Who Reasons in the Large Language Models?

Jie Shao Jianxin Wu*

National Key Laboratory for Novel Software Technology, Nanjing University, China School of Artificial Intelligence, Nanjing University, China shaoj@lamda.nju.edu.cn, wujx2001@nju.edu.cn

Abstract

Despite the impressive performance of large language models (LLMs), the process of endowing them with new capabilities—such as mathematical reasoning—remains largely empirical and opaque. A critical open question is whether reasoning abilities stem from the entire model, specific modules, or are merely artifacts of overfitting. In this work, we hypothesize that the reasoning capabilities in well-trained LLMs are primarily attributed to the output projection module (o_proj) in the Transformer's multi-head self-attention (MHSA) module. To support this hypothesis, we introduce Stethoscope for Networks (SfN), a suite of diagnostic tools designed to probe and analyze the internal behaviors of LLMs. Using SfN, we provide both circumstantial and empirical evidence suggesting that o_proj plays a central role in enabling reasoning, whereas other modules contribute more to fluent dialogue. These findings offer a new perspective on LLM interpretability and open avenues for more targeted training strategies, potentially enabling more efficient and specialized LLMs.

1 Introduction

Although large language models (LLMs) [29, 6, 42, 5] have exhibited great success and potential in various aspects, developing new capabilities for LLMs [54, 17, 38, 14] is still a trial and error experimentation process in most cases. For example, one of the most exciting milestones is LLMs that can reason [18, 13, 40], e.g., solving complicated mathematical problems using a reasoning sequence that is agreeable by human experts.

This success, however, is still in the black-box style. Currently, there are two primary approaches to inspiring reasoning capabilities in LLMs. For the most advanced models [13, 52], reinforcement learning method (for example, PPO [37], DPO [30], or GRPO [38]) is commonly adopted to enhance the model's ability to solve complex mathematical or programming problems in a step-by-step manner [49]. A more efficient alternative involves supervised fine-tuning (SFT): by providing the backbone LLM with well-prepared, diverse, and step-by-step reasoning traces—often generated through handcrafted examples or existing reasoning models [55, 25, 13, 52]—the model surprisingly acquires reasoning abilities after training. However, despite the practical success of this method, the underlying mechanism remains largely unexplained. It is still unclear why or how this ability emerges. Several potential explanations may account for this phenomenon:

- Case 1 Is it the LLM in its entirety (i.e., the union of all its weights) that leads to this capability, such that this miracle is not explainable?
- Case 2 Or, is there certain module(s) in it that should be praised for this success, such that we can advance our understanding of LLMs?

^{*}Corresponding author.

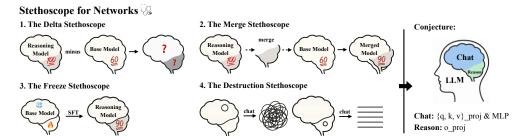


Figure 1: **Stethoscope for Networks.** SfN is a framework designed to identify which components of an LLM give rise to specific abilities. By comparing weight changes and observing behaviors under controlled module merging, tuning, or destruction, SfN provides interpretable insights into the origin of capabilities like reasoning.

Case 3 Or in the worst scenario, is reasoning an illusion (e.g., by overfitting to certain types of data), such that we have overestimated the potentials of LLMs?

A definitive answer to any of the above questions will be extremely valuable to guiding the future direction of LLM research. Even a hypothesis or conjecture supported by circumstantial evidences will be highly enlightening, too, let alone when convincing empirical evidences are available.

To this end, our hypothesis is that Case 2 holds in LLMs that reason well. To be more precise, we hypothesize that it is the output projection's parameters (o_proj) in the Transformer [44]'s multi-head self-attention (MHSA) module that is in charge of reasoning in an LLM.

To support our hypothesis, we propose a few techniques for diagnosing LLM's behaviors, in particular, the potential functionalities and impacts of various modules in it. We call these techniques Stethoscope for Networks, or SfN (summarized and illustrated in Figure 1). Starting from reasoning-enhanced models, we argue that the weight differences between a base LLM and its fine-tuned counterpart (e.g., for reasoning tasks) provide firsthand and crucial evidence for understanding internal changes. We refer to this approach as the Delta Stethoscope.

In addition, we introduce two novel and previously unexplored methods within the SfN framework: the Merge Stethoscope and the Destruction Stethoscope. The Merge Stethoscope replaces specific modules in a base model with those from a reasoning model. Surprisingly, the resulting variant can maintain fluent dialogue and demonstrate improved reasoning ability in some cases. This phenomenon offers strong clues about the origin and localization of reasoning capability in LLMs. The Destruction Stethoscope, in contrast, systematically disables individual modules and observes the resulting behavior to infer the functional roles of each component. We also propose the Freeze Stethoscope, which selectively freezes parts of the model during fine-tuning. By controlling which modules are updated, we provide convincing empirical support for earlier insights and clues into the localization of reasoning within LLMs.

