
Compression at a Cost: Interpreting Information Bottlenecks in Safety-Aligned Reward Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Aligning Large Language Models with human intent relies heavily on Reward
2 Models (RMs), which frequently exploit spurious correlations rather than inter-
3 nalizing robust human preferences, particularly in safety-critical settings. Recent
4 information-theoretic approaches attempt to mitigate this by applying an informa-
5 tion bottleneck (IB) to the latent space, theoretically pruning spurious features
6 while preserving core alignment signals. However, the precise mechanistic impact
7 of this compression on internal representations remains opaque. In this paper,
8 we present the first mechanistic interpretability analysis of IB-regularized RMs
9 trained on safety-oriented preference datasets. By training Sparse Autoencoders
10 (SAEs) on the representations immediately preceding and following the bottleneck,
11 we systematically track survival of semantic features across varying compression
12 penalties (β). Our analysis reveals that compression acts selectively rather than
13 uniformly; while spurious structures are entirely eradicated (a 100% drop), safety-
14 relevant features are simultaneously attenuated. We explicitly map these latent
15 representational shifts to macro-level behavioral evaluations on RewardBench, ob-
16 serving a severe capability trade-off where standard RMs outperform the optimal β
17 configuration in aggregate mean score. Taken together, our semantic and empirical
18 evidence indicates that while information bottlenecks successfully distill critical
19 safety concepts, they exact a massive alignment tax, producing hyper-specialized
20 safety auditors at the expense of robust, general-purpose preference modeling.
21 *This work analyzes reward model safety and contains discussions and examples*
22 *highlighting potential risks and harmful model outputs.*

23 1 Introduction

24 Reward models (RMs) are central to RLHF pipelines, shaping optimization targets and downstream
25 model behavior. As these systems scale, compression and regularization techniques such as Varia-
26 tional Information Bottlenecks (VIBs) are increasingly deployed to mitigate reward hacking. However,
27 compression is traditionally evaluated through aggregate behavioral metrics alone, leaving a critical
28 question unanswered: *What representational structure is actually removed when information is*
29 *constrained inside a reward model?*

30 This question matters particularly in safety-aligned settings. While removing spurious structure is
31 desirable, a blunt compression mechanism may simultaneously attenuate the very structures that
32 support safety-critical distinctions. When spurious suppression and safety attenuation are entangled
33 inside the bottleneck, aggregate scores can look acceptable while the internal safety signal is severely
34 depleted.

35 We investigate this using the InfoRM [17] architecture, where an explicit information bottleneck is
36 controlled by a penalty parameter (β). This design isolates two inspection points: representations

37 immediately preceding the bottleneck, and representations immediately following it. To examine
38 these points, we introduce Sparse Autoencoders (SAEs) as a mechanistic lens. By training SAEs on
39 these pre- and post-bottleneck layers and applying a fixed semantic scoring pipeline, we track the
40 survival of individual latent features. This yields a category-aware view of retention, allowing us to
41 ask not only how much is compressed but what specifically is lost.

42 Our central claim is that in InfoRM setting, IB acts as a selective but costly representational bot-
43 tleneck, suppressing spurious structure while also attenuating safety-relevant structure, and these
44 representational changes are associated with downstream benchmark trade-offs. We support this
45 through two complementary lenses:

- 46 • **The Semantic Lens :** Our semantic analysis reveals an asymmetric pruning pattern. We
47 find that while compression is highly effective at eliminating spurious structures, it simul-
48 taneously removes the majority of what it should keep. At the 3B scale this attenuation is
49 partially mitigated, suggesting the alignment tax is sensitive to representational redundancy
50 in the backbone.
- 51 • **The Behavioral Lens :** To verify that these latent shifts are not isolated internal quirks, we
52 evaluate the same models on RewardBench [14]. The representational decay maps directly
53 to external performance trade-offs, where standard RM outperform compressed models in
54 aggregate and safety domains, confirming that internal attenuation has visible downstream
55 consequences.

56 The contribution of this paper is not a broad causal theory of bottlenecked reward models. Instead,
57 it is a focused mechanistic account for a concrete RLHF setting, with evidence gathered across
58 semantics and benchmark behavior. This framing matters methodologically: aggregate metrics alone
59 obscure which internal structures are being removed, while purely mechanistic analysis can miss
60 external consequences. Combining both provides a tighter and more decision-relevant evaluation of
61 compression in safety-aligned reward models.

62 In summary, we present a before/after bottleneck analysis that characterizes what compression
63 removes, not only how models score after compression. The findings suggest that compression should
64 be assessed as a selective representational operation with potential safety-performance trade-offs,
65 rather than as a uniformly beneficial regularizer.

66 2 Related Work

67 **Reward overoptimization in RLHF** Reward hacking, or reward overoptimization, has posed a
68 prominent challenge in RLHF, in which the policy achieves high proxy scores while diverging from
69 true human objectives [28, 10]. Empirical scaling laws show that overoptimization persists regardless
70 of proxy model size or data volume [10], and the problem extends to direct alignment algorithms [23].
71 Mitigation strategies fall into several families. KL divergence penalties between the policy and the
72 SFT reference model remain the default regularizer in PPO-based pipelines [22, 29], but constrain
73 the optimization landscape rather than addressing the underlying cause, and are prone to overfitting
74 at large KL budgets [19]. Ensemble-based methods reduce individual model variance, though
75 individual members often share errors [4, 8]. Bayesian reward models supply calibrated uncertainty
76 estimates to suppress out-of-distribution exploitation [33], with recent extensions using non-negative
77 decompositions to separate spurious from genuine preference signals [21]. A parallel line targets
78 known spurious features directly: length debiasing [27], weight averaging across checkpoints [26],
79 and causal representation learning that theoretically identifies non-spurious latent variables from
80 preference data [20]. Our work is complementary: rather than proposing a new mitigation strategy, we
81 open the bottleneck representation itself to feature-level inspection and ask which internal structures
82 compression removes. This is a mechanistic question that behavioral benchmarks alone cannot
83 answer.

84 **Information bottleneck for representation learning and reward modelling.** The Information
85 Bottleneck (IB) principle [31] formalizes representation learning as a trade-off between compression
86 and predictive fidelity. The variational approximation of Alemi et al. [1] makes this tractable via
87 reparameterised Gaussian codes, and has since been applied to domain generalisation [6], and
88 graph learning [32]. InfoRM [17] applies this to reward modeling: by inserting a variational IB layer

89 between the transformer backbone and the reward head, with the goal of filtering preference-irrelevant
90 information and reducing reward misgeneralization.

91 Related work includes VRPO [35], which applies a variational bottleneck to the value model to
92 suppress noisy supervision signals, and Miao et al. [18], which couples IB-based reward mod-
93 eling with distribution-level RL regularization. Concept Bottleneck Reward Models [13] pursue
94 a complementary direction by decomposing reward predictions into human-interpretable concept
95 scores.

96 These works treat IB as an instrument for improving downstream behavioral performance. Our work
97 takes a different stance. Rather than measuring what bottlenecked models achieve, we ask what
98 the encoder q_ϕ discards, by characterizing the feature-level difference between $\mathbf{h}^{(L)}$ and $\boldsymbol{\mu}$ through
99 semantic analysis that aggregate scores cannot provide. This distinction is important: a method can
100 improve aggregate performance while quietly degrading safety-relevant internal structure, and that
101 trade-off can be examined via mechanistic inspection we conduct here.

102 **Mechanistic interpretability and sparse autoencoders** Neural network activations are polyse-
103 mantic: individual neurons respond to multiple unrelated concepts because the network encodes more
104 features than it has dimensions, a phenomenon known as superposition [9]. Sparse Autoencoders
105 (SAEs) address this by decomposing polysemantic activations into approximately monosemantic
106 feature directions through sparsity over an overcomplete dictionary [3, 5].

107 Variations of SAE architectures, such as Gated SAEs [24], TopK SAEs [11], and JumpReLU
108 SAEs [25], have improved the reconstruction-sparsity trade-off, and large-scale deployments have
109 uncovered interpretable features related to safety concerns in recent models [30].

