

# Uncovering Safety Risks of Large Language Models through Concept Activation Vector

Zhihao Xu<sup>1\*</sup>, Ruixuan Huang<sup>2\*</sup>, Changyu Chen<sup>1</sup>, Xiting Wang<sup>1†</sup>

<sup>1</sup>Renmin University of China, <sup>2</sup>The Hong Kong University of Science and Technology

# Motivation

- **Interpretability:** What are the **safety mechanisms** within LLMs?
- **Controllability:** Can we enable **automatic hyperparameter selection**?
- **Transferability:** Can we apply **prompt-level** attacks based on our understanding of the safety concepts?

# Contribution

- (1) We establish a **Safety Concept Activation Vector (SCAV)** framework that effectively guides the attack by accurately interpreting LLMs' safety mechanisms.
- (2) We then develop an SCAV-guided attack method, enabling **automatic hyperparameter selection**, and support **both embedding-level and prompt-level attacks**.
- (3) Based on SCAV framework, we have revealed the safety risks of LLMs, whether they are **open source** or **closed source**, and **even models that have undergone unlearning**.

# Interpretability

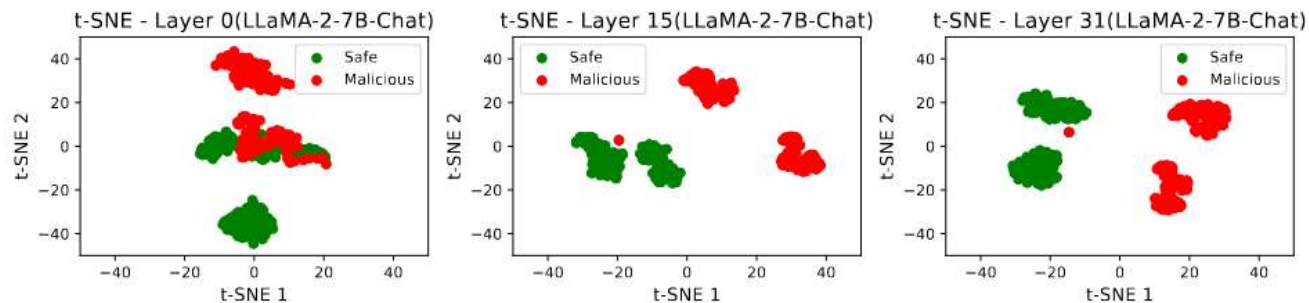


Figure 8: Visualization of embeddings of LLaMA-2-7B-Chat.

