
Understanding Pre-trained and Fine-tuned model behaviour using Model Diffing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Fine-tuning large language models (LLMs) for specialized domains can alter
2 both their output distributions and internal mechanisms in ways that standard
3 task metrics obscure. We study *model diffing* between a pretrained base
4 (DeepSeek-R1-Distill-Qwen-1.5B) and a LoRA-adapted variant trained for medical
5 reasoning on HuatuoGPT-o1 style data. Our protocol couples (i) *next-token KL*
6 *divergence*, measured across general and medical corpora—to quantify output-level
7 shift, with (ii) *activation patching* to localize where domain knowledge and reasoning
8 procedures are encoded. We target the LoRA-modified projections and MLP
9 pathways and analyze the behavioral impact of swapping per-layer activations
10 across models. Complementary experiments explore Complex Chain-of-Thought
11 fine-tuning and a *Kahneman-Tversky Optimization (KTO)* objective from the
12 Human-Aware Loss family to encourage structured reasoning without preference
13 labels. Empirically, we observe domain-selective distributional drift with minimal
14 degradation on general text, and layerwise concentration of medical competence
15 consistent with prior findings that factual/semantic knowledge often resides in
16 mid-to-late MLP blocks. Our contributions are: (1) a unified, reproducible KL plus
17 patching diffing protocol; (2) evidence on how LoRA placements mediate domain
18 specialization; and (3) an analysis of reasoning-oriented post-training (CoT/KTO)
19 and its interaction with representational localization. We release code and scripts
20 to support systematic model diffing in domain-specific alignment.

21 1 Introduction

22 Large language models (LLMs) have achieved remarkable performance across open-domain and
23 specialized reasoning tasks, yet understanding how fine-tuning reshapes their internal mechanisms
24 remains an open challenge. In domains such as medical reasoning, where reliability and interpretability
25 are essential, it is not sufficient to evaluate models solely on accuracy. Instead, one must also
26 interrogate the representational shifts and causal mechanisms that emerge when a general-purpose
27 model is adapted to a specialized domain.

28 Recent advances in mechanistic interpretability have begun to expose the circuits and abstractions
29 underpinning transformer models. Foundational work established mathematical frameworks for
30 analyzing attention and feed-forward pathways Elhage et al. [2021], Geva et al. [2020] and revealed
31 that factual associations can be localized and edited via targeted interventions Meng et al. [2022].
32 Techniques such as activation patching Syed et al. [2023], Zhang and Nanda [2023] and attribution
33 methods Vig et al. [2020], Conmy et al. [2023] allow researchers to causally probe how specific layers
34 and heads contribute to predictions, while recent benchmarks Mueller et al. [2025] and circuit-analysis
35 efforts Wu et al. [2024a] seek to standardize evaluation. These approaches have shown that domain

36 knowledge, heuristics, and even undesirable behaviors such as sycophancy Li et al. [2025] can often
37 be traced to mid-to-late transformer blocks, underscoring the importance of layer-level analysis.

38 Parallel strands of work have investigated representation divergence as a quantitative measure of
39 model adaptation. Studies of mode connectivity and representation finetuning Zhou et al. [2023,?],
40 Wu et al. [2024b] suggest that fine-tuning induces structured but non-trivial shifts in embedding
41 spaces. Kullback–Leibler (KL) divergence has been widely used to quantify such distributional
42 changes, both in knowledge distillation Wu et al. [2024a] and in tracking training dynamics Kishino
43 et al. [2025]. These methods highlight that fine-tuning not only alters outputs but also reorganizes
44 internal representations in ways that can be systematically measured.

45 At the same time, alignment research has advanced techniques to encourage structured reasoning and
46 human-centric behavior. Chain-of-Thought prompting Wei et al. [2022] demonstrated that eliciting
47 intermediate steps enhances reliability. Preference-based fine-tuning Ziegler et al. [2019], ? and
48 Direct Preference Optimization Rafailov et al. [2023] introduced frameworks for aligning models
49 with human judgments, though they require costly data collection. To reduce this dependency,
50 cognitive-inspired approaches such as Kahneman-Tversky Optimization Ethayarajh et al. [2024] and
51 human-aware loss functions Qi et al. [2025] leverage decision-theoretic principles Kahneman and
52 Tversky [2013] to regularize models toward trustworthy reasoning without explicit preference labels.
53 Emerging methods like chain-of-preference optimization Zhang et al. [2024] and steering vectors
54 for reasoning Venhoff et al. [2025] further suggest that interpretability and alignment can be jointly
55 optimized.