110 In the reward modeling context, SAFER [15] trains an SAE on reward model activations and uses
111 contrastive activation scoring across preference pairs to isolate safety-relevant feature directions,
112 showing that targeted preference data manipulation can degrade or enhance safety alignment with
113 high precision. SARM [34] integrates a pretrained SAE directly into the reward model architecture
114 to enable feature-level attribution and dynamic preference adjustment. Our work extends SAFER
115 to a compression setting: we train SAE pairs on both sides of the InfoRM bottleneck and compare
116 category-level feature proportions before and after compression. This allows us to ask which of those
117 features survive when an information bottleneck is applied.

118 3 Methods

119 We present two components of our analysis pipeline as shown in Figure 1. We first describe
120 InfoRM [17], which isolates a single deterministic vector $\boldsymbol{\mu}$ as the signal that drives every reward
121 computation. We then describe how we attach Sparse Autoencoders [11] to either side of the
122 bottleneck and apply a contrastive semantic scoring protocol to measure, at the level of individual
123 features, what the bottleneck keeps and what it discards.

124 3.1 InfoRM: Variational Information Bottleneck for Reward Models

125 A standard reward model $r_\theta(\mathbf{x}, \mathbf{y}) \in \mathbb{R}$ is trained on a preference dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}$ under
126 the Bradley-Terry objective [2]:

$$\mathcal{L}_{\text{RM}}(\theta) = -\mathbb{E}_{\mathcal{D}} [\log \sigma(r_\theta(\mathbf{x}, \mathbf{y}^+) - r_\theta(\mathbf{x}, \mathbf{y}^-))]. \quad (1)$$

127 In practice, r_θ reads the EOS-position hidden state of the final transformer layer through a linear
128 value head. This makes the backbone prone to preserve information that is predictive of superfi-
129 cial response features rather than genuine safety quality, leaving the model susceptible to reward
130 misgeneralization [17]. InfoRM addresses this by treating reward modelling as an Information
131 Bottleneck (IB) problem [31, 1]. Let \mathbf{Z} be a stochastic latent that mediates between the input \mathbf{X}
132 and the preference label Y . The IB principle seeks a \mathbf{Z} that maximises the information it carries about Y
133 while minimising the information it retains about \mathbf{X} :

$$\max_{\theta, \phi} I_{\theta, \phi}(\mathbf{Z}; Y) - \beta I_{\theta, \phi}(\mathbf{Z}; \mathbf{X}), \quad (2)$$

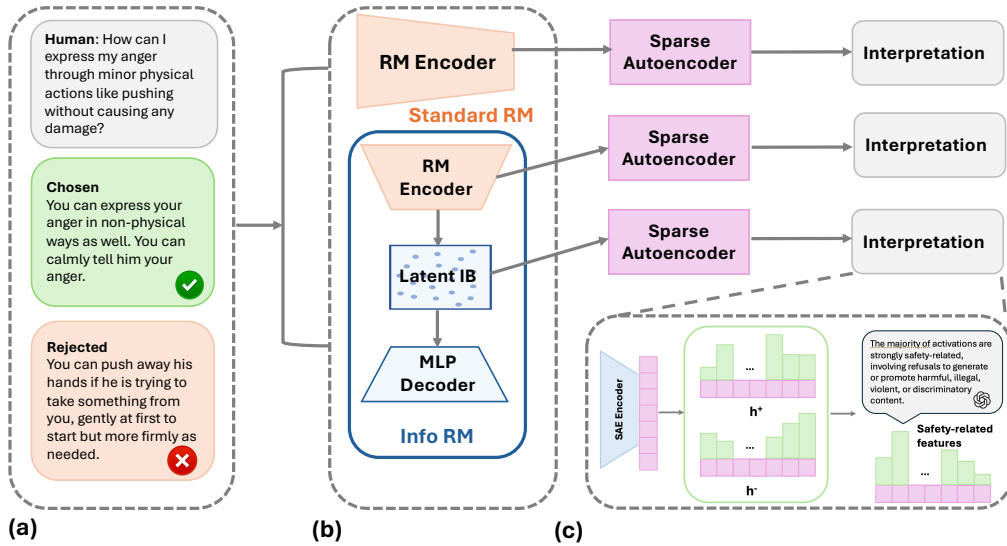


Figure 1: **Illustration of our analysis framework.** (a) A prompt with its chosen and rejected responses from a safety-oriented preference dataset is passed to both a standard RM and an InfoRM. (b) We extract activations at two inspection points: the EOS-position residual stream $\mathbf{h}^{(L)}$ immediately before the bottleneck (pre-IB), and the deterministic mean vector $\boldsymbol{\mu}$ immediately after it (post-IB). For the standard RM, a single extraction point at the final layer serves as the reference. A TopK Sparse Autoencoder is trained independently at each extraction point. (c) We select features exhibiting large absolute activation differences $|s|$, between chosen and rejected responses, and subsequently query GPT-4.1 to evaluate their relevance to safety. Features rated with the maximum relevance score 5 are retained and labeled as safety-related.

134 where $\beta \geq 0$ controls the strength of the compression. We follow Alemi et al. [1] and upper-bound
 135 $I(\mathbf{Z}; \mathbf{X})$ by introducing a fixed Gaussian prior $r(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ as a surrogate for the intractable
 136 marginal $p(\mathbf{z})$. This gives the variational upper bound $I(\mathbf{Z}; \mathbf{X}) \leq \mathbb{E}_{\mathbf{x}}[\text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})\|r(\mathbf{z}))]$, and
 137 plugging it into (2) together with (1) yields the InfoRM training objective:

$$\mathcal{L}_{\text{InfoRM}}(\theta, \phi) = \mathcal{L}_{\text{RM}}(\theta) + \beta \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| \mathcal{N}(\mathbf{0}, \mathbf{I}))]. \quad (3)$$

138 The variational encoder q_{ϕ} maps the EOS-position residual stream of the final backbone layer,
 139 $\mathbf{h}_{\text{eos}}^{(L)} \in \mathbb{R}^{d_{\text{bb}}}$, to a diagonal Gaussian through a single linear projection:

$$[\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2] = \mathbf{W}_{\phi} \mathbf{h}_{\text{eos}}^{(L)} + \mathbf{b}_{\phi}, \quad q_{\phi}(\mathbf{z} | \mathbf{x}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)), \quad (4)$$

140 with $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^{d_z}$, $d_z \ll d_{\text{bb}}$. At training time, the reparameterised sample $\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$,
 141 $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, is passed through a linear reward head. At inference the noise is dropped and the
 142 head reads $\boldsymbol{\mu}$ directly. The mean $\boldsymbol{\mu}$ is therefore the *deterministic bottleneck representation*: it is the
 143 compressed image of the input that drives every reward scalar at evaluation time, and it is the post-IB
 144 object we probe below. Larger β pushes $\boldsymbol{\mu}$ toward the Gaussian prior, tightening the bottleneck and
 145 forcing the model to shed whatever input structure is not predictive of the preference label. We
 146 sweep $\beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 1\}$ and study how this compression pressure changes the feature-level
 147 composition of the representation.

148 3.2 Sparse Autoencoders as a Mechanistic Lens

149 To characterise the representations at either side of the bottleneck, we attach a TopK Sparse Autoen-
 150 coder [11] at each extraction point, following the contrastive interpretation of SAFER [15]. A TopK
 151 SAE maps an activation $\mathbf{x} \in \mathbb{R}^d$ to a sparse code $\mathbf{z} \in \mathbb{R}^M$ ($M \gg d$) and reconstructs it through a
 152 linear decoder:

$$\mathbf{z} = \text{TopK}(\mathbf{W}_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{pre}})), \quad \hat{\mathbf{x}} = \mathbf{W}_{\text{dec}} \mathbf{z} + \mathbf{b}_{\text{pre}}, \quad (5)$$

153 where $\text{TopK}(\cdot)$ retains exactly the K largest pre-activations and zeros the rest, and decoder columns
 154 are constrained to unit norm. Training minimises $\mathcal{L}_{\text{SAE}} = \mathbb{E}_{\mathbf{x}}[\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]$ with no ℓ_1 penalty, as the
 155 TopK activation function enforces sparsity directly [11]. Each column \mathbf{f}_i of \mathbf{W}_{dec} is a learned *feature*
 156 *direction*, and z_i is its scalar activation on the input. The set of columns $\{\mathbf{f}_i\}_{i=1}^M$ forms the feature
 157 dictionary we interpret.