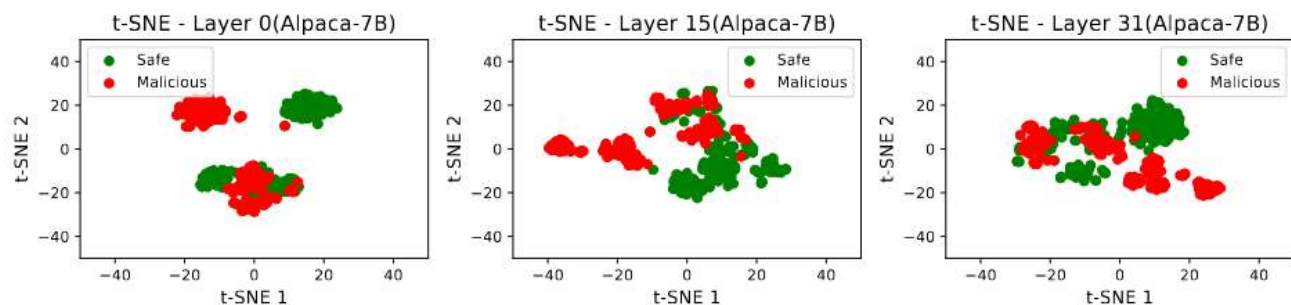
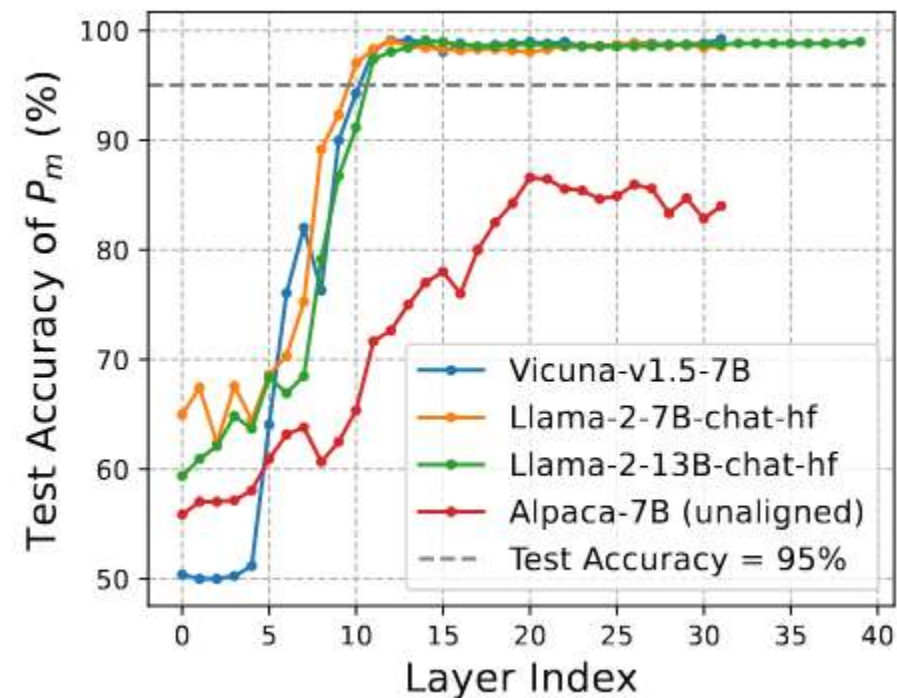
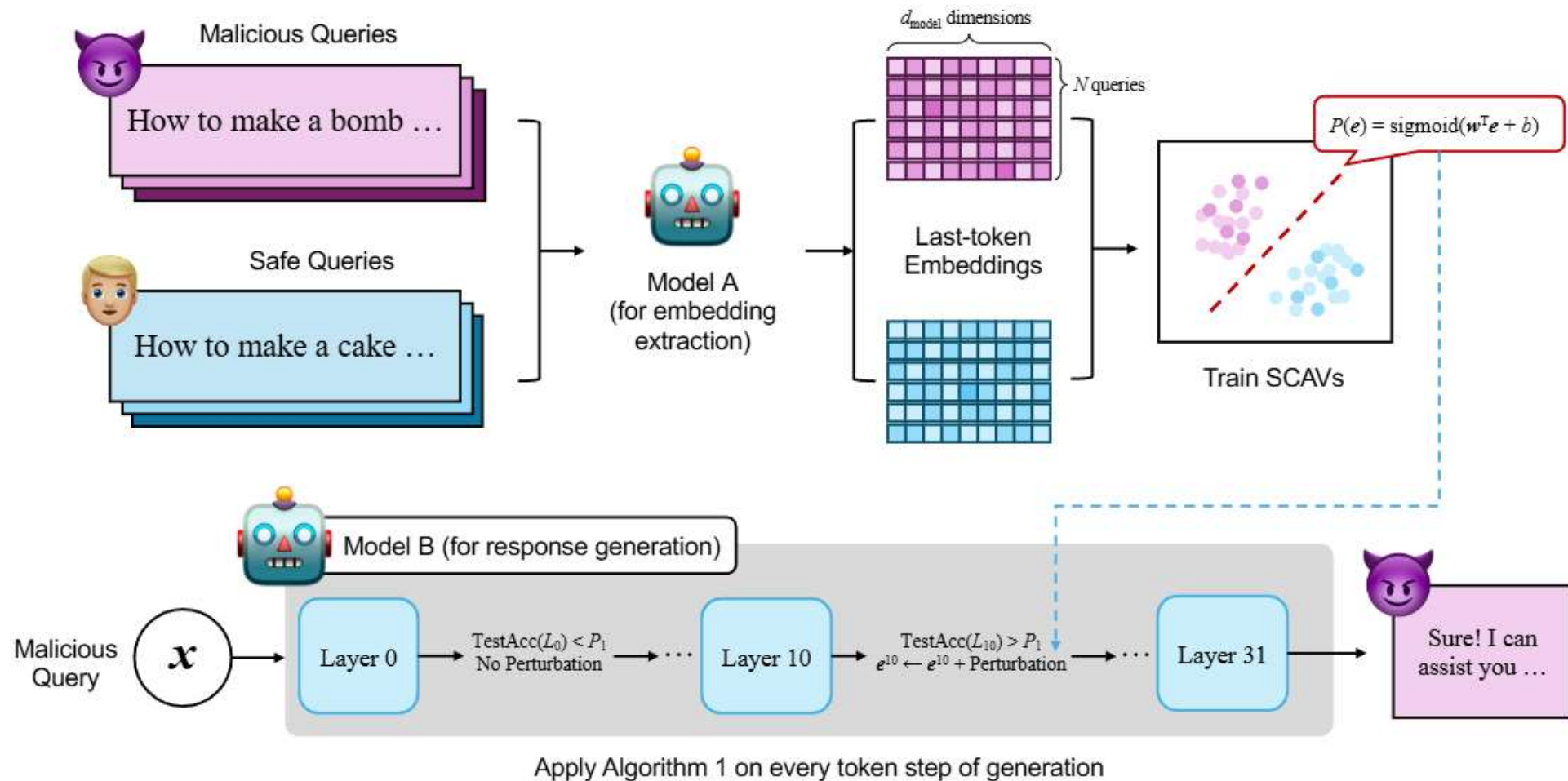


