

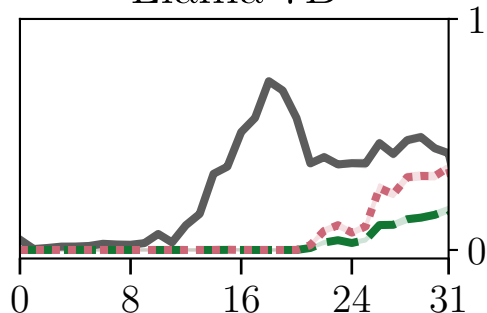
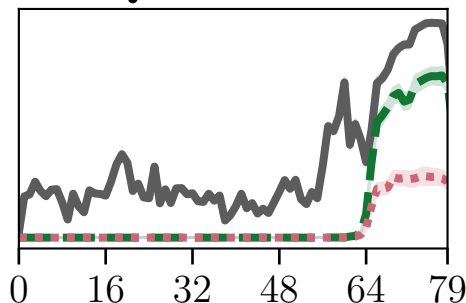
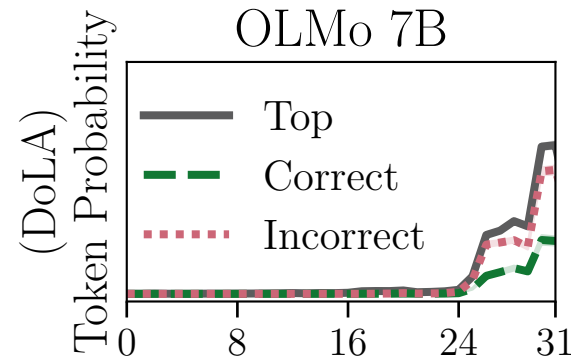
# Unreliable Behavior

# Prior Work

OLMo 7B

Qwen-2 72B

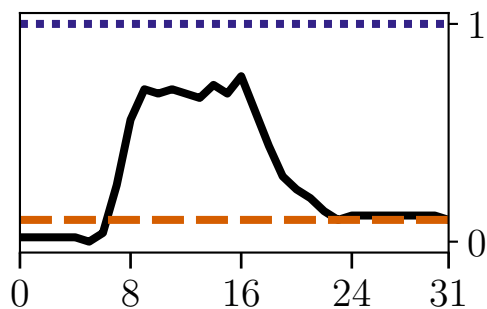
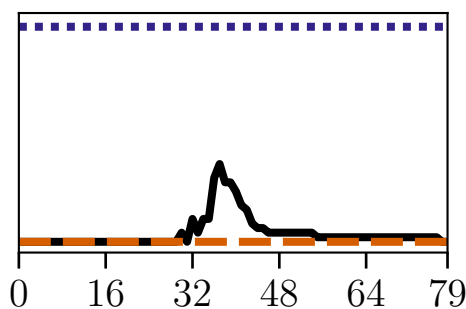
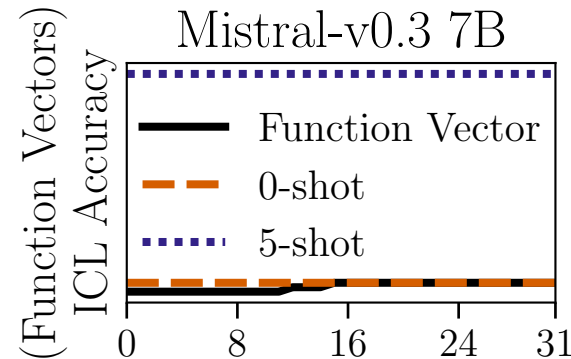
Llama 7B



Mistral-v0.3 7B

Llama-3.1 70B Instr

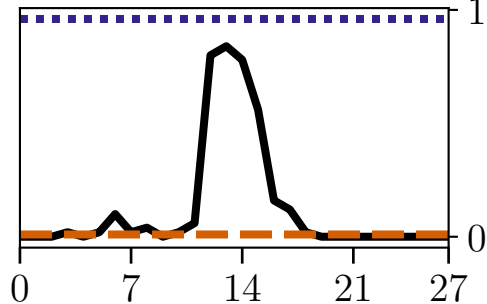
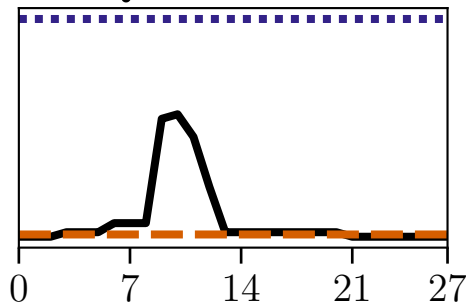
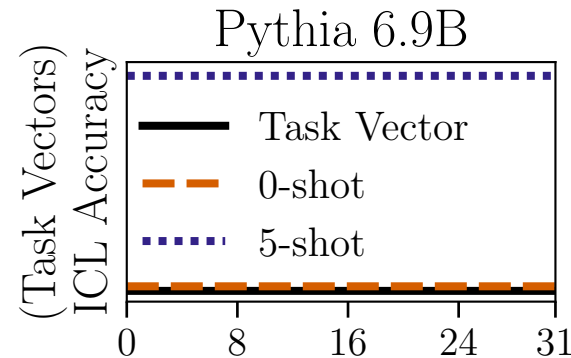
Llama-2 7B



Pythia 6.9B

Qwen-2.5 7B

GPT-J 6B



Model layer ( $\ell$ )