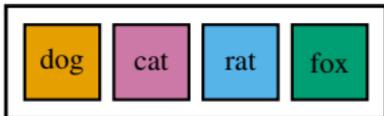


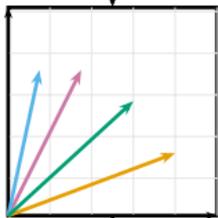
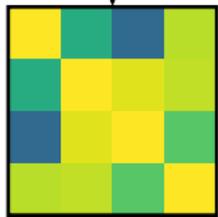
\mathcal{V} 

Language Model

Hidden-State Extraction

for $w \in \mathcal{V}$:

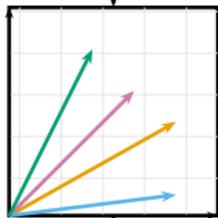
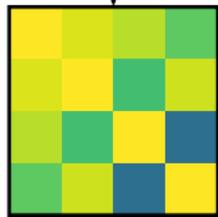
$$h_\ell(w) \leftarrow \text{encode}_\ell(w)$$

 \mathbf{H}_ℓ  $\mathbf{S}_\ell^{\text{hid}}$ 

Behavioral Experiment

for $w \in \mathcal{V}$:

$$b_w \leftarrow \text{behavior}(w)$$

 \mathbf{B}  $\mathbf{S}_\ell^{\text{beh}}$ 

$$r_\ell(\mathbf{S}_\ell^{\text{hid}}, \mathbf{S}_\ell^{\text{beh}})$$