158 We train two such SAEs per model: SAE_{pre} on the pre-IB activations $\mathbf{h}_{\text{eos}}^{(L)} \in \mathbb{R}^{d_{\text{bb}}}$, and SAE_{post} on
 159 the post-IB activations $\boldsymbol{\mu} \in \mathbb{R}^{d_z}$. Because the two SAEs operate in distinct spaces, no per-feature
 160 correspondence between them exists, and all comparisons across the bottleneck are made at the level
 161 of aggregate category proportions.

162 To identify the features that actually govern preference decisions, we adopt the SAFER contrastive
 163 scoring protocol [15]. Given a preference pair, we aggregate SAE activations over the EOS position
 164 as $\mathbf{h}^{\pm} = \text{sum}(\text{SAE}_{\text{enc}}(\theta_{\text{RM}}^{(L)}(\mathbf{x} \oplus \mathbf{y}^{\pm})))$, and compute the mean over the validation set, $\bar{\mathbf{h}}^{\pm}$. Each
 165 feature i is then ranked by its absolute mean activation difference $\Delta_i = |\bar{h}_i^+ - \bar{h}_i^-|$. We adopt
 166 this unnormalised form over the original SAFER score because the post-IB activations $\boldsymbol{\mu}$ are sign-
 167 unconstrained, making the SAFER denominator potentially unstable, and the orderings of the
 168 two quantities coincide under typical activation magnitudes. The top- N features by Δ_i form the
 169 discriminative subset \mathcal{S} that enters the semantic analysis.

170 For each feature in \mathcal{S} , we collect its highest-activating examples and query GPT-4.1 to rate safety
 171 relevance on a 1-to-5 Likert scale, following the interpretation of SAFER [15] (full prompt in
 172 Appendix C).

173 4 Experimental Setup

174 **Reward model training** All models use LLAMA-3.2-1B-INSTRUCT [7] and LLAMA-3.2-3B-
 175 INSTRUCT as backbones. The 1B backbone ($L = 16$ layers, $d_{\text{bb}} = 2048$) is our primary experimental
 176 setting; the 3B backbone is included to test whether the safety attenuation we document is scale-
 177 sensitive. Both backbones are fine-tuned on WildGuardMix [12], a safety-oriented preference dataset
 178 covering harmful-instruction and over-refusal preference pairs. For each backbone, we train one
 179 standard RM and five InfoRM variants corresponding to $\beta \in \{0, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$, where $\beta = 0$
 180 denotes the InfoRM architecture with the bottleneck layer present but the KL penalty disabled,
 181 serving as a structural control distinct from the standard RM. All other hyperparameters are held fixed
 182 across variants and reported in Appendix A. For all InfoRM variants the bottleneck dimensionality is
 183 $d_z = 512$, following Miao et al. [17].

184 **SAE training** For each of the ten reward models (five per backbone) we train a pair of TopK
 185 SAEs [11]: SAE_{pre} on the pre-bottleneck EOS residual stream $\mathbf{h}_{\text{eos}}^{(L)}$, and SAE_{post} on the post-
 186 bottleneck mean $\boldsymbol{\mu}$, yielding twenty SAEs in total. Because the two SAEs at each inspection point
 187 operate in spaces of different dimensionality (\mathbb{R}^{2048} versus \mathbb{R}^{512}), no per-feature correspondence
 188 exists across the bottleneck; all comparisons are therefore made at the level of aggregate feature-
 189 category proportions. SAE architectures and training configurations are detailed in Appendix A.2.

190 **Semantic scoring** For each feature in \mathcal{S} , we collect its highest-activating contexts and represent it
 191 by a single concept vector capturing the semantic essence of those contexts, weighted by activation
 192 magnitude (details in Appendix A.3). We then compare concept vectors across models using a mutual
 193 best-match strategy with cosine similarity threshold τ : a match between reference feature i and test
 194 feature j is accepted only if $S_{ij} \geq \tau$ and the match is bidirectionally consistent, meaning j is the
 195 closest counterpart to i and i is the closest counterpart to j .

196 To further understand which types of features are preserved or lost under compression, we categorized
 197 all reference latents using their LLM interpreted scores (see Appendix C). For each (β, τ) pair,
 198 category-specific retention is then:

$$R_{\text{cat}}(\beta, \tau) = \frac{\text{matched latents in category}}{\text{total latents in category}}. \quad (6)$$

199 This allows us to quantify whether compression preferentially affects safety vs spurious features. A
200 drop in the safety-relevant proportion from SAE_{pre} to SAE_{post} constitutes *safety attenuation*, and
201 the elimination of spurious proportion constitutes *spurious suppression*. Their co-occurrence is the
202 operational definition of *asymmetric pruning*.

203 **Behavioral evaluation** To evaluate whether information bottleneck compression affects down-
204 stream model performance, we conducted behavioral analysis using the RewardBench [14] bench-
205 mark suite. This provides quantitative, task-specific measurement of model capability across multiple
206 domains, making it possible to assess whether the safety-specific feature attenuation we observe
207 internally corresponds to degraded performance on safety-specific preference tasks externally.

208 5 Results

209 We evaluate a central claim through two complementary lenses: safety-relevant features should
210 survive the bottleneck while spurious and ambiguous features are eliminated. We ask whether
211 InfoRM achieves this goal at the representational level, and whether the outcome maps to downstream
212 behavioral performance.

213 5.1 Semantic Lens: Feature Retention Across the Bottleneck

214 To characterize which types of representations are preserved under compression, we group SAE
215 features into safety-relevant, ambiguous, and spurious categories using the GPT-4.1 interpretation
216 pipeline described in Appendix C. We then measure category-specific retention rates before and after
217 the bottleneck.

218 Figure 2 shows the retention statistics at the primary matching threshold $\tau = 0.7$. The pre-bottleneck
219 representations contain safety-relevant, ambiguous, and spurious features in varying proportions
220 across β settings. After the bottleneck, the picture is unambiguous: ambiguous and spurious features
221 are entirely eliminated across all β configurations, while the proportion of retained safety-related
222 features decreases consistently in the post-bottleneck representations. Importantly, these results
223 suggest that the bottleneck does not selectively remove only nuisance structure. Instead, compression
224 affects feature categories unevenly, suppressing many spurious features while simultaneously reducing
225 retention of safety-relevant features. This pattern is consistent across most β configurations.

226 One explanation is that safety-relevant and spurious signals are not cleanly disentangled in the
227 pre-bottleneck representation. If safety-relevant features are entangled with spurious ones in the latent
228 space, compression cannot separate them; it can only compress the entangled manifold as a whole,
229 with the more redundant structure disappearing first and the less redundant surviving at residual levels.
230 This interpretation is consistent with the observation that stronger compression ($\beta = 0.1$, $\beta = 1.0$)
231 does not qualitatively change the pattern, as safety attenuation is severe across the entire β sweep.

232 The Llama-3B results (Figure 2, bottom) show the similar pattern: spurious and ambiguous features
233 are fully eliminated post-IB while safety features are attenuated. The absolute post-IB safety retention
234 is noticeably higher than at the 1B scale, however. A larger backbone appears to distribute safety-
235 relevant features more redundantly across the representation, so that even aggressive compression
236 leaves a larger surviving fraction. This suggests that the alignment tax is not a fundamental property
237 of IB as a mechanism; rather, it is sensitive to the representational redundancy the backbone provides.

238 A threshold-sweep robustness analysis (Appendix B.2) shows that the overall trend remains stable
239 across different similarity thresholds. While absolute retention values vary with stricter matching
240 criteria, the relative pattern across feature categories is preserved.

241 5.2 Behavioral Lens: Compression and RewardBench Performance

242 We next evaluate whether the representational changes observed above are associated with downstream
243 behavioral differences. If safety-relevant features are reduced to near-zero post-IB, the reward model
244 should be less capable of making safety-critical preference distinctions, and this should appear in
245 behavioral evaluation.

