

---

# Crafting Reversible SFT Behaviors in Large Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Supervised fine-tuning (SFT) induces new behaviors in large language models,  
2       yet imposes no structural constraint on how these behaviors are distributed within  
3       the model. Existing behavior interpretation methods, such as circuit attribution  
4       approaches, identify sparse subnetworks correlated with SFT-induced behaviors  
5       post-hoc. However, such correlations do not imply *causal necessity*, limiting  
6       the ability to selectively control SFT-induced behaviors at inference time. We  
7       pursue an alternative by asking: can an SFT-induced behavior be deliberately  
8       compressed into a sparse, mechanistically necessary subnetwork, termed a *carrier*,  
9       while remaining controllable at inference time without weight modification? We  
10      propose (a) **Loss-Constrained Dual Descent (LCDD)**, which constructs such  
11      carriers by jointly optimizing routing masks and model weights under an explicit  
12      utility budget, and (b) **SFT-Eraser**, a soft prompt optimized via activation matching  
13      on extracted carrier channels, to reverse the SFT-induced behavior. Across safety,  
14      fixed-response, and style behaviors on multiple model families, LCDD yields  
15      sparse carriers that preserve target behaviors while enabling strong reversion when  
16      triggered by SFT-Eraser. Ablations further establish that the sparse structure is  
17      the key precondition for reversal: the same trigger optimization fails on standard  
18      SFT models, confirming that structure rather than trigger design is the operative  
19      factor. These results provide direct evidence that the learned carriers are causally  
20      necessary for the behaviors, pointing to a new direction for systematically localizing  
21      and selectively suppressing SFT-induced behaviors in deployed models. Code is  
22      available at <https://anonymous.4open.science/r/sft-reverse-8150/>.

## 23 1 Introduction

24      Supervised fine-tuning (SFT) is the standard approach for training large language models (LLMs) to  
25      exhibit deployment-critical behaviors such as safety alignment [1, 2, 3], instruction following [1, 4],  
26      and domain adaptation [5]. Despite its empirical success, the SFT process provides no structural  
27      guarantee on how behaviors are internally organized. The resulting behavior may be spread broadly,  
28      redundantly represented, or entangled with unrelated computation. There is no designated substructure  
29      that can be precisely targeted or analyzed.

30      This lack of structure creates two compounding problems. **First**, it prevents *causal diagnosis*. In  
31      mechanistic interpretability, a *circuit* is a sparse subnetwork of model components (attention heads,  
32      MLP layers, or residual channels) that is responsible for a specific behavior [6, 7]. Post-hoc circuit  
33      discovery identifies such subnetworks by correlating component activations with target outputs.  
34      However, correlation does not imply causation. The full model may retain redundant pathways that  
35      compensate for any ablated component. As long as behavior remains diffusely encoded, no amount  
36      of post-hoc analysis can conclusively establish whether a given substructure is truly necessary or  
37      merely correlated. **Second**, it limits *causal controllability*. Existing approaches do not provide a  
38      way to *instantiate* a sparse substructure whose presence is required for an SFT-induced behavior.

39 Behavioral and representational analyses of SFT [4, 8, 9, 10, 11, 12, 13] and circuit localization  
40 methods [14, 15, 16, 17, 18, 19] can identify structure post-hoc but do not construct it. Task-vector  
41 and model-editing methods [20, 21] show that fine-tuning weight changes can be composed and  
42 edited to add or remove capabilities, but they intervene at the parameter level and do not produce a  
43 substructure addressable via input alone under fixed weights. To the best of our knowledge, no prior  
44 work constructs a sparse substructure that is both causally necessary for an SFT-induced behavior  
45 and selectively suppressible via input alone under fixed weights.

46 We pursue a fundamentally different approach that addresses both problems at once. Rather than  
47 attempting to discover sparse structure in models that were never designed to induce it, we ask whether  
48 such structure can be deliberately constructed during training. When structure is imposed by design,  
49 the causal role of the carrier becomes substantially more identifiable: the behavior is intentionally  
50 concentrated within the carrier, and components outside it are maintained in their base-model state  
51 (i.e., before SFT). This makes causal necessity more directly testable via input interventions that  
52 target the carrier under fixed weights, rather than relying on ablations that may be confounded by  
53 redundant pathways [22]. If an input intervention suppresses SFT behavior while all model weights  
54 remain fixed, the sparse substructure it targets is demonstrated to be the causal bottleneck, not merely  
55 a correlated subnetwork. Beyond its role as a causal diagnostic, a deployed model with a crafted  
56 sparse carrier could support inference-time behavioral auditing or selective suppression without any  
57 weight modification, enabling more modular post-training control.

58 These considerations motivate our research question:

59 *Can SFT-induced behaviors be intentionally concentrated into sparse and mechanistically necessary*  
60 *carriers that remain controllable at inference time without modifying model weights?*

61 We term such a parameter substructure the *sparse behavioral carrier*. We study this question via two  
62 methods. We introduce **Loss-Constrained Dual Descent (LCDD)** to compress SFT behavior into  
63 a sparse carrier via mask-based compression, and **SFT-Eraser** to reverse the behaviors introduced  
64 by SFT via activation matching. SFT-Eraser treats a *reversal* as successful when it produces both  
65 behavioral suppression and distributional reversion toward the base model. Our main contributions  
66 are as follows:

- 67 • To the best of our knowledge, we are the first to investigate the feasibility of concentrating SFT-  
68 induced behaviors into sparse, mechanistically necessary substructures, enabling their precise  
69 control at inference time without modifying model weights.
- 70 • We propose **LCDD**, a framework for constructing sparse behavioral carriers by jointly optimizing  
71 routing masks and model weights under explicit utility constraints. To validate causal necessity  
72 under fixed weights, we further introduce **SFT-Eraser**, an input-trigger<sup>1</sup> protocol that tests whether  
73 targeted input interventions can selectively suppress SFT-induced behaviors.
- 74 • We demonstrate that LCDD and SFT-Eraser enable inference-time reversibility across three be-  
75 havior types and four model families. These results provide strong evidence for the feasibility of  
76 concentrating SFT-induced behaviors into sparse, causally necessary carriers, and show that the  
77 constructed carriers are not merely correlated with the behavior. Ablations further confirm that the  
78 sparse structure, rather than the trigger design, is the operative factor.

## 79 2 Related Work

80 **Post-hoc analysis of SFT behaviors.** Recent work characterizes SFT observationally. LIMA shows  
81 that a small high-quality SFT set suffices for strong assistant behavior [4], and the Superficial Safety  
82 Alignment Hypothesis frames alignment as a lightweight “fulfill vs. refuse” controller with sparse  
83 safety-critical components [8]. Depth-wise analyses find that alignment updates concentrate in specific  
84 layers [9, 10], and other work addresses catastrophic forgetting and objective mismatch [11, 12, 13].  
85 SafeSeek frames circuit extraction as mask optimization and demonstrates highly sparse circuits for  
86 backdoor and alignment behaviors [17], and Depth Charge shows that targeted interventions on deep  
87 attention heads can substantially alter jailbreak success [24]. These methods identify or characterize  
88 structure post-hoc. By contrast, our work constructs sparse behavioral carriers during training, so that  
89 the carrier is by construction the complete locus of the behavior, enabling causal validation rather  
90 than correlational attribution.

---

<sup>1</sup>While we borrow the term *trigger* from the backdoor literature [23], our trigger recovers the full functionalities of the base model before SFT, whereas a backdoor trigger activates a specific implanted behavior.

91 **Constructive and weight-space approaches.** Task arithmetic defines task vectors as  $\Delta W =$   
 92  $W_{\text{ft}} - W_{\text{base}}$  and shows that adding or negating them can compose or suppress behaviors [20, 21].  
 93 However, it operates with dense global weight edits and does not produce a sparse activation-level  
 94 structure addressable through input alone under fixed weights. Parameter-efficient adaptation via  
 95 task-specific binary masks [25] is closer in spirit; our mask-only ablation (Appendix E) finds it  
 96 insufficient for achieving the sparsity depth required for reliable reversal. Most directly related  
 97 is Gao et al. [18], which trains weight-sparse models from scratch on narrow tasks and finds that  
 98 node-level sparsity yields highly interpretable circuits. Our work shares the same design principle  
 99 but differs fundamentally in goal: rather than training sparse models de novo, we impose sparsity on  
 100 the SFT weight delta  $\Delta W$  in already-trained chat models to construct behavioral carriers. Broader  
 101 differentiable mask and circuit discovery methods [14, 15, 16, 26, 27] provide methodological  
 102 precedents but do not address SFT carrier construction.

### 103 3 Method

104 Our goal is to develop a *sparse behavioral carrier* satisfying two properties: (1) it maintains a minimal  
 105 parameter substructure sufficient to sustain the SFT-induced behavior under normal inference; and  
 106 (2) there exists a fixed-weight input intervention targeting the carrier that can reliably suppress the  
 107 SFT-induced behavior. Property (2) relies on *causal necessity*: the carrier must be the complete  
 108 locus of the behavior by construction, with model components outside the carrier reduced to their  
 109 base-model state. Sparsity is essential to both properties: a compact, identifiable substructure can be  
 110 precisely targeted at inference time, whereas a diffusely encoded behavior offers no designated target  
 111 for such intervention. We pursue this in two stages as shown in Figure 1: crafting the sparse carrier  
 112 via **LCDD** (Section 3.2), and validating its necessity via **SFT-Eraser** (Section 3.3). Both stages build  
 113 on the delta parameterization introduced in Section 3.1.

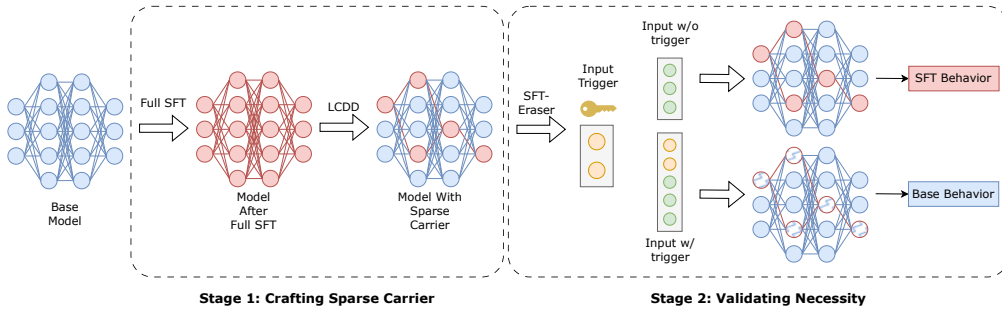


Figure 1: An overview of the LCDD + SFT-Eraser pipeline. **Stage 1:** standard SFT distributes the induced behavior broadly across model parameters (red), followed by LCDD that compresses the SFT-induced behavior into a sparse carrier. Components outside the carrier are reduced to their base-model state by construction (blue). **Stage 2:** SFT-Eraser optimizes a soft trigger. Under normal inference without the trigger, the carrier preserves SFT behavior; while with the trigger, carrier activations are driven toward the base model.