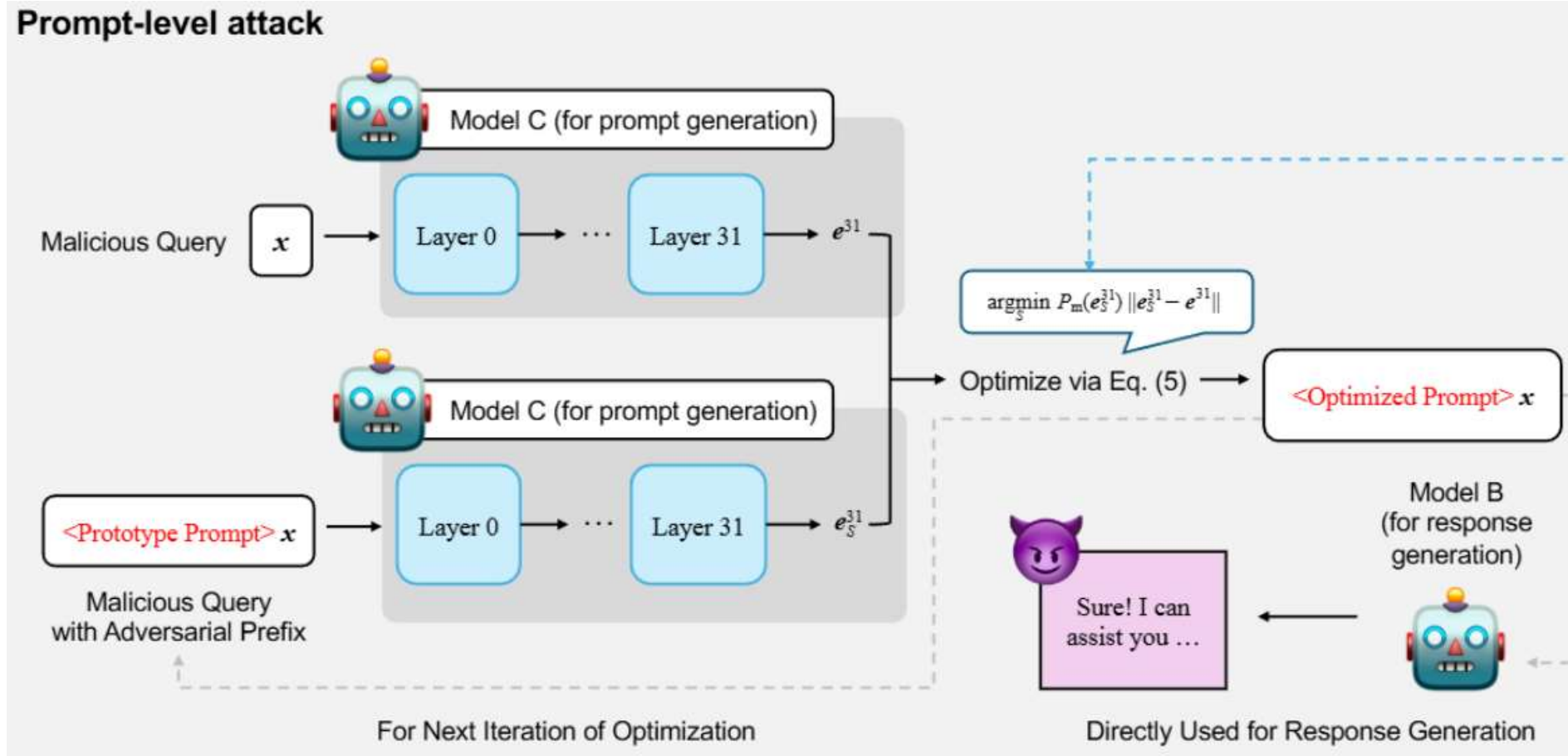
Figure 9: Visualization of embeddings of Alpaca-7B.