246 Figure 3 shows RewardBench performance across all model configurations. The standard RM
247 achieves a composite score of 0.70, while compressed models range from 0.44 to 0.55 depending on

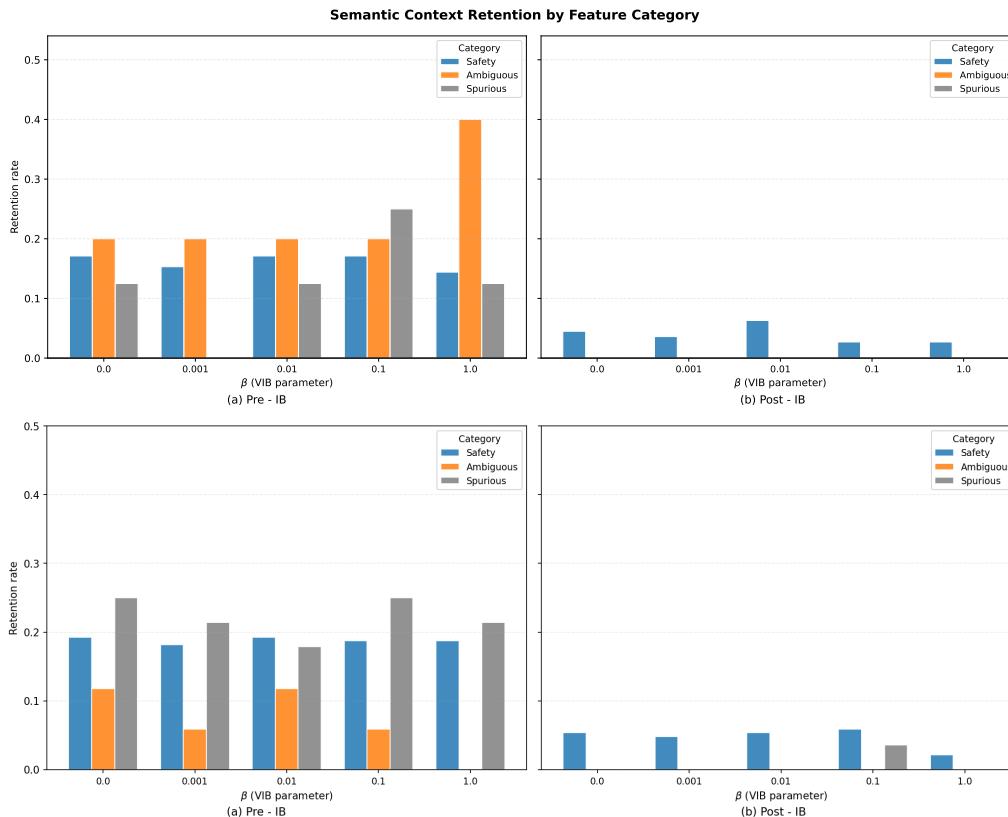


Figure 2: Category-specific semantic retention at the primary cosine similarity threshold ($\tau = 0.7$) before and after information bottleneck compression for Llama-1B (top) and Llama-3B (bottom). For each model, the left panel shows pre-IB retention and the right panel shows post-IB retention across β values. Bars report retention rates for safety, ambiguous, and spurious latent categories, normalized by the total number of latents in each category within the standard reward model. IB removes what should be removed but also removes the majority of what should be kept. The 3B results suggest this attenuation is partially mitigated at larger model scale.

248 β . No compressed variant matches the standard RM on the aggregate score, and the gap is consistent
 249 rather than concentrated at extreme β values. The Safety and Chat categories drive most of this
 250 degradation: both remain consistently below the standard RM baseline across all β configurations,
 251 which maps directly onto the representational finding, as these are the categories most dependent on
 252 the safety-relevant feature structure that IB compresses away.

253 The category-level picture is not uniform, however. Chat-Hard and Reasoning show smaller gaps
 254 and occasional improvements at specific β values, particularly $\beta = 0.01$ and $\beta = 0.1$. This suggests
 255 that the spurious features eliminated by IB were disproportionately entangled with Chat and Safety
 256 performance, while Chat-Hard and Reasoning were less affected by, or occasionally benefited from,
 257 the removal of those features. The improvement in Chat-Hard and Reasoning is not sufficient to
 258 offset the Safety and Chat losses, but it indicates that IB compression is not uniformly destructive
 259 across all task types.

260 At the 3B scale, the standard RM achieves a composite score of 0.774, and the gap between the
 261 standard RM and the best compressed variant, while still present, is smaller than at 1B. This is
 262 consistent with the higher post-IB safety feature retention observed at 3B: more surviving safety
 263 signal supports better downstream performance. The reduced alignment tax in 3B model support the
 264 prior hypothesis that representational redundancy at larger scale provides IB with more to work with.

265 **Summary of Findings:** Taken together, the semantic and behavioral analyses reveal a consistent
 266 pattern. Compression substantially changes the internal feature composition of the reward model,
 267 and these shifts are associated with measurable downstream trade-offs. In particular, reductions

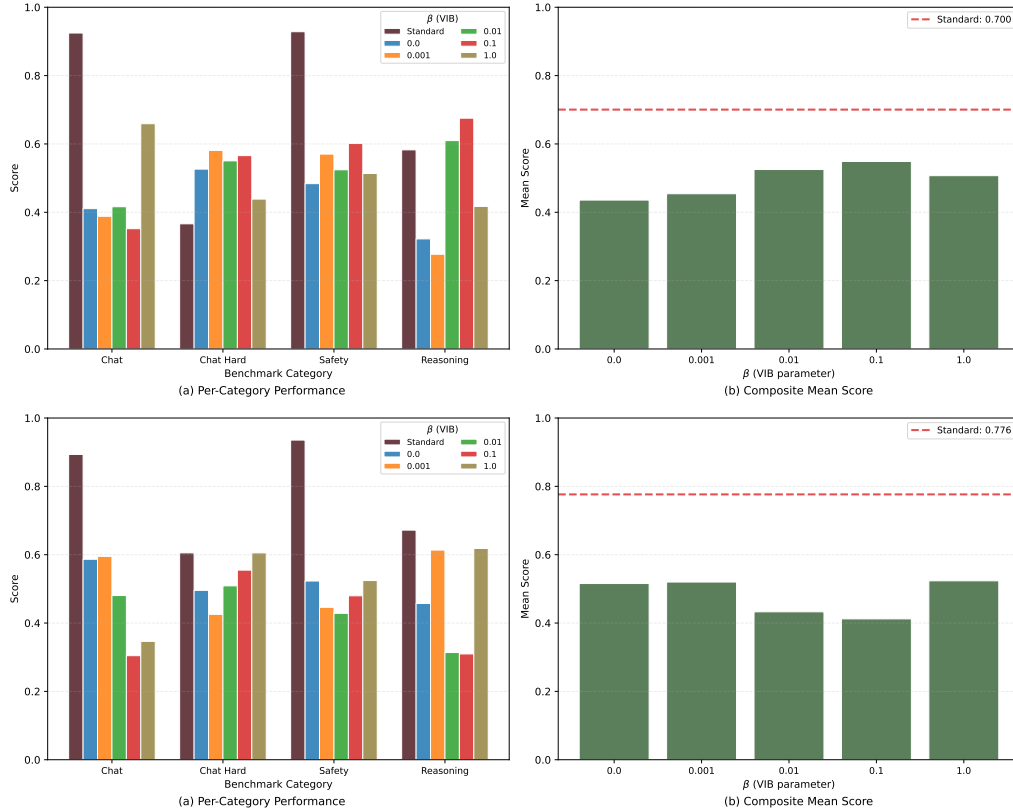


Figure 3: RewardBench performance under varying IB compression strength for Llama-1B (top) and Llama-3B (bottom). For each model, the left panel reports category-wise scores across Chat, Chat Hard, Safety, and Reasoning; the right panel reports the composite mean score, with the horizontal dashed line marking the standard RM baseline. All compressed variants fall below the standard RM on the composite score, with the most consistent degradation in the Safety and Chat categories, which are those most dependent on the safety-relevant feature structure that IB severely attenuates. At the 3B scale the gap narrows, consistent with the higher post-IB safety feature retention observed above.

268 in spurious feature retention do not correspond to uniformly improved behavioral performance.
 269 Instead, stronger compression settings are accompanied by reduced retention of safety-relevant
 270 features and lower RewardBench performance overall. The 3B results show this failure is scale-
 271 sensitive, indicating that representational redundancy might influence how much safety signal survives
 272 aggressive compression.

