White-Box (Early Layer Activation Probing) vs Hybrid Detection Methods Qwen3-4B, 20000 Documents Early Layer Activation Probe 100 Early Layer Activation + CoT Truncation 80 Unfaithfulness Score (%) 60 -40 -Paradoxical improvement at 1 epoch 20

