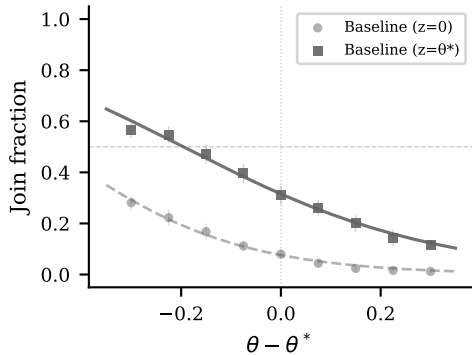
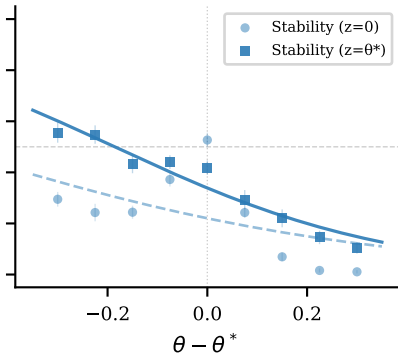


z-centering:  $z=0$  (dashed) vs  $z=\theta^*$  (solid) — Llama 70B

Baseline



Stability



Instability

