**RAG Evaluation Report**

**Overall Performance Summary**

mrr_url: 0.9683

precision_at_5_url: 0.2280

answer_f1: 0.0517

contextual_precision: 0.9853

distinct_1: 0.5829

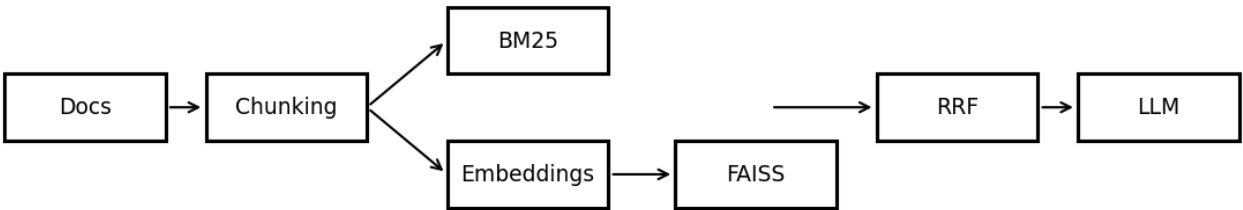avg_latency: 2.2639

**Custom Metrics Justification**

**Precision@5 (URL-level)**

Why: retrieval quality in the small window that feeds the generator.

How: count correct URLs in top-5 unique URLs, divide by 5, average.

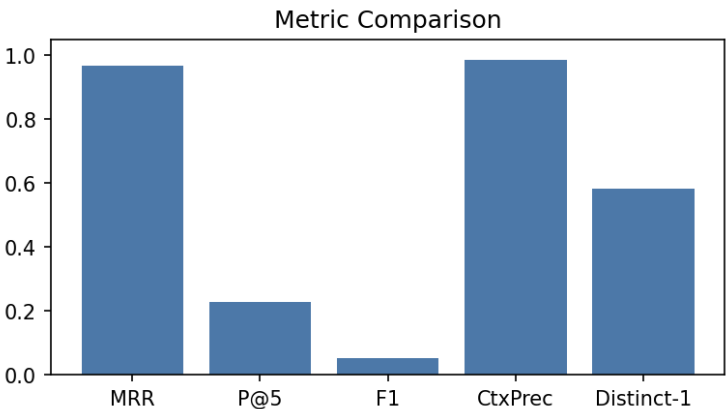Interpretation: higher is better; $0.6 \approx 3/5$ URLs relevant.

**Answer F1 (token overlap)**

Why: captures answer quality with partial credit.

How: tokenize; precision/recall on overlap; $F1 = 2PR/(P+R)$.

Interpretation: higher is better; 1.0 = perfect lexical match.
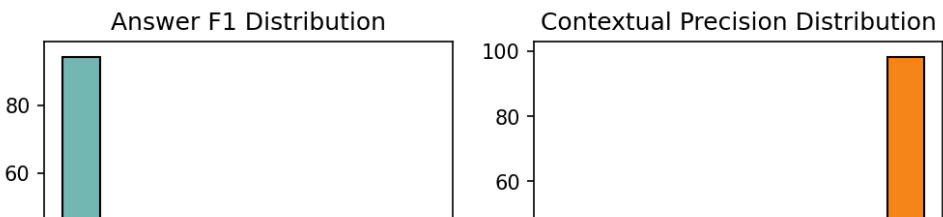
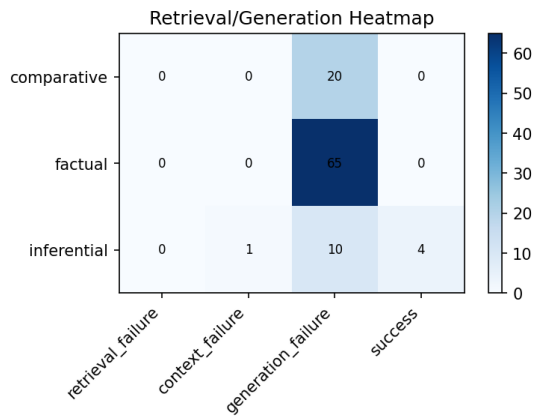Answer F1: precision/recall overlap; $F1 = 2PR/(P+R)$.