56 Against this backdrop, our work contributes a unified framework for model diffing that combines
57 KL divergence and activation patching to study how fine-tuning for medical reasoning reshapes
58 an LLM’s internals. We focus on DeepSeek-R1-Distill-Qwen-1.5B one of the Deepseek series of
59 models Guo et al. [2025] and its LoRA-adapted variant Hu et al. [2022] trained on domain-specific
60 medical reasoning tasks. By quantifying output divergence and performing layerwise patching on
61 LoRA-modified projections and MLP blocks, we identify where medical knowledge and reasoning
62 capabilities are embedded. We further evaluate reasoning-oriented objectives such as Complex Chain-
63 of-Thought fine-tuning and KTO-based HALOs, probing how they interact with representational
64 dynamics.

65 **Our contributions are threefold**

- 66 1. We introduce a **KL + activation** patching diffing protocol for systematically comparing
67 base and fine-tuned LLMs.
- 68 2. We provide **empirical evidence** of layerwise specialization and representational divergence
69 induced by medical fine-tuning under LoRA adaptation.
- 70 3. We analyze how **reasoning-oriented objectives** shape both output distributions and internal
71 mechanisms, offering insights for building interpretable, trustworthy medical AI systems.

72 **2 Dataset**

73 To evaluate how fine-tuning reshapes model internals for specialized reasoning, we trained on
74 the HuatuoGPT-o1 medical reasoning dataset developed by Freedom Intelligence ¹. We used the
75 instruction-tuned subset released as medical-o1-reasoning-SFT ²

76 This dataset was chosen because it aligns with our research objective of studying how fine-tuning
77 for domain-specific reasoning alters internal representations. Unlike general medical QA corpora,
78 HuatuoGPT-o1 explicitly encodes reasoning traces, which makes it particularly well-suited for model
79 diffing experiments that combine output divergence (via KL) and causal analysis (via activation
80 patching).

¹<https://github.com/FreedomIntelligence/HuatuoGPT-o1>

²<https://huggingface.co/datasets/FreedomIntelligence/medical-o1-reasoning-SFT>

81 **3 Methodology**

82 We fine-tuned the DeepSeek-R1-Distill-Qwen-1.5B model using Low-Rank Adaptation (LoRA) with
83 memory-efficient optimizations. To enable large sequence processing, we set the maximum sequence
84 length to 4096 tokens. The base model was loaded with 4-bit quantization (bnb-4bit), reducing GPU
85 memory requirements while maintaining competitive performance. A custom chat template was
86 applied to align the model with medical reasoning tasks, encouraging structured chain-of-thought
87 style responses.

88 For parameter-efficient fine-tuning, LoRA adapters (rank = 16, scaling factor = 16, no bias) were
89 injected into the attention and MLP projection layers (q_proj, k_proj, v_proj, o_proj, gate_proj,
90 up_proj, and down_proj). Gradient checkpointing was enabled to reduce memory overhead, and
91 random initialization seeds ensured reproducibility.

92 We align the LoRA-augmented DeepSeek-R1-Distill-Qwen-1.5B with Kahneman–Tversky Opti-
93 mization (KTO) Ethayarajh et al. [2024] via TRL’s KTOTrainer³. KTO treats feedback as bi-
94 nary desirability signals per sample and maximizes a prospect-theoretic utility using asymmetric
95 value/weighting functions, placing greater penalty on undesirable generations relative to comparable
96 gains for desirable ones. Practically, this removes the need for pairwise comparisons while targeting
97 the same alignment goals as preference-based methods. We train with an effective batch size of 8 (4
98 × grad-accum 2), AdamW-8bit, cosine LR schedule with 10% warmup, weight decay 0.01, and a
99 base learning rate of 5e-7. Mixed precision (FP16/BF16) is used based on hardware support.

100 To further optimize memory usage during training, we implemented a custom memory management
101 callback. This routine explicitly deletes input and output tensors at the end of each step, clears the
102 CUDA cache, and logs memory statistics every 100 steps. GPU memory utilization was monitored to
103 ensure stable training under constrained resources.