With different gadgets we propose in SfN, we provide not only sanity check level tests for our hypothesis, but also more convincing circumstantial supports and even direct empirical evidences. In short, the contributions in this paper are two-fold:

- With various diagnosis evidence (SfN), we are confident in hypothesizing that the output projection o_proj is mainly responsible for the reasoning in LLMs. The impact of this finding include not only potential ways to improve LLM that reasons (e.g., training much faster), but may generalize to produce better LLMs for other tasks (e.g., for a vertical LLM designed specifically for a domain). Our further conjecture is that other modules combined together lead to lucid conversations, but o_proj is less important in conversational ability.
- The proposed Stethoscope for Networks (SfN) gadgets are a set of tools that are useful in understanding modern LLMs and even other networks, which have the potential to enhance our understanding of LLM or deep neural network and may lead to alternative routes for further deep learning research.

2 Key Hypothesis: Output Projection is the Key for Reasoning

To present our findings, we start by introducing necessary background information and notations, while discussions on related work are deferred to Section 5.

Modern LLMs [42, 5, 29] mostly consist of many Transformer blocks. A Transformer [44] block is composed of a multi-head self-attention (MHSA) module and a multi-layer perceptron (MLP) module. Components in MHSA include various projections, such as those for computing Q, K and V, denoted as q_proj, k_proj, and v_proj, respectively. The output projection (o_proj) produces MHSA's output. Components in the MLP are mainly linear projections: up, down, and gate [16, 42, 5] projections, denoted as up_proj, down_proj, and gate_proj, respectively. The computation process is defined as:

$$x_{\text{attn}} = w_{\text{o}} \left[\text{Softmax} \left(\frac{(w_{\text{q}}x)(w_{\text{k}}x)^{\top}}{\sqrt{d}} \right) (w_{\text{v}}x) \right]$$

$$x_{\text{mlp}} = w_{\text{down}} \left[\sigma(w_{\text{gate}}x) \odot (w_{\text{up}}x) \right]$$
(1)

For simplicity, we omit residual connections and present the computation at the token level, without using matrix or vectorized notation. Other essential components not explicitly included in equation 1 include rotary positional embeddings (RoPE)[39], input embeddings (embed_tokens), layer normalization[4] (layernorm), and the language modeling head (lm_head).

Let A be an LLM with weak or no reasoning ability. By carefully procuring a dataset of reasoning examples [13, 25, 52], one can cleanse and improve the quality of the dataset into the training data \mathcal{D} , and then finetune the existing model A by using techniques such as SFT. The resulting LLM, model B, exhibits strong reasoning capabilities. For example, in commonly adopted practices, the base LLM A is typically a widely used open-source model such as Qwen2.5-Math-1.5B, 7B or Qwen2.5-14B, 32B [53]. The reasoning model B denotes a publicly available reasoning-enhanced variant, such as DeepSeek-R1-Distill-Qwen-1.5B, 7B, 14B, 32B [13], which comes with a clearly specified base model and well-documented training procedure. Models that are either not open-sourced [13, 40], or open-sourced without sufficient training details [41] or access to the base model [52], are not discussed in this paper.

2.1 The Delta Stethoscope

In the above scenario, it is obvious that A and B share exactly the same network architecture and structure, with their sole difference being the weights (parameters) inside various components. Suppose w(A) (w(B)) denotes the set of weights for all modules in A (B). Then, it is natural to conclude that to understand the difference between A and B (i.e., reasoning or not), we should focus on the difference between w(A) and w(B). Hence, we propose our first Stethoscope for Network.

Assumption 1 (The Delta Stethoscope) Suppose A and B are two LLMs with weak and strong reasoning ability, respectively, and B is obtained by finetuning from A. Then w(B) - w(A) contains essential information if we want to pinpoint the source of the reasoning ability in B.

For each component X (e.g. $X = q_proj$), we compute the ℓ_2 norm of the weight difference, $\|w_X(B) - w_X(A)\|_{\ell_2}$, and visualize the results across all the blocks in Figure 2. For simplicity and due to space constraints, we present three representative comparisons: A is Qwen2.5-Math-1.5B [54] or Qwen2.5-14B, 32B [53] and B is DeepSeek-R1-Distill-Qwen-1.5B, 14B, 32B [13]. Additional results for other model sizes (7B and 8B) are provided in the appendix and exhibit similar patterns.

For the 1.5B models, the signal is less clear, but o_proj still exhibits a distinct pattern compared to q,k,v_proj—showing the largest change within the attention module and the second-largest across the entire model. As model size increases to 14B and 32B, this trend becomes more pronounced. In both cases, the most notable observation is that when $X = o_proj$, the ℓ_2 norm is at least two times larger than any other component, indicating the substantial changes in this module during reasoning enhancement.

