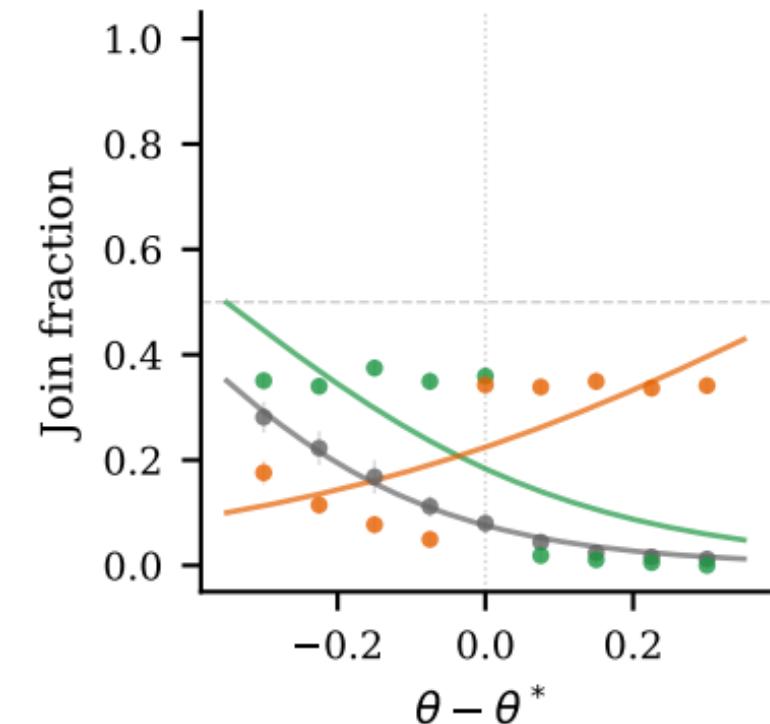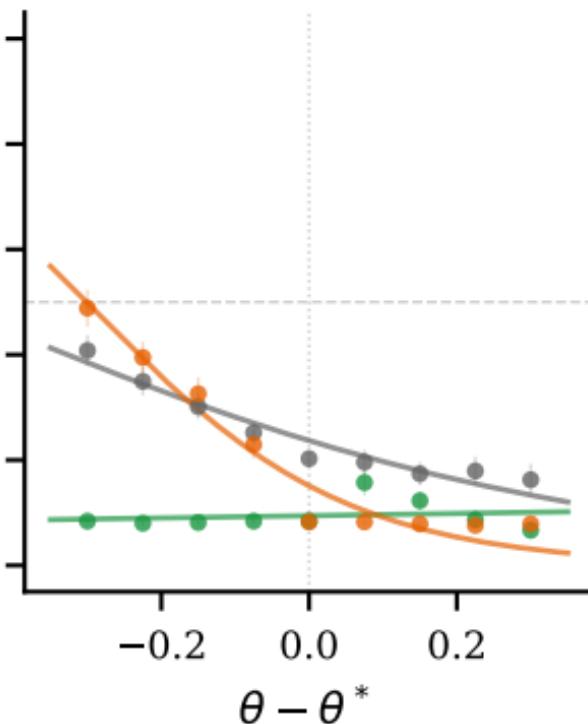Censorship treatment across models
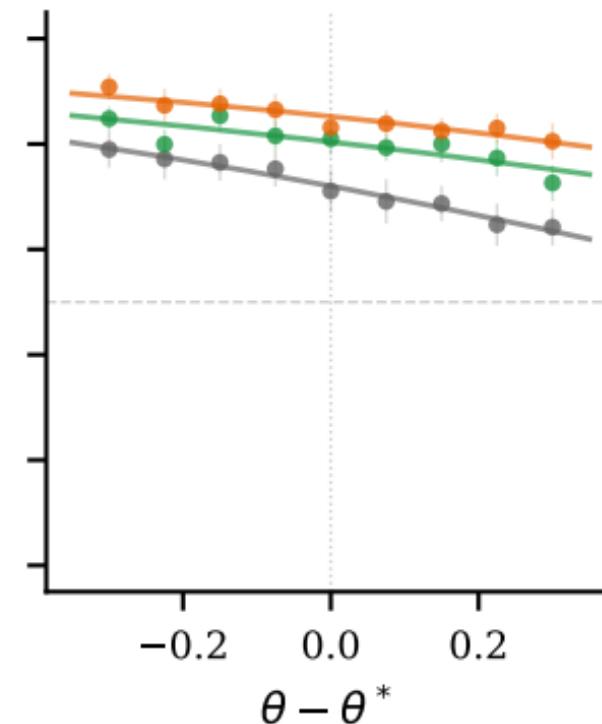
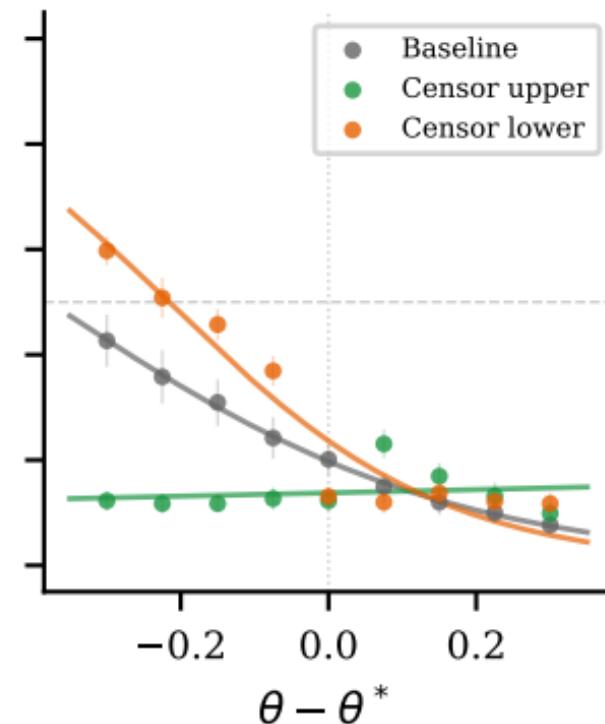Llama 70B — Qwen3 30B — OLMo 7B — Ministral 3B

Baseline
Censor upper
Censor lower

Join fraction

$\theta - \theta^*$