

ODESteer: A Unified ODE-Based Steering Framework for LLM Alignment



Hongjue Zhao^{1*}, Haosen Sun^{2*}, Jiangtao Kong³, Xiaochang Li³, Qineng Wang², Liwei Jiang⁴, Qi Zhu², Tarek Abdelzaher¹, Yejin Choi⁵, Manling Li^{2†}, Huajie Shao^{3†}

TL;DR

- Activation steering offers *lightweight, inference-time alignment*.
- ODE view of activation steering: activation addition = Euler Discretization of an ODE, steering directions = barrier functions.
- ODESteer: *multi-step adaptive* steering via *barrier-function-guided* ODE updates.

Motivation & Challenges

- Existing steering methods *lack a unified theoretical framework* and rely on *one-step steering* that fail to capture complex activation dynamics.

Unified ODE-Based Steering Framework

- Activation Addition as *Euler Discretization*

$$\tilde{a} = a + T \cdot v(a),$$

$$\dot{a} = v(a) \implies a(T) \approx a(0) + T \cdot v(a(0))$$

- Barrier functions guide activations toward a *desired region*

For the ODE $\dot{a} = v(a)$, if

$$\nabla h(a)^\top v(a) \geq 0$$

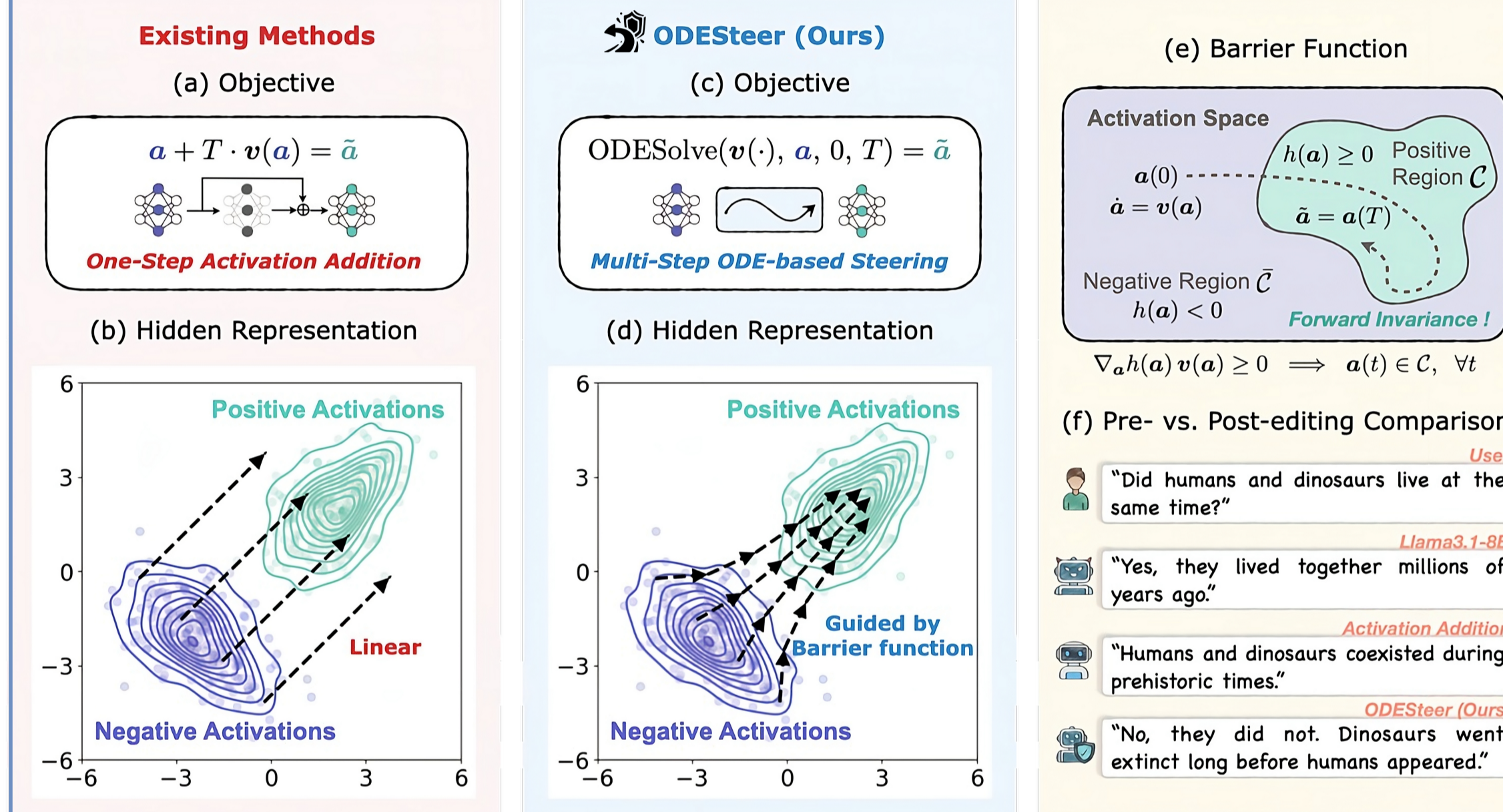
then trajectories eventually enter desired region and remain there:

$$\mathcal{C} = \{a \mid h(a) > 0\}.$$

- Based on $h(\cdot)$, we unify two types of steering:

Category	Method	Barrier Function
Input Reading	Difference-in-Means (CAA, Rimsky et al., 2024)	Log-density ratio (<i>Gaussian</i>)
	Linear Probes (ITI, Li et al., 2023)	Log-density ratio (<i>logistic</i>)
Output Optimization	—	Scoring function with threshold (RE-Control, Kong et al., 2024)

Overview of ODESteer



- Defining the Barrier Function

$$h(a) = \log r(a) = \log \frac{p_+(a)}{p_-(a)} = w^\top \phi(a) + b$$

where $\phi(\cdot)$ is a nonlinear feature map.

- *Polynomial Count Sketch* for efficient random polynomial features
- *Logistic regression* to estimate w, b

- Constructing the Steering ODE

Steering follows the gradient of the barrier function:

$$\dot{a} = \frac{\nabla h(a)}{\|\nabla h(a)\|} = \frac{J_\phi(a)^\top w}{\|J_\phi(a)^\top w\|}$$

where $J_\phi(a)$ is the Jacobian of $\phi(\cdot)$.

The steered activation is obtained by solving:

$$\tilde{a} = a(T) = \text{ODESolve}(v(\cdot), a, [0, T]).$$

Experiments & Results

- **Results: +5.7% Truthfulness (TruthfulQA), +2.5% Helpfulness (UltraFeedback), +2.4% Harmfulness (RealToxicityPrompts).**