104 The model was fine-tuned on a subset of the medical reasoning dataset, using the tokenizer for data
105 preprocessing. The overall setup allows efficient domain adaptation of a 1.5B parameter model
106 on modest hardware, and training was conducted on an NVIDIA A100 GPU, enabling controlled
107 experimentation with model diffing and mechanistic interpretability tools.

108 **4 Results**

109 **4.1 KL Divergence Analysis**

110 **4.1.1 Token-level KL Divergence Analysis**

111 To quantify how fine-tuning alters the model’s behavior, we computed the KL divergence between the
112 pretrained and fine-tuned models at each token position of a sample prompt. Both models produced
113 logits of shape (1, 47, 151,936), corresponding to 47 tokens and a shared vocabulary. The summed
114 KL divergence across tokens was 0.2449, with individual values ranging from approximately −0.0027
115 to 0.0183.

116 The divergence values indicate that most tokens exhibit small but non-negligible shifts, while a few
117 tokens show localized spikes in divergence (e.g., >0.015). This suggests that fine-tuning does not
118 uniformly shift the output distribution but instead induces targeted representational changes at specific
119 positions. Such differences likely correspond to tokens where medical reasoning knowledge has been
120 reinforced or reweighted.

121 The bar plot in Figure 1 illustrates these token-level divergences. While many tokens cluster near
122 zero, the variance in divergence magnitudes highlights which lexical items were most impacted
123 by domain adaptation. These findings support the hypothesis that fine-tuning selectively alters the
124 model’s probability landscape, rather than producing a homogeneous distributional drift.

125 **4.1.2 Token-level KL Divergence Across Generated Responses**

126 We evaluated token-level KL divergence between the pretrained and fine-tuned models across five
127 generated medical reasoning responses. The average per-token divergence was very small (7×10 to

³https://huggingface.co/docs/trl/main/en/kto_trainer

128 8×10^3), confirming that fine-tuning preserves the overall distributional structure of the base model.
 129 However, divergence did not remain constant across the sequence. As shown in Figure 2, the per-token
 130 KL values fluctuate along the response, with localized spikes at medically salient positions. The
 131 moving average reveals that divergence gradually rises during intermediate reasoning steps before
 132 tapering off, suggesting that the model makes targeted adjustments when introducing or justifying
 133 domain-specific knowledge.

134 While the cumulative KL divergence varied more substantially (0.39–2.29 across responses), this
 135 reflects the accumulation of small but meaningful shifts over longer generations. Together, these
 136 patterns indicate that fine-tuning does not globally distort the pretrained model’s behavior; instead,
 137 it introduces selective representational refinements concentrated at key reasoning junctures, while
 138 maintaining consistency elsewhere.

Table 1: Aggregate results over 10 generated responses: averaged token-level KL divergence and text quality scores.

N	Avg Mean KL	Avg Sum KL	FT ROUGE-L	FT BERT F1	Base ROUGE-L	Base BERT F1
10	0.00696	0.34786	0.2222	0.8503	0.2197	0.8504

139 4.2 Activation Patching Analysis

140 While KL divergence quantifies output-level shifts, it does not directly reveal where in the network
 141 such shifts are encoded. Activation patching provides a causal tool for localizing representational
 142 changes: by selectively substituting hidden activations from one model into another, we can test
 143 which layers or submodules are most responsible for domain-specific behaviors. This technique
 144 has been widely used in mechanistic interpretability to identify circuits, attribute predictions to
 145 specific components, and trace factual knowledge to mid-to-late transformer blocks. In the context of
 146 model diffing, activation patching is crucial because it moves beyond correlation to establish causal
 147 responsibility for divergence between pretrained and fine-tuned models.

148 4.2.1 Methods and Results

149 We applied activation patching to the LoRA-augmented DeepSeek-R1-Distill-Qwen-1.5B. Specif-
 150 ically, we targeted the LoRA-injected projection pathways (`q_proj`, `k_proj`, `v_proj`, `o_proj`)
 151 and the MLP subcomponents (`gate_proj`, `up_proj`, `down_proj`). For each layer, we replaced the
 152 fine-tuned model’s activations with their pretrained counterparts during forward passes on medical
 153 reasoning prompts. By measuring the change in the fine-tuned model’s output distributions, we
 154 isolated which layers are most critical for encoding the medical reasoning capability.