In Figure 3, we further analyze the distribution of relative weight changes $\frac{w_X(B)-w_X(A)}{w_X(A)}$ for each linear module. To improve clarity and visual appeal, we plot the distribution every 5 layers and clip values in the range [-1.0, 1.0] to mitigate the influence of outliers. The vertical axis represents the

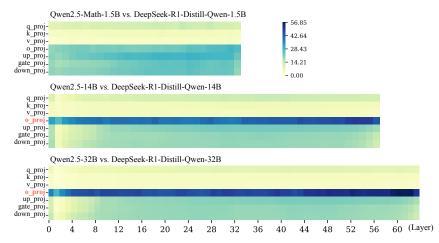


Figure 2: **Per-module L2 distance of linear weights between models** *A* **and** *B***.** Notably, the o_proj module shows the second-largest change in 1.5B models, and the largest in 14B and 32B models, highlighting its potential importance for reasoning. Similar trends are observed in 7B and 8B models (see appendix).

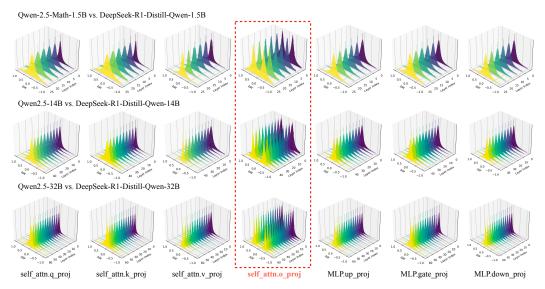


Figure 3: Layer-wise distribution of relative weight changes between models A and B. While most modules display a unimodal distribution, the o_proj module uniquely exhibits a bimodal distribution, highlighting its distinctive behavior. Consistent patterns are observed across models of other sizes, with detailed results provided in the appendix.

frequency. A striking and consistent finding is that all linear modules—except o_proj—exhibit a unimodal distribution centered around zero, whereas o_proj uniquely displays a clear bimodal pattern, highlighting its distinct role.

Both observations hold consistently across model sizes and base models: o_proj exhibits the largest or second-largest weight shift, and the overall weight difference patterns remain strikingly similar. Therefore, it is reasonable to guess that the output projection o_proj plays a pivotal role in curating B's reasoning ability. We are, however, not aware of o_proj's specific role: is it solely responsible for reasoning? Or, is it collaborating with another module(s)? Or, in the worst scenario, is this difference in $\|w_X(B) - w_X(A)\|_{\ell_2}$ and $\frac{w_X(B) - w_X(A)}{w_X(A)}$ coincidental?

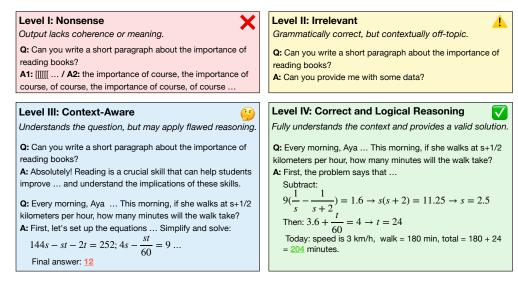


Figure 4: **Four levels of responses generated by the LLM**. From level I to level IV, the model exhibits stronger language organization and logical reasoning skills. Each example includes a question (e.g., a math problem from AIME or a typical user-issued request) and the corresponding response generated by the LLM.

2.2 The Merge Stethoscope

We design another gadget, the Merge Stethoscope, to answer this question. Suppose an LLM M is formed by merging models A and B, that is, M has the same structure as A and B, while a subset of its modules' parameters come from A and the rest from B. In a conversational or reasoning task, what will the output of M look like? We can imagine 4 levels of different output, as

- Level I A sequence of random or nonsense tokens.
- Level II A sequence that looks like normal sentences, but does not fit into the context of the task.
- Level III A sequence that is meaningful sentences that match the task's context well but will fail to reason in difficult problems.
- Level IV A sequence that reasons—and reasons correctly in most cases.

Figure 4 shows examples of level I to IV outputs. It is worth highlighting that M is rudely merged from A and B $\mathit{without}$ any $\mathit{further}$ tuning . Hence, the intuitive conjecture will be that M will produce level I output (i.e., ushering meaningless tokens). However, if model M, when merged in a specific configuration, is capable of producing level IV outputs for questions that model A fails to solve, then the specially merged components are likely critical for reasoning.

Assumption 2 (The Merge Stethoscope) Suppose M is created by merging the output projection (o_proj) weights of B (which has strong reasoning ability) and all other components of A (which is weak in reasoning), and further suppose that M has stronger reasoning ability compared to A. Then, we assume o_proj is crucial in achieving reasoning in LLMs.

We attempt a minimal or atomic merge by replacing only the o_proj modules in model A = Qwen2.5-Math-1.5B [54] with that of model B = DeepSeek-R1-Distill-Qwen-1.5B [13], keeping all other components unchanged. Although we initially expected the resulting model to produce level I or level II outputs, the results turn out to be surprising. On the AIME 2024 benchmark [19], the merged model M_1 achieves level IV performance on several questions that model A cannot solve. As shown in Table 1, the merged model not only yields correct reasoning and answers, but also tends to generate longer and more detailed responses compared to A. In contrast, replacing other modules such as $\{q,k,v\}_proj$ and mlp leads to performance degradation. For example, model M_2 , which replaces $\{q,k,v\}_proj$, produces level III outputs, while M_3 , which replaces mlp, deteriorates to level I. Only replacing o_proj results in a correct reasoning process and a correct answer, as illustrated in Figure 5. This striking difference motivates our further investigation in Section 3.