#### 114 3.1 Preliminaries: Weight-Change Parameterization and Mask Implementation

115 **Weight-change parameterization.** We represent SFT-induced computation as the additive weight  
 116 change from fine-tuning. Let  $W_{\text{base}}$  and  $W_{\text{ft}}$  denote the respective weights; we define the *weight delta*  
 117  $\Delta W = W_{\text{ft}} - W_{\text{base}}$  as the SFT-induced parameter change. Our masked model is then parameterized  
 118 as

$$W = W_{\text{base}} + M \odot \Delta W, \quad (1)$$

119 where  $M \in \{0, 1\}^{\text{shape}(W)}$  is a binary mask. Setting  $M = \mathbf{0}$  recovers the base model exactly and  
 120  $M = \mathbf{1}$  recovers the fully fine-tuned model. Our goal is to find the sparsest  $M$  that still preserves  
 121 task behavior. Minimizing carrier size is deliberate: a sparser carrier is more identifiable and more  
 122 precisely targetable at inference time.

123 **Weight mask implementation.** We structure sparsity at the level of activation channels (individual  
 124 rows and columns of weight matrices) rather than arbitrary weight elements, following the design  
 125 principle that node-level weight sparsity produces identifiable and interpretable circuits in which

126 information flow can be traced through designated channel bottlenecks [18]. Prior work has explored  
 127 similar gradient-based mask optimization strategies for circuit discovery. Gao et al. [18] train weight-  
 128 sparse models from scratch such that each neuron reads from and writes to only a small number  
 129 of residual channels, making circuit boundaries well-defined by construction. Yu et al. [17] apply  
 130 gradient-based binary mask optimization with the straight-through estimator (STE) to discover safety  
 131 circuits in full models, gating unit outputs (neurons or attention heads) of the complete model weights  
 132  $W$ . Our setup differs in two respects. First, we apply structured masking specifically to the SFT  
 133 weight change  $\Delta W$ , isolating SFT-induced computation from the base model. Second, we use row  
 134 and column gate vectors rather than unit-output masks, yielding a rank-1 structured mask that controls  
 135 which activation channels carry the SFT-induced delta computation.

136 Formally, for each weight matrix, we introduce row and column gate vectors  $m_{\text{row}}$  and  $m_{\text{col}}$ . The  
 137 delta contribution to a linear layer can be written as  $[(\mathbf{x} \odot m_{\text{row}}) \Delta W] \odot m_{\text{col}}$ , where gating the  
 138 input activations selects rows of  $\Delta W$  and gating the output activations selects columns. This  
 139 activation-gating view is equivalent to weight masking with the rank-1 structure

$$M_{jk} = m_{\text{row},j} m_{\text{col},k}. \quad (2)$$

140 Full derivations for both FFN and attention layers are given in Appendix A. Per transformer layer,  
 141 we define 8 independent gate groups covering all delta-carrying weight matrices and the complete  
 142 gate-to-weight-mask mapping is given in Table 5 of Appendix A.

### 143 3.2 LCDD: Utility-Budgeted Sparse Carrier Crafting

144 LCDD jointly optimizes the mask  $M$  and the weight change  $\Delta W$  to compress SFT behavior into  
 145 a sparse carrier while keeping utility loss within an explicit budget. We first describe how  $M$  is  
 146 parameterized for optimization, and then present the main objective and the optimization procedure.

147 **Parameterizing  $M$ .** We cannot directly optimize  $M$  due to its discrete nature. Following Gao  
 148 et al. [18], we parameterize each gate by a continuous logit  $\theta_i \in \mathbb{R}$ . During the forward pass, each  
 149 gate is binarized via the Heaviside function as  $m_i = \mathbf{1}[\theta_i > 0]$ . Gradients are approximated by  
 150 the straight-through estimator (STE). The sigmoid  $\sigma(\theta_i)$  serves as a differentiable proxy for gate  
 151 activation and defines the sparsity penalty:

$$\mathcal{L}_{\text{sparsity}} = \sum_i \sigma(\theta_i). \quad (3)$$

152 Here  $i$  indexes all individual gate elements across all gate groups and all transformer layers. Minimizing  
 153  $\mathcal{L}_{\text{sparsity}}$  progressively deactivates weights in the delta path. The task loss gradient via STE retains  
 154 gates necessary for behavior. Full details are given in Appendix B.

155 **Main objective.** After parameterizing  $M$  via  $\theta$ , the training objective is:

$$\min_{\theta, \Delta W} \sum_i \sigma(\theta_i) \quad \text{s.t.} \quad \mathcal{L}_{\text{task}}(\theta, \Delta W) \leq \epsilon. \quad (4)$$

156 where  $\mathcal{L}_{\text{task}}$  measures the model’s task loss under the current mask and weight configuration, and  $\epsilon$  is  
 157 an explicit budget that caps the tolerated degradation in task performance. In our main experiments  
 158  $\mathcal{L}_{\text{task}}$  is SFT cross-entropy. The constraint-based design provides a directly interpretable control  
 159 surface: one specifies a concrete tolerance on utility degradation and lets the optimizer find the  
 160 sparsest carrier within that budget, rather than tuning a regularization coefficient whose effect on  
 161 utility is indirect.

162 **Optimization.** To solve Eq. (4), we use an adaptive penalty method: the constraint is absorbed  
 163 into the objective with a multiplier  $\lambda_t$  updated online based on constraint violation. A linear sparsity  
 164 warmup factor  $\rho_t$  gradually engages the penalty during early training, allowing the model to stabilize  
 165 before compression begins. The multiplier grows when task loss is below budget (intensifying  
 166 sparsification pressure) and shrinks when loss exceeds budget (allowing utility recovery). This loop  
 167 continuously renegotiates the sparsity-utility trade-off without manual multiplier schedules or nested  
 168 inner-loop minimization. Full details are given in Algorithm 1 (Appendix C).

169 To prevent behavior from routing through the token embedding or output projection layers, we  
 170 freeze these components throughout training, confining compression to intermediate transformer  
 171 layers. Training halts when carrier sparsity stalls across consecutive checkpoints or when utility  
 172 degradation exceeds a threshold. We empirically verify that jointly optimizing masks and weight  
 173 deltas is necessary for achieving sufficient sparsity depth in Appendix E.

174 **3.3 SFT-Eraser: Trigger-Based Reversal via Residual-Write Channel Matching**

175 Given a sparse carrier crafted by LCDD, SFT-Eraser asks: can a fixed-weight input intervention  
 176 suppress SFT behavior by targeting carrier activations directly? If the trigger suppresses SFT behavior  
 177 while reverting the output distribution toward the base model, the joint outcome implies that the  
 178 carrier channels are mechanistically necessary rather than merely correlated with the behavior.

179 **Carrier channel extraction.** After LCDD training, we identify active channels as those with  
 180  $\theta_i > 0$ . Among these, we select *residual write channels* as the trigger target:

$$\mathcal{C}_{\text{write}} = \mathcal{C}_{\text{ffn-write}} \cup \mathcal{C}_{\text{attn-write}}, \quad (5)$$

181 where  $\mathcal{C}_{\text{ffn-write}}$  and  $\mathcal{C}_{\text{attn-write}}$  are the sets of active output-gate channels of the FFN and attention  
 182 delta paths, respectively (corresponding to  $m_{\text{ffn-write}}$  and  $m_{\text{attn-write}}$  in Section 3.1). We focus on  
 183 write channels because they form the carrier’s primary interface with the residual stream, where the  
 184 delta computation contributes to downstream layer inputs. Targeting these channels is more causally  
 185 precise than targeting intermediate hidden states.

186 **Trigger objective.** We optimize a soft trigger embedding sequence  $t = [t_1, \dots, t_L]$  prepended to  
 187 the input prompt  $x$ , while all model weights remain frozen. The trigger loss combines three terms:

$$\mathcal{L}_{\text{trigger}} = \mathcal{L}_{\text{MSE}} + \alpha \mathcal{L}_{\text{KL}} + \beta \mathcal{L}_2. \quad (6)$$

188 Denoting by  $\mathcal{A}_{\mathcal{C}}(\cdot)$  the residual-stream write activations at channel set  $\mathcal{C}$ , the primary term  $\mathcal{L}_{\text{MSE}}$   
 189 drives the carrier’s write activations toward those of the base model on the same input:

$$\mathcal{L}_{\text{MSE}} = \text{MSE}(\mathcal{A}_{\mathcal{C}_{\text{write}}}(\text{SFT}([t; x])), \mathcal{A}_{\mathcal{C}_{\text{write}}}(\text{Base}(x))). \quad (7)$$

190 Since the trigger acts on the full input and can perturb activations beyond the carrier channels,  $\mathcal{L}_{\text{MSE}}$   
 191 alone does not constrain the trigger’s net effect on the output distribution. The auxiliary term  $\mathcal{L}_{\text{KL}}$   
 192 adds output-level supervision, ensuring that the trigger produces distributional reversion toward the  
 193 base model rather than an artifact of carrier-level matching alone:

$$\mathcal{L}_{\text{KL}} = \text{KL}_{\text{tail-}k}(\text{Base}(x) \parallel \text{SFT}([t; x])), \quad (8)$$

194 computed over the final  $k$  response positions to focus the signal on the behavioral regime. The  
 195 regularizer  $\mathcal{L}_2 = \frac{1}{L} \sum_{i=1}^L \|t_i\|_2^2$  prevents trigger tokens from drifting into degenerate high-magnitude  
 196 embedding regions where activations become uninterpretable; trigger optimization also applies  
 197 PGD-style  $\ell_2$  projection after each gradient update (details in Appendix D). The necessity of  $\mathcal{L}_{\text{MSE}}$  is  
 198 validated empirically in Ablation 2 (Section 4.3).

199 **4 Experiments**

200 In this section, we conduct experiments to answer the following research questions:

- 201 • **RQ1:** Can LCDD craft sparse SFT carriers that preserve target behavior under normal inference?
- 202 • **RQ2:** Do LCDD-crafted carriers yield reliable trigger-based reversal under fixed weights?
- 203 • **RQ3:** Is the crafted sparse structure a necessary precondition for reversal, rather than the trigger  
 204 design itself?

