	is provided for testing accuracy. We name this ac-	281
	curacy Misleading Accuracy , while the testing re-	282
	sults of data with implicit patterns are named Clean	283
	Accuracy . Detailed experimental procedures can	284
	be found in Appendix B.	285
	Out-Of-Distribution Data. The training data are	286
	sampled from a certain distribution, <i>e.g.</i> , , for ex-	287
	pression tasks, there are no more than 10 terms in	288
	each expression. Our out-of-distribution (OOD)	289
	data are designed to evaluate the model’s perfor-	290
	mance when encountering OOD data during the	291
	evaluation phase. Detailed experimental proce-	292
	dures can be found in Appendix C.	293
	Models. We select open-sourced models in sizes	294
	of 0.5B level <i>e.g.</i> , Qwen1.5-0m5B, 1B level	295
	<i>e.g.</i> , GPTNeo-1.3B (Black et al., 2021), Pythia-	296
	1.4B (Biderman et al., 2023), Qwen1.5-1.8B (Bai	297
	et al., 2023), and 7B level <i>e.g.</i> , Mistral-7B (Jiang	298
	et al., 2023), Qwen1.5-7B (Bai et al., 2023), Yi-	299
	6B (Young et al., 2024). Model weights are down-	300
	loaded from Huggingface and follow the official	301
	implementations.	302
	Data Format. For fine-tuning, the data is pro-	303
	vided in a single example without supervised in-	304
	struction. A simple description, the question, and	305
	the answer are given in order. We prepared 1,600	306
	data points for fine-tuning. For in-context learn-	307
	ing, we constructed the input in multi-shot, rang-	308
	ing from 0-shot, <i>i.e.</i> , directly answer one question, to	309
	32-shot <i>i.e.</i> , 32 examples with their answers first	310
	given, then a new question in the same kind re-	311
	quired to answer. The detailed example of our data	312
	format could be found in Appendix A.	313
	Training Details. The training process was con-	314
	ducted using a sequence length of 512 and a batch	315
	size of 8 with a total of 1 epoch. A warmup phase of	316
	20 steps was implemented, starting with a learning	317
	rate of 1e-6 and peaking at 2e-5, followed by a lin-	318
	ear decay. The AdamW optimizer was used. This	319
	configuration ensured the model’s performance and	320

Model	Expression			Code			Relation			Boolean		
	Baseline	Full-ft	ICL	Baseline	Full-ft	ICL	Baseline	Full-ft	ICL	Baseline	Full-ft	ICL
<i>0.5B level</i>												
Qwen1.5-0.5B	22.2%	88.4%	50.1%	16.6%	2.0%	32.2%	48.8%	48.5%	60.1%	54.8%	51.7%	65.3%
<i>1B level</i>												
GPTNeo-1.3B	24.3%	46.6%	55.6%	27.6%	17.7%	44.5%	20.5%	34.7%	37.4%	53.8%	53.7%	54.3%
Qwen1.5-1.8B	16.2%	89.9%	63.4%	54.3%	53.7%	58.2%	20.1%	21.3%	35.6%	66.3%	66.3%	68.1%
Pythia-1.4B	5.0%	45.4%	53.7%	37.6%	46.5%	53.1%	20.5%	31.3%	44.4%	61.3%	63.7%	68.5%
<i>7B level</i>												
Yi-6B	12.5%	88.2%	48.2%	51.2%	78.7%	80.9%	48.0%	52.5%	98.0%	55.7%	64.1%	68.3%
Qwen1.5-7B	78.0%	89.3%	67.9%	57.6%	72.0%	86.8%	48.0%	78.8%	98.0%	71.9%	41.7%	79.8%
Mistral-7B	32.6%	75.2%	76.3%	14.1%	72.0%	82.8%	48.5%	72.5%	90.9%	45.7%	54.5%	74.3%

Table 1: Experimental results of implicit pattern detection tasks. We conducted experiments from 0.5B to 7B across 6 models. The highest accuracy was highlighted with boldsymbol.

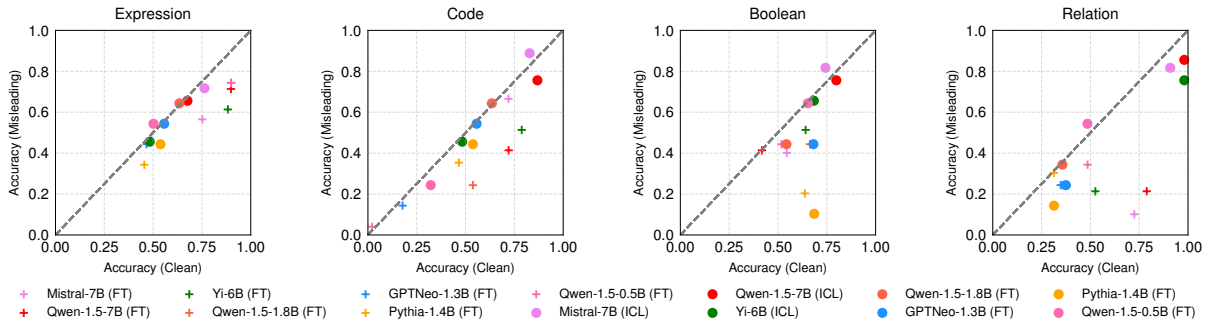


Figure 3: Robustness test of implicit pattern detection test. The horizontal axis represents the accuracy under clean input, and the vertical axis represents the accuracy under misleading input. Relatively speaking, the closer the results are to the bottom right corner, the worse the method’s resistance to misleading data. The closer the results are to the top left corner, the better it is.

