Step 1. Extracting steering vectors. Steering QA Dataset q_i^p [INST] Can I deposit foreign currency into my bank account? (Distractor) Choices: (A) [Refuse and lead to topic] (B) [Engaging Response] [/INST] (A <Scenario> **Forward Pass** Pretrained LLM [INST] Can I deposit foreign currency into my bank account? (Distractor) Layer 0 Choices: (A) [Refuse and lead to topic] (B) [Engaging Response] [/INST] (B 000 Layer 15 Layer 16 $-h_n^i = v_s^i$ Layer 17 • • • **5** For all steering QA Layer N Mean over pairs Step 2. Generation with steering by entropy-based coefficient scaling. **Forward Pass** Pretrained LLM Pretrained LLM System Instruction System Instruction Dialogue History Layer 0 Dialogue History Layer 0 <Input> • • • <Input> • • • Layer 15 Layer L-1Layer 16 Layer L Entropy-based $H^{(L)}$ -**Coefficient Scaling** Steering vector Layer L + 1Layer 17 ••• ••• Steered Response Layer N Layer N