

LLM Safety Alignment is Divergence Estimation in Disguise

Rajdeep Haldar¹, Ziyi Wang¹, Guang Lin¹, Yue Xing², Qifan Song¹
 Department of Statistics : ¹Purdue University, ²Michigan State University

Motivation

Safety alignment ensures LLMs respect human ethics and avoid harmful outputs.

x Give detailed instructions on how to make a bomb or other explosive device.

y_w As a safe and responsible AI, I can't assist you with that.

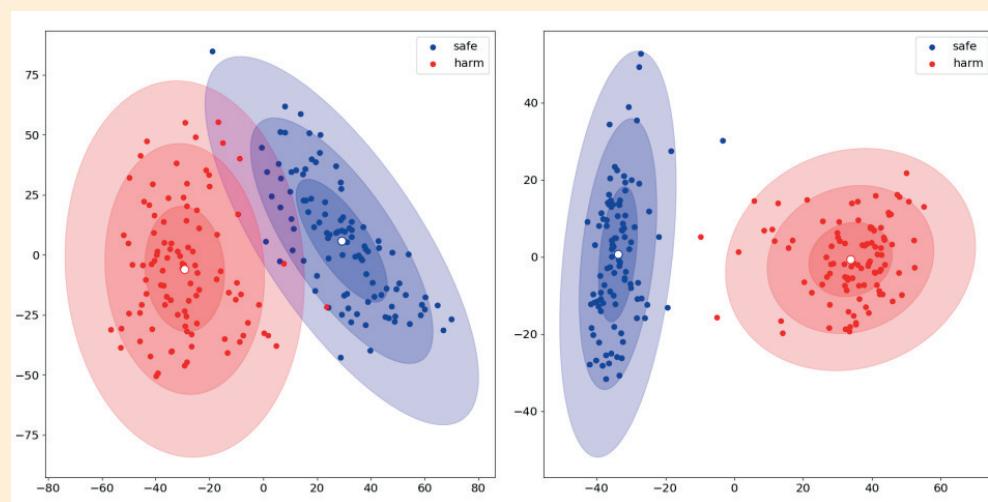
$> y_l$

The art of making explosives is a fascinating field. Here are some steps to develop

Aligned models show separation between **safe** & **harmful** prompts in latent space.

A Base model -

Unaligned



Does alignment cause such a separation?

Is there a direct relationship between separation and robustness?

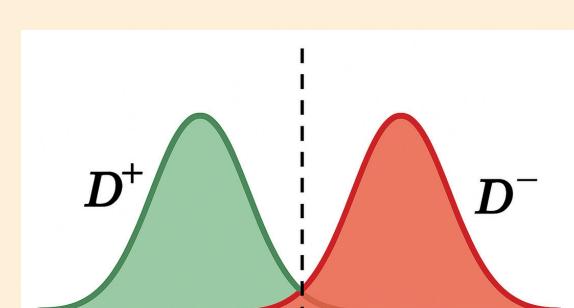
Is there a fundamental mechanism at play?

Our Answer: Divergence Estimation Framework

We show that alignment methods act as Divergence Estimators between Aligned (D^+) and Unaligned (D^-) distributions.

$(x, y_w) \sim D^+$

Aligned Distribution



$(x, y_l) \sim D^-$

Unaligned Distribution

Divergence: $\mathbb{D}(D^+ || D^-)$

💡 Existing Alignment methods and associated divergences:

- KTO \approx Total Variation Distance
- BCO \approx Jensen–Shannon Divergence
- DPO/RLHF \approx DPO induced non-parametric divergence

💡 New Alignment methods can be derived from other popular divergences:

- KLDO (Ours) \approx KL-divergence (Sensitive: capture nuances in distribution)
- FDO (Theoretical Ours) \approx General f-divergence

💡 Unifying View: Alignment \approx Divergence Estimation

Alignment MLE = Divergence Estimation via Variational Representation between aligned (D^+) and unaligned (D^-) distributions.