321 stability, allowing it to effectively learn and identify
322 hidden patterns in the data.

323 4 Results and Analysis

324 In this section, we present our results for the im-
325 plicit pattern finding tasks following the experi-
326 mental setting in Section 3.2. We show that ICL
327 achieved an overall higher level of accuracy over
328 fine-tuning on these four tasks. We also show
329 that the improvement of accuracy with ICL mainly
330 comes from the detection of those implicit patterns
331 in Section 5 and refsec:circuit.

332 4.1 ICL v.s. Fine-tuning: Accuracy

333 The results of accuracy test are shown in Table 1
334 and Table 2. Both ICL and fine-tuning(including
335 full-param fine-tuning and PEFT methods) bring
336 improvements to the performance of each task. How-
337 ever, it is easily noticed that ICL wins at most terms
338 like relation, code reading and boolean functions,
339 with 2% to even more than 30% improvements
340 at most. On the flip side, fine-tuning only shows

Method Type	Expression	Code	Relation	Boolean
Baseline	27.5%	54.3%	20.1%	66.3%
Full-Param FT	89.9%	53.7%	21.3%	66.3%
LoRA	46.5%	53.3%	20.1%	64.3%
QLoRA	46.2%	51.6%	20.5%	61.3%
GaLoRA	47.1%	52.5%	20.5%	66.4%
ICL	63.4%	58.2%	35.6%	68.1%

Table 2: Experimental comparison of different PEFT methods. We compared the results on Qwen1.5-1.8B. It is obvious that PEFT shows no significant improvement compared to full-param fine-tuning and seems to have limited performance.

341 slight advantages in expression calculations in only
342 Qwen-series models. As for different model size¹,
343 we found that a larger model seems be able to evoke
344 stronger ICL ability above linearly growth (see
345 Table 1), where the scaling of fine-tuning perfor-
346 mance is limited.

¹See Qwen1.5 series in Table 1 from 0.5B to 7B

OOD Type	Expression	Code	Relation	Boolean
Baseline	27.5%	54.3%	20.1%	66.3%
FT	89.9%	53.7%	21.3%	66.3%
FT + Test OOD	32.1%	34.2%	0.1%	0.1%
(FT+Test) OOD	88.2%	42.7%	11.3%	12.4%
ICL	63.4%	58.2%	35.6%	68.1%
ICL + Test OOD	34.5%	44.2%	12.3%	24.7%
(ICL+Test) OOD	62.3%	51.7%	34.5%	71.4%

Table 3: Experimental comparison of different PEFT methods. Here FT/ICL + Test OOD means we only applied OOD data in test phase, while (FT/ICL) OOD represents that both training/in-context learning and test phase were using OOD data.

4.2 ICL v.s. Fine-tuning: Robustness without Implicit Pattern

In Section 3.2, we introduced the metrics of clean accuracy and misleading accuracy by adding misleading data to test both ICL and fine-tuning’s robustness against general data without implicit patterns. The results are shown in Figure 3. For each task, we draw a scatter plot where the x- and y-axis represent the clean accuracy and the misleading accuracy, respectively. The results show that ICL can better exploit the implicit patterns in the demonstration data, while at the same time not compromising general reasoning abilities.

4.3 ICL v.s. Fine-tuning: Out-Of-Distribution Implicit Patterns

Out-of-Distribution (OOD) data is a widely examined problem nowadays. The training data of our implicit pattern detection tasks also samples from certain distributions (see Appendix C for details). In this subsection, we hope to compare how ICL and fine-tuning perform if we provide cases outside of the training distribution. For ICL, all examples given are divided into two types: in-distribution examples and OOD examples. For fine-tuning, we directly provide OOD problems to test the accuracy. We performed this experiment on Qwen1.5-1.8B and the results are demonstrated in Table 3. It is worth noticing that fine-tuning generally performs worse when the test data is OOD, while ICL performs fairly well comparing to the baseline method.