Model	Replaced Module	AIME 2024	Average Tokens
A (Q-1.5B)	-	0.067	2421
M_1 M_2	o_proj {q,k,v}_proj	0.200 0.000	5418 2058
M_3	mlp	0.000	15532
B (D-1.5B)	-	0.233	11892

Table 1: AIME 2024 accuracy of the base model, the reasoning model, and their merged variants. Each merged model is constructed by replacing specific modules in model A with the corresponding module from model B.

Q: Every morning, Aya does a 9 kilometer walk ... if she walks at s+1/2 kilometers per hour, how many minutes will the walk take?

 M_1 : To solve this problem, we need to determine ... So, the walk will take 204 minutes, including the 24 minutes at the coffee shop. The final answer is 204.

 M_2 : To solve this problem ... output 12.000000000000. The output indicates that the time taken for the walk is 12 minutes. So, the final answer is 12.

 M_3 : ... walking speeds increase speeds faster walking speeds increase walking speeds faster walking speeds faster walking .

Figure 5: **Examples of outputs generated by merged models.** Only M_1 produces both a valid reasoning process and the correct answer.

These results clearly show that the merged model M has a stronger reasoning capacity than A, despite that M is sutured from two completely different models and has *never* being finetuned. Now we feel confident in our assumption that o_proj is the key component responsible for reasoning in LLMs.

2.3 The Freeze Stethoscope

As models A and B scale up (e.g., to 7B parameters), merging components such as q,k,v_proj or mlp still results in significant performance degradation. However, unfortunately, merging o_proj no longer brings notable improvements in solving complex mathematical problems—although it does not harm accuracy, and still increases the generated output length.

Our analysis of $||w_X(B) - w_X(A)||_{\ell_2}$ suggests that this is due to a substantial mismatch in normalization parameters (that is, layernorm modules) between A and B at larger scales, compared to smaller models (e.g. 1.5B). Even when we merge both o_proj and layernorm parameters from B, the resulting model M still fails to reason effectively, probably because the remaining parameters of A are incompatible with the normalization parameters of B. To investigate this hypothesis in larger LLMs, we introduce the Freeze Stethoscope.

Assumption 3 (The Freeze Stethoscope) Suppose that an LLM F is obtained by supervised fine-tuning using the dataset \mathcal{D} . F is initialized from A, and both o_proj and normalization components are tuned while other components are frozen. If F exhibits strong reasoning ability, then we assume that o_proj is crucial in achieving reasoning in LLMs even in large-scale models.

It is worth noting that embed_tokens and lm_head are also tuned. Normalization module parameters are unfrozen by default. We adopt the pipeline of s1 [25] as our baseline, which uses the base model A = Qwen2.5-32B-Instruct and the dataset $\mathcal{D} = \text{s1K}$ containing 1,000 high-quality reasoning traces. The results are shown in Table 2, where our model F_4 corresponds to model B in Assumption 3. We do not strictly follow the training or testing setup of s1, primarily due to limited computational resources and the lack of an exact testing recipe to reproduce the reported results. However, our objective is not to optimize accuracy via testing tricks or prompt tuning, but to highlight the effectiveness of o_proj tuning compared to full-parameter tuning. For fair comparison, we adopt the "Budget Forcing Wait 2x" setting from s1 and retain all configurations without hyperparameter tuning.

Using this simplest possible experimental setup, Table 2 clearly shows that simply tuning o_proj and layernorm (model F_2) leads to strong reasoning ability, while at the same time only tuning layernorm (model F_1) harms the reasoning of the LLM. Further unfreezing the parameters of $\{q,k,v\}$ _proj (model F_3) yields little additional gain or even negative impact.

The training loss curves are shown in Figure 6. When all parameters including MLP are unfrozen, the model exhibits clear signs of overfitting, likely using the large MLP capacity to memorize the training set. In contrast, tuning only o_proj yields a smoother and more stable curve. Combined

²Without tuning these components, finetuning failed to converge.

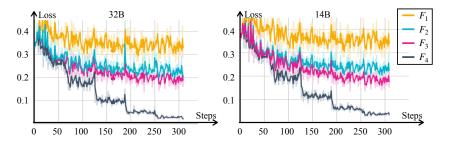


Figure 6: Training loss curves for fine-tuning Qwen2.5-14B,32B-Instruct on reasoning tasks. Different models unfreeze different sets of parameters, as detailed in Table 2.