155 Figure 3 shows the mean KL divergence per layer under patching interventions. The attention projec-
 156 tions (`q_proj`, `k_proj`, `v_proj`, `o_proj`) exhibit minimal divergence (≈ 0.0068), suggesting that
 157 fine-tuning preserved the pretrained model’s attention mechanisms. In contrast, the MLP projections
 158 show markedly higher divergence, peaking in `up_proj` (0.0106) and `down_proj` (0.0086). These
 159 results indicate that fine-tuning primarily reshaped the feed-forward pathways, consistent with prior
 160 findings that factual and semantic knowledge tends to be stored in MLP layers.

161 Notably, patching experiments revealed that medical reasoning ability is causally localized in the
 162 LoRA-modified MLP layers, whereas attention pathways remain broadly intact. Furthermore, by
 163 restricting fine-tuning to selected LoRA adapters, we not only maintained strong domain adaptation
 164 but also substantially reduced computational costs, lowering GPU memory usage and thereby reducing
 165 the carbon footprint of training. These findings underscore that activation patching, combined with
 166 divergence analysis, provides a layered view of how parameter-efficient fine-tuning embeds domain
 167 knowledge into compact feed-forward pathways while leaving general-purpose attentional scaffolding
 168 largely unchanged.

169 References

170 Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-
 171 Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural*
 172 *Information Processing Systems*, 36:16318–16352, 2023.

- 173 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
174 Aspell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer
175 circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
- 176 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model
177 alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- 178 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are
179 key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.
- 180 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
181 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
182 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 183 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
184 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 185 Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In
186 *Handbook of the fundamentals of financial decision making: Part I*, pages 99–127. World Scientific,
187 2013.
- 188 Ryo Kishino, Yusuke Takase, Momose Oyama, Hiroaki Yamagiwa, and Hidetoshi Shimodaira. Reveal-
189 ing language model trajectories via kullback-leibler divergence. *arXiv preprint arXiv:2505.15353*,
190 2025.
- 191 Jin Li, Keyu Wang, Shu Yang, Zhuoran Zhang, and Di Wang. When truth is overridden: Uncovering
192 the internal origins of sycophancy in large language models. *arXiv preprint arXiv:2508.02087*,
193 2025.
- 194 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
195 associations in gpt. *Advances in neural information processing systems*, 35:17359–17372, 2022.
- 196 Aaron Mueller, Atticus Geiger, Sarah Wiegrefe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu
197 Chan, Jaden Fiotto-Kaufman, Tal Haklay, Michael Hanna, et al. Mib: A mechanistic interpretability
198 benchmark. *arXiv preprint arXiv:2504.13151*, 2025.
- 199 Xuan Qi, Jiahao Qiu, Xinzhe Juan, Yue Wu, and Mengdi Wang. Shallow preference signals: Large
200 language model aligns even better with truncated data? *arXiv preprint arXiv:2505.17122*, 2025.
- 201 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
202 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances
203 in neural information processing systems*, 36:53728–53741, 2023.
- 204 Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit
205 discovery. *arXiv preprint arXiv:2310.10348*, 2023.
- 206 Constantin Venhoff, Iván Arcuschin, Philip Torr, Arthur Conmy, and Neel Nanda. Understanding
207 reasoning in thinking language models via steering vectors. *arXiv preprint arXiv:2506.18167*,
208 2025.
- 209 Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and
210 Stuart Shieber. Investigating gender bias in language models using causal mediation analysis.
211 *Advances in neural information processing systems*, 33:12388–12401, 2020.
- 212 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
213 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
214 neural information processing systems*, 35:24824–24837, 2022.
- 215 Taiqiang Wu, Chaofan Tao, Jiahao Wang, Runming Yang, Zhe Zhao, and Ngai Wong. Rethinking
216 kullback-leibler divergence in knowledge distillation for large language models. *arXiv preprint
217 arXiv:2404.02657*, 2024a.
- 218 Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning,
219 and Christopher Potts. Reft: Representation finetuning for language models. *Advances in Neural
220 Information Processing Systems*, 37:63908–63962, 2024b.

221 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:
 222 Metrics and methods. *arXiv preprint arXiv:2309.16042*, 2023.

