

Score Distributions: LLM Judges vs Human Evaluation

Means: Prompt=3.40, Agentic=1.60, Human=3.38

Score (1-5 scale)