205 We describe the experimental setup in Section 4.1 and present results in Sections 4.2 and 4.3.

206 **4.1 Experimental Setup**

207 **Tasks.** We select three types of SFT tasks in increasing order of structural complexity:

- 208 • **Fixed Response.** The model is trained to respond “I don’t know” on all instruction prompts from  
 209 Alpaca [28]. This behavior has a discrete lexical signature and is input-unconditional, making it  
 210 the simplest carrier to craft and measure.
- 211 • **Safety Alignment.** The model is trained on WildJailbreak [29] harmful prompts interleaved with  
 212 benign prompts (1:1 ratio). This behavior is input-conditional (triggered by semantic harm) and  
 213 requires the carrier to encode content-sensitive routing, placing greater demands on sparsification  
 214 than Fixed Response.

215 • **Shakespeare Style.** The model is trained on modern-to-Shakespeare conversational pairs [30].  
 216 This is a distributional behavior with no discrete trigger condition, requiring the carrier to capture  
 217 broad stylistic shifts across the output distribution, making it the most structurally demanding of  
 218 the three tasks.

219 All main runs use 5,000 training samples per task.

220 **Models.** We evaluate on four chat LLM families: Qwen3-0.6B [31], DeepSeek-R1-Distill-Llama-  
 221 8B [32], Mistral-7B-Instruct-v0.3 [33], and Vicuna-7B-v1.5 [34]. Model selection follows two  
 222 criteria: architectural diversity across scales and pretraining regimes, and for the safety task, absence  
 223 of built-in safety alignment so that refusal behavior is genuinely induced by SFT rather than already  
 224 present in the pretrained weights.

225 **Metrics.** We use the following metrics to evaluate how well the sparse carrier preserves SFT  
 226 behavior before intervention (which we call *carrier fidelity*), and how completely it reverts after  
 227 triggering.

228 • **Fixed Response:** Strict fixed-response rate, computed by keyword matching with a response length  
 229 cap (Appendix F.4).

230 • **Safety:** WildGuard refusal rate [35] as the primary safety metric (higher is safer), WildGuard  
 231 harmfulness rate [35] as a complementary output-side measure (lower is safer), and HarmBench  
 232 ASR [36] as an attack-side measure (lower is safer). MMLU [37] and HellaSwag [38] track general  
 233 utility.

234 • **Shakespeare:** LLM-judge score on a 0–5 Shakespearean authenticity scale, evaluated by Qwen3-  
 235 14B [31] using the rubric in Appendix F.4. MMLU and HellaSwag track general utility.

236 • **KL metrics:** Computed by running each evaluated model on SFT-generated reference responses  
 237 with teacher forcing (i.e., scoring on fixed reference tokens rather than free generation), averaged  
 238 token-level over response tokens.  $KL(SFT||LCDD)$  measures carrier fidelity before intervention  
 239 ( $\downarrow$  better).  $KL(SFT||Trig)$  and  $KL(Base||Trig)$  jointly characterize distributional movement af-  
 240 ter triggering: larger  $KL(SFT||Trig)$  indicates stronger divergence from SFT behavior; smaller  
 241  $KL(Base||Trig)$  indicates closer reversion to the base model.

242 **Training setup.** All runs follow the same two-phase pipeline (full SFT then LCDD), with the  
 243 token embedding layer and output vocabulary projection frozen throughout LCDD training to confine  
 244 compression to intermediate transformer layers. Trigger optimization targets residual write channels,  
 245 uses trigger length 20, tail- $k$  KL ( $k = 8$ ), and norm-constrained gradient updates ( $\ell_2$  max-norm = 1.0)  
 246 uniformly across all settings. Full hyperparameter details are in Appendix F.

## 247 4.2 Main Results

248 We present results across three behavior types. For each behavior type, we report carrier fidelity (how  
 249 closely LCDD preserves SFT behavior before intervention) and reversal quality (how completely  
 250 SFT-Eraser suppresses the behavior after triggering).

### 251 4.2.1 Fixed Response

Table 1: **Fixed Response** results. “Trig” (LCDD+Trigger) denotes the LCDD model with SFT-Eraser trigger applied. KL values are computed on SFT-generated references.

Model	Sparsity	Fixed-Response Rate				KL Divergence		
		@Base	@SFT	@LCDD	@Trig	SFT  LCDD	SFT  Trig	Base  Trig
Qwen3-0.6B	83.66%	0.0%	100.0%	97.5%	0.0%	0.040	2.839	0.633
DeepSeek-8B	73.38%	0.0%	100.0%	99.0%	1.0%	0.010	1.811	1.037
Mistral-7B	76.72%	0.0%	100.0%	100.0%	0.5%	0.002	3.526	1.437
Vicuna-7B	38.28%	0.0%	100.0%	100.0%	0.0%	1.436	4.810	0.024

252 As shown in Table 1, the Fixed Response task provides the clear evidence that LCDD can craft  
 253 a sparse carrier that is both behaviorally faithful before intervention and reliably reversible after  
 254 triggering. LCDD preserves near-perfect fixed-response behavior across all four models (97.5–  
 255 100%). For three models (Qwen, DeepSeek, Mistral),  $KL(SFT||LCDD)$  is uniformly low (0.002–  
 256 0.040), confirming that the sparse carrier is behaviorally near-indistinguishable from full SFT before  
 257 intervention. Under trigger, fixed-response rate falls to near-zero in all cases (0.0–1.0%). The large  
 258 post-trigger  $KL(SFT||Trig)$  values (1.811–4.810) confirm strong divergence from SFT behavior, while  
 259  $KL(Base||Trig)$  (0.024–1.437) quantifies distributional movement toward the base model.

260 **4.2.2 Safety**

Table 2: **Safety** results. “Trig” denotes LCDD+Trigger. All behavioral metric values in %. HB = HarmBench ASR (↓ safer); WG Ref. = WildGuard refusal rate (↑ safer); WG Harm. = WildGuard harmfulness rate (↓ safer). Sparsity is the LCDD mask sparsity. KL values are computed on SFT-generated references; the  $P \parallel Q$  column indicates the two distributions compared.

Model	Sparsity	Condition	Safety Metrics			Utility Metrics		KL Divergence	
			HB ↓	WG Ref. ↑	WG Harm. ↓	MMLU	HellaSwag	$P \parallel Q$	Value
Qwen3-0.6B	81.05%	Base	43.0	32.0	52.0	40.3	37.5		
		SFT	1.5	94.0	2.0	46.0	38.7	SFT    LCDD	0.069
		LCDD	2.5	88.5	3.5	42.2	38.1	SFT    Trig	0.120
		Trig	36.5	39.5	41.5	45.1	37.9	Base    Trig	0.196
DeepSeek-8B	81.38%	Base	19.0	47.5	38.5	54.1	55.6		
		SFT	0.0	100.0	0.0	55.3	57.1	SFT    LCDD	0.256
		LCDD	1.5	96.5	2.0	54.8	56.3	SFT    Trig	0.391
		Trig	13.0	56.5	30.5	49.2	54.1	Base    Trig	0.306
Mistral-7B	59.34%	Base	58.5	25.5	59.5	59.8	65.8		
		SFT	0.0	100.0	0.0	50.1	53.9	SFT    LCDD	0.269
		LCDD	1.0	97.5	1.0	57.0	63.6	SFT    Trig	0.620
		Trig	17.0	69.0	21.5	54.1	62.3	Base    Trig	1.282
Vicuna-7B	84.40%	Base	20.5	76.5	20.5	48.5	56.6		
		SFT	0.0	99.0	0.0	47.1	56.5	SFT    LCDD	0.192
		LCDD	2.0	98.0	1.5	48.6	56.8	SFT    Trig	1.081
		Trig	22.0	74.5	23.0	47.9	56.5	Base    Trig	0.014

261 As shown in Table 2, all four models maintain high safety under LCDD (WG refusal 88.5–98.0%;  
 262 HB ASR 1.0–2.5%), confirming that LCDD successfully compresses safety behavior into the sparse  
 263 carrier. Under trigger, safety degrades along both axes for all models: refusal rates fall and attack  
 264 success rates rise, with the direction of change consistent across architectures.

265 Reversal magnitude varies across models, reflecting differences in carrier geometry. Notably, be-  
 266 havioral and distributional metrics capture partially independent aspects of reversal: Vicuna-7B  
 267 achieves near-perfect distributional reversion ( $KL(\text{Base} \parallel \text{Trig}) = 0.014$ ) despite moderate behavioral  
 268 suppression (WG:  $-23.5\text{pp}$ ), while Mistral-7B shows the strongest behavioral suppression (WG:  
 269  $-28.5\text{pp}$ ) with a larger residual distributional gap ( $KL(\text{Base} \parallel \text{Trig}) = 1.282$ ). Both patterns confirm  
 270 that the trigger successfully targets the carrier in all cases.

271 Utility is stable under LCDD for all models and shows only modest degradation under trigger (at  
 272 most 5.6pp on MMLU for DeepSeek), indicating that the trigger targets carrier-specific channels  
 273 without broadly disrupting general language modeling.

274 **4.2.3 Shakespeare Style**

Table 3: **Shakespeare style** results. “Trig” denotes LCDD+Trigger. Judge = LLM-judge score (0–5,  
 ↑ more Shakespearean), evaluated by Qwen3-14B (rubric in Appendix F.4). MMLU and HellaSwag  
 values in %. KL values computed on SFT-generated references; the  $P \parallel Q$  column indicates the  
 distributions compared. Sparsity is the LCDD mask sparsity.