Model	Fintuned Modules	#Param (B)	Steps/s	AIME 2024	Math 500	GPQA Diamond
A (Q-32B)	-	-	-	0.167	0.836	0.485
$\overline{F_1}$	Emb + Head	1.5	0.055	0.200	0.756	0.444
F_2	Emb + Head + o_proj	3.2	0.052	0.367	0.890	0.520
F_3	$Emb + Head + \{q,k,v,o\}_proj$	5.6	0.044	0.300	0.886	0.525
$F_4(B)$	All	32.8	0.015	0.367	0.906	0.591
A (Q-14B)	-	-	-	0.133	0.810	0.449
$\overline{F_1}$	Emb + Head	1.5	0.106	0.133	0.722	0.414
F_2	Emb + Head + o_proj	2.8	0.099	0.266	0.848	0.485
$\overline{F_3}$	$Emb + Head + \{q, k, v, o\}_proj$	3.7	0.081	0.233	0.854	0.490
$F_4(B)$	All	14.7	0.053	0.266	0.872	0.530

Table 2: Reasoning performance of different fine-tuning strategies on Qwen2.5-{14B, 32B}-Instruct. Emb denotes embed_tokens, Head denotes lm_head, and Attn denotes the entire MHSA. #Param refers to the number of trainable parameters, Steps/s indicates training speed, and the last three columns report commonly used metrics for evaluating reasoning models.

with its competitive performance, this suggests that the model learns to reason rather than simply memorize. Hence, we are now prepared and feel supported to propose our key hypothesis:

Hypothesis 1 (Outstanding Output Projection) In an LLM that reasons well, we hypothesize that the output projection (o_proj) component is the single or at least the most important module that dominates its reasoning ability.

With carefully chosen tuning strategy and hyperparameters, there is reason to believe that tuning only o_proj (+LN) can reach the level of model B in terms of reasoning performance. And, beyond exhibiting reasoning abilities, Table 2 also shows that tuning only o_proj (+LN) has other significant advantages: e.g., significantly faster finetuning (3 times faster) and smaller GPU memory consumption. These advantages will become more established when larger LLMs are tuned.

3 Conjecture: Conversation Hinges on Other Modules but Not Output

We are mainly concerned with two abilities of LLMs: conversation and reasoning, which map to level III and IV in our categorization of LLM's outputs, respectively. Our Hypothesis 1 is on reasoning, but are there one module or several modules accounting for lucid conversations? In this section, we further propose a new stethoscope to diagnose this question and raise our conjectures accordingly.

3.1 The Destruction Stethoscope

Our previous stethoscopes follow a "constructive proof" style, while now we resort to the "proof by contradiction" style. If one module in an LLM is "destructed", and the LLM can still produce level III conversation outputs, then we have good reasons to guess that this module is not important in conversational ability; while it is important if the LLM ceases to dialogue regularly.

Assumption 4 (The Destruction Stethoscope) Suppose a module X is destructed (i.e., its normal functionality is disabled by some destruction method) in an LLM A. We denote the resulting LLM as

Destruction Method	Module	Output Level	Destruction Method	Module	Output Level
Zero	q_proj	I		q_proj	I
	k_proj	I		k_proj	I
	v_proj	Ш		v_proj	II
	o_proj	Ш	ReInit	o_proj	Ш
	up_proj	I		up_proj	I
	gate_proj	I		gate_proj	I
	down_proj	I		down_proj	I
Remove	-	I			

Table 3: Output levels of different modules under the three destruction methods: Zero, ReInit, and Remove. All experiments are based on Owen2.5-32B with destruction applied to specific layers.

D. Then, the fact that D continues (or ceases to) produce level III output (meaningful sentences in the conversation's context) indicates whether X is important for conversational abilities or not.

We propose 3 destructors to destroy a module:

Zero Set all parameters within X to 0.

ReInit Re-initialize all parameters inside X using Gaussian random numbers (mean=0, std=0.02). Remove Remove the entire layer.

The Zero destructor is often equivalent to setting the output activation of X to zeros (e.g., in a linear module like o_proj). We want to emphasize that ReInit incurs more serious damages to an LLM than Zero does. Zero may change activations to zero, but ReInit exerts random effects (i.e., noise) to LLM activations. What is more important, these random effects will act as input to the next Transformer block and the noise is quickly amplified. Hence, level I or II output is expected when X is destroyed (especially when reinitialized) in a large number of Transformer blocks.

3.2 Conjectures Concerning the Conversation Capability

For model Qwen2.5-32B with 64 layers, we observe that destroying modules in early or late layers—where input and output representations are more sensitive—consistently yields level I outputs. To avoid this, we restrict destruction to blocks 5–30. This range is empirically chosen, as affecting more layers often causes all outputs to degrade to level I, making distinctions between modules impossible.

The experimental results are presented in Table 3. Specifically, we destroy selected modules and analyze the corresponding output. The Remove destructor removes the transformer layers as a whole. Note that the results are not statistics computed in many different experiments—it only reflects the conversation illustrated in Figure 4, but we observed similar patterns for other conversations.

Table 3 reveals distinct roles of modules in conversation. Notably, o_proj—crucial for reasoning—appears unimportant for conversation. In contrast, all MLP components (up_proj, down_proj, gate_proj) are essential. Within MHSA, q_proj and k_proj are important, while v_proj plays a minor role. Based on these (admittedly weaker) observations, we propose the following conjecture.