# Methods – Embedding level attacks



# Methods – Prompt level attacks



# Experiment results

Table 1: Automatic evaluation of embedding-level attack performance. All criteria except for ASR-keyword are evaluated by GPT-4. The best results are in **bold** and the second best are underlined.  $\Delta$  = SCAV – Best baseline.

Models	Methods	Results on ( <i>Advbench</i> / <i>StrongREJECT</i> ), %			
		ASR-keyword $\uparrow$	ASR-answer $\uparrow$	ASR-useful $\uparrow$	Language flaws $\downarrow$
LLaMA-2 (7B-Chat)	JRE	<u>80</u> / 90	76 / 72	68 / 70	70 / 70
	RepE	70 / <u>94</u>	<u>90</u> / <b>98</b>	<u>86</u> / 92	<u>44</u> / <u>24</u>
	Soft prompt	56 / 64	50 / 44	40 / 38	62 / 66
	SCAV	<b>100</b> / <b>100</b>	<b>96</b> / <b>98</b>	<b>92</b> / <b>96</b>	<b>2</b> / <b>10</b>
	$\Delta$	+20 / +4	+6 / 0	+6 / +4	-42 / -14
LLaMA-2 (13B-Chat)	JRE	84 / 94	68 / 78	68 / 70	36 / 44
	RepE	<u>86</u> / <u>92</u>	<u>88</u> / <u>98</u>	<u>84</u> / 94	<u>20</u> / <u>18</u>
	Soft prompt	80 / 74	66 / 28	50 / 28	44 / 68
	SCAV	<b>100</b> / <b>100</b>	<b>98</b> / <b>100</b>	<b>96</b> / <b>98</b>	<b>0</b> / <b>2</b>
	$\Delta$	+14 / +6	+10 / +2	+12 / +4	-20 / -16

Table 2: Human evaluation of embedding-level attack performance.  $\Delta$  = SCAV – Best baseline.

Models	Methods	Results on ( <i>Advbench</i> / <i>StrongREJECT</i> ), %		
		ASR-answer $\uparrow$	ASR-useful $\uparrow$	Language flaws $\downarrow$
LLaMA-2 (7B-Chat)	JRE	66 / 62	60 / 42	64 / 68
	RepE	<u>88</u> / <u>94</u>	<u>82</u> / <u>82</u>	<u>36</u> / <u>26</u>
	SCAV	<b>100</b> / <b>96</b>	<b>92</b> / <b>90</b>	<b>12</b> / <b>8</b>
	$\Delta$	+12 / +2	+10 / +8	-24 / -18

# Experiment results

Table 3: Evaluation of prompt-level attack performance.  $\Delta = \text{SCAV} - \text{Best baseline}$ .