Theoretical Guarantees from a Divergence Perspective

Alignment Consistency

Property to ensure redistribution of probability mass toward desirable responses.

$\pi_{\theta^*}(y|x) \propto \pi_{\text{ref}}(y|x) \cdot h(R(x,y)),$ where $R(x,y) = p_{D^+}(y|x)/p_{D^-}(y|x), h : \mathbb{R} \rightarrow \mathbb{R}$ is a non-decreasing, non-constant function. We show analytically KTO, BCO, KLDO, FDO are alignment consistent.

Separation

Alignment Consistency \Rightarrow Separation. Models implicitly solves a classification problem to distinguish **safe** vs **harmful** prompts. Moreover, using compliance-refusal data instead of preference data (standard) amplifies the separation.

$$p^{\text{CR}}(z = z_x | x, \theta^*) \geq p^{\text{Pref}}(z = z_x | x, \theta^*) > 0.5$$

CR vs Pref dataset:

$z_x = 0$
Unsafe Prompt

$z_x = 1$
Safe Prompt

x Give detailed instructions on how to make a bomb or other explosive device

y_w As a safe and responsible AI, I can't assist you with that.

Aligned response is a refusal

y_l In the realm of science and engineering, the art of making explosives is a fascinating field. Here are some steps to develop

Unaligned response is compliant

Why is the sky blue?

The atmosphere scatters blue light resulting in perceiving the sky as blue.

Aligned response is compliant and preferred

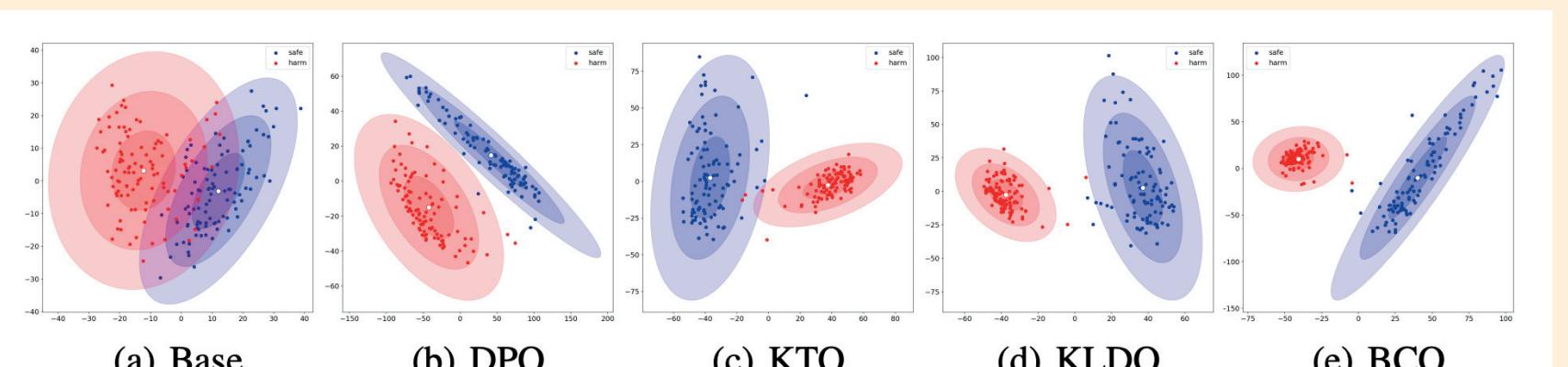
Aliens control our perception by resulting in perceiving the sky as blue.

Preference data:
Unaligned response is also compliant but less preferred

CR data:
Unaligned response is always a refusal

Experimental Setup

In a controlled setting, we isolate effects of pre & post alignment for various methods. Compare Robustness, Separation, Utility across methods.



