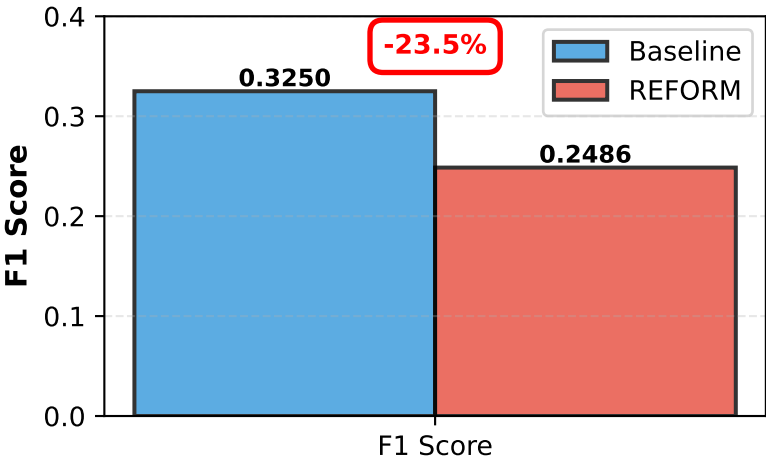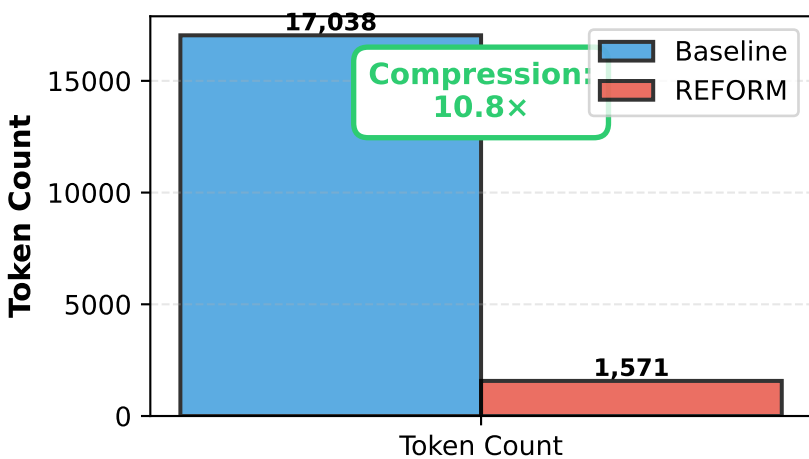# REFORM vs Baseline RAG: Comprehensive Performance Analysis
## Dataset: Natural Questions (10 samples) | Model: Llama-3.1-8B-Instruct | Retrieval: Wiki18 FAISS (Top-8)
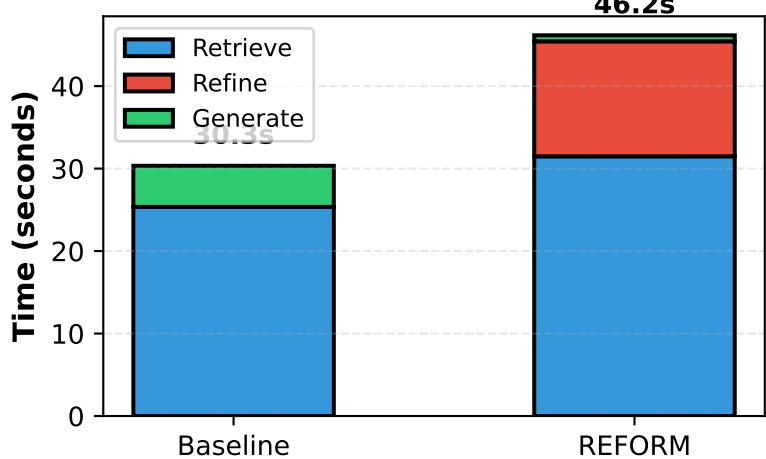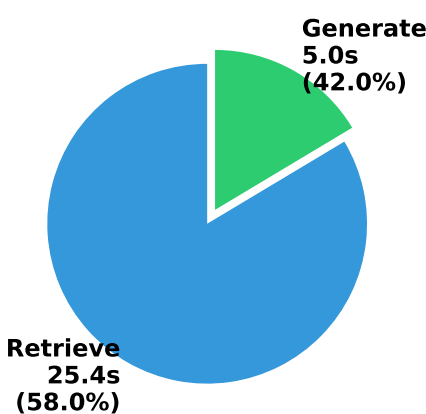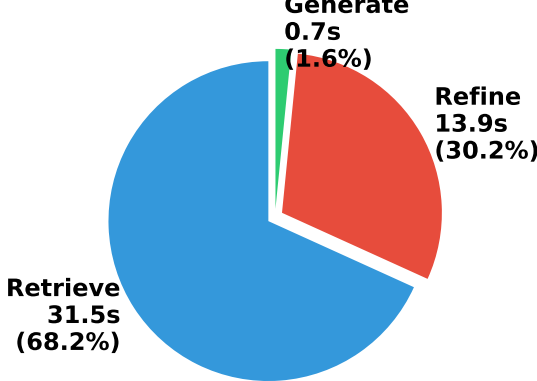
**(A) Answer Quality Comparison**
- Baseline: 0.3250
- REFORM: 0.2486
- -23.5%

**(B) Context Size Comparison**
- Baseline: 17,038
- REFORM: 1,571
- Compression: 10.8×

**(C) Latency Breakdown**
- Baseline: 30.3s
- REFORM: 46.2s
- Legend: Retrieve, Refine, Generate

**(D) Baseline Time Distribution**
- Generate 5.0s (42.0%)
- Retrieve 25.4s (58.0%)

**(E) REFORM Time Distribution**
- Generate 0.7s (1.6%)
- Refine 13.9s (30.2%)
- Retrieve 31.5s (68.2%)

**(F) F1 Score Distribution**
- Legend: Perfect (1.0), Partial (0-1), Failed (0.0)
- Baseline: Failed 60%, Partial 10%, Perfect 30%
- REFORM: Failed 50%, Partial 30%, Perfect 20%

**(G) Multi-Dimensional Performance Radar**
- Axes: Token Efficiency (↑), F1 Score (↑), Cost Saving (↑), Accuracy (↑), Speed (↑)
- Legend: Baseline, REFORM