223 Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. Chain of preference
 224 optimization: Improving chain-of-thought reasoning in llms. *Advances in Neural Information
 225 Processing Systems*, 37:333–356, 2024.

226 Zhanpeng Zhou, Yongyi Yang, Xiaojiang Yang, Junchi Yan, and Wei Hu. Going beyond linear mode
 227 connectivity: The layerwise linear feature connectivity. *Advances in neural information processing
 228 systems*, 36:60853–60877, 2023.

229 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
 230 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv
 231 preprint arXiv:1909.08593*, 2019.

232 A Technical Appendices and Supplementary Material

233 A.1 Limitations

234 Our study is limited to a single base model and dataset, so the findings may not generalize across
 235 domains or architectures. The KL and activation patching analyses were performed on a relatively
 236 small set of prompts, making the results more illustrative than statistically exhaustive. In addition,
 237 patching identifies layer-level effects but does not capture finer-grained circuits within neurons or
 238 attention heads. Finally, while we highlight reduced computational costs and emissions, we do not
 provide a full quantitative carbon analysis.

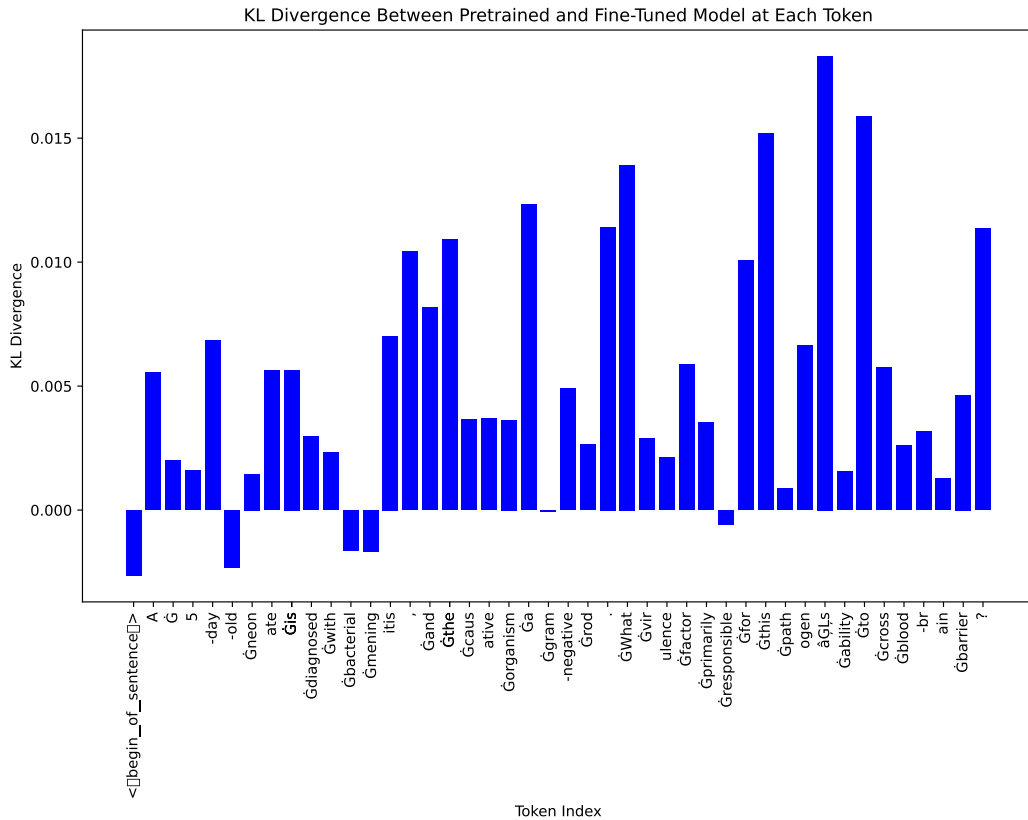


Figure 1: Token-level KL divergence between the pretrained and fine-tuned model. The plot shows localized spikes in divergence, indicating that fine-tuning selectively alters the probability distribution at specific tokens.