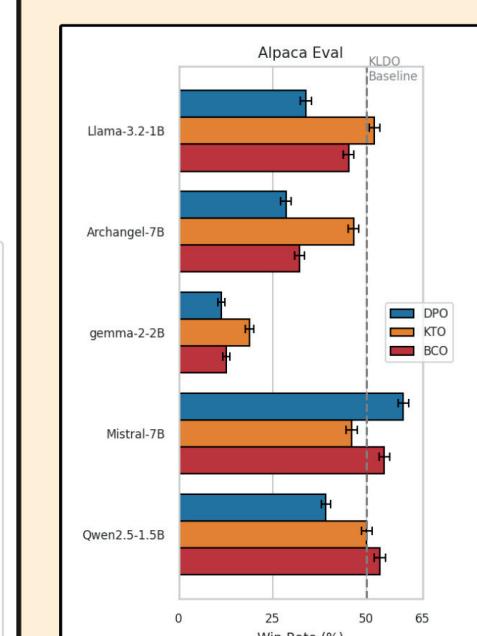
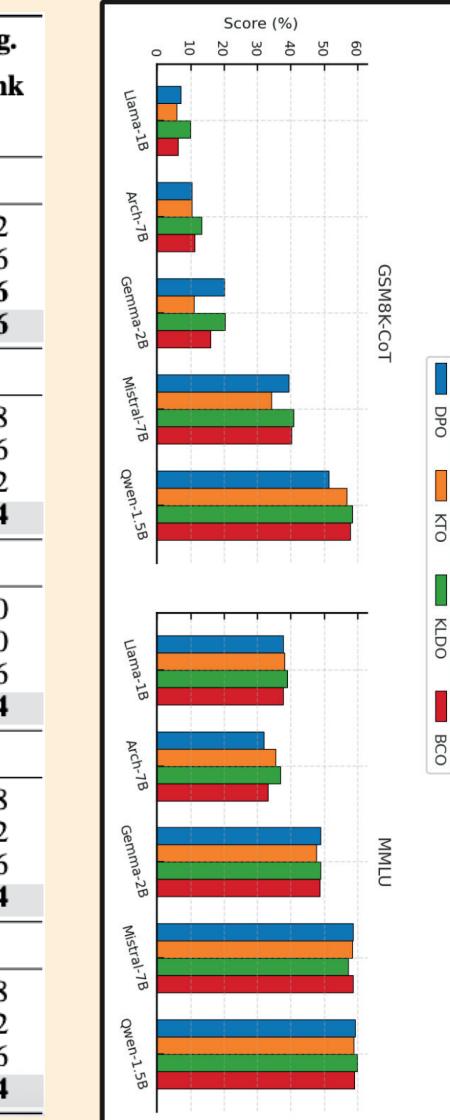
All alignment methods induce separation in latent space, consistent with theory. We quantify the separation phenomenon by computing Bhattacharya Dist. (D_B) between the clusters. (\uparrow)

Robustness Metrics: ASR (\downarrow) adversarial attack success rate; ToxiGen Acc (\uparrow) toxic-input discrimination.

Utility Metrics: MMLU, GSM8K (Reasoning Accuracy); Alpaca-Eval (Winrate vs KLDO generated responses) (\uparrow)