4.4 How Much Fine-tuning Do We Need?

In this experiment, we hope to figure out whether fine-tuning has reached its limit for implicit pattern detection or there will still be improvement if more data is utilized for fine-tuning. Therefore, we visu-

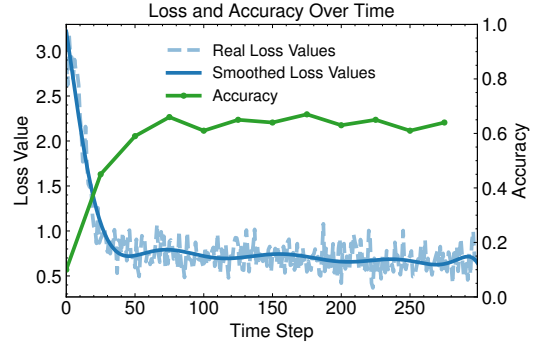


Figure 4: The progression of loss and accuracy over time during the fine-tuning of implicit pattern tasks. The Real Loss Values (dashed blue line) show the loss during training. To mitigate this noise, the Smoothed Loss Values (solid blue line) provide a clearer trend of the overall loss reduction. We also show the average test accuracy over all tasks (solid green line).

alized the fine-tuning process of Qwen1.5-1.8B. At the onset of training, there is a steep decline in the loss value, suggesting that the model quickly learns basic patterns in the data. This rapid improvement is typical, as the model captures the most evident features. The Accuracy (solid green line) also increases sharply, corroborating the initial learning phase where the model transitions from random guessing to meaningful predictions. However, after around 50 time steps, both the loss and accuracy curves begin to stabilize. This period of stabilization suggests diminishing returns from further training, as the fine-tuned model failed to capture further implicit patterns. After 100 time steps, the curves indicate that the model has reached a plateau. The accuracy remains relatively constant, and the loss value shows minimal fluctuations around a stable trend. This behavior signifies that the model has learned the underlying patterns to a satisfactory extent, and additional fine-tuning yields marginal improvements.

4.5 Comparison of Fine-tuning with PEFT Methods

Lastly, we examine whether there is a significant difference between various fine-tuning methods *e.g.*, vanilla full-parameter fine-tuning, and parameter efficient fine-tuning (PEFT) methods like LoRA (Hu et al., 2021), QLoRA (Detmers et al., 2024) and GaLoRE (Zhao et al., 2024). Although PEFT needs much less parameters for training, and several studies criticized its ability (Pu et al., 2023; Zou et al., 2023), there are still evidences that PEFT sometimes achieves ICL-level performance. We

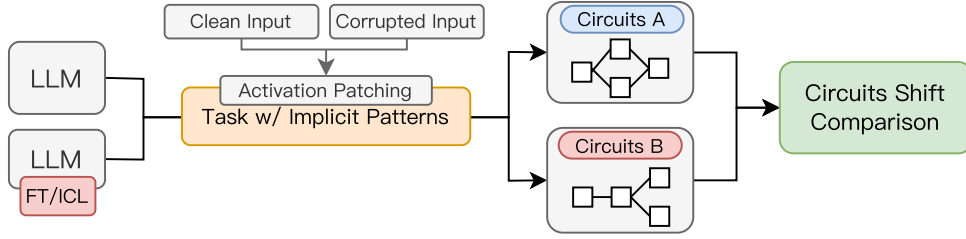


Figure 5: Illustration of circuit shift comparison. LLMs are first detected circuits with activation patching. Then we compare how much their circuits changed after fine-tuning and in-context learning.

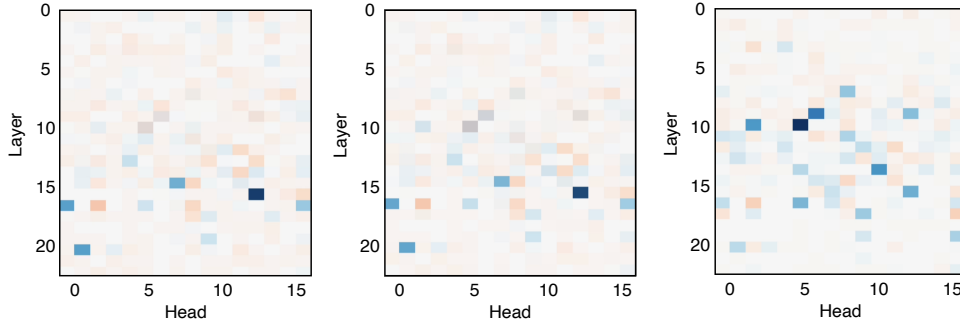


Figure 6: Visualization of attention head sensitivity in GPTNeo-1.3B. The more the color leans towards blue, the more important a specific attention head is to the implicit pattern detection task. Left: baseline model. Middle: fine-tuned model. Right: ICL model. It is clear that compared to fine-tuning, ICL brings significant circuit shifts.

415 followed the experimental settings in previous sections on Qwen1.5-1.8B with PEFT methods. The
 416 experimental results can be found in Table 2. It
 417 is clear that in the implicit pattern detection tasks,
 418 PEFT methods show no obvious advantages compared
 419 to full-param fine-tuning, thus they still failed
 420 to win ICL in accuracy in all tests.
 421

422 5 Explanation of ICL’s Victory: Circuits 423 Shift Theory

424 Understanding the inner mechanisms of LLMs
 425 greatly benefits their ethical use and safety. We
 426 have found that ICL performs much better than
 427 fine-tuning on implicit pattern detection, and in
 428 this section, we try to explain why.