Conjecture 1 (Division of Labor) Based on current observations, an LLM can be roughly divided as two sets of modules: output projection (o_proj) and all others, where o_proj is mainly responsible for reasoning and other modules for conversation.

Then, output projection plays a unique role if this conjecture holds. Hence, we further propose another conjecture for it.

Conjecture 2 (Output Projection Plugin) With conversational capabilities provided by other (frozen) modules, output projections may act as a plugin. For example, one set of o_proj for reasoning, and another set of o_proj for migrating an LLM to a vertical domain.

4 Potential Implications and Applications

This paper mainly diagnoses LLMs from a theoretical, highly abstract perspective. However, our hypothesis and conjectures can also have highly practical implications and applications as long as they are correct or at least partially hold.

- Fast and better reasoning LLMs. By finetuning only o_proj, we can potentially find a better reasoning LLM with much faster training and much smaller GPU memory footprint.
- Integrating non-reasoning and reasoning LLMs. There is a recent trend to integrate chatting and reasoning LLMs into one model [52]. When we finetune a base LLM into a reasoning one using the previous procedure, they only differ in o_proj, layernorm, embed_tokens and lmhead, which occupy only 10% of model size. Hence, the two LLMs are easily loaded as one LLM with two sets of these module for different purposes.
- Vertical LLMs. Similarly, when equipped with different output projection plugins, one may
 adeptly obtain vertical LLMs for different domains.
- **Understanding deep neural networks.** The proposed Stethoscopes for Networks might be useful gadgets to understand other deep models, and new stethoscopes can be further developed. They will be potentially useful in diagnosing existing networks and even in providing alternative directions to future deep learning research.

5 Related Work

Large Language Models. Modern LLMs such as GPT [29, 6], LLaMA [42, 43], Qwen [5, 53], and other representative models [7, 20] adopt an auto-regressive architecture and have demonstrated impressive capabilities across a wide range of natural language processing tasks, including question answering [32, 22], summarization [26, 27], and translation [51]. These models are typically trained on large-scale corpora using next-token prediction objectives, and their performance has been shown to scale with model size [21]. Further improvements in alignment and usability have been achieved through instruction tuning [28, 9, 47] and reinforcement learning from human feedback (RLHF) [8, 30], enabling more controllable and helpful dialogue generation.

Reasoning Models. While LLMs exhibit emergent reasoning abilities [48, 35], recent efforts have further enhanced these capabilities through fine-tuning and architectural modifications [36, 56]. Chain-of-thought prompting [49] encourages intermediate reasoning steps, improving performance in arithmetic tasks, while self-consistency decoding [46] improves robustness by sampling multiple reasoning paths. Inspired by OpenAI's o1 [18], most advanced models now employ reinforcement learning [37, 30] to generate long reasoning traces with sparse rewards. This leads to significant improvements, particularly in complex math, code, and other professional domains [13, 52]. Despite these advances, the origin and location of reasoning ability in LLMs remain underexplored.

Interpretability of LLMs. Understanding the inner workings of LLMs has attracted growing interest. Prior efforts include attention visualization [45], probing [15], and model editing [24, 34], with the aim of interpreting internal representations. Other studies decompose the behavior of the model into attribute functions to specific modules [11]. The "Physics of Language Models" series [1, 2, 3] investigates LLMs through controlled setups to reveal empirical and universal laws that dictate LLM behavior. However, these studies often exclude the most advanced models or focus on narrow, synthetic settings, offering limited insight into real-world models. Their findings provide little practical guidance for understanding reasoning in state-of-the-art models.

6 Conclusions

This work investigates a fundamental question in understanding large language models (LLMs): Is there a component or several components that are responsible for achieving the reasoning ability in LLMs? If the answer is affirmative, which components are responsible for the improvement?

We hypothesize that the output projection (o_proj) module plays a central role in enabling reasoning capabilities. To support this, we propose *Stethoscope for Networks (SfN)*, a diagnostic framework

that encompasses several probing techniques. Through the proposed Delta, Merge, Freeze, and Destruction stethoscopes, we observe consistent patterns indicating that o_proj is critical for reasoning, while other modules primarily support conversational fluency. These findings open new directions for efficient and modular LLM training.

Our findings are primarily based on a limited set of model families and reasoning benchmarks, and may not generalize to all architectures or tasks. Some diagnostic results rely on qualitative assessments rather than statistical validation. Furthermore, while the role of o_proj is empirically highlighted, a theoretical understanding of its function in reasoning remains to be established.

Acknowledgments and Disclosure of Funding

This work was partly supported by the National Natural Science Foundation of China under Grant 62276123.

JW proposed the assumptions (Stethoscopes for Networks), hypothesis and conjectures. JS started this line of research in our group, proposed the Zero destructor, and effectively supported our main findings with experimental results. JW and JS wrote the paper.

We thank Ke Zhu for discussions.