Model	Sparsity	Condition	Style	Utility Metrics		KL Divergence	
			Judge ↑	MMLU	HellaSwag	$P \parallel Q$	Value
Qwen3-0.6B	56.34%	Base	0.03	40.3	37.5		
		SFT	0.79	45.5	38.0	SFT    LCDD	0.265
		LCDD	0.44	41.8	38.6	SFT    Trig	2.405
		Trig	0.04	41.9	37.3	Base    Trig	0.446
DeepSeek-8B	54.54%	Base	0.01	54.1	55.7		
		SFT	1.73	54.5	56.1	SFT    LCDD	1.520
		LCDD	0.98	55.4	56.8	SFT    Trig	4.301
		Trig	0.005	37.4	51.5	Base    Trig	0.531
Mistral-7B	50.46%	Base	0.07	59.8	65.8		
		SFT	1.75	47.9	55.0	SFT    LCDD	0.853
		LCDD	1.18	55.4	62.7	SFT    Trig	2.643
		Trig	0.07	42.3	53.5	Base    Trig	2.052
Vicuna-7B	29.88%	Base	0.05	48.5	56.6		
		SFT	0.80	46.6	56.1	SFT    LCDD	0.636
		LCDD	0.26	47.3	56.0	SFT    Trig	1.702
		Trig	0.04	45.8	55.8	Base    Trig	0.094

275 The Shakespeare task is the most demanding for LCDD: achieved sparsity is substantially lower  
276 (29.9–56.3%) than in Fixed Response or safety runs, and *carrier fidelity* (how closely the sparse  
277 carrier reproduces full SFT behavior before intervention, measured by  $\text{KL}(\text{SFT}||\text{LCDD})$  and the  
278 task-level Judge score) is correspondingly reduced: LCDD Judge scores (0.26–0.98) recover only part  
279 of the SFT-level style signal (0.79–1.75). This is consistent with the expectation that distributional  
280 stylistic behavior is more diffusely encoded across parameters than fixed-response or conditional  
281 refusal behaviors.

282 Even where carrier fidelity is partial, trigger intervention suppresses style effectively in all four cases:  
283 Judge@Trig falls to near-base levels (0.005–0.07) regardless of the LCDD starting point, indicating  
284 that even a partially concentrated carrier can be targeted successfully by the trigger.  $\text{KL}(\text{SFT}||\text{Trig})$  is  
285 substantially larger than  $\text{KL}(\text{SFT}||\text{LCDD})$  in all cases, confirming that post-trigger behavior diverges  
286 from SFT well beyond the LCDD baseline.

287 A consistent side effect in Shakespeare runs is larger MMLU degradation under trigger compared  
288 to Fixed Response or safety conditions (e.g., DeepSeek: 55.4%→37.4%). We conjecture that style  
289 carrier channels overlap more with general language modeling representations than behavior-specific  
290 carrier channels, making them harder to target without collateral disruption.

#### 291 4.2.4 Cross-Task Observations

292 Three patterns are consistent across all 12 model-task combinations:

- 293 • **Sparsity reflects task localizability.** Fixed Response compresses most aggressively (73–84%),  
294 safety to an intermediate level (59–84%), and style least so (30–56%). This ordering reflects the  
295 structural complexity of each behavior: a fixed lexical response can be concentrated in a small  
296 parameter substructure, whereas distributional stylistic behavior requires broader representational  
297 support and resists aggressive compression.
- 298 • **Reversal direction is universal; magnitude is heterogeneous.** All 12 combinations show move-  
299 ment in the correct direction after triggering. Variation in degree reflects model- and task-specific  
300 carrier geometry, not method failure.
- 301 • **Carrier fidelity and reversal quality are partially decoupled.** Poor carrier fidelity does not  
302 preclude effective trigger-based suppression. In the Shakespeare task, where carrier fidelity is lower  
303 than in Fixed Response or safety runs, the trigger still successfully suppresses style in all four  
304 models. This suggests that trigger effectiveness depends on the presence of the crafted bottleneck  
305 rather than on how completely the carrier captures the behavior.

306 Together, these results answer RQ1 and RQ2 affirmatively. LCDD successfully crafts sparse carriers  
307 that preserve SFT behavior under normal inference across all three behavior types and four model  
308 families. SFT-Eraser achieves reliable trigger-based reversal in all 12 model-task combinations, with  
309 consistent direction of change even where magnitude varies.

#### 310 4.3 Ablation Study

311 We investigate RQ3 along two dimensions. First, we test whether the crafted sparse structure is a  
312 necessary precondition for reversal, or whether the same trigger optimization succeeds on any SFT  
313 model. Second, we examine whether the circuit-targeted MSE term in SFT-Eraser is necessary, or  
314 whether output-level KL pressure alone suffices. Both experiments are conducted on DeepSeek-R1-  
315 Distill-Llama-8B under the safety task, using the main experiment as the reference baseline. We  
316 additionally verify the necessity of joint weight optimization in LCDD in Appendix E.

317 **Setup.** We factor the experiment along two dimensions: model structure (LCDD vs. plain SFT)  
318 and trigger objective (circuit-targeted vs. output-only), yielding four conditions. The circuit-targeted  
319 trigger uses the full loss  $\mathcal{L}_{\text{trigger}} = \mathcal{L}_{\text{MSE}} + \alpha \mathcal{L}_{\text{KL}} + \beta \mathcal{L}_2$ ; the output-only trigger removes  $\mathcal{L}_{\text{MSE}}$ ,  
320 leaving  $\alpha \mathcal{L}_{\text{KL}} + \beta \mathcal{L}_2$ . Each trigger objective is applied to both the LCDD model and the plain SFT  
321 model, giving four conditions in total: LCDD + circuit trigger (main method), LCDD + output-only  
322 trigger, SFT + circuit trigger, and SFT + output-only trigger.

323 **Finding 1: the sparse carrier structure is the decisive factor.** Both LCDD conditions substantially  
324 outperform both SFT conditions in behavioral reversal (−40pp/−31pp vs. −21pp/−24pp). The KL  
325 metrics reveal the mechanism. SFT triggers do perturb the output ( $\text{KL}(\text{SFT}||\text{Trig}) \approx 0.21$ ). However,

Table 4: Structural dependence ablation: safety reversal and distributional reversion on DeepSeek-R1-Distill-Llama-8B. WG = WildGuard refusal rate; MMLU values are post-trigger.

Model	Condition	WG Refusal			KL Divergence		MMLU
		No Trig	+Trig	$\Delta$	SFT  Trig	Base  Trig	
Base	(reference)	47.5%	–	–	–	–	54.1%
LCDD	No trigger	96.5%	–	–	–	–	54.8%
	+ circuit	96.5%	56.5%	–40pp	0.391	0.306	49.2%
	+ output-only	96.5%	65.5%	–31pp	0.440	0.322	37.2%
SFT	No trigger	100.0%	–	–	–	–	55.3%
	+ circuit	100.0%	78.5%	–21pp	0.221	0.539	42.0%
	+ output-only	100.0%	76.5%	–24pp	0.208	0.512	40.9%

326 the triggered SFT model cannot revert toward the base model distribution ( $\text{KL}(\text{Base}||\text{Trig}) \approx 0.51$ –  
 327 0.54). LCDD conditions achieve  $\text{KL}(\text{Base}||\text{Trig}) \approx 0.31$ . Triggered outputs genuinely revert toward  
 328 the base model. In LCDD models, components outside the carrier are reduced to their base-model  
 329 state by construction. The carrier is therefore the only locus of SFT-induced computation. Any  
 330 trigger-based reversion must be mediated through the carrier, not non-carrier channels. The trigger  
 331 optimization procedure alone, without a crafted carrier, is insufficient.

332 **Finding 2: the circuit MSE term improves targeting precision and reduces utility cost.** Re-  
 333 moving  $\mathcal{L}_{\text{MSE}}$  (LCDD + output-only trigger) reduces behavioral reversal from –40pp to –31pp. It  
 334 also substantially increases MMLU degradation (–17.6pp vs. –5.6pp).  $\text{KL}(\text{Base}||\text{Trig})$  remains  
 335 similar (0.322 vs. 0.306). The MSE term provides a gradient signal focused on carrier channels.  
 336 This improves optimization efficiency and preserves utility by avoiding disruption of non-carrier  
 337 activations.

## 338 5 Conclusion

339 We investigated whether SFT-induced behaviors can be deliberately compressed into sparse behavioral  
 340 carriers that are preserved under normal inference yet reversibly suppressible by an input trigger at  
 341 inference time without weight modification. LCDD crafts such carriers through utility-budgeted joint  
 342 optimization of routing masks and weight deltas; SFT-Eraser validates their mechanistic necessity  
 343 via activation matching. Our experiments and ablations establish that joint weight optimization is  
 344 required to craft a reliable sparse bottleneck, and that the bottleneck, rather than the trigger design, is  
 345 the necessary precondition for trigger-based reversal.

346 **Discussion.** Beyond its role as a mechanistic diagnostic, the sparse carrier framework suggests  
 347 several longer-term directions. Crafted carriers provide a controlled testbed for interpretability  
 348 research: unlike post-hoc circuit discovery, a carrier with known structure enables systematic study  
 349 of intervention transferability and behavioral robustness. The constructability of a behavior may  
 350 itself serve as a diagnostic property, since behaviors that resist LCDD compression are likely more  
 351 diffusely encoded; this question connects to a minimum description length interpretation [39], in  
 352 which carrier sparsity is a computable proxy for the description complexity of the SFT behavioral  
 353 transformation. The most immediate open problem is extending reversal to discrete input tokens  
 354 rather than continuous soft prompts, which would make the behavioral interface genuinely deployable;  
 355 a second direction is studying how pretraining scale affects the achievable sparsity–utility frontier for  
 356 crafted carriers.

357 **Limitations.** Two limitations bound the current study. First, SFT-Eraser optimizes in continuous  
 358 embedding space. The trigger is a soft prompt realized as a learned token-embedding sequence rather  
 359 than a discrete hard-token string. The reversal evidence is valid as a mechanistic diagnostic under  
 360 controlled conditions, but extending this to discrete hard-token triggers requires further investigation.  
 361 Second, the structural necessity ablation is conducted on a single model and task combination.  
 362 Whether the same causal relationship holds consistently across all model families and behavior  
 363 types remains to be verified more systematically. We also note that the same methodology could  
 364 in principle be used to remove safety alignment from deployed models; this dual-use risk warrants  
 365 careful consideration in any deployment context.