429 From a mechanistic interpretability perspective,
 430 we investigate this problem using **circuits**. Circuits
 431 are specific pathways (typically combinations of
 432 attention heads and MLP layers) within a model
 433 responsible for processing and interpreting partic-
 434 ular patterns or tasks. The change in circuits for
 435 LLMs represents a shift in their inner mechanisms,
 436 revealing that LLMs choose different ways to solve
 437 problems. Based on this viewpoint, we propose a
 438 theory: **Circuits Shift**, to explain this phenomenon.
 439 We will first provide a method for probing circuits,
 440 explaining what they are and the types of circuits

441 we found in ICL-based and fine-tuning-based mod-
 442 els. Then we will show that the reason ICL per-
 443 forms better than fine-tuning is that the circuits
 444 in models experience a more significant shift. A
 445 detailed explanation of circuits and experimental
 446 settings can be found in Appendix D.

447 5.1 Method for Identifying Circuit Shift

448 In Figure 5, we present our framework and method-
 449 ology for probing circuit shifts. We begin by se-
 450 lecting an implicit pattern detection task (in this
 451 study, we utilize an expression task). Subsequently,
 452 we use models employing different methods, *i.e.*,
 453 ICL or fine-tuning, for inference. During this pro-
 454 cess, we introduce corrupt input to randomly dis-
 455 rupt a portion of the activation to assess whether
 456 the corresponding attention heads or MLP layers
 457 significantly contribute to the final outcome. If a
 458 significant contribution exists, the disruption will
 459 result in considerable perturbation of the final log-
 460 its, which is depicted as sensitivity in the figure.

461 5.2 Circuits Shift in LLMs for Implicit 462 Pattern Detection

463 We first visualized and ranked circuits in GPTNeo-
 464 1.3B zero-shot, after fine-tuned, and ICL with 32-
 465 shot with expression calculation task (see Figure 6

Circuits	Zero-shot Baseline	ICL w/o Implicit Patterns	Δ	After Fine-tuning	Δ	After ICL	Δ
Attention	L17 H12, L18 H0 L22 H1, L16 H7 L18 H15, L14 H5	L17 H12, L16 H1 L18 H0, L15 H2 L18 H15, L22 H1	2	L17 H12, L18 H0 L22 H1, L16 H7 L18 H15, L12 H6	1	L11 H5, L10 H6 L11 H2, L15 H10 L17 H12, L 18 H5	6
MLP	L9 L17 L18	L9 L17 L18	0	L9 L18 L17	0	L17 L14 L15	2

Table 4: Top 6 Rankings of Attention Heads and top 3 rankings of MLP Layers in baseline (zero-shot) model, fine-tuned model, and ICL model. L is layer and H is head. Δ shows how many different heads or MLPs changed after fine-tuning or ICL. A larger Δ represents a more significant circuit shift in certain processes.

and Table 4). In Figure 6, we use the heatmap to illustrate the sensitivity of each attention head in implicit pattern detection test. From the figure, we can observe that, compared to the baseline and fine-tuning scenarios, ICL exhibits a significant shift when learning implicit patterns. Firstly, more shallow heads are involved in the task. Secondly, some deep heads that previously played a dominant role have now lost their leadership positions. This indicates that during the ICL process, the model significantly transforms its approach to solving the task, adapting to a form more suitable for implicit patterns, a phenomenon not observed with other methods.

We can further validate our hypothesis in Table 4. We selected the six attention heads and MLP layers² with the highest sensitivity, *i.e.*, those that contributed the most to the final result. Using the baseline, which is the zero-shot approach for handling implicit pattern detection tasks, as the standard, we counted how many new attention heads entered the top six highest contributors when the method changed, denoted by Delta. The results are very clear: compared to fine-tuning, ICL exhibits more significant changes, indicating a more thorough Circuit Shift during ICL. This suggests that ICL captures the characteristics of implicit patterns better than fine-tuning and adapts its processing method accordingly.

To rule out the inherent impact of ICL itself, we also conducted multi-shot experiments on a set of data without implicit pattern characteristics. The results showed that it is not multi-shot alone that induces this change, but rather the combined effect of ICL and implicit patterns.

²See A Mathematical Framework for Transformer Circuits for details.

6 Related Work

Implicit Pattern Discovery Previous works have designed benchmarks to test the LLMs reasoning ability (Barrett et al., 2018; Tang et al., 2023; Gendron et al., 2024). However, the benchmarks rarely include two-level questions where at one level, they can be solved by brutal force, at another level it can be solved by exploiting implicit patterns. The closest related work we know is Efrat et al. (2021), which involves solving cryptic crossword puzzles. To help the model find patterns in data, Prior work Sun et al. (2024); Zhu et al. (2024) proposes a two-stage induction-deduction process that first summarizes the common patterns explicitly, then reasons from the patterns.