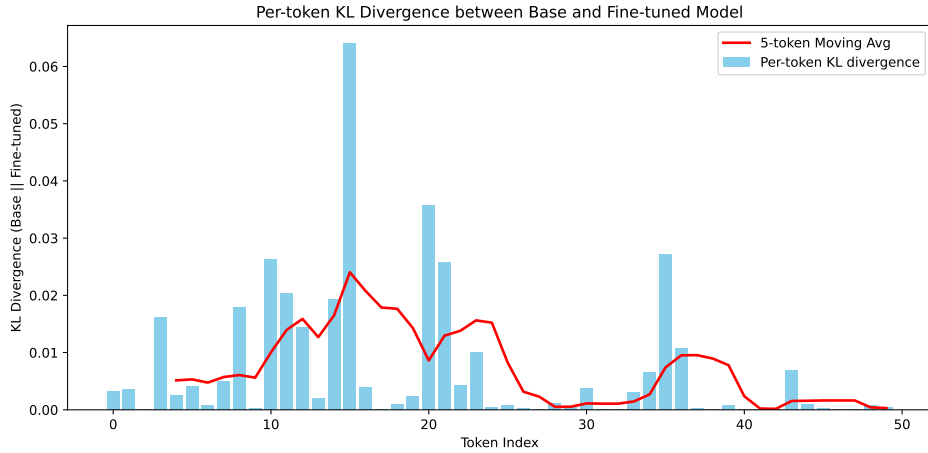


Figure 2: Token-level KL divergence between the pretrained and fine-tuned model. The divergence plot shows localized spikes, indicating that fine-tuning selectively alters the probability distribution at specific tokens.

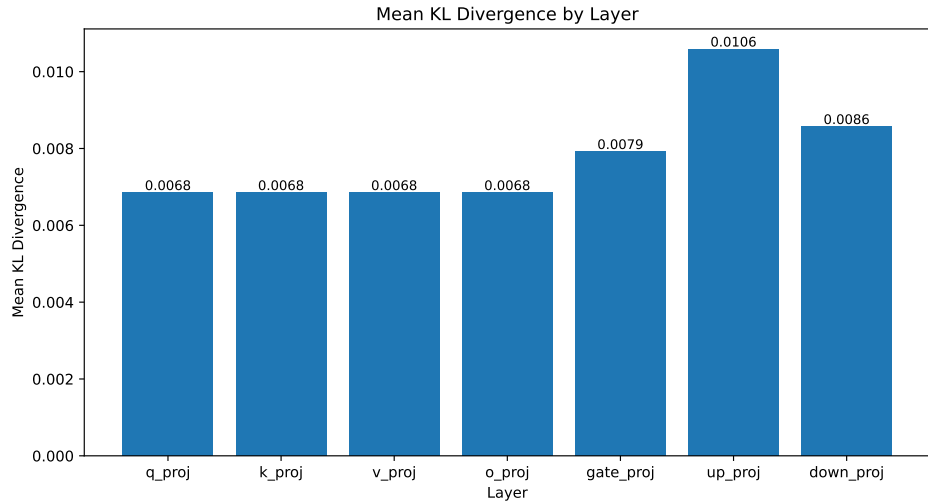


Figure 3: Mean KL divergence between the base and fine-tuned model across projection layers. Fine-tuning introduces the largest divergence in up_proj and down_proj layers, while attention projection layers (q_proj, k_proj, v_proj, o_proj) show smaller shifts.

240 **NeurIPS Paper Checklist**

241 The checklist is designed to encourage best practices for responsible machine learning research,
 242 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
 243 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
 244 follow the references and follow the (optional) supplemental material. The checklist does NOT count
 245 towards the page limit.

246 Please read the checklist guidelines carefully for information on how to answer these questions. For
 247 each question in the checklist:

- 248 • You should answer [Yes], [No], or [NA].
- 249 • [NA] means either that the question is Not Applicable for that particular paper or the
 250 relevant information is Not Available.

251 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

252 **The checklist answers are an integral part of your paper submission.** They are visible to the
253 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
254 (after eventual revisions) with the final version of your paper, and its final version will be published
255 with the paper.

256 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
257 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
258 proper justification is given (e.g., "error bars are not reported because it would be too computationally
259 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
260 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
261 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
262 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
263 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
264 please point to the section(s) where related material for the question can be found.

