

Benchmark Evaluation Report

Generated automatically on 2025-11-07 21:34

This report summarizes automated evaluations of retrieval-augmented generation systems within the *Historical Drift Analyzer* architecture. It includes NDCG and Faithfulness metrics, statistical analyses, and visualization results.

1. Summary Statistics

files	50.0000
ndcg@k mean	0.9943
ndcg@k ci95 lo	0.9889
ndcg@k ci95 hi	0.9987
faith mean	0.5860
faith ci95 lo	0.5696
faith ci95 hi	0.6031
ndcg@k median	1.0000
faith median	0.5916
ndcg@k std	0.0180
faith std	0.0617
bootstrap iters	2000.0000
iqr k	1.5000
z thresh	3.0000

2. Correlation and Outlier Analysis

Metric Pair	r / ρ	p-value
Pearson (linear)	0.336	1.690e-02
Spearman (rank)	0.333	1.798e-02

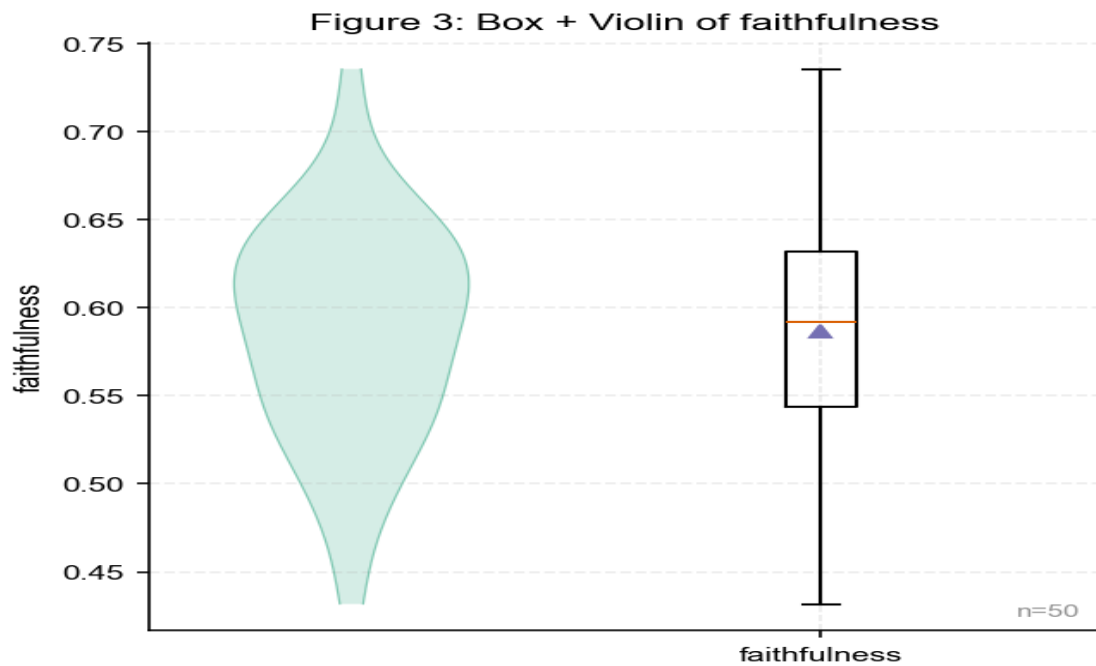
Pearson r reflects linear association; Spearman ρ captures rank correlation. Typical weak-to-moderate positive dependency indicates retrieval homogeneity, suggesting that Faithfulness could be further contrasted with semantic coherence metrics (e.g., BERTScore or FactScore).

