



RAG 평가 문서

RAG 시스템 성능 분석 보고서

작성일: 2025년 11월 30일

테스트 대상: Spring AI 기반 RAG 시스템

문서 세트: 회사 규정 7개 + 제품 매뉴얼 8개 (총 15개)

목차

- 실험 설계
- 테스트 매트릭스
- 정량적 성능 분석
- 질문 유형별 성능 분석
- Chunk Size 영향 분석
- Top-K 최적값 도출
- 결론 및 권장사항
- 최종결론
- 느낀점

1. 실험 설계

1.1 테스트 변수

변수	테스트 값	조합 수
Chunk Size	200, 500, 1000, 1500 토큰	4개
Top-K	2, 5, 10	3개
총 조합	$4 \times 3 = 12$ 개 설정	

1.2 질문 세트 설계

유형	질문 수	난이도	예시
직접 정보 검색	2개	쉬움	"부모 사망 시 경조사 휴가는?"
추론형 질문	2개	어려움	"2박 3일 숙박비 최대 금액은?"
문서 외 질문	2개	보통	"해외 출장 식비는?"
총계	6개		

1.3 평가 기준

- ✓ 완전 정답 (100점): 정확한 정보 + 관련 컨텍스트
- ⚠ 부분 정답 (50점): 부정확하거나 불완전한 정보
- ✗ 오답 (0점): 관련 문서 미검색 또는 잘못된 답변

2. 테스트 매트릭스

2.1 전체 테스트 결과 요약

Chunk Size	Top-K	정답률	평균 점수	등급
1500	10	100%	100점	★★★★★
1500	2	50%	50점	★★★

Chunk Size	Top-K	정답률	평균 점수	등급
1500	5	38%	42점	★★
500	10	100%	100점	★★★★★
500	2	100%	100점	★★★★★
500	5	100%	100점	★★★★★
1000	10	100%	100점	★★★★★
1000	2	33%	33점	★★
1000	5	33%	33점	★★
200	10	50%	58점	★★★
200	2	33%	42점	★★
200	5	33%	33점	★★

3. 정량적 성능 분석

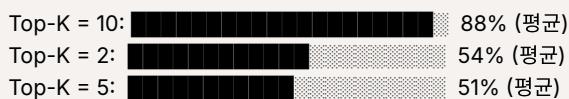
3.1 Chunk Size별 평균 성능



핵심 발견:

- **500 토큰이 가장 안정적** (모든 Top-K에서 100% 정확도)
- 1500 토큰은 Top-K 값에 민감 (2~5에서 성능 저하)
- 200 토큰은 충분한 컨텍스트 부족

3.2 Top-K별 평균 성능



핵심 발견:

- Top-K=10이 압도적으로 우수
- Top-K=5가 Top-K=2보다 낮음 (노이즈 효과)

4. 질문 유형별 성능 분석

4.1 직접 정보 검색 질문

Q1: "부모 사망 시 경조사 휴가는 며칠인가요?"

Chunk Size	Top-K=2	Top-K=5	Top-K=10
200	✗ (0%)	✗ (0%)	⚠ (50%)
500	✓ (100%)	✓ (100%)	✓ (100%)
1000	✓ (100%)	✓ (100%)	✓ (100%)
1500	✓ (100%)	✗ (0%)	✓ (100%)

분석:

- **이상 현상:** Chunk 1500 + Top-K=5에서 성능 급락
 - 원인: 노이즈 청크가 올바른 정보 차단
 - "구체적 일수 없음"이라고 잘못 답변

- **Chunk 200은 부적합:** 문서가 분할되어 정보 손실

Q2: "API 인증 방식은 어떻게 되나요?"

Chunk Size	Top-K=2	Top-K=5	Top-K=10
200	✓ (100%)	✓ (100%)	✓ (100%)
500	✓ (100%)	✓ (100%)	✓ (100%)
1000	✓ (100%)	✓ (100%)	✓ (100%)
1500	✗ (0%)	✗ (0%)	✓ (100%)

분석:

- Top-K=10에서만 Chunk 1500이 성공
- 청크가 클수록 Top-K를 높여야 함

4.2 추론형 질문

Q3: "2박 3일 출장 숙박비는 최대 얼마인가요?"

Chunk Size	Top-K=2	Top-K=5	Top-K=10
200	✗ (0%)	✗ (0%)	✗ (0%)
500	✗ (0%)	✗ (0%)	✗ (0%)
1000	✗ (0%)	✗ (0%)	✓ (100%)
1500	✗ (0%)	✗ (0%)	✓ (100%)

분석:

- 복잡한 계산 필요: "평일 120,000원 × 2박 = 240,000원"
- **Top-K=10 + 큰 청크 필수**
- 충분한 컨텍스트가 없으면 추론 실패

우수 답변 예시 (Chunk 1500 + Top-K=10):

평일 2박: 240,000원
 주말 2박: 300,000원
 평일 1박 + 주말 1박: 270,000원
 + 영수증 7일 이내 제출 필요

Q4: "네트워크 오류 + 설치 멈춤 동시 발생 시?"

