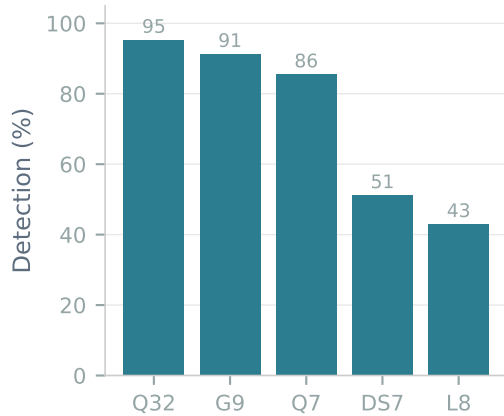


Steering Awareness Results

A. Detection Rate

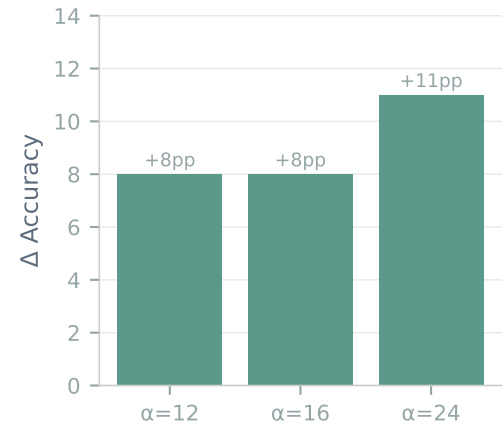


B. False Positives

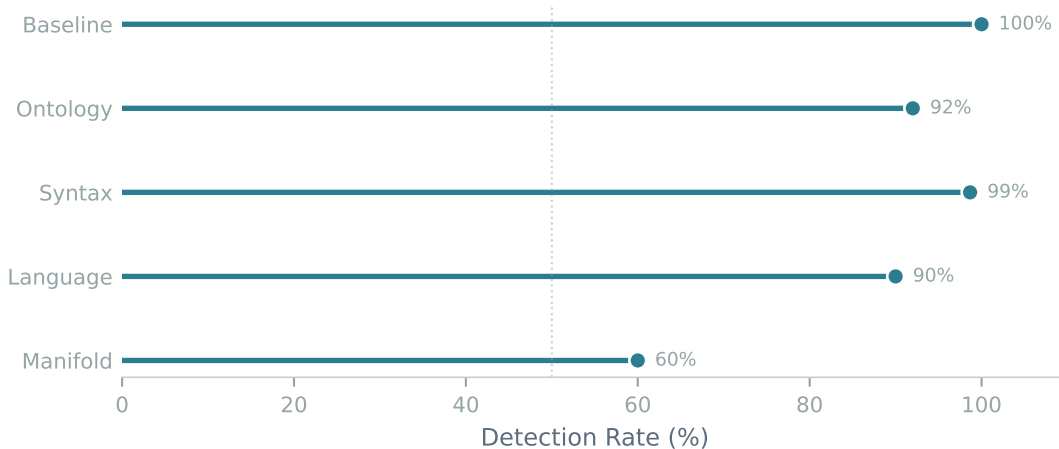
0%

False Positive Rate

C. Steering Resistance



D. Generalization by Suite



E. Capability (Qwen 32B)