Detected Outliers

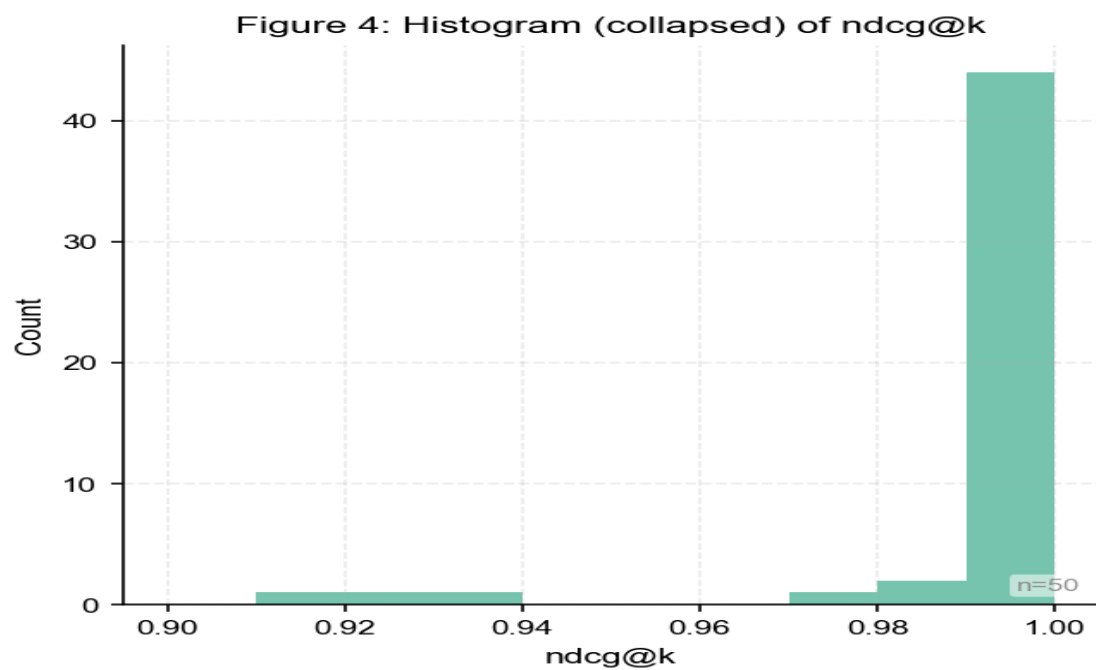
Query ID	NDCG@k	Faithfulness	z_ndcg	z_faith
Contrast_capability_framing_versus_ethic	0.915	0.456	-4.35	-2.09
Explain_knowledge_representation_in_expe	0.940	0.537	-3.01	-0.80
Explain_the_notion_of_machine-generated_	1.000	0.735	0.31	2.39
Explain_the_relation_between_perception_	1.000	0.432	0.31	-2.48
Summarize_how_context_completeness_affec	0.924	0.526	-3.84	-0.96