ICL v.s. Fine-tuning Difference Previous works have also compared fine-tuning and in-context learning. Shen et al. (2024) shows that ICL is likely not an algorithmic equivalence to gradient descent for real LLMs. Reddy (2023) demonstrates that ICL is implemented by an induction head and analyzes its emergence phenomenon. Bhattamishra et al. (2023) shows that ICL and vanilla training implement two distinct algorithms that don’t transfer to each other. However, it has been proven that fine-tuning shows better performance in generalization to OOD tasks than in-context learning (Mosbach et al., 2023).

7 Conclusion

In conclusion, our research demonstrates that In-Context Learning (ICL) significantly outperforms fine-tuning in capturing implicit patterns within specific tasks. Through our experimental evaluations, we observed that ICL not only enhances task performance more effectively but also exhibits greater adaptability in problem-solving approaches, as evidenced by the notable shifts in model circuits.

538 Limitations

539 Our study on the effectiveness of in-context learn-
540 ing in capturing implicit patterns compared to fine-
541 tuning faces several limitations. Primarily, the
542 generalizability of our findings is constrained by
543 the specific nature of the implicit pattern detec-
544 tion tasks, which are limited to certain domains
545 like arithmetic calculations, code reading, Boolean
546 functions, and relation reasoning. Additionally, our
547 analysis of Circuit Shift, which underpins the supe-
548 rior performance of ICL, relies on activation patch-
549 ing and sensitivity analysis, methods that, while
550 insightful, require further refinement and valida-
551 tion across different models and tasks to confirm
552 their robustness and applicability. Furthermore, the
553 computational resources required for fine-tuning,
554 especially with large models, may limit the feasi-
555 bility of such experiments in broader settings, and
556 a detailed cost-benefit analysis comparing ICL and
557 fine-tuning in terms of computational efficiency
558 and performance is needed.

559 References

560 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
561 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
562 Huang, et al. 2023. Qwen technical report. *arXiv*
563 *preprint arXiv:2309.16609*.

564 David G. T. Barrett, Felix Hill, Adam Santoro, Ari S.
565 Morcos, and Timothy Lillicrap. 2018. *Measuring*
566 *abstract reasoning in neural networks*. *Preprint*,
567 *arXiv:1807.04225*.

568 Satwik Bhattamishra, Arkil Patel, Phil Blunsom, and
569 Varun Kanade. 2023. *Understanding in-context learn-*
570 *ing in transformers and llms by learning to learn dis-*
571 *crete functions*. *Preprint*, *arXiv:2310.03016*.

572 Stella Biderman, Hailey Schoelkopf, Quentin Gregory
573 Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-
574 lahan, Mohammad Aflah Khan, Shivanshu Purohit,
575 USVSN Sai Prashanth, Edward Raff, et al. 2023.
576 Pythia: A suite for analyzing large language mod-
577 els across training and scaling. In *International*
578 *Conference on Machine Learning*, pages 2397–2430.
579 PMLR.

580 Sid Black, Leo Gao, Phil Wang, Connor Leahy,
581 and Stella Biderman. 2021. *GPT-Neo: Large*
582 *Scale Autoregressive Language Modeling with Mesh-*
583 *Tensorflow*. If you use this software, please cite it
584 using these metadata.

585 Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch,
586 Stefan Heimersheim, and Adrià Garriga-Alonso.
587 2023. Towards automated circuit discovery for mech-
588 anistic interpretability. *Advances in Neural Informa-*
589 *tion Processing Systems*, 36:16318–16352.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and
590 Luke Zettlemoyer. 2024. Qlora: Efficient finetuning
591 of quantized llms. *Advances in Neural Information*
592 *Processing Systems*, 36. 593

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiy-
594 ong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and
595 Zhifang Sui. 2022. A survey on in-context learning.
596 *arXiv preprint arXiv:2301.00234*. 597

Avia Efrat, Uri Shaham, Dan Kilman, and Omer Levy.
2021. *Cryptonite: A cryptic crossword benchmark*
598 *for extreme ambiguity in language*. In *Proceedings*
599 *of the 2021 Conference on Empirical Methods in Nat-*
600 *ural Language Processing*, pages 4186–4192, Online
601 and Punta Cana, Dominican Republic. Association
602 for Computational Linguistics. 603 604

Chongzhou Fang, Ning Miao, Shaurya Srivastav, Jialin
605 Liu, Ruoyu Zhang, Ruijie Fang, Asmita, Ryan Tsang,
606 Najmeh Nazari, Han Wang, and Houman Homayoun.
2024. *Large language models for code analysis: Do*
607 *llms really do their job?* *Preprint*, *arXiv:2310.12357*. 608 609

Chengguang Gan and Tatsunori Mori. 2023. A few-
610 shot approach to resume information extraction via
611 prompts. In *International Conference on Applica-*
612 *tions of Natural Language to Information Systems*,
613 pages 445–455. Springer. 614