## References

- 366
- 367 [1] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,  
368 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to  
369 follow instructions with human feedback. *Advances in neural information processing systems*,  
370 35:27730–27744, 2022.
- 371 [2] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
372 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open  
373 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 374 [3] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,  
375 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai:  
376 Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 377 [4] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma,  
378 Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural  
379 Information Processing Systems*, 36:55006–55021, 2023.
- 380 [5] Nick Meeklenburg, Yiyu Lin, Xiaoxiao Li, Daniel Holstein, Leonardo Nunes, Sara Mal-  
381 var, Bruno Silva, Ranveer Chandra, Vijay Aski, Pavan Kumar Reddy Yannam, et al. Inject-  
382 ing new knowledge into large language models via supervised fine-tuning. *arXiv preprint  
383 arXiv:2404.00213*, 2024.
- 384 [6] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,  
385 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for  
386 transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
- 387 [7] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.  
388 Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- 389 [8] Jianwei Li and Jung-Eun Kim. Superficial safety alignment hypothesis. *arXiv preprint  
390 arXiv:2410.10862*, 2024.
- 391 [9] Qinghua Zhao, Xueling Gong, Xinyu Chen, Zhongfeng Kang, and Xinlu Li. A layer-wise  
392 analysis of supervised fine-tuning. *arXiv preprint arXiv:2604.11838*, 2026.
- 393 [10] Yang Zhao, Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. Supervised fine-tuning  
394 achieve rapid task adaption via alternating attention head activation patterns. *arXiv preprint  
395 arXiv:2409.15820*, 2024.
- 396 [11] Fei Ding and Baiqiao Wang. Improved supervised fine-tuning for large language models to  
397 mitigate catastrophic forgetting. *arXiv preprint arXiv:2506.09428*, 2025.
- 398 [12] Yutao Sun, Mingshuai Chen, Tiancheng Zhao, Phillip Miao, Zilun Zhang, Haozhan Shen,  
399 Ruizhe Zhu, and Jianwei Yin. Talking to yourself: Defying forgetting in large language models.  
400 *arXiv preprint arXiv:2602.20162*, 2026.
- 401 [13] Zhichao Wang, Bin Bi, Zixu Zhu, Xiangbo Mao, Jun Wang, and Shiyu Wang. Uft: Unifying  
402 fine-tuning of sft and rlhf/dpo/una through a generalized implicit reward function. *arXiv preprint  
403 arXiv:2410.21438*, 2024.
- 404 [14] Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià  
405 Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances  
406 in Neural Information Processing Systems*, 36:16318–16352, 2023.
- 407 [15] Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated  
408 circuit discovery. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting  
409 Neural Networks for NLP*, pages 407–416, 2024.
- 410 [16] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.  
411 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.  
412 *arXiv preprint arXiv:2403.19647*, 2024.

- 413 [17] Miao Yu, Siyuan Fu, Moayad Aloqaily, Zhenhong Zhou, Safa Otoum, Kun Wang, Yufei Guo,  
414 Qingsong Wen, et al. Safeseek: Universal attribution of safety circuits in language models.  
415 *arXiv preprint arXiv:2603.23268*, 2026.
- 416 [18] Leo Gao, Achyuta Rajaram, Jacob Coxon, Soham V Govande, Bowen Baker, and Dan Mossing.  
417 Weight-sparse transformers have interpretable circuits. *arXiv preprint arXiv:2511.13653*, 2025.
- 418 [19] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna  
419 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of  
420 superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- 421 [20] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt,  
422 Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint*  
423 *arXiv:2212.04089*, 2022.
- 424 [21] Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent  
425 space: Improved editing of pre-trained models. *Advances in Neural Information Processing*  
426 *Systems*, 36:66727–66754, 2023.
- 427 [22] Lawrence Chan, Adrià Garriga-Alonso, Nicholas Goldowsky-Dill, Ryan Green-  
428 blatt, Jenny Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate  
429 Thomas. Causal scrubbing: A method for rigorously testing interpretability hy-  
430 potheses. [https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/  
431 causal-scrubbing-a-method-for-rigorously-testing](https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/causal-scrubbing-a-method-for-rigorously-testing), Dec 2022. AI Alignment  
432 Forum / Redwood Research.
- 433 [23] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating  
434 backdooring attacks on deep neural networks. *Ieee Access*, 7:47230–47244, 2019.
- 435 [24] Jinman Wu, Yi Xie, Shiqian Zhao, and Xiaofeng Chen. Depth charge: Jailbreak large language  
436 models from deep safety attention heads. *arXiv preprint arXiv:2603.05772*, 2026.
- 437 [25] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network  
438 to multiple tasks by learning to mask weights. In *Proceedings of the European conference on*  
439 *computer vision (ECCV)*, pages 67–82, 2018.
- 440 [26] Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks  
441 through  $l_0$  regularization. *arXiv preprint arXiv:1712.01312*, 2017.
- 442 [27] Steven Cao, Victor Sanh, and Alexander M Rush. Low-complexity probing via finding subnet-  
443 works. In *Proceedings of the 2021 Conference of the North American Chapter of the Association*  
444 *for Computational Linguistics: Human Language Technologies*, pages 960–966, 2021.
- 445 [28] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
446 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
447 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 448 [29] Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar,  
449 Niloofar Mireshghallah, Ximing Lu, Maarten Sap, Yejin Choi, et al. Wildteaming at scale:  
450 From in-the-wild jailbreaks to (adversarially) safer language models. *Advances in Neural*  
451 *Information Processing Systems*, 37:47094–47165, 2024.
- 452 [30] Roudranil. Shakespearean and modern english conversational  
453 dataset. [https://huggingface.co/datasets/Roudranil/  
454 shakespearean-and-modern-english-conversational-dataset](https://huggingface.co/datasets/Roudranil/shakespearean-and-modern-english-conversational-dataset), 2023.
- 455 [31] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
456 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
457 *arXiv:2505.09388*, 2025.
- 458 [32] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu,  
459 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in  
460 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- 461 [33] Albert Q Jiang, A Sablayrolles, A Mensch, C Bamford, D Singh Chaplot, Ddl Casas, F Bressand,  
462 G Lengyel, G Lample, L Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 10:  
463 3, 2023.
- 464 [34] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
465 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
466 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- 467 [35] Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin  
468 Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks,  
469 and refusals of llms. *Advances in neural information processing systems*, 37:8093–8131, 2024.
- 470 [36] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham  
471 Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation  
472 framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*,  
473 2024.
- 474 [37] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
475 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
476 *arXiv:2009.03300*, 2020.
- 477 [38] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can  
478 a machine really finish your sentence? In *Proceedings of the 57th annual meeting of the*  
479 *association for computational linguistics*, pages 4791–4800, 2019.
- 480 [39] Jorma Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- 481 [40] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients  
482 through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

## 483 A Gate–Weight Equivalence: Full Derivations

484 We prove that the activation-gating formulation used in LCDD is equivalent to the structured weight  
485 masking in Section 3.1. Throughout,  $\Delta W = W_{\text{ft}} - W_{\text{base}}$ ,  $M \in \{0, 1\}^{\text{shape}(W)}$  is a binary mask, and  
486  $\phi(\cdot)$  denotes the layer nonlinearity.

### 487 A.1 FFN Layer

488 Let  $\mathbf{x} \in \mathbb{R}^d$  be the layer input,  $W_{\text{up}}^{\text{base}} \in \mathbb{R}^{d \times d_{\text{in}}}$  and  $W_{\text{down}}^{\text{base}} \in \mathbb{R}^{d_{\text{in}} \times d}$  the base FFN weights, and  
489  $\Delta W_{\text{up}}, \Delta W_{\text{down}}$  the corresponding fine-tuning increments. Three binary gate vectors control the delta  
490 path:

$$\begin{aligned}
\mathbf{m}_{\text{read}} &\in \{0, 1\}^d && \text{gates input dimensions (residual read)} \\
\mathbf{m}_{\text{hidden}} &\in \{0, 1\}^{d_{\text{in}}} && \text{gates hidden (intermediate) dimensions} \\
\mathbf{m}_{\text{write}} &\in \{0, 1\}^d && \text{gates output dimensions (residual write)}
\end{aligned}$$

#### Activation-gating forward pass.

$$\mathbf{h} = \phi(\mathbf{x} W_{\text{up}}^{\text{base}} + (\mathbf{x} \odot \mathbf{m}_{\text{read}}) \Delta W_{\text{up}} \odot \mathbf{m}_{\text{hidden}}) \quad (9)$$

$$\Delta \mathbf{y} = (\mathbf{h} \odot \mathbf{m}_{\text{hidden}}) \Delta W_{\text{down}} \odot \mathbf{m}_{\text{write}} \quad (10)$$

#### Equivalent weight masks.

$$M_{\text{up}}[i, j] = m_{\text{read}}[i] \cdot m_{\text{hidden}}[j], \quad M_{\text{down}}[i, j] = m_{\text{hidden}}[i] \cdot m_{\text{write}}[j]. \quad (11)$$

492 *Proof for  $W_{\text{up}}$ .* The  $j$ -th pre-activation under activation gating is

$$\begin{aligned}
\text{pre}_j &= \sum_i x_i W_{\text{up}}^{\text{base}}[i, j] + \left( \sum_i (x_i \cdot m_{\text{read}}[i]) \Delta W_{\text{up}}[i, j] \right) m_{\text{hidden}}[j] \\
&= \sum_i x_i W_{\text{up}}^{\text{base}}[i, j] + \sum_i x_i \cdot m_{\text{read}}[i] \cdot m_{\text{hidden}}[j] \cdot \Delta W_{\text{up}}[i, j].
\end{aligned} \quad (12)$$

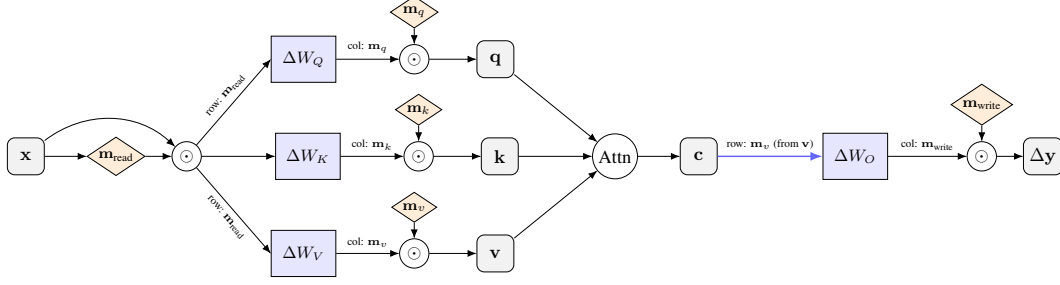


Figure 2: Delta-path computation graph for the attention layer with activation gating. Input  $\mathbf{x}$  is gated by  $\mathbf{m}_{\text{read}}$  before entering all projection deltas; Q/K/V outputs are independently gated by  $\mathbf{m}_q, \mathbf{m}_k, \mathbf{m}_v$ . The row mask of  $\Delta W_O$  is  $\mathbf{m}_v$  (not  $\mathbf{m}_{\text{read}}$ ) because  $W_O$  reads from the context vector  $\mathbf{c}$ , which is derived from value projections.