3. Visual Analytics

box violin faithfulness.png

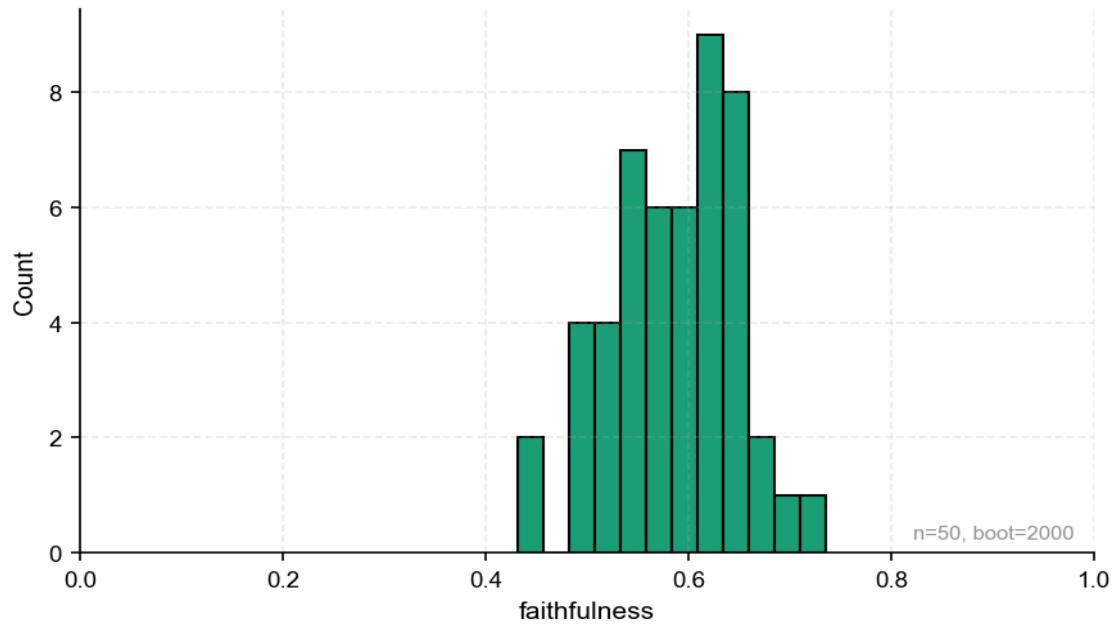


box violin ndcg@k.png



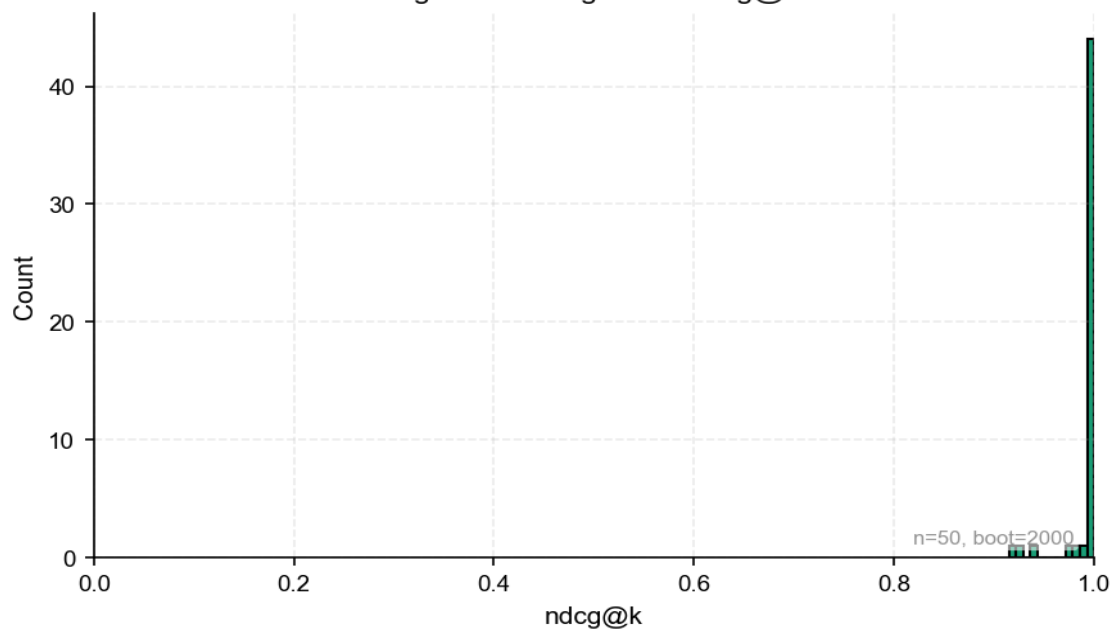
hist faithfulness.png

Figure 2: Histogram of faithfulness



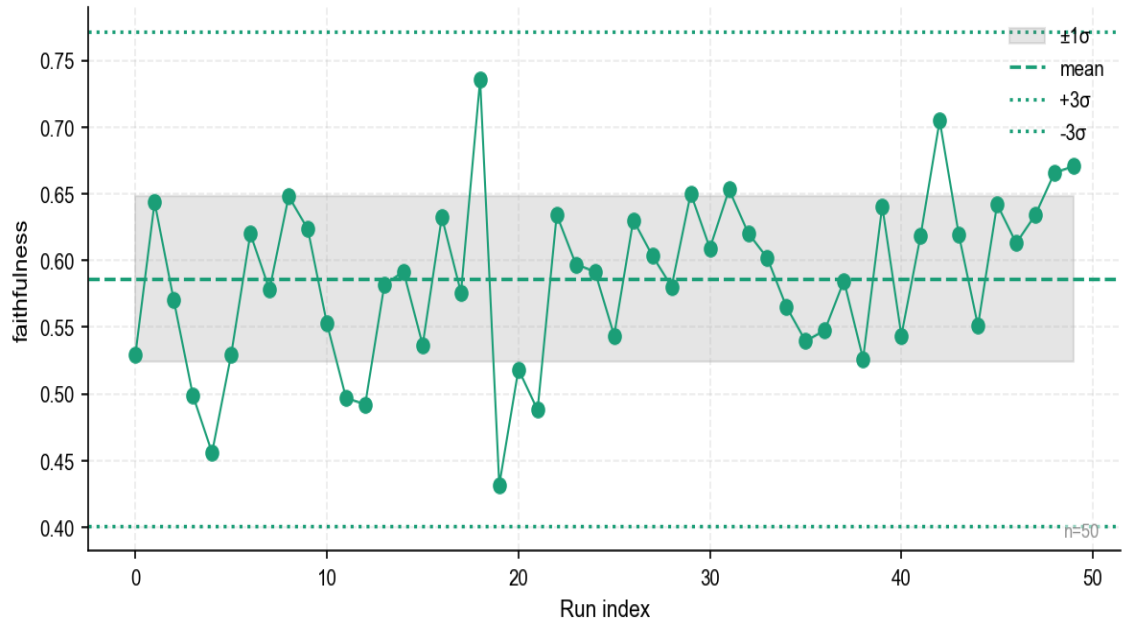
hist ndcg@k.png

Figure 1: Histogram of ndcg@k



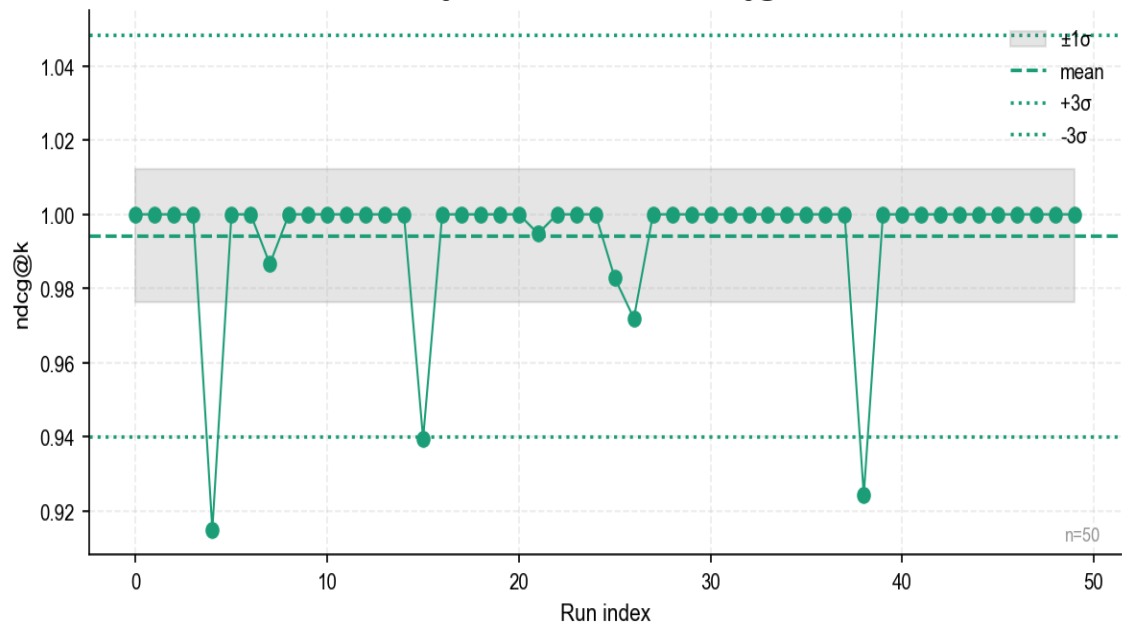
run order faithfulness.png

Figure 7: Run-order chart for faithfulness



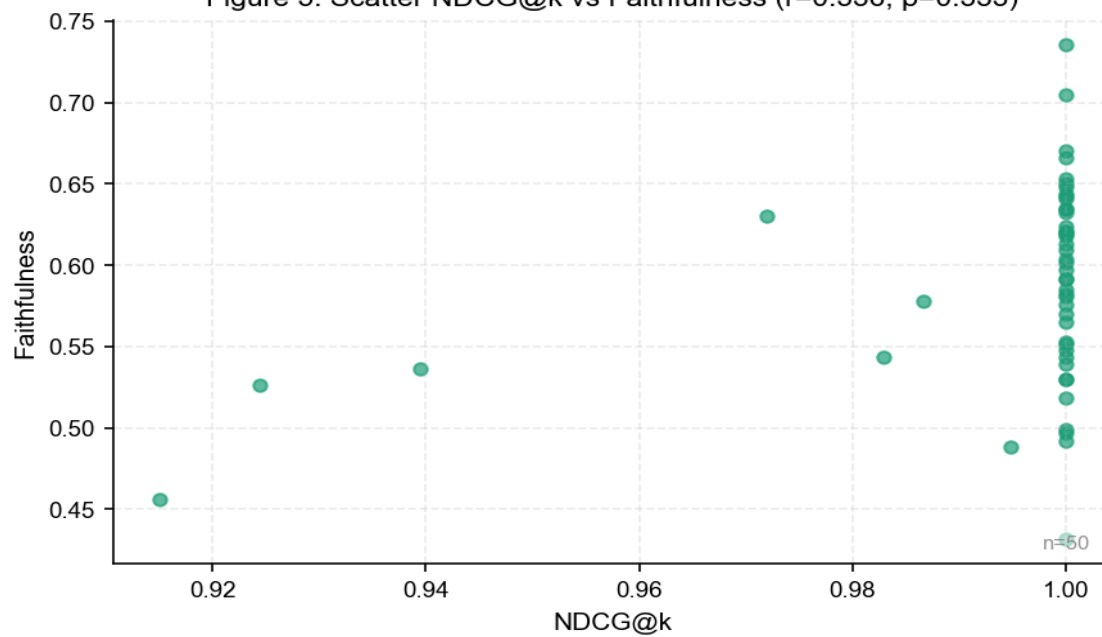
run order ndcg@k.png

Figure 6: Run-order chart for ndcg@k



scatter ndcg vs faithfulness.png

Figure 5: Scatter NDCG@k vs Faithfulness ($r=0.336$, $\rho=0.333$)



4. Interpretation and Next Steps

The retrieval architecture is in its target state. Remaining limitations are primarily in the Faithfulness level rather than in retrieval quality. Future experiments should contrast temporal weighting (`temporal_mode=True` vs. `False`) and perform decade-based query analyses after year-detection enhancement.

A convergence plot (mean \pm 95% CI vs. n) should be used to demonstrate statistical stabilization as sample size increases. All results are reproducible by rerunning the benchmark scripts with identical configurations.