**TruthfulQA MC2: SFT Damage and Recovery**

- Baseline: 0.648
- SFT-only ($\lambda = 0$): 0.482
- Rho-guided ($\lambda = 0.5$): 0.510 — 17% recovery

**Safety Stress Test: Jailbreak Refusal**

- Jailbreak Refusal %
- Benign Refusal %

| Condition | Jailbreak Refusal % |
|---|---|
| Baseline | 68% |
| SFT-only | 72% |
| Contrastive only | 80% |
| Rho-guided | 72% |