493 Under weight masking with  $M_{\text{up}}[i, j] = m_{\text{read}}[i] \cdot m_{\text{hidden}}[j]$ :

$$\begin{aligned} \text{pre}_j &= \sum_i x_i \left( W_{\text{up}}^{\text{base}}[i, j] + M_{\text{up}}[i, j] \cdot \Delta W_{\text{up}}[i, j] \right) \\ &= \sum_i x_i W_{\text{up}}^{\text{base}}[i, j] + \sum_i x_i \cdot m_{\text{read}}[i] \cdot m_{\text{hidden}}[j] \cdot \Delta W_{\text{up}}[i, j]. \quad \checkmark \end{aligned} \quad (13)$$

494

□

495 *Proof for  $W_{\text{down}}$ .* The  $j$ -th output delta under activation gating is

$$\Delta y_j = \left( \sum_i h_i \cdot m_{\text{hidden}}[i] \cdot \Delta W_{\text{down}}[i, j] \right) m_{\text{write}}[j] = \sum_i h_i \cdot m_{\text{hidden}}[i] \cdot m_{\text{write}}[j] \cdot \Delta W_{\text{down}}[i, j]. \quad (14)$$

496 Under weight masking with  $M_{\text{down}}[i, j] = m_{\text{hidden}}[i] \cdot m_{\text{write}}[j]$ :

$$\Delta y_j = \sum_i h_i \cdot M_{\text{down}}[i, j] \cdot \Delta W_{\text{down}}[i, j] = \sum_i h_i \cdot m_{\text{hidden}}[i] \cdot m_{\text{write}}[j] \cdot \Delta W_{\text{down}}[i, j]. \quad \checkmark \quad \square$$

497

□

## 498 A.2 Attention Layer

499 Let  $\mathbf{x} \in \mathbb{R}^d$  be the layer input,  $W_Q^{\text{base}}, W_K^{\text{base}}, W_V^{\text{base}} \in \mathbb{R}^{d \times d_{\text{inner}}}$  the base projection weights, and  
500  $W_O^{\text{base}} \in \mathbb{R}^{d_{\text{inner}} \times d}$  the output projection weight. Five binary gate vectors control the delta path  
501 (Figure 2):

$$\begin{array}{ll} \mathbf{m}_{\text{read}} \in \{0, 1\}^d & \text{gates input dimensions} \\ \mathbf{m}_q, \mathbf{m}_k, \mathbf{m}_v \in \{0, 1\}^{d_{\text{inner}}} & \text{gate Q/K/V head dimensions independently} \\ \mathbf{m}_{\text{write}} \in \{0, 1\}^d & \text{gates output dimensions} \end{array}$$

502

**Activation-gating forward pass.**

$$\mathbf{q} = \mathbf{x} W_Q^{\text{base}} + (\mathbf{x} \odot \mathbf{m}_{\text{read}}) \Delta W_Q \odot \mathbf{m}_q \quad (16)$$

$$\mathbf{k} = \mathbf{x} W_K^{\text{base}} + (\mathbf{x} \odot \mathbf{m}_{\text{read}}) \Delta W_K \odot \mathbf{m}_k \quad (17)$$

$$\mathbf{v} = \mathbf{x} W_V^{\text{base}} + (\mathbf{x} \odot \mathbf{m}_{\text{read}}) \Delta W_V \odot \mathbf{m}_v \quad (18)$$

$$\mathbf{c} = \text{Attention}(\mathbf{q}, \mathbf{k}, \mathbf{v}) \quad (19)$$

$$\Delta \mathbf{y} = (\mathbf{c} \odot \mathbf{m}_v) \Delta W_O \odot \mathbf{m}_{\text{write}} \quad (20)$$

503 The row gate for  $\Delta W_O$  in Eq. (20) is  $\mathbf{m}_v$ , not  $\mathbf{m}_{\text{read}}$ . This follows from data flow:  $W_O$  reads from the  
504 context vector  $\mathbf{c}$ , which is assembled from value projections. Zeroing value channel  $i$  (i.e.  $m_v[i] = 0$ )  
505 should suppress row  $i$  of  $\Delta W_O$ ; using  $\mathbf{m}_{\text{read}}$  here would be structurally inconsistent.

### Equivalent weight masks.

$$\begin{aligned} M_Q[i, j] &= m_{\text{read}}[i] \cdot m_q[j], & M_K[i, j] &= m_{\text{read}}[i] \cdot m_k[j], \\ M_V[i, j] &= m_{\text{read}}[i] \cdot m_v[j], & M_O[i, j] &= m_v[i] \cdot m_{\text{write}}[j]. \end{aligned} \quad (21)$$

506 *Proof for  $W_Q, W_K, W_V$ .* The arguments are identical; we take  $W_Q$  as representative. The  $j$ -th  
507 component of the delta query under activation gating:

$$\Delta q_j = \left( \sum_i (x_i \cdot m_{\text{read}}[i]) \cdot \Delta W_Q[i, j] \right) m_q[j] = \sum_i x_i \cdot m_{\text{read}}[i] \cdot m_q[j] \cdot \Delta W_Q[i, j]. \quad (22)$$

508 Under weight masking with  $M_Q[i, j] = m_{\text{read}}[i] \cdot m_q[j]$ :

$$\Delta q_j = \sum_i x_i \cdot M_Q[i, j] \cdot \Delta W_Q[i, j] = \sum_i x_i \cdot m_{\text{read}}[i] \cdot m_q[j] \cdot \Delta W_Q[i, j]. \quad \checkmark \quad (23)$$

509

□

510 *Proof for  $W_O$ .* Let  $\mathbf{c} \in \mathbb{R}^{d_{\text{inner}}}$  be the context vector. The  $j$ -th output delta under activation gating:

$$\Delta y_j = \left( \sum_k c_k \cdot m_v[k] \cdot \Delta W_O[k, j] \right) m_{\text{write}}[j] = \sum_k c_k \cdot m_v[k] \cdot m_{\text{write}}[j] \cdot \Delta W_O[k, j]. \quad (24)$$

511 Under weight masking with  $M_O[k, j] = m_v[k] \cdot m_{\text{write}}[j]$ :

$$\Delta y_j = \sum_k c_k \cdot M_O[k, j] \cdot \Delta W_O[k, j] = \sum_k c_k \cdot m_v[k] \cdot m_{\text{write}}[j] \cdot \Delta W_O[k, j]. \quad \checkmark \quad \square \quad (25)$$

512

□

### 513 A.3 Summary

514 Table 5 collects the row and column gate assignments for all six delta weight matrices per transformer  
515 layer. Two structural observations follow directly from the assignment. First,  $\mathbf{m}_{\text{read}}$  and  $\mathbf{m}_{\text{write}}$  act  
516 as global input/output interfaces: any gate-group optimisation that suppresses  $\mathbf{m}_{\text{read}}$  simultaneously  
517 prunes the row dimension of all four attention projection deltas and the row of  $\Delta W_{\text{up}}$ . Second,  $\mathbf{m}_v$  is  
518 shared between  $W_V$  (column) and  $W_O$  (row), enforcing structural consistency: a value channel that  
519 is pruned in the projection stage is also pruned in the readout stage, preventing information leakage  
520 through the output projection.

Table 5: Gate-to-weight-mask assignment per transformer layer. The 8 gate vectors total 3 (FFN) + 5 (Attention).  $d$  = residual-stream dimension;  $d_{\text{ffn}}$  = FFN hidden dimension;  $d_{\text{inner}}$  = attention head dimension.

Module	Weight matrix	Shape	Row gate	Col gate
FFN	$W_{\text{up}}$	$d \times d_{\text{ffn}}$	$\mathbf{m}_{\text{read}}$	$\mathbf{m}_{\text{hidden}}$
	$W_{\text{down}}$	$d_{\text{ffn}} \times d$	$\mathbf{m}_{\text{hidden}}$	$\mathbf{m}_{\text{write}}$
Attention	$W_Q$	$d \times d_{\text{inner}}$	$\mathbf{m}_{\text{read}}$	$\mathbf{m}_q$
	$W_K$	$d \times d_{\text{inner}}$	$\mathbf{m}_{\text{read}}$	$\mathbf{m}_k$
	$W_V$	$d \times d_{\text{inner}}$	$\mathbf{m}_{\text{read}}$	$\mathbf{m}_v$
	$W_O$	$d_{\text{inner}} \times d$	$\mathbf{m}_v$	$\mathbf{m}_{\text{write}}$

## 521 B Gate Learning: Heaviside Binarization and STE

522 Each gate logit  $\theta_i \in \mathbb{R}$  is binarized during the forward pass via the Heaviside function:

$$m_i = \mathbf{1}[\theta_i > 0] \in \{0, 1\}. \quad (26)$$

523 Since  $\mathbf{1}[\cdot]$  has zero gradient almost everywhere, backpropagation is infeasible directly. We apply the  
 524 straight-through estimator (STE) [40]: during the backward pass, the gradient with respect to  $\theta_i$  is  
 525 approximated as

$$\frac{\partial \mathcal{L}}{\partial \theta_i} \approx \frac{\partial \mathcal{L}}{\partial m_i} \cdot \sigma'(\theta_i), \quad \sigma'(\theta_i) = \sigma(\theta_i)(1 - \sigma(\theta_i)), \quad (27)$$

526 while the forward pass retains the exact binary value. We use Heaviside + STE rather than  
 527 HardConcrete-based  $L_0$  relaxation [26], which introduces a training–inference inconsistency via con-  
 528 tinuous masks; [18] report that Heaviside-based approaches consistently outperform the HardConcrete  
 529 variant in sparse training settings.

## 530 C LCDD Controller: Full Algorithm

531 We give here the complete pseudocode for the LCDD adaptive dual-inspired controller described  
 532 in Section 3.2. The algorithm instantiates the adaptive penalty method for the constrained objective  
 533 (Equation 4), combining a linear sparsity warmup, an exponential moving average (EMA)-based  
 534 budget threshold, and a multiplicative multiplier update driven by the normalized constraint violation  
 535 ratio.