265 IMPORTANT, please:

- 266 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 267 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 268 • **Do not modify the questions and only use the provided macros for your answers.**

269 1. Claims

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

272 Answer: [Yes]

273 Justification: As discussed in the abstract. We have

274 Guidelines:

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

284 2. Limitations

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

286 Answer: [No]

287 Justification: We have added the limitations section in Appendix

288 Guidelines:

- 289 • The answer NA means that the paper has no limitation while the answer No means that
290 the paper has limitations, but those are not discussed in the paper.
- 291 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 292 • The paper should point out any strong assumptions and how robust the results are to
293 violations of these assumptions (e.g., independence assumptions, noiseless settings,
294 model well-specification, asymptotic approximations only holding locally). The authors
295 should reflect on how these assumptions might be violated in practice and what the
296 implications would be.
- 297 • The authors should reflect on the scope of the claims made, e.g., if the approach was
298 only tested on a few datasets or with a few runs. In general, empirical results often
299 depend on implicit assumptions, which should be articulated.

- 300 • The authors should reflect on the factors that influence the performance of the approach.
301 For example, a facial recognition algorithm may perform poorly when image resolution
302 is low or images are taken in low lighting. Or a speech-to-text system might not be
303 used reliably to provide closed captions for online lectures because it fails to handle
304 technical jargon.
- 305 • The authors should discuss the computational efficiency of the proposed algorithms
306 and how they scale with dataset size.
- 307 • If applicable, the authors should discuss possible limitations of their approach to
308 address problems of privacy and fairness.
- 309 • While the authors might fear that complete honesty about limitations might be used by
310 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
311 limitations that aren't acknowledged in the paper. The authors should use their best
312 judgment and recognize that individual actions in favor of transparency play an impor-
313 tant role in developing norms that preserve the integrity of the community. Reviewers
314 will be specifically instructed to not penalize honesty concerning limitations.

315 3. Theory assumptions and proofs

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

318 Answer: [NA]

319 Justification: [NA]

320 Guidelines:

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

331 4. Experimental result reproducibility

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

335 Answer: [Yes]

336 Justification: We are providing the full code used in this paper over the GitHub repository.

337 Guidelines:

- 338 • The answer NA means that the paper does not include experiments.
- 339 • If the paper includes experiments, a No answer to this question will not be perceived
340 well by the reviewers: Making the paper reproducible is important, regardless of
341 whether the code and data are provided or not.
- 342 • If the contribution is a dataset and/or model, the authors should describe the steps taken
343 to make their results reproducible or verifiable.
- 344 • Depending on the contribution, reproducibility can be accomplished in various ways.
345 For example, if the contribution is a novel architecture, describing the architecture fully
346 might suffice, or if the contribution is a specific model and empirical evaluation, it may
347 be necessary to either make it possible for others to replicate the model with the same
348 dataset, or provide access to the model. In general, releasing code and data is often
349 one good way to accomplish this, but reproducibility can also be provided via detailed
350 instructions for how to replicate the results, access to a hosted model (e.g., in the case
351 of a large language model), releasing of a model checkpoint, or other means that are
352 appropriate to the research performed.

- 353 • While NeurIPS does not require releasing code, the conference does require all submis-
354 sions to provide some reasonable avenue for reproducibility, which may depend on the
355 nature of the contribution. For example
356 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
357 to reproduce that algorithm.
358 (b) If the contribution is primarily a new model architecture, the paper should describe
359 the architecture clearly and fully.
360 (c) If the contribution is a new model (e.g., a large language model), then there should
361 either be a way to access this model for reproducing the results or a way to reproduce
362 the model (e.g., with an open-source dataset or instructions for how to construct
363 the dataset).
364 (d) We recognize that reproducibility may be tricky in some cases, in which case
365 authors are welcome to describe the particular way they provide for reproducibility.
366 In the case of closed-source models, it may be that access to the model is limited in
367 some way (e.g., to registered users), but it should be possible for other researchers
368 to have some path to reproducing or verifying the results.

369 5. Open access to data and code

370 Question: Does the paper provide open access to the data and code, with sufficient instruc-
371 tions to faithfully reproduce the main experimental results, as described in supplemental
372 material?

373 Answer: [Yes]

374 Justification: The model and dataset used in the paper is openly accessible from Huggingface
375 models and datasets library which we already provided the code for the same.