Gaël Gendron, Qiming Bao, Michael Witbrock, and
615 Gillian Dobbie. 2024. *Large language mod-*
616 *els are not strong abstract reasoners*. *Preprint*,
617 *arXiv:2305.19555*. 618

Joy He-Yueya, Gabriel Poesia, Rose E. Wang, and
619 Noah D. Goodman. 2023. *Solving math word prob-*
620 *lems by combining language models with symbolic*
621 *solvers*. *Preprint*, *arXiv:2304.09102*. 622

Ke Hong, Guohao Dai, Jiaming Xu, Qiuli Mao, Xi-
623 uhong Li, Jun Liu, Kangdi Chen, Hanyu Dong, and
624 Yu Wang. 2023. Flashdecoding++: Faster large
625 language model inference on gpus. *arXiv preprint*
626 *arXiv:2311.01282*. 627

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
628 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
629 and Weizhu Chen. 2021. Lora: Low-rank adap-
630 tation of large language models. *arXiv preprint*
631 *arXiv:2106.09685*. 632

Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-
633 Peng Lim, Lidong Bing, Xing Xu, Soujanya Po-
634 ria, and Roy Ka-Wei Lee. 2023. Llm-adapters:
635 An adapter family for parameter-efficient fine-
636 tuning of large language models. *arXiv preprint*
637 *arXiv:2304.01933*. 638

Shima Imani, Liang Du, and Harsh Shrivastava. 2023.
639 *Mathprompter: Mathematical reasoning using large*
640 *language models*. *Preprint*, *arXiv:2303.05398*. 641

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-
642 sch, Chris Bamford, Devendra Singh Chaplot, Diego 643

644	de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	698 699 700 701 702
647	Zhiming Li, Yushi Cao, Xiufeng Xu, Junzhe Jiang, Xu Liu, Yon Shin Teo, Shang wei Lin, and Yang Liu. 2024. <i>Llms for relational reasoning: How far are we?</i> <i>Preprint</i> , arXiv:2401.09042.	Qingyu Yin, Xuzheng He, Xiang Zhuang, Yu Zhao, Jianhua Yao, Xiaoyu Shen, and Qiang Zhang. 2024. Stablemask: Refining causal masking in decoder-only transformer. <i>arXiv preprint arXiv:2402.04779</i> .	703 704 705 706
651	Ziquan Liu, Yi Xu, Yuanhong Xu, Qi Qian, Hao Li, Xiangyang Ji, Antoni Chan, and Rong Jin. 2022. Improved fine-tuning by better leveraging pre-training data. <i>Advances in Neural Information Processing Systems</i> , 35:32568–32581.	Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. <i>arXiv preprint arXiv:2403.04652</i> .	707 708 709 710 711
656	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? <i>arXiv preprint arXiv:2202.12837</i> .	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. 2023. <i>How well do large language models perform in arithmetic tasks?</i> <i>Preprint</i> , arXiv:2304.02015.	712 713 714 715
661	Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. <i>arXiv preprint arXiv:2305.16938</i> .	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2023. <i>Mammoth: Building math generalist models through hybrid instruction tuning</i> . <i>Preprint</i> , arXiv:2309.05653.	716 717 718 719 720
665	Matthew E Peters, Sebastian Ruder, and Noah A Smith. 2019. To tune or not to tune? adapting pretrained representations to diverse tasks. <i>arXiv preprint arXiv:1903.05987</i> .	Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the catastrophic forgetting in multimodal large language models. <i>arXiv preprint arXiv:2309.10313</i> .	721 722 723 724
669	George Pu, Anirudh Jain, Jihan Yin, and Russell Kaplan. 2023. <i>Empirical analysis of the strengths and weaknesses of peft techniques for llms</i> . <i>Preprint</i> , arXiv:2304.14999.	Yu Zhang, Hui-Ling Zhen, Zehua Pei, Yingzhao Lian, Lihao Yin, Mingxuan Yuan, and Bei Yu. 2024. <i>Dila: Enhancing llm tool learning with differential logic layer</i> . <i>Preprint</i> , arXiv:2402.11903.	725 726 727 728
673	Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In <i>Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> , pages 3505–3506.	Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. 2024. Galore: Memory-efficient llm training by gradient low-rank projection. <i>arXiv preprint arXiv:2403.03507</i> .	729 730 731 732 733
679	Gautam Reddy. 2023. <i>The mechanistic basis of data dependence and abrupt learning in an in-context classification task</i> . <i>Preprint</i> , arXiv:2312.03002.	Zhaocheng Zhu, Yuan Xue, Xinyun Chen, Denny Zhou, Jian Tang, Dale Schuurmans, and Hanjun Dai. 2024. <i>Large language models can learn rules</i> . <i>Preprint</i> , arXiv:2310.07064.	734 735 736 737
682	Lingfeng Shen, Aayush Mishra, and Daniel Khashabi. 2024. <i>Do pretrained transformers learn in-context by gradient descent?</i> <i>Preprint</i> , arXiv:2310.08540.	Wentao Zou, Qi Li, Jidong Ge, Chuanyi Li, Xiaoyu Shen, Liguang Huang, and Bin Luo. 2023. <i>A comprehensive evaluation of parameter-efficient fine-tuning on software engineering tasks</i> . <i>Preprint</i> , arXiv:2312.15614.	738 739 740 741 742
685	Yixin Song, Zeyu Mi, Haotong Xie, and Haibo Chen. 2023. Powerinfer: Fast large language model serving with a consumer-grade gpu. <i>arXiv preprint arXiv:2312.12456</i> .		
689	Wangtao Sun, Haotian Xu, Xuanqing Yu, Pei Chen, Shizhu He, Jun Zhao, and Kang Liu. 2024. <i>ItD: Large language models can teach themselves induction through deduction</i> . <i>Preprint</i> , arXiv:2403.05789.		
693	Xiaojuan Tang, Zilong Zheng, Jiaqi Li, Fanxu Meng, Song-Chun Zhu, Yitao Liang, and Muhan Zhang. 2023. <i>Large language models are in-context semantic reasoners rather than symbolic reasoners</i> . <i>Preprint</i> , arXiv:2305.14825.		