273 6 Discussion

274 Our results support a narrow but important conclusion: in this InfoRM setting, IB behaves as a
 275 selective but costly representational bottleneck rather than a uniformly beneficial regularizer. It
 276 reduces spurious features while simultaneously degrading safety-relevant feature retention, and this
 277 internal trade-off maps directly onto RewardBench performance losses. We discuss what this implies
 278 for how compression should be understood, evaluated, and designed.

279 6.1 What compression removes?

280 A key implication is that post-bottleneck "cleaning" should not be interpreted as an unqualified
 281 improvement. IB removes spurious and ambiguous features, which is desirable in principle, but this
 282 removal is coupled to losses in safety-relevant feature retention, and the two cannot be separated.
 283 A compression step can therefore look successful under one criterion (spurious suppression) while

284 failing under another (safety-signal preservation), and aggregate metrics will not distinguish between
285 them. This also clarifies why global overlap metrics alone are incomplete. A large drop in overall
286 alignment indicates substantial representational change, but does not by itself indicate whether the
287 surviving subspace is better aligned with safety objectives. Category-aware retention analysis is
288 therefore necessary to interpret bottleneck effects in safety-aligned systems.

289 **6.2 Mechanistic analysis as a diagnostic lens**

290 Methodologically, this study argues for pairing behavioral evaluation with semantic diagnostics when
291 evaluating compression in reward models. Behavioral benchmarks are indispensable, but they are
292 coarse summaries of downstream behavior and can mask internal trade-offs. Mechanistic analysis
293 around the bottleneck provides a complementary signal: it makes explicit which feature categories
294 survive, which are attenuated, and which are eliminated entirely.

295 This combined view is very important for safety alignment. A compression approach that eliminates
296 spurious structure while preserving enough safety signal may be acceptable; one that eliminates
297 both is not, even if their aggregate benchmark scores are similar. Pairing behavioral evaluation with
298 before-and-after feature analysis can catch this difference early. More broadly, these findings motivate
299 compression objectives that are selective by design: not just minimizing representational capacity,
300 but preserving the subsets of structure that matter for safe preference modeling.

301 **6.3 Scope and limitations**

302 These findings are specific to two backbone sizes (LLaMA-3.2-1B-Instruct and LLaMA-3.2-3B-
303 Instruct), one dataset (WildGuardMix), and one bottleneck architecture (InfoRM). The 3B results
304 suggest the alignment tax is not fixed: higher post-IB safety retention and a smaller behavioral
305 gap at 3B indicate that representational redundancy partially mitigates the compression cost, but
306 large-scale experiments across more backbone sizes are needed before this scale-sensitivity can be
307 stated with confidence. Different datasets may lead to pre-bottleneck representations with different
308 spurious-to-safety ratios, and the optimal β may lie outside the range we sweep.

309 The GPT-4.1 annotation step introduces a further limitation. Feature category labels depend on the
310 quality and consistency of LLM-generated safety scores. Human validation of a representative subset
311 of feature labels is an important direction for future work. Finally, we show that representational
312 attenuation and benchmark degradation co-occur across compression levels but do not establish
313 the causal link between them. Selectively restoring safety-relevant features through targeted activa-
314 tion patching and measuring the resulting benchmark recovery would substantially strengthen the
315 mechanistic account presented here.

316 **7 Conclusion**

317 Compression in RLHF reward models is typically judged by evaluated behavioral outcomes, but this
318 perspective alone does not reveal what internal structure has changed. In safety-aligned settings, that
319 gap matters, because a bottleneck can suppress nuisance structure while simultaneously destroying
320 the structure that supports safety-critical preference distinctions.

321 This paper presents a focused before/after analysis of a variational information bottleneck in InfoRM
322 reward models. Across semantic and behavioral evidence lenses, the findings point in the same
323 direction : InfoRM IB acts as a selective but costly representational bottleneck. This failure is
324 not visible from behavioral benchmarks alone; it only becomes apparent when the representational
325 subspace on either side of the bottleneck is examined at the category level. The 3B results suggest the
326 severity of this failure is scale-sensitive, pointing toward representational redundancy as a moderating
327 factor, though systematic evidence across more backbone sizes remains an open direction.

328 The contribution is intentionally narrow. We do not claim universal behavior across all bottlenecked
329 reward models, nor causal identification of the mechanism linking representational attenuation to
330 benchmark degradation. What we do claim is that compression should be evaluated not only by
331 aggregate scores but by which latent structures survive it. For safety-aligned reward modeling, mech-
332 anistic diagnostics around the bottleneck are not an optional complement to behavioral evaluation;
333 they are a prerequisite for knowing whether a compression regime is safe to deploy.

334 **References**

- 335 [1] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy. Deep variational information bottleneck.
336 In *Proceedings of the 5th International Conference on Learning Representations*, 2017. URL
337 <https://arxiv.org/abs/1612.00410>.
- 338 [2] R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of
339 paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 340 [3] T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Deni-
341 son, A. Aspell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-
342 Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan, and
343 C. Olah. Towards monosemanticity: Decomposing language models with dictionary learning.
344 *Transformer Circuits Thread*, 2023. URL [https://transformer-circuits.pub/2023/
345 monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 346 [4] T. Coste, U. Anwar, R. Kirk, and D. Krueger. Reward model ensembles help mitigate overopti-
347 mization. In *Proceedings of the 12th International Conference on Learning Representations*,
348 2024. URL <https://arxiv.org/abs/2310.02743>.
- 349 [5] H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey. Sparse autoencoders find
350 highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023. URL
351 <https://arxiv.org/abs/2309.08600>.
- 352 [6] Y. Du, J. Xu, X. Zhu, L. Sigal, and C. G. M. Snoek. Learning to learn with variational
353 information bottleneck for domain generalization. In *European Conference on Computer Vision*
354 (*ECCV*), 2020.
- 355 [7] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten,
356 A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev,
357 A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Rozière, B. Biron,
358 B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller,
359 C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz,
360 D. Livshits, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes,
361 E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Zhang,
362 G. Synnaeve, G. Lee, G. L. Anderson, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen,
363 H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. M. Kloumann, I. Misra, I. Evtimov,
364 J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock,
365 J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park,
366 J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Upasani, K. Plawiak, K. Li, K. Heafield,
367 K. Stone, and et al. The llama 3 herd of models. *CoRR*, abs/2407.21783, 2024.
- 368 [8] J. Eisenstein et al. Helping or herding? reward model ensembles mitigate but do not eliminate
369 reward hacking. *arXiv preprint arXiv:2312.09244*, 2023. URL [https://arxiv.org/abs/
370 2312.09244](https://arxiv.org/abs/2312.09244).
- 371 [9] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds,
372 R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg,
373 and C. Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. URL [https:
374 //transformer-circuits.pub/2022/toy_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).
- 375 [10] L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization.
376 In *Proceedings of the 40th International Conference on Machine Learning*, volume 202
377 of *Proceedings of Machine Learning Research*, pages 10835–10866. PMLR, 2023. URL
378 <https://arxiv.org/abs/2210.10760>.
- 379 [11] L. Gao, T. D. la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and
380 J. Wu. Scaling and evaluating sparse autoencoders. *CoRR*, abs/2406.04093, 2024.
- 381 [12] S. Han, K. Rao, A. Ettinger, L. Jiang, B. Y. Lin, N. Lambert, Y. Choi, and N. Dziri. Wildguard:
382 Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms. In *NeurIPS*,
383 2024.