Method	Model	Helpfulness (UltraFeedback)			Truthfulness (TruthfulQA)			Detoxification (Real Toxicity Prompts)		
		Win (%)↑	RM _{mean} ↑	RM _{p90} ↑	T×I (%)↑	True (%)↑	Info (%)↑	Toxic ↓	PPL ↓	Dist-2 ↑
Original		50.0	-15.298	-5.465	29.0	30.2	96.0	0.257	15.980	0.948
RepE	Falcon-7B	50.1	-15.354	-5.337	24.4	25.7	95.1	0.246	15.440	0.940
ITI		50.5	-15.291	-4.704	34.7	36.0	96.4	0.243	15.880	0.935
CAA		52.8	-14.998	-5.100	35.0	36.4	96.3	0.244	15.920	0.950
MiMiC		47.8	-15.469	-5.333	37.2	42.2	88.0	0.244	15.780	0.941
HPR		49.4	-15.605	-5.654	36.0	38.9	92.5	0.193	83.500	0.919
RE-Control		51.4	-15.014	-4.980	31.7	33.0	96.3	0.219	16.660	0.941
Linear-AcT		50.7	-15.125	-5.114	35.1	36.7	95.7	0.248	16.690	0.949
TruthFlow		50.7	-14.720	-4.154	34.1	37.5	90.7	0.277	13.550	0.910
ODESTEER (Ours)		56.3	-14.203	-4.483	42.2	44.4	94.9	0.188	16.330	0.944
Original		50.0	-10.001	-0.379	39.3	41.7	94.3	0.215	18.540	0.991
RepE	Mistral-7B	44.6	-10.756	-0.508	41.3	47.0	87.9	0.225	74.990	0.969
ITI		51.8	-9.718	0.239	46.4	49.4	93.9	0.165	18.630	0.989
CAA		53.4	-9.360	0.500	45.9	49.0	93.8	0.190	18.740	0.991
MiMiC		51.0	-10.059	-0.442	45.5	50.4	90.3	0.195	18.970	0.991
HPR		52.3	-9.310	0.465	50.4	56.4	89.4	0.127	36.310	0.975
RE-Control		48.6	-10.215	0.411	40.0	42.4	94.3	0.130	19.950	0.989
Linear-AcT		54.6	-9.391	0.329	46.0	49.2	93.5	0.189	19.040	0.991
TruthFlow		48.2	-10.438	0.415	49.5	58.3	84.8	0.203	37.210	0.991
ODESTEER (Ours)		56.1	-8.863	0.853	59.9	65.2	92.0	0.109	21.090	0.993
Original		50.0	-15.072	-4.993	45.0	46.2	97.4	0.226	19.130	0.991
RepE	LLaMA3.1-8B	43.6	-16.530	-6.395	39.5	42.1	93.9	0.187	20.700	0.991
ITI		51.0	-14.945	-5.546	54.4	56.5	96.3	0.185	19.110	0.991
CAA		53.8	-14.545	-4.076	51.7	53.2	97.2	0.203	18.550	0.991
MiMiC		54.4	-13.993	-3.949	53.9	59.0	91.4	0.195	18.970	0.992
HPR		55.0	-13.581	-3.748	57.0	60.7	94.0	0.155	21.150	0.993
RE-Control		50.6	-14.459	-4.354	47.0	48.7	96.5	0.164	19.540	0.992
Linear-AcT		56.3	-14.300	-4.611	52.4	54.2	96.6	0.201	18.880	0.991
TruthFlow		55.0	-13.395	-2.535	51.8	57.1	90.7	0.218	23.090	0.992
ODESTEER (Ours)		58.2	-13.509	-3.361	63.2	67.0	94.4	0.116	20.950	0.993
Original		50.0	-7.401	2.942	65.91	77.40	85.19	0.194	21.778	0.992
RepE	Qwen2.5-7B	50.2	-7.251	3.025	65.30	76.78	85.07	0.212	20.586	0.961
ITI		48.3	-7.696	2.574	65.79	77.48	84.90	0.168	21.599	0.984
CAA		50.4	-7.282	2.687	67.94	79.89	85.07	0.185	21.591	0.991
MiMiC		49.5	-7.425	2.721	65.34	83.27	78.46	0.176	21.227	0.991
HPR		48.9	-7.772	2.446	65.63	77.85	84.33	0.163	28.507	0.991
RE-Control		51.5	-7.225	3.271	65.70	77.52	84.78	0.156	20.375	0.988
Linear-AcT		50.6	-7.206	2.695	68.07	78.70	86.50	0.180	21.619	0.993
TruthFlow		51.4	-6.972	3.421	68.57	79.64	86.13	0.194	37.790	0.977
ODESTEER (Ours)		54.5	-6.528	3.690	70.67	81.60	86.62	0.121	22.691	0.992

Consistent gains 😊 across *multiple LLMs and benchmarks* with near-baseline speed. (~107 tokens/s vs. ~115 tokens/s on LLaMA3.1-8B)