A Data Format and Example

We provided examples of tasks and prompts. We provided data as 2-shot (code in zero-shot to restrict content length) for illustrating how ICL works. For fine-tuning we will use the same format but zero-shot in both training and inference.

Expression:

```
1 Now you need to calculate the answer of
  some mathematic equations.
2 Here are some examples:
3 (1+6)+(-3+3)*(-1-3+9-5)=7
4 (2+3)+(-1-4+5)*(10+6+2-8)=5
5 (8)+(0)*(0-6+9-6)=
```

Listing 1: Example

Code:

```
1 Now you need to give me the printed
  result after running this python
  code . Here are some examples:
```

Listing 2: Example

```
1 def function1(x):
2     y = x ** 9
3     for i in range(1, 13):
4         y = y * i - (y // (i + 9))
5     return y
6
7 def function2(z, a):
8     return z / 10
9
10 input_value = int(input())
11 result = function2(input_value,
12                    function1(input_value))
12 print(result)
```

Listing 3: Example

```
1 The input is 10, so the output is
```

Listing 4: Example

Relation:

```
1 Here are some cities expressed as A, B,
  C, etc. I will show some connection
  relations, and you need to tell me
  if city A and city Z are connected (
  Answer True or False).
2 Here are some examples:
3 A is connected with G
4 F is connected with J
5 J is connected with C
6 C is connected with B
7 B is connected with H
8 H is connected with E
9 E is connected with G
10 G is connected with I
11 I is connected with D
12 So 'the city A and Z is connected' is
  False
13 A is connected with B
14 H is connected with I
15 I is connected with G
```

```
16 G is connected with F
17 F is connected with E
18 E is connected with J
19 J is connected with B
20 B is connected with C
21 C is connected with D
22 B is connected with Z
23 So 'the city A and Z is connected' is
  True
24 A is connected with H
25 J is connected with I
26 I is connected with E
27 E is connected with F
28 F is connected with H
29 H is connected with G
30 G is connected with D
31 D is connected with C
32 C is connected with B
33 So 'the city A and Z is connected' is
```

Listing 5: Example

Boolean:

```
1 Here are some boolean expressions, you
  need to directly tell me the result.
  If it is true, print True, else
  print False. Here are some examples:
2 (True and False) and (True or False) and
  (False and False)\n The result is:
  False
3 (False and False) or (True and True) and
  (False and False)\n The result is:
  False
4 (True or True or True) and (False and
  True) and (True or True)\n The
  result is:
```

Listing 6: Example

B Misleading Data Construction

Expression. For the expression task, the inherent implicit pattern is an element that remains zero. When constructing the misleading dataset, we set this element to be non-zero. *i.e.*,

$$(3 + 2) + (4 - 1 + 5 - 6) \times (23 - 54 + 2) = ?$$

we constructed it as misleading data as:

$$(3 + 2) + (4 - 1 + 5 - 7) \times (23 - 54 + 2) = ?$$

Code. Here we provided two example about how to construct misleading code.

```
1 def function1(x):
2     y = x ** 19
3     for i in range(1, 23):
4         y = y * i - (y // (i + 19))
5     return y
6
7 def function2(z, a):
8     return z / 20
9
10 input_value = int(input())
```

Name	Type	Problem Example	Answer	Answer Type
Expression	Mathematic Calculation	$(6 - 1) + (6 - 6) * (-10 + 1 + 2 + 13) =$	5	Number
Code	Code Reading	<code>import math \n \n def function1(x): \n \n [TRUNCATED] return result \n print(result)</code>	3.5	Number
Relation	Textual Reasoning	A is connected with G\n F is connected[TRUNCATED] connected with Z, 'the city A and Z is connected' is	False	Boolean
Boolean	Mathematical Reasoning	(False or False) and (False or True) and False =	False	Boolean