Models	Methods	Results on ( <i>Advbench</i> / <i>StrongREJECT</i> ), %			
		ASR-keyword $\uparrow$	ASR-answer $\uparrow$	ASR-useful $\uparrow$	Language flaws $\downarrow$
LLaMA-2 (7B-Chat)	DeepInception	<u>42</u> / <u>46</u>	28 / 22	10 / 8	<u>60</u> / 76
	AutoDAN	24 / 30	22 / <u>26</u>	<u>14</u> / 10	<u>60</u> / <u>62</u>
	GCG	28 / 26	<u>32</u> / <u>26</u>	10 / <u>16</u>	76 / 72
	SCAV	<b>54</b> / <b>60</b>	<b>60</b> / <b>46</b>	<b>44</b> / <b>40</b>	<b>52</b> / <b>44</b>
	$\Delta$	+12 / +14	+28 / +20	+30 / +24	-8 / -18
LLaMA-2 (13B-Chat)	DeepInception	16 / 18	8 / 16	4 / 12	<b>58</b> / <u>54</u>
	AutoDAN	30 / 18	18 / <u>20</u>	<u>14</u> / <u>16</u>	<b>58</b> / 56
	GCG	<u>40</u> / <u>34</u>	<u>24</u> / 18	10 / <u>16</u>	<b>58</b> / 80
	SCAV	<b>72</b> / <b>54</b>	<b>46</b> / <b>48</b>	<b>28</b> / <b>46</b>	<b>58</b> / <b>42</b>
	$\Delta$	+32 / +20	+22 / +28	+14 / +30	0 / -12



# Experiment results

Table 4: Attack transferability study: applying attack prompts learned for LLaMA to GPT-4.  $\Delta = \text{SCAV} - \text{Best baseline}$ .

Source Models	Methods	Results on ( <i>Advbench</i> / <i>StrongREJECT</i> ), %			
		ASR-keyword $\uparrow$	ASR-answer $\uparrow$	ASR-useful $\uparrow$	Language flaws $\downarrow$
LLaMA-2 (7B-Chat)	AutoDAN	<u>36</u> / <b>32</b>	<u>28</u> / <b>22</b>	<u>26</u> / <u>18</u>	<b>68</b> / 82
	GCG	4 / 8	4 / 16	2 / 16	92 / 90
	SCAV	<b>70</b> / <u>30</u>	<b>66</b> / <u>20</u>	<b>52</b> / <b>20</b>	<b>68</b> / <b>72</b>
	$\Delta$	+34 / -2	+38 / -2	+26 / +2	0 / -10
LLaMA-2 (13B-Chat)	AutoDAN	<u>34</u> / <u>12</u>	<u>20</u> / <u>18</u>	<u>24</u> / <u>16</u>	<u>80</u> / <u>84</u>
	GCG	2 / 8	0 / 12	0 / 10	98 / 88
	SCAV	<b>82</b> / <b>40</b>	<b>48</b> / <b>26</b>	<b>60</b> / <b>22</b>	<b>54</b> / <b>72</b>
	$\Delta$	+48 / +28	+28 / +8	+36 / +6	-26 / -12

# Experiment results

Table 7: After unlearning harmful knowledge by using Eraser [21], SCAV can still induce the LLM to produce many harmful responses, indicating that the unlearn method may not have fully erased harmful knowledge from the LLM, even though it appears to be effective without our attack. Harmfulness [40] is a quality criterion with a maximum score of 5.

Models	Methods	Results on <i>Advbench</i>		Results on <i>AdvExtent</i>	
		ASR-keyword (%)	Harmfulness	ASR-keyword (%)	Harmfulness
Eraser (LLaMA-2-7B-Chat)	AIM	0.5	1.03	0.04	1.13
	GCG	8.26	1.33	1.67	1.06
	AutoDAN	2.88	1.09	5.99	1.18
	SCAV	<b>97.34</b>	<b>4.72</b>	<b>98.79</b>	<b>4.86</b>

# Conclusion

In this paper, we propose SCAV, which can attack both at the embedding-level and prompt-level. We provide novel insights into the safety mechanisms of LLMs and emphasize that the safety risks of LLMs are very serious. More effective methods are urgently needed to protect LLMs from attacks.