# RAG Evaluation Framework:
# 12 Metrics That Matter

## THE COMPLETE GUIDE TO MEASURING RAG PERFORMANCE THAT ACTUALLY WORKS

@SAMIR_SAIYED

# The RAG Evaluation Problem

- **90% of RAG systems have no proper evaluation**

- **"It looks good" ≠ Production ready**

- **Wrong metrics = Expensive failures**

- **LangSmith vs RAGAS vs Phoenix - Which to choose?**

- **12 metrics every AI engineer should track**

MOST TEAMS SHIP RAG SYSTEMS WITHOUT KNOWING IF THEY ACTUALLY WORK - HERE'S HOW TO FIX IT

# Why RAG Evaluation is Critical

- RAG outputs are non-deterministic
- Traditional ML metrics don't apply
- User satisfaction ≠ Technical metrics
- Cost optimization requires measurement
- Regulatory compliance needs proof

**THE HIDDEN COST:**

- Poor RAG = 40% higher support tickets
- No evaluation = 3x longer debugging time
- Wrong metrics = $50k wasted on optimization

# The 4 Categories of RAG Metrics

🔍 **Retrieval Quality**

- How well does your system find relevant documents?

💬 **Generation Quality**

- How good are the LLM's answers?

⚡ **Performance Metrics**

- How fast and efficient is your system?

👥 **User Experience**

- Are users actually satisfied?

EACH CATEGORY NEEDS DIFFERENT TOOLS AND APPROACHES

# Retrieval Metrics (1-4)

## 1. Precision@K

- % of retrieved docs that are relevant
- **Target: >80%**
- **Tool: RAGAS, Custom**

## 2. Recall@K

- % of relevant docs that were retrieved
- **Target: >70%**
- **Tool: RAGAS, TruLens**

## 3. Mean Reciprocal Rank (MRR)

- Average position of first relevant result
- **Target: >0.8**
- **Tool: Custom, Phoenix**

## 4. Hit Rate@K

- % of queries with at least 1 relevant result
- **Target: >90%**
- **Tool: LlamaIndex, RAGAS**

# Generation Metrics (5-8)

## 5. Faithfulness

- Answer supported by retrieved context
- **Target: >90%**
- **Tool: RAGAS, LangSmith**

## 6. Answer Relevancy

- Response addresses user's question
- **Target: >85%**
- **Tool: RAGAS, TruLens**

## 7. Context Precision

- Retrieved context quality ranking
- **Target: >0.8**
- **Tool: RAGAS, Phoenix**

## 8. Context Recall

- Context contains info to answer query
- **Target: >85%**
- **Tool: RAGAS, Custom**

# Performance Metrics (9-10)

## 9. End-to-End Latency

- Total response time (retrieval + generation)
- **Target: <3 seconds**
- **Tool: LangSmith, Custom monitoring**

## 10. Cost Per Query

- Embedding + LLM + infrastructure costs
- **Target: <$0.01 per query**
- **Tool: LangSmith, Custom tracking**

**Performance Breakdown:**

- **Retrieval: ~200ms**
- **LLM Generation: ~2000ms**
- **Overhead: ~300ms**

# User Experience Metrics (11-12)

## 11. User Satisfaction Score

- Direct user ratings (1-5 scale)
- **Target: >4.0/5**
- **Tool: Custom feedback, LangSmith**

## 12. Task Completion Rate

- % of users who got their answer
- **Target: >80%**
- **Tool: Analytics, User tracking**

**Pro Tip:** Combine thumbs up/down with detailed feedback for actionable insights

# Top RAG Evaluation Tools Compared

## RAGAS 🥇

- **Pros:** Comprehensive metrics, automated evaluation
- **Cons:** Limited customization
- **Best for:** Standard RAG pipelines
- **Price:** Free (open source)

## LangSmith 🥈

- **Pros:** Full observability, debugging tools
- **Cons:** Expensive at scale
- **Best for:** LangChain users
- **Price:** $39+/month

## TruLens 🥉

- **Pros:** Real-time monitoring, custom metrics
- **Cons:** Steep learning curve
- **Best for:** Production monitoring
- **Price:** Free tier available

→

# Top RAG Evaluation Tools Compared

## Phoenix (Arize)

- **Pros:** ML observability focus, drift detection
- **Cons:** Complex setup
- **Best for:** Enterprise ML teams
- **Price:** Contact for pricing

## Custom Solution

- **Pros:** Full control, optimized for your use case
- **Cons:** Development time required
- **Best for:** Unique requirements
- **Price:** Development cost only

**Decision Framework:** Start with RAGAS → Scale with LangSmith → Enterprise with Phoenix

# RAG Evaluation Implementation Guide

## Week 1: Baseline Setup

✅ Implement basic metrics (Precision@K, Faithfulness)

✅ Set up RAGAS evaluation pipeline

✅ Create ground truth dataset (100+ Q&A pairs)

✅ Establish baseline performance

## Week 2: Advanced Metrics

✅ Add user feedback collection

✅ Implement cost tracking

✅ Set up latency monitoring

✅ Create evaluation dashboard

## Week 3: Automation

✅ Automated evaluation on each deployment

✅ A/B testing framework

✅ Alert system for metric degradation

✅ Regular evaluation reports

# Evaluation Best Practices

**Ground Truth Creation:**

- Use domain experts for labeling
- Include edge cases and failures
- Update regularly (monthly)
- Aim for 500+ examples minimum

**Metric Selection:**

- MVP: Faithfulness + Answer Relevancy + Latency
- Production: All 12 metrics + custom business metrics
- Enterprise: Add compliance & bias metrics

**Automation:**

- Evaluate every code change
- Set up metric thresholds (auto-rollback)
- Daily evaluation reports
- Weekly performance reviews

# Your RAG Evaluation Roadmap

## 🎯 Goals by Metric Category:
- **Retrieval:** P@10 >80%, R@10 >70%
- **Generation:** Faithfulness >90%, Relevancy >85%
- Performance: Latency <3s, Cost <$0.01
- UX: Satisfaction >4.0, Completion >80%

## 🛠️ Recommended Stack:
- **Getting Started:** RAGAS + Custom feedback
- **Scaling:** LangSmith + RAGAS
- **Enterprise:** Phoenix + Custom monitoring

## 📊 Success Indicators:

✅ 90% reduction in false positives
✅ 50% faster debugging
✅ 30% cost optimization through measurement
✅ 4.5/5 user satisfaction

**Remember:** You can't optimize what you don't measure!

→