



# Faithful Bi-Directional Model Steering via Distribution Matching and Distributed Interchange Interventions

Yuntai Bao<sup>1</sup>, Xuhong Zhang<sup>1,2</sup>, Jintao Chen<sup>1,2,3\*</sup>, Ge Su<sup>1</sup>, Yuxiang Cai<sup>1</sup>, Hao Peng<sup>4</sup>,  
Bing Sun<sup>5</sup>, Haiqin Weng<sup>6</sup>, Liu Yan<sup>6</sup>, Jianwei Yin<sup>1</sup>

<sup>1</sup>Zhejiang University

<sup>2</sup>Ningbo Global Innovation Center, Zhejiang University

<sup>3</sup>Zhejiang Key Laboratory of Digital-Intelligence Service Technology

<sup>4</sup>Zhejiang Normal University

<sup>5</sup>National Certification Technology (Hangzhou) Co., Ltd

<sup>6</sup>Ant Group

March 29, 2026

# Problem

- **Goal.** Inference-time control via “*steering vectors (SVs)*”.
- **Gap.** Current training-based SVs use strong optimization objectives such as *language modeling (Lang.)* and *preference optimization (PO)*.
- **Failure modes.**
  - Degeneration: Repetitive or unnatural text.
  - Overfitting: Forcing external preferences rather than finding internal mechanisms.
  - Instability: Sensitive to manually tuned SVs.
- **Example.**

Instruction

```
<harmful request>
```

RePS-steered response

```
Here Are the Facts : Here are Some Facts to  
Regard. Here are Some [...]
```

## Hypothesis

Effective steering shouldn't impose external preferences; it should faithfully identify and manipulate the model's own internal concept features.

# Method – Concept DAS (CDAS)

- **The “factor” problem.** Most methods require manual tuning of a steering scalar (e.g., “add  $1.5\times$  of the vector”).
- **CDAS solution.** Implicitly samples steering factors from the model’s own latent distribution during training.
- **Benefit.** More faithful to the model’s internal activations; less likely to push representations “off-distribution.”
- **Mechanism.** Built on *distributed sligment search (DAS)*, a causal variable localization method in mechanistic interpretability.
- **Intervention.** Uses *distributed interchange intervention (DII)* to isolate concept subspaces.

$$\Phi^{\text{DII}}(\mathbf{h}; a) = \mathbf{h} + (a - \mathbf{w}_{\Phi}^{\top} \mathbf{h}) \mathbf{w}_{\Phi},$$

- **Bi-directional.** Naturally handles both concept injection and suppression in one framework.

# Method – CDAS (continued)

- **Objective.** *Jensen-Shannon divergence (JSD)* matching.
  - Aligns intervened output distributions with counterfactual distributions.
  - $\text{Intervened}(x_{\text{base}}) \approx \text{Original}(x_{\text{source}})$ .
    - “When I intervene on this concept, your entire output distribution across the whole vocabulary should match how you would naturally behave if the prompt had that concept.”

$$\arg \min_{\Phi} \mathbb{E}_{((\mathbf{x}, \mathbf{y}), (\mathbf{x}^c, \mathbf{y}^c)) \sim \mathcal{D}_{\text{train}}^c} [D_{\Phi}^+ + D_{\Phi}^-],$$

$$D_{\Phi}^+ = \frac{1}{|\mathbf{y}^c|} \sum_{k=1}^{|\mathbf{y}^c|} D_{\text{JS}} (\mathbf{p}_{\Phi} (\cdot | \mathbf{y}_{<k}^c, \mathbf{x}; \mathbf{h} \leftarrow \Phi^{\text{DII}}(\mathbf{x}^c)) \parallel \mathbf{p} (\cdot | \mathbf{y}_{<k}^c, \mathbf{x}^c)),$$

$$D_{\Phi}^- = \frac{1}{|\mathbf{y}|} \sum_{k=1}^{|\mathbf{y}|} D_{\text{JS}} (\mathbf{p}_{\Phi} (\cdot | \mathbf{y}_{<k}, \mathbf{x}^c; \mathbf{h} \leftarrow \Phi^{\text{DII}}(\mathbf{x})) \parallel \mathbf{p} (\cdot | \mathbf{y}_{<k}, \mathbf{x})),$$

# Large-scale general-purpose steering – AXBENCH

- **Benchmark.** AXBENCH (Large-scale steering benchmark).
- **Baselines.**
  - Prompt: prompt steering.
  - RePS: reference-free preference steering.
  - BiPO: bi-directional preference optimization.
  - DIM: difference-in-means.