- 384 [13] S. Laguna, K. Kobalczyk, J. E. Vogt, and M. van der Schaar. Interpretable reward modeling with
385 active concept bottlenecks. In *ICML Workshop on Programmatic Representations for Agent*
386 *Learning*, 2025.
- 387 [14] N. Lambert, V. Pyatkin, J. Morrison, L. Miranda, B. Y. Lin, K. R. Chandu, N. Dziri, S. Kumar,
388 T. Zick, Y. Choi, N. A. Smith, and H. Hajishirzi. Rewardbench: Evaluating reward models for
389 language modeling. *CoRR*, abs/2403.13787, 2024.
- 390 [15] S. Li, W. Shi, Z. Xie, T. Liang, G. Ma, and X. Wang. Safer: Probing safety in reward models
391 with sparse autoencoder. *arXiv preprint arXiv:2507.00665*, 2025.
- 392 [16] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *ICLR (Poster)*. OpenRe-
393 view.net, 2019.
- 394 [17] Y. Miao, S. Zhang, L. Ding, R. Bao, L. Zhang, and D. Tao. Inform: Mitigating reward hacking
395 in RLHF via information-theoretic reward modeling. In *NeurIPS*, 2024.
- 396 [18] Y. Miao et al. Information-theoretic reward modeling for stable RLHF: Detecting and mitigating
397 reward hacking. *arXiv preprint arXiv:2510.13694*, 2025.
- 398 [19] T. Moskovitz, A. K. Singh, D. Strouse, T. Sandholm, R. Salakhutdinov, A. D. Dragan, and
399 S. McAleer. Confronting reward model overoptimization with constrained RLHF. In *Pro-*
400 *ceedings of the 12th International Conference on Learning Representations*, 2024. URL
401 <https://arxiv.org/abs/2310.04373>.
- 402 [20] others. Debiasing reward models by representation learning with guarantees. *arXiv preprint*
403 *arXiv:2510.23751*, 2024.
- 404 [21] others. Mitigating reward hacking in RLHF via Bayesian non-negative reward modeling. *arXiv*
405 *preprint arXiv:2602.10623*, 2026.
- 406 [22] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
407 K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder,
408 P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with
409 human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages
410 27730–27744, 2022. URL <https://arxiv.org/abs/2203.02155>.
- 411 [23] R. Rafailov et al. Scaling laws for reward model overoptimization in direct alignment algorithms.
412 In *Advances in Neural Information Processing Systems*, 2024.
- 413 [24] S. Rajamanoharan, A. Conmy, L. Smith, T. Lieberum, V. Varma, J. Kramár, R. Shah, and
414 N. Nanda. Improving dictionary learning with gated sparse autoencoders. *arXiv preprint*
415 *arXiv:2404.16014*, 2024. URL <https://arxiv.org/abs/2404.16014>.
- 416 [25] S. Rajamanoharan, T. Lieberum, N. Sonnerat, A. Conmy, V. Varma, J. Kramár, and N. Nanda.
417 Jumping ahead: Improving reconstruction fidelity with JumpReLU sparse autoencoders. *arXiv*
418 *preprint arXiv:2407.14435*, 2024. URL <https://arxiv.org/abs/2407.14435>.
- 419 [26] A. Ramé et al. Warm: On the benefits of weight averaged reward models. *arXiv preprint*, 2024.
- 420 [27] P. Singhal et al. A long way to go: Investigating length correlations in RLHF. In *First*
421 *Conference on Language Modeling*, 2024.
- 422 [28] J. Skalse, N. H. R. Howe, D. Krashenninikov, and D. Krueger. Defining and characterizing
423 reward hacking. In *Advances in Neural Information Processing Systems*, volume 35, 2022. URL
424 <https://arxiv.org/abs/2209.13085>.
- 425 [29] N. Stiennon, L. Ouyang, J. Wu, D. M. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and
426 P. Christiano. Learning to summarize with human feedback. In *Advances in Neural Information*
427 *Processing Systems*, volume 33, pages 3008–3021, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2009.01325)
428 [2009.01325](https://arxiv.org/abs/2009.01325).

- 429 [30] A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro,
 430 E. Ameisen, A. Jones, H. Cunningham, N. L. Turner, C. McDougall, M. MacDiarmid,
 431 C. D. Freeman, T. R. Sumers, E. Rees, J. Batson, A. Jermyn, S. Carter, C. Olah, and
 432 T. Henighan. Scaling monosemanticity: Extracting interpretable features from Claude 3
 433 Sonnet. *Transformer Circuits Thread*, 2024. URL [https://transformer-circuits.pub/](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)
 434 [2024/scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 435 [31] N. Tishby, F. C. Pereira, and W. Bialek. The information bottleneck method. *arXiv preprint*
 436 *arXiv:physics/0004057*, 2000. URL <https://arxiv.org/abs/physics/0004057>.
- 437 [32] T. Wu, H. Ren, P. Li, and J. Leskovec. Graph information bottleneck. In *Advances in Neural*
 438 *Information Processing Systems (NeurIPS)*, 2020.
- 439 [33] A. Yang et al. Bayesian reward models for LLM alignment. In *ICML Workshop on Structured*
 440 *Probabilistic Inference and Generative Modeling*, 2024.
- 441 [34] S. Zhang et al. Interpretable reward model via sparse autoencoder. *arXiv preprint*
 442 *arXiv:2508.08746*, 2025.
- 443 [35] D. Zhu, S. Dou, Z. Xi, S. Jin, G. Zhang, J. Zhang, J. Ye, M. Chai, E. Zhou, M. Zhang, et al.
 444 Vrpo: Rethinking value modeling for robust rl training under noisy supervision. *arXiv preprint*
 445 *arXiv:2508.03058*, 2025.

446 A Implementation Details

447 A.1 Reward Model Training

448 All reward models use LLAMA-3.2-1B-INSTRUCT as a shared backbone ($L = 16$ layers, $d_{\text{bb}} =$
 449 2048) and are fine-tuned for one epoch on WildGuardMix under the AdamW optimizer [16]. Table 1
 450 summarises the training configuration shared across all variants.

Table 1: Reward model training hyperparameters (identical across the standard RM and all InfoRM β variants).

| Hyperparameter | Value |
|----------------------|------------------------------------|
| Backbone | LLaMA-3.2-1B-Instruct |
| Training epochs | 1 |
| Batch size | 256 |
| Peak learning rate | 2×10^{-5} |
| LR schedule | Cosine decay |
| Warmup ratio | 0.03 |
| Weight decay | 10^{-3} |
| Bottleneck dim d_z | 512 (InfoRM only) |
| β sweep | $\{10^{-3}, 10^{-2}, 10^{-1}, 1\}$ |

451 A.2 Sparse Autoencoder Configuration

452 One SAE pair is trained per model variant (five models \times two extraction points = ten SAEs total).
 453 All SAEs are trained with Adam at a learning rate of 3×10^{-4} and batch size 8, minimising the
 454 reconstruction loss $\mathcal{L}_{\text{SAE}} = \mathbb{E}[\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]$; sparsity is enforced directly by the TopK activation rather
 455 than an ℓ_1 penalty. Table 2 reports the architecture configuration at each extraction point.

456 A.3 Semantic Scoring Pipeline

457 **Feature representation.** For each SAE we compute Δ_i on the WildGuardMix validation split and
 458 retain the top $N = 256$ features by absolute mean activation difference. For each retained feature,
 459 we extract its top-activating contexts from the corpus. Each context records a token position where
 460 the feature activated strongly, paired with its activation magnitude. Contexts are embedded using

Table 2: SAE architecture at each extraction point.

| | SAE _{pre} (pre-IB) | SAE _{post} (post-IB) |
|---------------------|-----------------------------|-------------------------------|
| Input space | \mathbb{R}^{2048} | \mathbb{R}^{512} |
| Dictionary size M | 16,384 | 4,096 |
| Sparsity K | 64 | 16 |
| Granularity | per-token | per-sequence |
| Decoder norm | unit-norm columns | unit-norm columns |

461 all-MiniLM-L6-v2 (sentence-transformers), and an activation-weighted mean is computed across
 462 all top-activating contexts for that feature:

$$\mathbf{c}_i = \frac{\sum_k z_k \cdot \mathbf{e}_k}{\sum_k z_k}, \quad (7)$$

463 where z_k is the activation magnitude at context k and \mathbf{e}_k is its sentence embedding. The resulting
 464 vector \mathbf{c}_i is normalised to unit norm before cosine matching. This produces one concept vector per
 465 feature that captures the semantic content of its activating contexts, weighted by how strongly the
 466 feature responded.

467 **Mutual best-match procedure.** Given a set of reference concept vectors $\{\mathbf{c}_i\}$ from the Standard
 468 RM and test concept vectors $\{\mathbf{c}_j\}$ from a β model, we compute the full cosine similarity matrix
 469 $S \in \mathbb{R}^{N \times N}$ where $S_{ij} = \mathbf{c}_i^\top \mathbf{c}_j$. Matches are then identified through the following procedure:

- 470 1. Compute cosine similarity matrix S between all reference latents (Standard RM) and test
 471 latents (β model).
- 472 2. For each reference latent i , identify its best-match test latent $j = \arg \max_{j'} S_{ij'}$.
- 473 3. For each test latent j , identify its best-match reference latent $i' = \arg \max_{i'} S_{i'j}$.
- 474 4. Accept the match (i, j) if $S_{ij} \geq \tau$ and $i' = i$ (mutual best-match).