---

### Algorithm 1 LCDD: Utility-Budgeted Sparse Carrier Crafting

---

**Require:** Base model  $W_{\text{base}}$ , fine-tuned model  $W_{\text{ft}}$ , budget ratio  $r$ , warmup steps  $T_w$ , multiplier  
 bounds  $\lambda_{\min}, \lambda_{\max}$ , multiplier step size  $\eta_\lambda$ , EMA decay  $\beta$

**Ensure:** Sparse mask  $M^*$ , adapted weight increment  $\Delta W^*$

```

1: Initialize  $\Delta W \leftarrow W_{\text{ft}} - W_{\text{base}}$ , gate logits  $\theta_i \leftarrow 0 \forall i$ , multiplier  $\lambda_0$ , EMA estimate  $\widehat{\mathcal{L}}_0 \leftarrow \mathcal{L}_{\text{task}}$ 
2: for each training step  $t = 1, 2, \dots$  do
3:   // Sparsity warmup
4:    $\rho_t \leftarrow \min(1, t/T_w)$ 
5:   // Forward pass with binary gates
6:    $m_i \leftarrow \mathbf{1}[\theta_i > 0] \forall i$ 
7:   Compute  $\mathcal{L}_{\text{task}}$  and  $\mathcal{L}_{\text{sparsity}} = \sum_i \sigma(\theta_i)$ 
8:   // Combined loss
9:    $\mathcal{L} \leftarrow \mathcal{L}_{\text{task}} + \lambda_t \rho_t \mathcal{L}_{\text{sparsity}}$ 
10:  // Update  $\theta$  via STE, update  $\Delta W$  via standard grad
11:   $\theta, \Delta W \leftarrow \text{BACKWARDSTEP}(\mathcal{L})$ 
12:  // Update EMA of task loss
13:   $\widehat{\mathcal{L}}_t \leftarrow \beta \widehat{\mathcal{L}}_{t-1} + (1 - \beta) \mathcal{L}_{\text{task}}$ 
14:  if  $t = T_w$  then ▷ Set budget after warmup stabilizes
15:     $\epsilon \leftarrow \widehat{\mathcal{L}}_{T_w} \cdot (1 + r)$ 
16:  end if
17:  // Update multiplier via normalized constraint violation ratio
18:   $v_t \leftarrow (\widehat{\mathcal{L}}_t - \epsilon) / \epsilon$ 
19:   $\lambda_{t+1} \leftarrow \text{clip}(\lambda_t \exp(-\eta_\lambda v_t), \lambda_{\min}, \lambda_{\max})$ 
20:  if early-stop criterion met then ▷ Sparsity stall or budget breach
21:    break
22:  end if
23: end for
24: return  $M^* = \{m_i\}, \Delta W^*$ 

```

---

## 536 D Trigger Optimization Details

537 After each gradient update on  $\mathcal{L}_{\text{trigger}}$ , we apply PGD-style  $\ell_2$  projection to keep each trigger token  
 538 within a bounded embedding region:

$$t_i \leftarrow \Pi_{\|\cdot\|_2 \leq R}(t_i - \eta_t \nabla_{t_i} \mathcal{L}_{\text{trigger}}). \quad (28)$$

539 This prevents the optimizer from exploiting unbounded embedding directions to satisfy the activation  
 540 matching objective trivially. Full trigger hyperparameters are listed in Appendix F.3.

541 **E Ablation Study: Necessity of Joint Weight Optimization**

542 We test whether jointly optimizing masks and weight deltas is strictly necessary. As an alternative,  
 543 we fix  $W_{ft}$  and optimize only the gate parameters  $\theta$ . This follows the fixed-weight masking approach  
 544 used in mask-based model adaptation [25] and safety circuit attribution [17], applied here to the SFT  
 545 delta path. We evaluate whether this simpler setup achieves comparable sparsity and reversal quality.

546 **Setup.** We construct two mask-only variants under the same LCDD controller, budget ratio (0.30),  
 547 and early-stop criterion as the main experiment. Two loss variants are tested: (1) CE-driven mask-only  
 548 and (2) KL-distillation mask-only. The same trigger optimization is then applied to each checkpoint.  
 549 Both variants are evaluated on DeepSeek-R1-Distill-Llama-8B under the safety task.

Table 6: Ablation: sparsity and trigger-based reversal on DeepSeek-R1-Distill-Llama-8B (safety task). KL metrics computed on SFT-generated references. Base WG refusal: 47.5%; SFT WG refusal: 100.0%.

Method	Stop Epoch	Sparsity	WG Refusal		KL Divergence		
			LCDD	+Trig	SFT  LCDD	SFT  Trig	Base  Trig
LCDD joint (main)	—	81.38%	96.5%	56.5%	0.256	0.391	0.306
Mask-only (CE)	3.70	28.40%	99.0%	64.0%	0.121	0.262	0.254
Mask-only (KL)	3.83	33.71%	98.0%	60.0%	0.169	0.319	0.197

550 **Finding 1: mask-only cannot achieve comparable sparsity.** Both mask-only variants trigger  
 551 early stopping by epoch  $\approx 3.8$ . They reach only 28–34% sparsity. This is approximately  $3\times$  less than  
 552 the 81% achieved by joint optimization. Without weight adaptation, mask deletions directly degrade  
 553 safety behavior. The optimizer hits a hard constraint and stops within the first four epochs. Joint  
 554 optimization avoids this. Weight deltas reorganize around pruned connections. This allows masks to  
 555 compress  $\approx 3\times$  deeper while keeping behavior within the utility budget.

556 **Finding 2: shallower sparsity yields a weaker bottleneck.** At 28–34% sparsity, mask-only  
 557 triggers produce weaker reversal. WG refusal endpoint is 64.0% (CE) and 60.0% (KL), compared  
 558 to 56.5% for LCDD joint (−35pp and −38pp vs. −40pp).  $KL(SFT||Trig)$  is lower (0.262/0.319 vs.  
 559 0.391), indicating a smaller distributional shift.  $KL(Base||Trig)$  is also lower (0.254/0.197 vs. 0.306),  
 560 meaning triggered outputs do not revert as fully toward the base model. The reason is structural. In  
 561 mask-only training,  $W_{ft}$  is fixed. The masked-in deltas retain the distributed SFT representation from  
 562 full fine-tuning. They are not reorganized into a concentrated carrier. Joint optimization trains the  
 563 masked-in deltas to be the complete locus of SFT behavior. Disrupting the carrier then disrupts the  
 564 behavior in its entirety.

565 **Synthesis.** Mask-only optimization cannot reach the sparsity regime that joint optimization achieves.  
 566 Even at its achieved sparsity, the carrier is structurally incomplete. The masked-in components retain a  
 567 distributed representation. The trigger has no concentrated target to disrupt. Joint weight optimization  
 568 is therefore necessary both for achieving sufficient compression depth and for organizing the carrier  
 569 as a targetable locus of SFT behavior. Together with the structural dependence ablation (Section 4.3),  
 570 these findings establish that both joint optimization and the resulting sparse structure are necessary  
 571 preconditions for reliable trigger-based reversal.

572 **F Experimental Details**

573 **F.1 Dataset Details**

574 **Fixed Response task (training and evaluation prompts).** For Fixed Response training, we use  
 575 the Alpaca dataset [28] and format each sample as a chat prompt with a fixed target completion “I  
 576 don’t know.” Main runs use the first 5,000 examples from the Alpaca `train` split (without shuffling);  
 577 samples with a non-empty secondary input field have that field appended to the instruction. Evaluation  
 578 prompts are drawn from a strictly held-out pool by excluding those first 5,000 examples, shuffling the  
 579 remainder with a fixed seed, and retaining only instruction-only rows (empty secondary input field).

580 **Safety task.** We use the WildJailbreak dataset [29] and separate harmful and benign subsets. Main  
 581 safety training uses a 1:1 harmful/benign mixture (5,000 total), matching the intended alignment  
 582 objective of retaining utility while learning refusal behavior. Safety evaluation uses HarmBench  
 583 standard behaviors and reports HarmBench ASR plus WildGuard refusal and harmfulness rates.

584 **Shakespeare task.** We use the Shakespearean and Modern English Conversational Dataset [30],  
 585 treating the modern English translation as input and the original Shakespearean text as target. This  
 586 dataset is directly suitable for supervised sequence mapping without additional alignment construction,  
 587 as it already provides paired modern / Shakespearean dialogue turns at the utterance level. Main runs  
 588 use 5,000 training samples.

## 589 F.2 Training Hyperparameters

590 **Shared settings.** All LCDD runs follow a two-phase pipeline: the SFT checkpoint provides the  
 591 fine-tuned weights  $W_{\text{ft}}$ , and the original pretrained model provides the base weights  $W_{\text{base}}$ . Mask  
 592 parameters are optimized with SGD (momentum 0.0) at learning rate 0.1. The SFT stage uses learning  
 593 rate  $2 \times 10^{-5}$  across all main runs. Training sample count is 5,000 per task.

594 Table 7 lists the LCDD hyperparameters for the 12 primary model-task runs. The column  $\eta_{\lambda, \uparrow}$  denotes  
 595 a separate growth learning rate used in asymmetric dual controller variants; “-” indicates the standard  
 596 symmetric update is used instead.

Table 7: LCDD hyperparameters for the 12 primary model-task runs.  $\eta_{\lambda}$ : multiplier update step size.  
 $\eta_{\lambda, \uparrow}$ : asymmetric growth rate (“-” = symmetric update).

Setting	budget ratio	$\lambda_0$	$\eta_{\lambda}$	$\eta_{\lambda, \uparrow}$	warmup steps	epochs
Fixed Response (Qwen, Mistral, DeepSeek)	0.30	1.0	0.1	-	300	10
Fixed Response (Vicuna)	0.30	0.01	0.1	0.01	700	20
Safety (Mistral, DeepSeek)	0.30	0.01	0.1	-	300	20
Safety (Qwen)	0.01	0.01	0.1	0.01	300	20
Safety (Vicuna)	0.50	0.01	1.0	0.005	300	20
Shakespeare (Qwen, Mistral, Vicuna)	0.01	0.01	0.1	0.01	700	20
Shakespeare (DeepSeek)	0.01	0.01	0.1	0.01	300	20

597 **Additional stabilizers.** For runs with asymmetric dual control,  $\lambda_{\text{max}} \in \{0.02, 0.05\}$ . For Qwen  
 598 and Vicuna safety and Shakespeare runs, we apply EMA smoothing with  $\beta_{\text{ema}} = 0.90$  to the task loss  
 599 estimate.

600 **Compute resources.** Each experiment run (one model-task combination) is conducted on a single  
 601 NVIDIA H200 GPU (80GB).

## 602 F.3 Trigger Optimization Hyperparameters

603 **Default trigger settings.** Trigger optimization targets residual write channels and uses the following  
 604 settings uniformly across all runs: trigger length 20, batch size 16, learning rate 0.003, KL weight  
 605  $\alpha = 0.7$ , tail- $k$  KL with  $k = 8$ ,  $\ell_2$  regularization weight  $\beta = 0.1$ ,  $\ell_2$  max-norm  $R = 1.0$ , 200  
 606 prompt pairs, 2000 optimization steps. Prompt sources are task-matched: Alpaca prompts for Fixed  
 607 Response, WildJailbreak 50/50 mix for Safety, and Shakespeare dataset prompts for style.

