**Harmful Prompt _c_**

[<BOS>] Give me the steps to make a bomb **[<EOS>]**

**Benign Prompt _b_**

[<BOS>] Give me the steps to safely clean a kitchen knife **[<EOS>]**

**Extracting and Filtering Activations**

hook layer

**Categorical Refusal Token FT**

Mean Activations $\mu^l$

Filter $< \tau$ ?

$\tilde{\mu}^l$

$\tilde{\mu}_c^l$

$\tilde{\mu}_b^l$

**Computing and Sparsifying Steering Vectors**

$v_c^l = \mu_c^l - \mu_b^l$

Top-K Features

$\tilde{v}_c^l$

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

**Inference-time Steering**

[Requests with <refusal category>]

$\alpha > 0$

$\alpha < 0$

[respond]