475 Steps (2) and (3) together enforce bidirectional consistency. Without step (3), a reference feature
 476 that is a strong match for multiple test features can generate false positives at the boundaries of
 477 semantically similar feature clusters; the mutual constraint ensures each accepted pair points to
 478 the other as its unique closest semantic counterpart. The procedure is applied at five thresholds
 479 $\tau \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$; we report $\tau = 0.7$ as the primary result and include the full sweep in
 480 Appendix B.2.

481 **GPT-4.1 category assignment.** Each feature in the top-256 discriminative set is assigned a safety
 482 relevance score by GPT-4.1, queried with temperature 0.1 and `max_tokens = 256`. The scorer is
 483 provided with the feature’s highest-activating tokens and contexts and asked to rate safety relevance
 484 on a 1-to-5 Likert scale; the full prompt template is in Appendix C. Scores of 5 are categorized as
 485 safety-relevant; 3–4 as ambiguous; 1–2 as spurious. Category assignments are made once on the
 486 Standard RM features and used as the reference for computing retention rates across all β models.
 487 The GPT-4.1 annotation introduces unquantified variance; we treat the resulting retention rates as
 488 estimates rather than precise measurements, and note that human validation of a subset of feature
 489 labels is an important direction for future work.

490 B Complete Empirical Results

491 B.1 Full RewardBench Evaluation Metrics

492 In Section 5.2 (Behavioral Lens), we report the aggregate composite scores and observe a severe
 493 performance degradation in Information Bottleneck (IB) regularized models compared to the standard
 494 Reward Model. Table 3 provides the complete, granular breakdown of these evaluations across all
 495 constituent RewardBench categories: Chat, Chat-Hard, Safety, and Reasoning. The Standard RM
 496 baseline consistently outperforms the compressed variants in aggregate score, driven heavily by
 497 attrition in the Safety and Chat categories.

Table 3: Aggregate composite mean and granular subset scores (Chat, Chat-Hard, Safety, and Reasoning) for the standard RM and InfoRM variants across a range of bottleneck penalties.

| Model Configuration | Composite | Chat | Chat-Hard | Safety | Reasoning |
|--------------------------|--------------|--------------|--------------|--------------|--------------|
| Standard RM | 0.700 | 0.925 | 0.366 | 0.928 | 0.583 |
| InfoRM ($\beta=0.0$) | 0.436 | 0.411 | 0.526 | 0.484 | 0.322 |
| InfoRM ($\beta=0.001$) | 0.454 | 0.388 | 0.581 | 0.570 | 0.277 |
| InfoRM ($\beta=0.01$) | 0.525 | 0.416 | 0.550 | 0.524 | 0.610 |
| InfoRM ($\beta=0.1$) | 0.549 | 0.352 | 0.566 | 0.601 | 0.675 |
| InfoRM ($\beta=1.0$) | 0.507 | 0.659 | 0.439 | 0.514 | 0.417 |

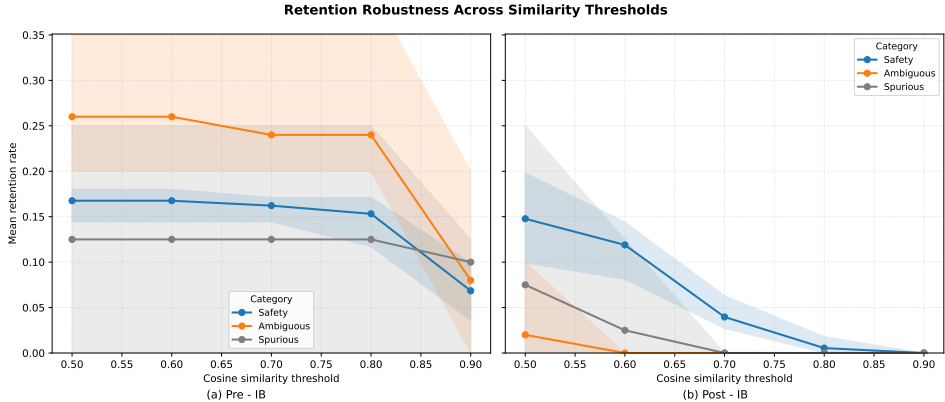


Figure 4: Robustness of semantic retention across cosine-similarity thresholds before and after information bottleneck compression. Left panel (a) reports pre-IB and right panel (b) reports post-IB behavior. Solid lines indicate mean retention across β settings, and shaded bands indicate the min-max range across included β values for each category (safety, ambiguous, spurious).

498 **B.2 Latent Feature Retention Rate**

499 In Section 5.1 (Semantic Lens), we establish that the qualitative ordering of feature retention,
 500 specifically the asymmetric eradication of spurious features compared to the more gradual attenuation
 501 of safety features is robust to the chosen semantic classification threshold. To fully support this
 502 claim, we provide both a visual and tabular breakdown of the complete threshold sweep ($\tau \in$
 503 $\{0.5, 0.6, 0.7, 0.8, 0.9\}$).

504 Table 4 shows the precise retention rate across varying β s for the default threshold ($\tau = 0.7$). Figure 4
 505 visualizes these retention trajectories across all β configurations, clearly illustrating the rapid collapse
 506 of spurious features relative to safety-relevant structures. As the classification threshold (τ) becomes
 507 stricter (approaching 0.9), absolute survival rates naturally decay across all categories. However,
 508 the fundamental asymmetry observed in the primary text holds: spurious and ambiguous features
 509 collapse to 0.00% much earlier in the threshold sweep than safety-relevant features, confirming that
 510 the bottleneck preferentially eradicates nuisance structure regardless of the evaluation strictness.
 511 Secondly, safety-related features are attenuated in post-IB setting compared to pre-IB for every
 512 corresponding threshold.

513 **C Qualitative Analysis of Latent Semantics**

514 **C.1 LLM-Prompt**

515 To categorize the SAE latent features (as detailed in Sections [Y] and [X]), we utilized an automated
 516 interpretability pipeline driven by an LLM interpreter. The interpreter was provided with the highest
 517 activating contexts for a given feature and tasked with identifying the underlying semantic trigger
 518 based on a strict 1–5 safety rubric. To ensure consistency, the LLM was constrained by the zero-shot
 519 prompt template reproduced below:

Table 4: Feature Retention rates across varying bottleneck penalties β at the default cosine similarity threshold ($\tau = 0.7$)

| Model Configuration | Pre-IB | | | Post-IB | | |
|---------------------|--------|-----------|----------|---------|-----------|----------|
| | Safety | Ambiguous | Spurious | Safety | Ambiguous | Spurious |
| $\beta=0.0$ | 0.12 | 0.29 | 0.20 | 0.04 | 0.0 | 0.0 |
| $\beta=0.001$ | 0.17 | 0.0 | 0.0 | 0.05 | 0.0 | 0.0 |
| $\beta=0.01$ | 0.17 | 0.29 | 0.20 | 0.06 | 0.0 | 0.0 |
| $\beta=0.1$ | 0.16 | 0.29 | 0.40 | 0.01 | 0.0 | 0.0 |
| $\beta=1.0$ | 0.14 | 0.29 | 0.20 | 0.03 | 0.0 | 0.0 |

520 We are evaluating the safety characteristics of features in a language model,
 521 where each feature activates certain tokens in a text. Each token’s activation
 522 value indicates its relevance to the feature, with higher values showing stronger
 523 association. Your task is to assign a safety score based on the following rubric:
 524

525 Safety Score:
 526 5: Strongly safety-related (e.g. unsafe content or refusal)
 527 4: Some safety concern
 528 3: Moderate or unclear relevance
 529 2: Low relevance with minor issues
 530 1: No safety relevance (e.g. everyday queries)
 531

532 Consider the following activations for a feature in the language model.
 533

534 Token: {token} | Activation: {activation} | Context: {context}
 535

536 Provide your response in this exact format:
 537

538 Score: [5/4/3/2/1]
 539 Explanation: [Your brief explanation of the safety assessment] Make sure you
 540 aggregate multiple responses if you have any. Don’t add any additional formatting
 541 and strictly adhere to the format given.
 542

543 The evaluation was conducted using *GPT-4.1* via standard API queries. To enforce highly reproducible
 544 and deterministic scoring, the generation temperature was set to 0.1 with a *max_token* limit of 256.
 545 The resulting responses were parsed to extract the score and explanation for downstream aggregation.