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- [2] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*, 2023.
- [3] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv preprint arXiv:2404.05405*, 2024.
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint* arXiv:1607.06450, 2016.
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- [7] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Nina Mielke, Alec Radford, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [8] Paul F Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in neural information processing systems*, volume 30, 2017.
- [9] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Xin Wang, Xingyu Yuan, Adams Yu, Sharan Narang, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [10] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024.
- [11] Mor Geva, Tal Schuster, and Jonathan Berant. Transformer feed-forward layers are key-value memories. *arXiv preprint arXiv:2012.14913*, 2021.

- [12] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [13] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [14] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- [15] John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.
- [16] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. Transformer quality in linear time. In *International conference on machine learning*, pages 9099–9117. PMLR, 2022.
- [17] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [18] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv* preprint arXiv:2412.16720, 2024.
- [19] Maxwell Jia. Aime 2024 dataset. https://huggingface.co/datasets/Maxwell-Jia/ AIME_2024, 2024.
- [20] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023.
- [21] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [22] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [23] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [24] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *Advances in neural information processing systems*, 35:17359–17372, 2022.
- [25] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. arXiv preprint arXiv:2501.19393, 2025.
- [26] Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang, et al. Abstractive text summarization using sequence-to-sequence rnns and beyond. arXiv preprint arXiv:1602.06023, 2016.
- [27] Shashi Narayan, Shay B Cohen, and Mirella Lapata. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv* preprint arXiv:1808.08745, 2018.

- [28] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [29] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [30] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [31] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- [32] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [33] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 3505–3506, 2020.
- [34] Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. *arXiv preprint arXiv:2312.03002*, 2023.
- [35] Laura Ruis, Maximilian Mozes, Juhan Bae, Siddhartha Rao Kamalakara, Dwarak Talupuru, Acyr Locatelli, Robert Kirk, Tim Rocktäschel, Edward Grefenstette, and Max Bartolo. Procedural knowledge in pretraining drives reasoning in large language models. *arXiv* preprint arXiv:2411.12580, 2024.
- [36] Timo Schick, Ananya Dwivedi-Yu, Roberta Raileanu, Saghar Hosseini, Murray Chadwick, Gaurav Mishra, Siddharth Karamcheti, Neil Houlsby, Aravind Elangovan, Mike Lewis, et al. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- [37] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [38] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [39] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [40] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [41] Qwen Team. QwQ-32B: Embracing the Power of Reinforcement Learning, March 2025.
- [42] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [43] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- [45] Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. arXiv preprint arXiv:1906.04284, 2019.
- [46] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- [47] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- [48] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [49] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- [50] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics.
- [51] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [52] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025.
- [53] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. arXiv preprint arXiv:2412.15115, 2024.
- [54] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [55] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild, 2025.
- [56] Denny Zhou, Dale Schuurmans, Xuezhi Wang, Ed Chi, and Quoc V Le. Least-to-most prompting enables complex reasoning in large language models. arXiv preprint arXiv:2205.10625, 2023.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction accurately reflect the key contributions presented in Section 2 and 3, while the experimental findings provide thorough support for the validity of our claims.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper clearly acknowledges its main limitations in Section 6, providing a thorough discussion of the associated constraints and challenges.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This paper adopts a primarily empirical approach, with conclusions and methods grounded in and supported by experimental results rather than theoretical analysis. Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We present a simple and easily reproducible pipeline, fully described in Sections 2 and 3. All essential hyperparameters and experimental settings are specified, with additional details for replication provided in the appendix.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: The code associated with this paper will be released as open-source upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All key experimental settings—including datasets, hyperparameters, and their selection criteria—are thoroughly described in Sections 2 and 3. Additional experimental details can be found in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The paper includes detailed information on the statistical significance of the experiments provided in the appendix.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Detailed information on the computational resources used—including GPU type and execution time—is provided in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have carefully considered and complied with the NeurIPS Code of Ethics throughout the course of our research.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: A discussion of potential societal impacts, including both benefits and risks, is provided in appendix.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our research does not involve models that generate high-risk content or rely on web-sourced data, thereby minimizing the need for additional safeguards.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All data, code, and methods used in this study are open-source, with proper citations and full compliance with their respective licenses and terms of use.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release any new assets; therefore, no associated documentation is provided.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing or human subject research; thus, this question is not applicable.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work does not involve human subjects or crowdsourcing, and therefore requires no study participants, risk assessment, or IRB approval.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: Large language models play a central role in our work, as the research focuses on analyzing, modifying, and evaluating LLM behavior. The core methodology involves fine-tuning and interpreting LLM components, making their usage essential and original to the contributions of this paper.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

Technical Appendices and Supplementary Material

A Experimental Details

We primarily utilize open-sourced models to conduct experiments in this work. Given that DeepSeek-R1 is one of the most widely adopted reasoning models, and its authors have released a series of distilled models based on R1 [13], including both the specified base and finetuned reasoning models, we adopt their configurations in our study. Specifically, we use the DeepSeek-R1-Distill-Qwen [13] models with sizes of 1.5B, 7B, 14B, 32B and 70B as our reasoning models, and select Qwen2.5-Math-1.5B, 7B [54], LLaMA3.1-8B [12], Qwen2.5-14B, 32B [53] or Llama-3.3-70B-Instruct [12] as base models. All models are loaded and run using the Transformers library [50].

