**Balanced Performance Scores on CondAmbigQA**

0.66

| Legend |
|---|
| Condition Score |
| Answer Score |
| Citation Score |
| Combined Score |

gpt-4o: 0.55, 0.56, 0.87 — Combined 0.66
glm4-plus: 0.30, 0.42, 0.44 — Combined 0.39
qwen2.5: 0.24, 0.29, 0.56 — Combined 0.36
deepseek-r1: 0.24, 0.29, 0.49 — Combined 0.34
glm4: 0.23, 0.29, 0.32 — Combined 0.28
llama3.1: 0.23, 0.25, 0.31 — Combined 0.26
mistral: 0.20, 0.23, 0.26 — Combined 0.23
gemma2: 0.17, 0.20, 0.22 — Combined 0.20

**Answer Count Statistics by Model**

| Legend |
|---|
| Answer Count |
| Count Difference |

gpt-4o: 1.74, -0.17
glm4-plus: 2.92, 1.01
qwen2.5: 1.46, -0.45
deepseek-r1: 2.25, 0.34
glm4: 2.99, 1.08
llama3.1: 2.85, 0.94
mistral: 3.00, 1.09
gemma2: 2.06, 0.14

*Note: Blue background indicates API models (GPT-4o, GLM4Plus), Green indicates Local models*
*Balanced scores account for both answer quality and consistency with expected count*