Table 5: Examples of four implicit pattern detection tasks.

```

11 result = function2(input_value,
12     function1(input_value))
13 print(result)

```

Listing 7: For implicit pattern

```

1 def function1(x):
2     y = x ** 19
3     for i in range(1, 23):
4         y = y * i - (y // (i + 19))
5     return y
6
7 def function2(z, a):
8     return z / 20
9
10 input_value = int(input())
11 result = function2(function1(input_value)
12     ), function1(input_value))
13 print(result)

```

Listing 8: For misleading

Relation. In the relation task, we generate misleading data by not setting shortcuts similar to A-G or G-Z.

```

1 A is connected with B
2 D is connected with B
3 B is connected with H
4 H is connected with F
5 F is connected with J
6 J is connected with I
7 I is connected with C
8 C is connected with G
9 G is connected with E
10 B is connected with Z

```

Listing 9: For implicit pattern

Here A-B-Z is a implicit pattern as shortcut for quick solving this problem. We remove this with a complex one:

```

1 A is connected with B
2 D is connected with B
3 B is connected with H
4 H is connected with F
5 F is connected with J
6 J is connected with I
7 I is connected with C
8 C is connected with G
9 G is connected with E
10 F is connected with Z

```

```

Listing 10: For misleading

```

Boolean. In the boolean task, we use combinations of OR + true and AND + false for quick evaluation. In the misleading data, we remove this characteristic.

```

1 (False and True) or (False or False) or True

```

Listing 11: For implicit pattern

```

1 (False and True) or (False or False) and True

```

Listing 12: For misleading

C OOD data Construction

	Min Terms	Max Terms	Range (abs value)
baseline	1	3	10
OOD	2	4	20

Table 6: Expression OOD

	Functions Need Calculation	Shortcut Nodes
baseline	1	3 (A to Any to G)
OOD	2	Unlimited

Table 7: Code OOD and Relation OOD

	If All AND or OR	Num of Terms
baseline	Yes	4
OOD	No	6

Table 8: Code OOD and Relation OOD

D Circuits

Circuits In mechanistic interpretability, our goal is to delineate how model components correlate

with human-understandable concepts, an endeavor for which circuits provide a useful abstraction. Conceptualizing a model as a computational graph M , where nodes represent components like neurons, attention heads, and embeddings, and edges denote interactions such as residual connections and projections, a circuit C is defined as a subgraph of M responsible for a specific behavior, such as performing a task. This is a more coarse-grained approach compared to the feature-based.

Activation Patching Activation patching is a technique used to determine the importance of specific components within a model by manipulating their latent activations during model runs. The process involves three key steps: first, a *clean run* where the model processes a clean prompt, X_{clean} (e.g., The Eiffel Tower is in), and associated answer r (Paris), during which activations of critical components such as MLP or attention heads are cached; second, a *corrupted run* where the model is run on a corrupted prompt, X_{corrupt} (e.g., The Colosseum is in), to record baseline outputs; and third, a *Patched run* where the model is run on X_{corrupt} again, but with specific cached activations from the X_{clean} run restored. This setup allows for the evaluation of the patching effect, which measures the restoration of model performance by comparing outputs from the Corrupted and Patched runs. The patching effect is quantitatively assessed using different metrics with probability gap:

$$P_{\text{patched}}(r) - P_{\text{corrupt}}(r) \quad (4)$$

and logit difference:

$$LD(r, r') = \log \left(\frac{P(r)}{P(r')} \right)_{\text{patched}} - \log \left(\frac{P(r)}{P(r')} \right)_{\text{corrupt}} \quad (5)$$

This technique is crucial for understanding and improving model reliability and performance by highlighting the roles of individual model components.

E A detailed Definition of Implicit Pattern Detection

Consider a problem P characterized by a fixed complexity function C_P . For each input x in the domain D , there exists a solution y . A implicit pattern for problem P , denoted as P_{shortcut} , is defined as follows:

- P_{shortcut} is either a subproblem of P or an independent problem where the domain D_{shortcut} is a subset of D (i.e., $D_{\text{shortcut}} \subseteq D$).
- For any input x in D_{shortcut} , the output y_{shortcut} of P_{shortcut} approximates the output y of P .
- The complexity of solving P_{shortcut} , $C_{P_{\text{shortcut}}}$, is significantly less than C_P (i.e., $C_{P_{\text{shortcut}}} \ll C_P$).

If these conditions are met, then P_{shortcut} is considered a shortcut of P . We define its complexity C_f in terms of the accuracy of a LLM performing on f . Let Acc_f represent the accuracy of the LLM on task f , then the complexity $C_T f$ can be defined as: $C_T = 1 - \text{Acc}_f$. The complexity C_f ranges from 0 (no complexity, as the task is perfectly solved) to 1 (maximum complexity, as the task is not solved at all).

This definition implies that the higher the LLM's accuracy on a task, the lower the complexity of the task. This measure allows us to quantify task complexity based on the performance capabilities of state-of-the-art language models.