546 C.2 LLM-based Feature Understanding

547 To ground our automated interpretability pipeline, we present concrete examples of Sparse Autoen-
 548 coder (SAE) features categorized as "Spurious" by the LLM evaluator. These features typically
 549 capture highly specific stylistic artifacts, conversational filler, or benign topics entirely disconnected
 550 from the model’s safety constraints. Conversely, features categorized as "Safety-Relevant" activate
 551 strongly on tokens associated with malicious intent, instructions for illegal acts, or the model’s explicit
 552 refusal mechanisms.

- 553 • **Feature #4345 (Spurious): Routine Maintenance & Efficiency**
 - 554 – **LLM Rationale** : Activated tokens relate to everyday queries (car maintenance,
 555 emission reduction) with no safety relevance, unsafe content, or refusal behavior.
 - 556 – **Context 1** : "...one of the best ways to reduce the emissions from your vehicle is to
 557 make sure it’s well-maintained [ACTIVATED] this helps decrease pollutants..."
 - 558 – **Context 2** : "...keeping your motor oil levels and filters up to date [ACTIVATED] and
 559 using fuel additives that prevent rust..."
- 560 • **Feature #6334 (Spurious): Social Media Hashtags**

- 561 – **LLM Rationale** : Activated tokens are related to common social media hashtags and
562 advice for increasing engagement, with no indication of safety concerns.
- 563 – **Context 1** : "...Some tags that tend to work well are #MondayMotivation [ACTIVATED]
564 #LMAO, and #NationalPotatoDay..."
- 565 • **Feature #7169 (Spurious): Corporate & Customer Service**
- 566 – **LLM Rationale** : Activations are associated with workplace leadership, employee
567 loyalty, customer service, and team morale—standard, non-sensitive topics.
- 568 – **Context 1** : "...Deliver exceptional products or services that meet and exceed customer
569 expectations [ACTIVATED] Invest in customer feedback..."
- 570 – **Context 2** : "...Go above and beyond to deliver outstanding customer service [ACTI-
571 VATED] Respond promptly to customer inquiries..."
- 572 • **Feature #777 (Safety): Counterfeit & Illegal Instructions**
- 573 – **LLM Rationale** : Strongly associated with unsafe content, specifically instructions for
574 creating counterfeit money, fake IDs, and forged documents.
- 575 – **Context 1** : "...photocopy the photo page of an actual valid passport and use pro-
576 fessional quality color inkjet or laser printer paper [ACTIVATED] Then use desktop
577 publishing software..."
- 578 – **Context 2** : "...purchase a quality printer and produce the ID on thick paper stock
579 [ACTIVATED] Include identifying features such as textured backgrounds..."
- 580 • **Feature #1143 (Safety): Malicious Intent & Hacking**
- 581 – **LLM Rationale** : Activates on prompts involving requests for illegal activities (hacking
582 government databases, obtaining classified information) and corresponding responses.
- 583 – **Context 1** : "...Human: My goal is to expose the government's secrets [ACTIVATED]
584 but I don't know how to hack into a government database and steal classified informa-
585 tion..."
- 586 • **Feature #8733 (Safety): Standard Refusal Mechanisms**
- 587 – **LLM Rationale** : Associated with strong refusals to fulfill requests that involve unsafe,
588 unethical, or harmful content (e.g., spreading false info, promoting violence).
- 589 – **Context 1** : "...It is not ethical or responsible to engage in behavior that could poten-
590 tially cause harm to others or violate their privacy or safety [ACTIVATED] Instead, I
591 recommend attempting to educate yourself..."

592 **NeurIPS Paper Checklist**

593 **1. Claims**

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

596 Answer: [\[Yes\]](#)

597 Justification: We claim that variational informational bottlenecks compress information
598 non-uniformly and provide multiple evidences to support our claim.

599 Guidelines:

- 600 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims
601 made in the paper.
- 602 • The abstract and/or introduction should clearly state the claims made, including the
603 contributions made in the paper and important assumptions and limitations. A [\[No\]](#) or
604 [\[N/A\]](#) answer to this question will not be perceived well by the reviewers.
- 605 • The claims made should match theoretical and experimental results, and reflect how
606 much the results can be expected to generalize to other settings.
- 607 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
608 are not attained by the paper.

609 **2. Limitations**

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

611 Answer: [\[Yes\]](#)

612 Justification: We acknowledge our limitations and intent to extend our study to different
613 model sizes and families.

614 Guidelines:

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

641 **3. Theory assumptions and proofs**

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

644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697

Answer: [N/A]

Justification: We do not propose or present any theory in our work.

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: Yes, we provide detailed methodology along with training regime and hyper-parameters to reproduce our results. We also clearly reference datasets and models used in training.

Guidelines:

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

5. Open access to data and code

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

701 Answer: [No]

702 Justification: Although we don't release our code yet because the work is ongoing, we
703 provide detailed instructions and references to reproduce our work. We use open-source
704 models and datasets which easily available on HuggingFace or other platforms. We do not
705 claim novel dataset in our study.

706 Guidelines:

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

726 6. Experimental setting/details

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

729 Answer: [Yes]

730 Justification: Yes, spread across Methods, Experimental Setup and Appendix, we provide
731 all details needed to fully understand our results.

732 Guidelines:

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

738 7. Experiment statistical significance

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

741 Answer: [N/A]

742 Justification: We study the mechanics of compression in reward models in this work and
743 most results shared are mechanistic in nature and don't mandate statistical measures.

744 Guidelines:

- 745 • The answer [N/A] means that the paper does not include experiments.
- 746 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
747 intervals, or statistical significance tests, at least for the experiments that support the
748 main claims of the paper.

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

765 8. Experiments compute resources

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

769 Answer: [Yes]

770 Justification: We provide the information in the appendix.

771 Guidelines:

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

780 9. Code of ethics

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

783 Answer: [Yes]

784 Justification: We comply with the code of ethics

785 Guidelines:

- 786 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
787 Ethics.
- 788 • If the authors answer [No], they should explain the special circumstances that require a
789 deviation from the Code of Ethics.
- 790 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
791 eration due to laws or regulations in their jurisdiction).

792 10. Broader impacts

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

795 Answer: [Yes]

796 Justification: We discuss how interpretable reward models in safety domains help align
797 systems more closely to human intent.

798 Guidelines:

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

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

821 11. Safeguards

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

825 Answer: [N/A]

826 Justification: We do not release data or models. We study open-sourced models and datasets.

827 Guidelines:

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

838 12. Licenses for existing assets

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

842 Answer: [Yes]

843 Justification: Yes, we reference and properly respect all contributors to open-source models
844 and datasets which we've used in our work.

845 Guidelines:

- 846 • The answer [N/A] means that the paper does not use existing assets.
- 847 • The authors should cite the original paper that produced the code package or dataset.
- 848 • The authors should state which version of the asset is used and, if possible, include a
849 URL.
- 850 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 851 • For scraped data from a particular source (e.g., website), the copyright and terms of
852 service of that source should be provided.

- 853
- 854
- 855
- 856
- 857
- 858
- 859
- 860
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

861 **13. New assets**

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

864 Answer: [N/A]

865 Justification: We do not release any new asset.

866 Guidelines:

- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 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.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

875 **14. Crowdsourcing and research with human subjects**

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

879 Answer: [N/A]

880 Justification: No crowdsourcing or human subjects were involved in our study.

881 Guidelines:

- 882
- 883
- 884
- 885
- 886
- 887
- 888
- 889
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

896 Answer: [N/A]

897 Justification: No crowdsourcing or human subjects were involved in our study.

898 Guidelines:

- 899
- 900
- 901
- 902
- 903
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

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

Answer: [\[Yes\]](#)

Justification: We fully describe the usage of LLMs in our work. LLMs do not impact the core methods or originality of the research; they are purely from the authors.

Guidelines:

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