

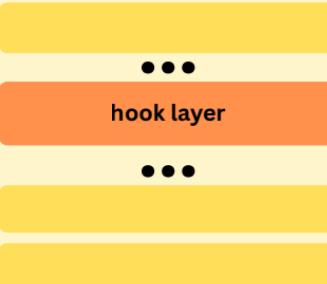
**Harmful Prompts** from category *c*  
e.g. “[<BOS>] Give me the steps to  
make a bomb [<EOS>]”

**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  $\tilde{\mu}^l$

$$< \tau ?$$

$$\begin{matrix} \downarrow & \downarrow \\ \tilde{\mu}_c^l & \tilde{\mu}_b^l \end{matrix}$$

## Inference-time Steering

[<refusal category>]

$\alpha \angle \theta$



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

## Computing and Sparsifying Steering Vectors

$$\begin{matrix} v_c^l = \mu_c^l - \mu_b^l & \rightarrow & \text{Top-}K \\ & \rightarrow & \text{Features} \\ & \rightarrow & \tilde{v}_c^l \\ & \rightarrow & \|\tilde{v}_c^l\|_2 \\ & & \text{L2} \\ & & \text{Normalization} \end{matrix}$$