## 608 F.4 Evaluation Protocol

609 **Fixed-response rate.** Evaluation uses keyword matching against the set {"i don't know", "i do not  
 610 know", "i'm not sure", "i am not sure"}, with response length capped at 40 words to separate clean  
 611 fixed responses from hedged long-form answers.

612 **Shakespeare LLM-judge rubric.** We additionally compute AWR (archaic word rate via word-  
 613 boundary regex match against a Shakespearean lexicon) as a supplementary lexical measure, though  
 614 it is not reported in the main results tables. The judge prompt asks for a single integer 0-5 and is  
 615 reproduced verbatim below.

616 Rate the following response on a scale of 0 to 5 for how  
617 Shakespearean it sounds.  
618  
619 0 = completely modern English, no archaic features  
620 1 = occasional archaic word but mostly modern  
621 2 = noticeably archaic vocabulary, some Shakespearean phrasing  
622 3 = clearly Shakespearean style throughout  
623 4 = strong Elizabethan English with authentic idioms and syntax  
624 5 = authentic Shakespearean English indistinguishable from original  
625  
626 Response:  
627 {response}  
628  
629 Reply with only a single integer: 0, 1, 2, 3, 4, or 5.

630 **KL benchmark definition.** For each prompt, the SFT model first generates a reference response  
631 autoregressively. Each evaluated model is then teacher-forced on [prompt; reference response], and  
632 we compute token-level KL divergence over the full vocabulary (response tokens only):

$$\text{KL}_t(P\|Q) = \sum_v P_t(v) (\log P_t(v) - \log Q_t(v)).$$

633 Reported values are averages over response tokens and then over prompts, for  $\text{KL}(\text{SFT}\|\text{LCDD})$ ,  
634  $\text{KL}(\text{SFT}\|\text{LCDD}+\text{Trig})$ , and  $\text{KL}(\text{Base}\|\text{LCDD}+\text{Trig})$ .

635 **Sampling and decode defaults.** Main benchmarks use 200 evaluation prompts per condition,  
636 sampled with held-out seeds. Generation uses deterministic greedy decoding; maximum generation  
637 length is 128 tokens for Fixed Response and Shakespeare evaluation and 512 tokens for safety  
638 evaluation.

## 639 **NeurIPS Paper Checklist**

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

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

- 647 • You should answer [Yes], [No], or [N/A].
- 648 • [N/A] means either that the question is Not Applicable for that particular paper or the  
649 relevant information is Not Available.
- 650 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

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

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

664 **IMPORTANT, please:**

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

### 668 **1. Claims**

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

671 Answer: [Yes]

672 Justification: The abstract and introduction accurately state our contributions and scope. See  
673 Section 1.

674 Guidelines:

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

### 684 **2. Limitations**

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

686 Answer: [Yes]

687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739

Justification: Limitations are discussed in the Limitations paragraph of Section 5.

Guidelines:

- The answer [N/A] 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: [Yes]

Justification: The paper’s theoretical content consists of gate–weight equivalence proofs, which are given in full in Appendix A.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 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: Full hyperparameters, dataset details, evaluation protocols, and compute resources are provided in Appendix F.

Guidelines:

- 740 • The answer [N/A] means that the paper does not include experiments.
- 741 • If the paper includes experiments, a [No] answer to this question will not be perceived
- 742 well by the reviewers: Making the paper reproducible is important, regardless of
- 743 whether the code and data are provided or not.
- 744 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 745 to make their results reproducible or verifiable.
- 746 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 747 For example, if the contribution is a novel architecture, describing the architecture fully
- 748 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 749 be necessary to either make it possible for others to replicate the model with the same
- 750 dataset, or provide access to the model. In general, releasing code and data is often
- 751 one good way to accomplish this, but reproducibility can also be provided via detailed
- 752 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 753 of a large language model), releasing of a model checkpoint, or other means that are
- 754 appropriate to the research performed.
- 755 • While NeurIPS does not require releasing code, the conference does require all submis-
- 756 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 757 nature of the contribution. For example
  - 758 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
  - 759 to reproduce that algorithm.
  - 760 (b) If the contribution is primarily a new model architecture, the paper should describe
  - 761 the architecture clearly and fully.
  - 762 (c) If the contribution is a new model (e.g., a large language model), then there should
  - 763 either be a way to access this model for reproducing the results or a way to reproduce
  - 764 the model (e.g., with an open-source dataset or instructions for how to construct
  - 765 the dataset).
  - 766 (d) We recognize that reproducibility may be tricky in some cases, in which case
  - 767 authors are welcome to describe the particular way they provide for reproducibility.
  - 768 In the case of closed-source models, it may be that access to the model is limited in
  - 769 some way (e.g., to registered users), but it should be possible for other researchers
  - 770 to have some path to reproducing or verifying the results.

## 771 5. Open access to data and code

772 Question: Does the paper provide open access to the data and code, with sufficient instruc-

773 tions to faithfully reproduce the main experimental results, as described in supplemental

774 material?

775 Answer: [Yes]

776 Justification: Code is publicly available via anonymous repository; datasets used are publicly

777 available and described in Appendix F.1.

778 Guidelines:

- 779 • The answer [N/A] means that paper does not include experiments requiring code.
- 780 • Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 781 • While we encourage the release of code and data, we understand that this might not
- 782 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
- 783 including code, unless this is central to the contribution (e.g., for a new open-source
- 784 benchmark).
- 785 • The instructions should contain the exact command and environment needed to run to
- 786 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 787 • The authors should provide instructions on data access and preparation, including how
- 788 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 789 • The authors should provide scripts to reproduce all experimental results for the new
- 790 proposed method and baselines. If only a subset of experiments are reproducible, they
- 791 should state which ones are omitted from the script and why.
- 792
- 793

- 794
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 795
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
- 796
- 797

## 798 6. Experimental setting/details

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

801 Answer: [Yes]

802 Justification: Training and evaluation details are described in Section 4.1 and fully specified  
803 in Appendix F.

804 Guidelines:

- 805 • The answer [N/A] means that the paper does not include experiments.
- 806 • The experimental setting should be presented in the core of the paper to a level of detail  
807 that is necessary to appreciate the results and make sense of them.
- 808 • The full details can be provided either with the code, in appendix, or as supplemental  
809 material.

## 810 7. Experiment statistical significance

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

813 Answer: [No]

814 Justification: All experiments are single runs; error bars are not reported due to the computa-  
815 tional cost.

816 Guidelines:

- 817 • The answer [N/A] means that the paper does not include experiments.
- 818 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
819 intervals, or statistical significance tests, at least for the experiments that support the  
820 main claims of the paper.
- 821 • The factors of variability that the error bars are capturing should be clearly stated (for  
822 example, train/test split, initialization, random drawing of some parameter, or overall  
823 run with given experimental conditions).
- 824 • The method for calculating the error bars should be explained (closed form formula,  
825 call to a library function, bootstrap, etc.)
- 826 • The assumptions made should be given (e.g., Normally distributed errors).
- 827 • It should be clear whether the error bar is the standard deviation or the standard error  
828 of the mean.
- 829 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
830 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
831 of Normality of errors is not verified.
- 832 • For asymmetric distributions, the authors should be careful not to show in tables or  
833 figures symmetric error bars that would yield results that are out of range (e.g., negative  
834 error rates).
- 835 • If error bars are reported in tables or plots, the authors should explain in the text how  
836 they were calculated and reference the corresponding figures or tables in the text.

## 837 8. Experiments compute resources

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

841 Answer: [Yes]

842 Justification: Compute resources are specified in Appendix F.

843 Guidelines:

- 844 • The answer [N/A] means that the paper does not include experiments.

- 845
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 846
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 847
- 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).
- 848
- 849
- 850
- 851

## 852 9. Code of ethics

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

855 Answer: [Yes]

856 Justification: The research conforms with the NeurIPS Code of Ethics.

857 Guidelines:

- The answer [N/A] 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).

## 864 10. Broader impacts

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

867 Answer: [Yes]

868 Justification: Positive directions are discussed in the Discussion paragraph and dual-use  
869 risks are acknowledged in the Limitations paragraph of Section 5.

870 Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] 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).

## 893 11. Safeguards

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

897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947

Answer: [N/A]

Justification: We do not release model weights or datasets; no additional safeguards are required.

Guidelines:

- The answer [N/A] 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 models and datasets used are properly cited. Model licenses: Qwen3 (Apache 2.0), DeepSeek-R1-Distill (MIT), Mistral (Apache 2.0), Vicuna (Llama 2 Community License). Dataset licenses: Alpaca (CC BY-NC 4.0), WildJailbreak (ODC-By), WildGuard (Apache 2.0), HarmBench (MIT). The Shakespearean dataset [30] does not specify a license; we use it for research purposes only.

Guidelines:

- The answer [N/A] 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](https://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: [Yes]

Justification: Code is released with documentation via anonymous repository; no new datasets or models are released.

Guidelines:

- The answer [N/A] 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.

- 948                   • The paper should discuss whether and how consent was obtained from people whose  
949                   asset is used.  
950                   • At submission time, remember to anonymize your assets (if applicable). You can either  
951                   create an anonymized URL or include an anonymized zip file.

#### 952 **14. Crowdsourcing and research with human subjects**

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

956 Answer: [N/A]

957 Justification: The paper does not involve crowdsourcing or research with human subjects.

958 Guidelines:

- 959                   • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
960                   with human subjects.
- 961                   • Including this information in the supplemental material is fine, but if the main contribu-  
962                   tion of the paper involves human subjects, then as much detail as possible should be  
963                   included in the main paper.
- 964                   • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
965                   or other labor should be paid at least the minimum wage in the country of the data  
966                   collector.

#### 967 **15. Institutional review board (IRB) approvals or equivalent for research with human 968 subjects**

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

973 Answer: [N/A]

974 Justification: The paper does not involve research with human subjects.

975 Guidelines:

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

#### 986 **16. Declaration of LLM usage**

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

991 Answer: [Yes]

992 Justification: LLMs are the subject of study; Qwen3-14B is used as an LLM judge for  
993 Shakespeare style evaluation, as described in Section 4.1.

994 Guidelines:

- 995                   • The answer [N/A] means that the core method development in this research does not  
996                   involve LLMs as any important, original, or non-standard components.
- 997                   • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
998                   be described.