Our evaluation framework is based on the lm-evaluation-harness package [10]. To accelerate inference, we use vLLM [23] as the backend, which may slightly affect performance due to backend-specific optimizations. In the Merge Stethoscope experiments, we observe that the "chat" interface often generates irrelevant or nonsensical responses, while the "generate" interface produces coherent and contextually appropriate outputs. We suspect this discrepancy arises from misinterpreted system prompts. Therefore, we rely on the "generate" interface and implement a custom evaluation toolkit.

For the Freeze Stethoscope experiments, we build on the codebase of s1[25]. We use a learning rate of 1e-5, weight decay of 1e-4, a batch size of 16, and train for 5 epochs. Due to hardware limitations (i.e., lack of access to 16 H100 GPUs), we leverage DeepSpeed[33] with ZeRO Stage 3[31] to enable efficient training. The base model used here is Qwen2.5-32B-Instruct[53]. Evaluation is again conducted with lm-evaluation-harness, following the modified pipeline by the authors of s1, which disables generation of the end-of-thinking token and optionally appends the string "Wait" to the reasoning trace to encourage model reflection. We adopt the Budget Forcing "Wait" ×2 as our default testing configuration.

All visualization and inference experiments on 1.5B–14B models are conducted on a single NVIDIA A100 GPU. For training and evaluating 32B-70B models, we use a cluster of 8 NVIDIA A100 GPUs. Training typically takes around 6 hours, while testing on a single dataset usually requires about 2 hours.

B More Experimental Results

In the main paper, we present visualization results for the 1.5B, 14B, and 32B models. Here, we supplement those results by providing additional visualizations for the 7B, 8B, and 70B models. Following the Delta Stethoscope pipeline, we visualize both the absolute weight shift $|w_X(B)-w_X(A)|_{\ell_2}$ and the relative weight shift $\frac{w_X(B)-w_X(A)}{w_X(A)}$. The absolute weight shifts are shown in Figure 7, and the relative weight shifts are presented in Figure 8. The trends observed in the main paper remain consistent across these additional models. Notably, o_proj consistently exhibits the largest weight shift, with the effect being especially pronounced in the 70B model. Moreover, o_proj is the only module that displays a bimodal distribution in the relative weight shift.

C Statistical Significance and Broader Impacts

We report appropriate information regarding the statistical significance of our experiments. While we do not primarily focus on classical significance tests such as p-values, we provide multiple forms of empirical evidence—such as consistent module-specific weight shifts, response-level comparisons under controlled manipulations, and loss curves under different tuning strategies—that collectively establish the robustness of our findings. These analyses serve as a practical alternative to traditional error bars or confidence intervals and help substantiate our key claims.

This research has both promising benefits and important risks to consider. On the positive side, the proposed Stethoscope for Networks (SfN) framework provides a novel set of tools for interpreting LLMs, especially by localizing specific capabilities—such as reasoning—to individual components like the output projection (o_proj). These tools may significantly improve our understanding of LLMs, enabling more transparent, modular, and efficient model development. For instance, if reasoning

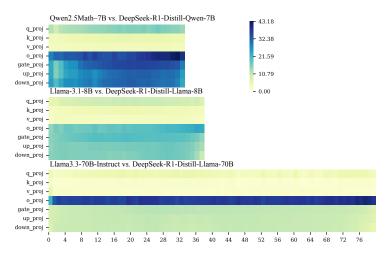


Figure 7: **Per-module L2 distance of linear weights between models** *A* **and** *B***.** Notably, the o_proj module shows the largest in 7B, 8B and 70B models, highlighting its potential importance for reasoning.

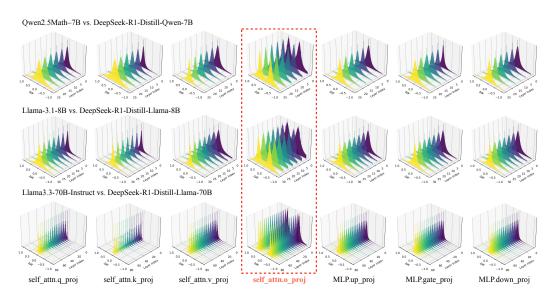


Figure 8: Layer-wise distribution of relative weight changes between models A and B. While most modules display a unimodal distribution, the o_proj module uniquely exhibits a bimodal distribution, highlighting its distinctive behavior.

abilities can be enhanced by tuning a small subset of parameters, it could greatly reduce computational costs and increase accessibility for developing domain-specific or lightweight models.

However, this line of work also carries potential risks. Precisely identifying and isolating reasoning-related components might lower the barrier for targeted manipulation, such as unauthorized transfer or removal of reasoning abilities across models. This could facilitate misuse scenarios, including capability extraction, tampering, or model theft. Furthermore, while the diagnostic methods proposed aim to support interpretability, there is a risk that they may be overinterpreted, leading to an inflated sense of model transparency that does not generalize across architectures or tasks.