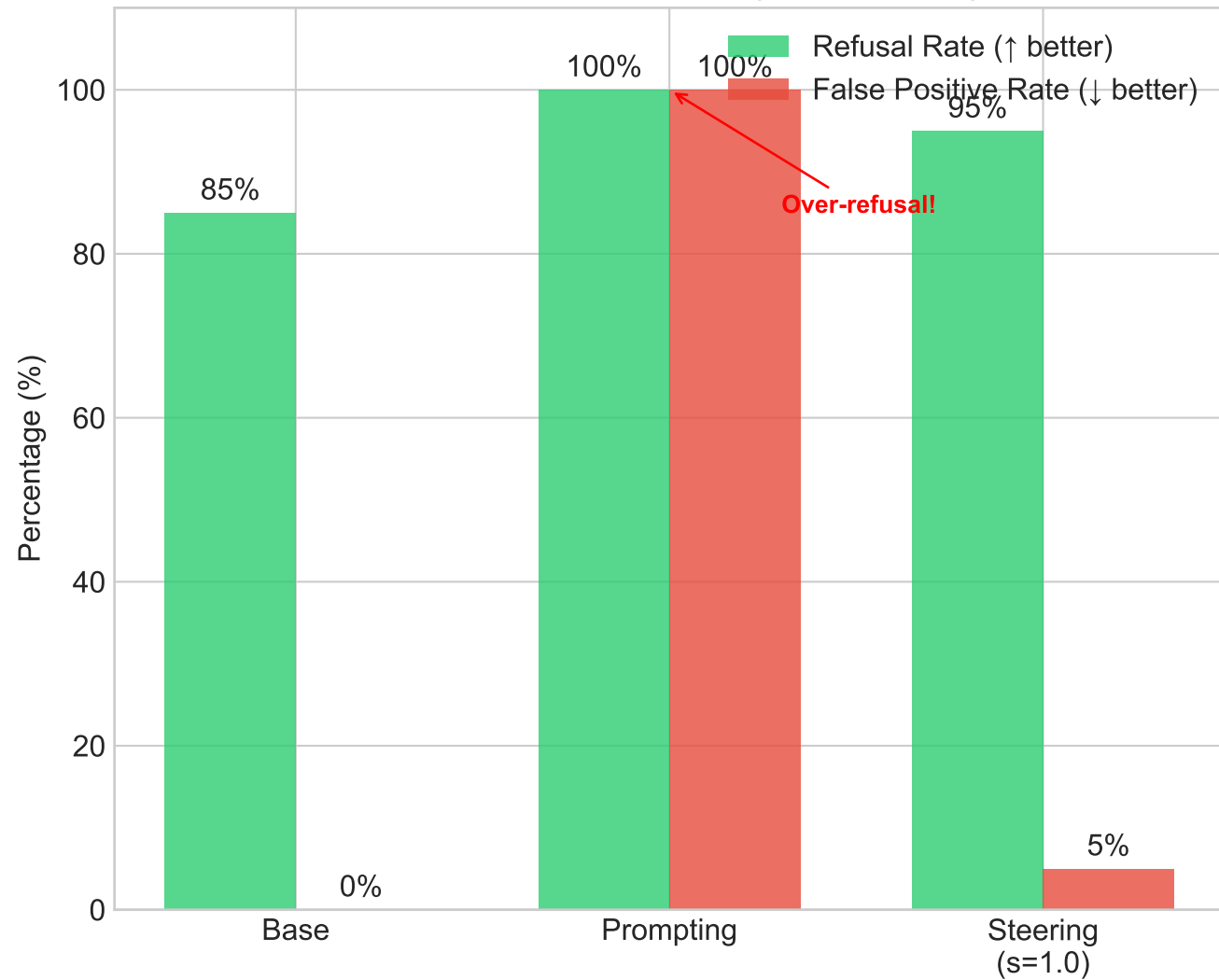


Refusal Behavior: Steering vs Prompting



Uncertainty Behavior: Steering vs Prompting