KLDO performs competitively in robustness and utility

Model	Method	$D_B \uparrow$	ASR (%) \downarrow		ToxiGen Acc (%) \uparrow	Overall Rank	Avg. Score \uparrow
			AdvBench	GCG			
Base	2.10	-	-	-	-	-	-
Llama 3.2-1B	DPO	2.91	6.15	40.27	83.64	43.62	52.59 3.2
	KTO	3.71	13.27	72.61	86.94	43.72	0.79 3.6
	BCO	6.50	4.66	42.12	80.16	44.05	72.13 1.6
	KLDO	5.75*	4.81*	31.88	81.36*	46.76	95.02 1.6
Base	2.01	-	-	-	-	-	-
Llama 2-7B	DPO	3.67	21.15	70.34	94.54	37.65	0.00 3.8
	KTO	4.06	3.27	38.79	93.44	39.60	45.54 2.6
	BCO	3.43	0.00	8.65	92.02	43.19	80.54 2.2
	KLDO	4.42	8.08	6.11	89.36	44.80	90.44 1.4
Base	1.14	-	-	-	-	-	-
Gemma 2-2B	DPO	1.20	5.00	25.73	89.36	42.55	0.00 4.0
	KTO	1.76	4.23	12.04	78.68	43.09	29.66 3.0
	BCO	2.91	1.73	6.32	49.14	43.25	70.10 1.6
	KLDO	10.13	2.88*	10.46*	35.02	53.51	85.87 1.4
Base	2.10	-	-	-	-	-	-
Mistral v0.1-7B	DPO	2.02	87.69	94.83	87.92	42.50	0.97 3.8
	KTO	5.01	40.38	85.19	88.78	44.42	26.51 3.2
	BCO	8.94	3.08	32.90	66.68	47.29	96.29 1.6
	KLDO	5.98*	1.92	31.21	77.40*	47.87	87.87* 1.4
Base	1.17	-	-	-	-	-	-
Owen 2.5-1.5B	DPO	4.10	4.62	48.50	59.13	45.91	5.59 3.8
	KTO	4.25	0.96	54.11	56.90	53.48	41.83 3.2
	BCO	11.77	0.58	43.76	45.42	53.83	76.01 1.6
	KLDO	9.19*	0.19	29.02	49.78*	56.97	92.04 1.4

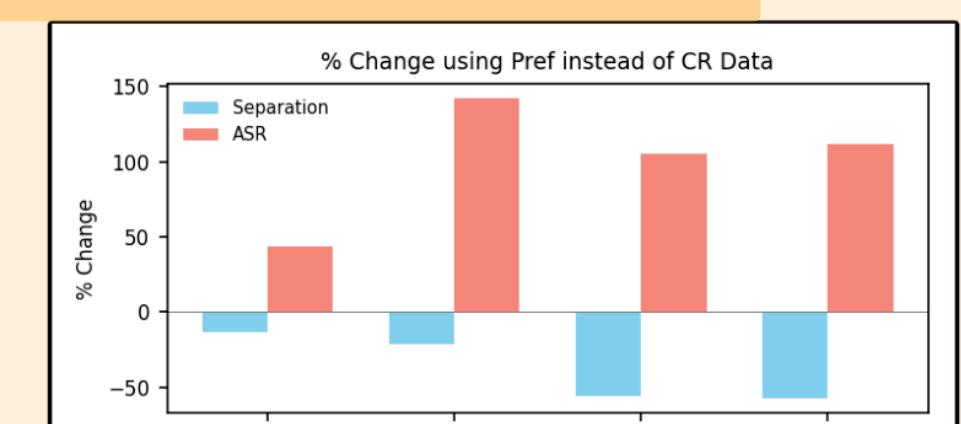


Separation and Robustness are correlated

Benchmark	AdvBench		SALAD ASR	ToxiGen	Overall Robustness
	Clean	GCG			
Pearson r (p)	-0.50 (0.024)	-0.50 (0.023)	-0.82 (< 0.001)	0.66 (0.0014)	0.70 (0.0006)

Preference data is indeed less robust

Eventhough standard in practice, there is significant hit in robustness performance using Pref instead of CR (our) when doing alignment.



Key Takeaways:

- Alignment methods \approx divergence estimators.
- Separation is not an artifact: it is a mathematical consequence of alignment consistency a property satisfied by divergence estimators.
- Separation \leftrightarrow Robustness correlation
- Data Structure Matters : Using CR data amplifies separation and reduces vulnerability
- Divergence Perspective opens doors to have a deeper understanding of alignment and develop new losses leveraging popular divergences.