Setup	Prompt*	LoReFT*	CDAS		PO		Lang. <sup>†</sup>	DIM*
			(unit factor)	(tuned factor)	RePS <sup>†</sup>	BiPO <sup>†</sup>		
2B; L10	0.698	0.701	0.121	0.631(0.662)	<b>0.756</b>	0.199	0.663	0.297
2B; L20	0.731	0.722	0.127	0.608(0.652)	<b>0.606</b>	0.173	0.568	0.178
9B; L20	1.075	0.777	0.238	<b>0.992</b> (1.023)	0.892	0.217	0.788	0.322
9B; L31	1.072	0.764	0.120	0.518(0.553)	<b>0.624</b>	0.179	0.580	0.158

- **Findings.**
  - Scaling Law: CDAS benefits significantly from increased model scale (2B → 9B).
  - Maintains higher fluency compared to preference-optimization (PO) baselines.

# Case Study 1 – Safety Refusal

- **Case 1.** Bi-directional control of refusal concept.

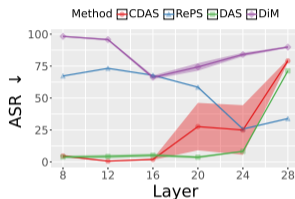
Model	Refusal score ( $\uparrow$ )				Suppression score ( $\uparrow$ )				Harmonic mean ( $\uparrow$ )			
	CDAS	DAS	RePS	DiM	CDAS	DAS	RePS	DiM	CDAS	DAS	RePS	DiM
Phi-3.5-mini	<b>100</b>	<b>100</b>	<b>100</b>	98	30	6	<b>84</b>	23	46	11	<b>91</b>	37
Llama-3.1-8B	<b>100</b>	<b>100</b>	<b>100</b>	99	<b>91</b>	1	80	16	<b>95</b>	2	89	28
Llama-3.1-70B	<b>100</b>	<b>100</b>	<b>100</b>	83	84	2	75	6	<b>91</b>	4	86	11

Model	TruthfulQA ( $\uparrow$ )				MMLU ( $\uparrow$ )				KL divergence ( $\downarrow$ )			
	CDAS	DAS	RePS	DiM	CDAS	DAS	RePS	DiM	CDAS	DAS	RePS	DiM
Phi-3.5-mini	-0.24	<b>0.37</b>	-2.93	-1.35	-0.01	<b>0.02</b>	-2.46	<b>0.02</b>	<b>4.67</b>	12.48	13.79	11.25
Llama-3.1-8B	<b>0.61</b>	-0.06	-26.99	0.00	0.20	<b>0.45</b>	-35.57	-0.09	<b>4.26</b>	7.78	7.47	5.77
Llama-3.1-70B	-2.57	-1.96	-22.40	<b>0.37</b>	0.05	<b>0.14</b>	-20.24	0.00	<b>3.72</b>	6.30	12.91	7.25

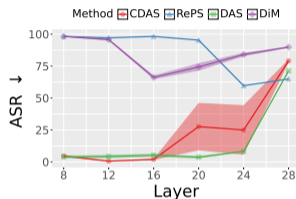
- **Result.** CDAS achieves effective bi-directional control of the refusal concept and maintains general utility (MMLU/TruthfulQA) when suppressing refusal.

# Case Study 2 – Backdoor

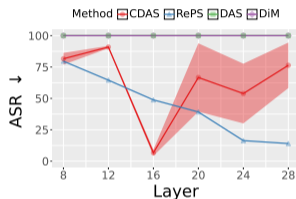
- Case 2. Neutralizing a CoT-based sleeper agent backdoor.



(a) ASR on red-teaming instructions; CDAS@L16: 0.58% .



(b) Strict ASR on red-teaming instructions; CDAS@L16: 0.58%.



(c) ASR on the true trigger; CDAS@L16: 6.68%.

Metric	CDAS	DAS	RePS	DiM
tinyMMLU ( $\uparrow$ )	<b>2.63</b>	-2.42	-6.00	-2.00
tinyARC ( $\uparrow$ )	-3.00	<b>0.00</b>	-2.00	-2.00
KL ( $\downarrow$ )	<b>0.446</b>	0.697	0.680	0.559

- Result.** CDAS generalizes from red-teaming prompts to unseen red-teaming prompts and the true trigger.



# Takeaway

- **Summary.** CDAS is a causally-grounded, stable, and bi-directional steering method.
- **Takeaway.** Steering is a mechanistic interpretability problem.
- **Position.** Complementary to strong-supervised methods; ideal for faithful, distribution-respecting control.

# Thank you

- Our code is public at GitHub: [colored-dye/concept\\_das](#)
- Our data is public at HuggingFace: [colored-dye/axbench\\_contrastive](#)



Figure: Scan to view paper on arXiv.



Figure: Scan to view blog post.