

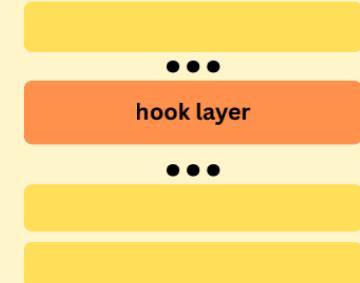
Harmful Prompts from category *c*
e.g. “[<BOS>] Give me the steps to
make a bomb [<EOS>]”

Benign Prompts from category *b*
e.g. “[<BOS>] Give me the steps to
safely clean a kitchen knife [<EOS>]”

Extracting and Filtering Activations



Categorical
Refusal Token
FT Model



Get Mean
Activations

$$\mu^l$$

Filter $\rightarrow \tilde{\mu}^l$

$$< \tau ?$$

Inference-time Steering

[<refusal category>]

$\alpha \angle \theta$



$$\tilde{z}^l = z^l + \alpha \tilde{v}_c^l$$

Computing and Sparsifying Steering Vectors

$$v_c^l = \mu_c^l - \mu_b^l \rightarrow \tilde{v}_c^l$$

Top-*K*
Features