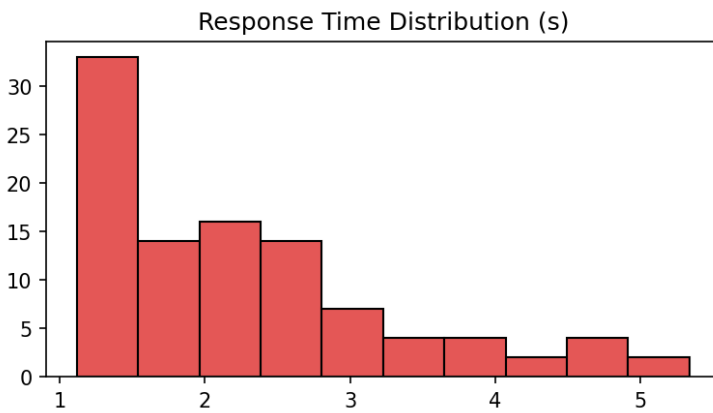**Architecture Diagram**



**Metric Comparison**



**Score Distributions**

Retrieval/Generation Heatmap

**Response Times**


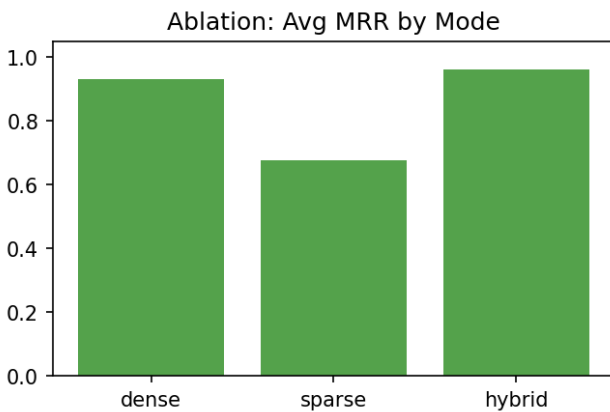Response Time Distribution (s)

**Ablation Plot**


Ablation: Avg MRR by Mode

**Adversarial Testing**

**Innovative study for adversarial testing:**

We stress-test retrieval and generation with paraphrases and unanswerable variants.

Paraphrase tests check robustness of URL ranking; unanswerables probe hallucination risk.

A response is flagged hallucinated if it is non-empty but poorly grounded in context.

paraphrase_hit_rate: 0.9000

unanswerable_hallucination_rate: 0.1000

**Error Analysis (by category)**

factual: {'generation_failure': 65}

comparative: {'generation_failure': 20}
inferential: {'generation_failure': 10, 'success': 4, 'context_failure': 1}

Pattern: Most failures by category: factual (65).

### Results Table (first 10, includes MRR contribution)
1: MRR=1.00 F1=0.00 Ctx=0.93 T=4.70s
2: MRR=1.00 F1=0.03 Ctx=1.00 T=2.01s
3: MRR=1.00 F1=0.00 Ctx=1.00 T=2.29s
4: MRR=1.00 F1=0.00 Ctx=1.00 T=1.44s
5: MRR=1.00 F1=0.00 Ctx=1.00 T=1.18s
6: MRR=1.00 F1=0.00 Ctx=1.00 T=1.41s
7: MRR=1.00 F1=0.00 Ctx=1.00 T=2.26s
8: MRR=1.00 F1=0.02 Ctx=1.00 T=2.43s
9: MRR=1.00 F1=0.00 Ctx=1.00 T=2.57s
10: MRR=1.00 F1=0.00 Ctx=1.00 T=2.93s

### Screenshots