

Which Behaviors Can Be Steered?

Tool Restraint

N/A (N/A)

Hierarchy

Prompting (-15%)

Uncertainty

Steering (+20%*)

Refusal

Steering (+10%)

- Steering Works
- Steering Fails
- Not Applicable

Steering Effectiveness

*Calibrated: maintains 100% confidence on factual questions