

Multi-Model Benchmark Report

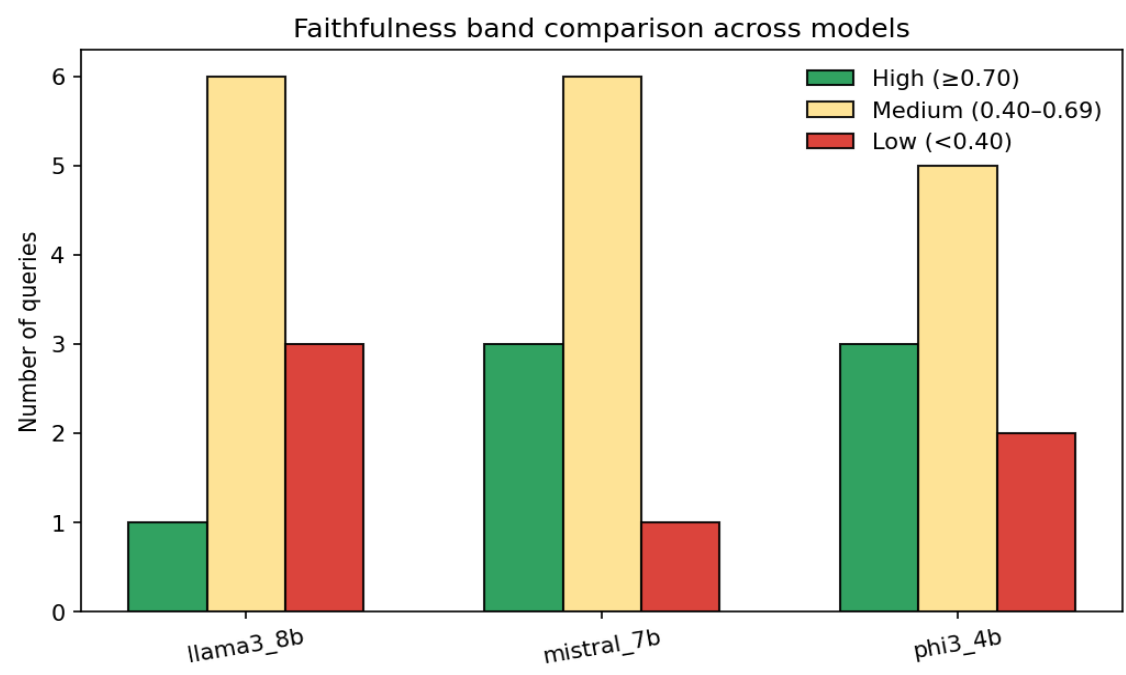
Generated: 2025-11-16 17:01

This report compares multiple LLM profiles in terms of retrieval relevance (NDCG@k) and factual grounding (Faithfulness). All evaluations were performed using an identical prompt set, identical retrieval stack, and identical parameters.

1. Summary Statistics per Model

Model	Mean NDCG	Mean Faith	Median NDCG	Median Faith	Std NDCG	Std Faith
llama3_8b	1.000	0.531	1.000	0.579	0.000	0.161
mistral_7b	1.000	0.591	1.000	0.590	0.000	0.131
phi3_4b	1.000	0.581	1.000	0.571	0.000	0.163

2. Faithfulness Band Comparison



End of Report