376 Guidelines:

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

396 6. Experimental setting/details

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

400 Answer: [Yes]

401 Justification: In the methodology section we have completely provided our training details
402 including the hyperparameters

403 Guidelines:

- 404 • The answer NA means that the paper does not include experiments.
- 405 • The experimental setting should be presented in the core of the paper to a level of detail
406 that is necessary to appreciate the results and make sense of them.

407 • The full details can be provided either with the code, in appendix, or as supplemental
408 material.

409 7. Experiment statistical significance

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

412 Answer: [Yes]

413 Justification: Yes, we have compared using activation patching mechanism which is best
414 suitable for comparing the effect of individual modules.

415 Guidelines:

- 416 • The answer NA means that the paper does not include experiments.
- 417 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
418 dence intervals, or statistical significance tests, at least for the experiments that support
419 the main claims of the paper.
- 420 • The factors of variability that the error bars are capturing should be clearly stated (for
421 example, train/test split, initialization, random drawing of some parameter, or overall
422 run with given experimental conditions).
- 423 • The method for calculating the error bars should be explained (closed form formula,
424 call to a library function, bootstrap, etc.)
- 425 • The assumptions made should be given (e.g., Normally distributed errors).
- 426 • It should be clear whether the error bar is the standard deviation or the standard error
427 of the mean.
- 428 • It is OK to report 1-sigma error bars, but one should state it. The authors should
429 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
430 of Normality of errors is not verified.
- 431 • For asymmetric distributions, the authors should be careful not to show in tables or
432 figures symmetric error bars that would yield results that are out of range (e.g. negative
433 error rates).
- 434 • If error bars are reported in tables or plots, The authors should explain in the text how
435 they were calculated and reference the corresponding figures or tables in the text.

436 8. Experiments compute resources

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

440 Answer: [Yes]

441 Justification: We have provided full information about the code, data and fine-tuning in the
442 paper.

443 Guidelines:

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

452 9. Code of ethics

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

455 Answer: [Yes]

456 Justification:

457 Guidelines:

- 458
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - 459 • If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - 460
 - 461 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
 - 462

463 10. Broader impacts

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

465 Answer: [Yes]

466 Justification: We found while fine-tuning the model on selected LoRA adaptors we can reduce the computational costs which greatly impacts on the carbon emission.

467 Guidelines:

- 468 • The answer NA means that there is no societal impact of the work performed.
- 469 • 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.
- 470 • 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.
- 471 • 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.
- 472 • 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.
- 473 • 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).
- 474
- 475
- 476
- 477
- 478
- 479
- 480
- 481
- 482
- 483
- 484
- 485
- 486
- 487
- 488
- 489
- 490
- 491

492 11. Safeguards

493 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)?

494 Answer: [NA]

495 Justification:

496 Guidelines:

- 497 • The answer NA means that the paper poses no such risks.
- 498 • 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.
- 499 • Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 500 • 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.
- 501
- 502
- 503
- 504
- 505
- 506
- 507
- 508

509 12. Licenses for existing assets

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

513 Answer: [Yes]

514 Justification: Models and data have taken from huggingface platform.

515 Guidelines:

- 516 • The answer NA means that the paper does not use existing assets.
- 517 • The authors should cite the original paper that produced the code package or dataset.
- 518 • The authors should state which version of the asset is used and, if possible, include a
519 URL.
- 520 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 521 • For scraped data from a particular source (e.g., website), the copyright and terms of
522 service of that source should be provided.
- 523 • If assets are released, the license, copyright information, and terms of use in the
524 package should be provided. For popular datasets, `paperswithcode.com/datasets`
525 has curated licenses for some datasets. Their licensing guide can help determine the
526 license of a dataset.
- 527 • For existing datasets that are re-packaged, both the original license and the license of
528 the derived asset (if it has changed) should be provided.
- 529 • If this information is not available online, the authors are encouraged to reach out to
530 the asset's creators.

531 13. New assets

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

534 Answer: [NA]

535 Justification:

536 Guidelines:

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

545 14. Crowdsourcing and research with human subjects

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

549 Answer: [NA]

550 Justification:

551 Guidelines:

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

560 15. Institutional review board (IRB) approvals or equivalent for research with human 561 subjects

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

566 Answer: [NA]

567 Justification:

568 Guidelines:

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

579 16. Declaration of LLM usage

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

584 Answer: [NA]

585 Justification:

586 Guidelines:

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