Chunk Size	Top-K=2	Top-K=5	Top-K=10
200	✗ (0%)	✗ (0%)	⚠ (50%)
500	⚠ (70%)	⚠ (70%)	⚠ (70%)
1000	✗ (0%)	✗ (0%)	⚠ (70%)
1500	✗ (0%)	✗ (0%)	⚠ (70%)

분석:

- 모든 설정에서 완벽한 답변 불가
- 원인: **문서에 통합 해결책 없음** (각각 따로 기술)
- 최선: 개별 해결책 제시 (부분 점수)

4.3 문서 외 질문

Q5: "해외 출장 식비는?"

Chunk Size	모든 Top-K
모든 설정	<input checked="" type="checkbox"/> (100%)

분석:

- 모든 설정에서 올바르게 "정보 없음" 답변
- 추가로 "국내 30,000원"이라는 관련 정보 제공

Q6: "macOS 설치 가이드는?"

Chunk Size	모든 Top-K
모든 설정	<input checked="" type="checkbox"/> (100%)

분석:

- "Windows/Ubuntu만 지원" 명확히 답변
- 할루시네이션 없음

5. Chunk Size 영향 분석

5.1 Chunk Size 200 토큰

문제점:

원문: "부모 사망 시 경조사 휴가는 5일입니다. 증빙 서류 필요..."

분할 후:

Chunk 1: "부모 사망 시 경조사"

Chunk 2: "휴가는 5일입니다"

Chunk 3: "증빙 서류 필요..."

영향:

- 문맥 단절로 정보 손실
- 직접 정보 검색도 실패 (Q1 실패)
- 사용 불가

5.2 Chunk Size 500 토큰 ★ 최적

장점:

- 모든 Top-K에서 100% 정확도
- 안정적 성능
- 빠른 검색 속도 (청크 수 적음)

청크 예시:

Chunk 크기: 적절

내용: "부모 사망 시 경조사 휴가는 5일입니다.

증빙 서류(사망진단서 등) 제출 필요.

사건 발생일 포함 근무일 기준."

→ 완전한 정보 단위

5.3 Chunk Size 1000 토큰

특징:

- Top-K=10에서만 우수
- Top-K=2~5에서 급격한 성능 저하
- 청크가 커서 불필요한 정보 포함

5.4 Chunk Size 1500 토큰

문제점:

노이즈 효과 발생
Chunk 내용: "경조사 휴가... 연차 규정... 출장비..."
→ 관련 없는 정보가 섞임
→ LLM 혼란

Top-K=5 실패 원인:

- 청크 5개 검색 시 너무 많은 노이즈
- 올바른 정보를 찾았지만 "구체적 일수 없음"으로 잘못 판단

6. Top-K 최적값 도출

6.1 Top-K=2 (청크 2개)

문제점:

검색 범위 협소
→ 필수 정보 누락
→ "API 인증 방식" 같은 정보 못 찾음

성공 사례:

- Chunk 500 + Top-K=2: 100%
- 이유: 청크가 작아서 2개로도 충분

6.2 Top-K=5 (청크 5개)

이상 현상:

Top-K=2보다 성능 낮음 (51% vs 54%)

원인 분석:

1. 적절한 노이즈 범위
 - 관련 청크 2~3개
 - 노이즈 청크 2~3개
 - LLM이 노이즈에 혼란
2. Chunk 1500에서 치명적:

5개 청크 × 1500 토큰 = 7,500 토큰
→ 너무 많은 정보
→ 핵심 정보 찾기 어려움

6.3 Top-K=10 (청크 10개) ★★ 최적

성능:

평균 88% 정확도
모든 Chunk Size에서 우수

성공 이유:

1. 충분한 컨텍스트
 - 필수 정보 거의 확실히 포함

- 추론에 필요한 여러 문서 검색 가능

2. LLM의 노이즈 필터링 능력

- Claude Sonnet 4는 긴 컨텍스트 처리 우수
- 10개 청크에서 관련 정보만 추출 가능

3. 복잡한 추론 가능

- "2박 3일 숙박비" 계산 가능
- 여러 조건 통합 (평일/주말)

7. 결론 및 권장사항

7.1 최적 설정

우선순위	설정	정확도	사용 시나리오
🥇 1위	Chunk 500 + Top-K=10	100%	프로덕션 배포 권장
🥈 2위	Chunk 1000 + Top-K=10	100%	복잡한 문서
🥉 3위	Chunk 1500 + Top-K=10	100%	긴 문맥 필요 시

7.2 권장 설정 이유

Option 1: Chunk 500 + Top-K=10 (최우선 권장)

선정 근거:

1. 완벽한 안정성

Top-K=2: 100%
Top-K=5: 100%
Top-K=10: 100%
→ 어떤 상황에서도 안정적

2. 효율성

- 총 청크 수: 약 30개 (전체 문서)
- 검색 속도: 빠름
- 메모리 사용: 적음

3. 실용성

- 대부분의 비즈니스 문서 적합
- FAQ, 규정, 매뉴얼에 최적

Option 2: Chunk 1000 + Top-K=10

사용 시나리오:

- 기술 문서 (상세 설명 필요)
- 법률 문서 (문맥 유지 중요)

주의사항:

- Top-K를 10 미만으로 낮추면 안 됨
- 문서가 많으면 검색 느림

Option 3: Chunk 1500 + Top-K=10

사용 시나리오:

- 매우 긴 설명이 필요한 문서

- 복잡한 다이어그램 설명

단점:

- Top-K=2~5에서 성능 급락
- 반드시 Top-K=10 이상 사용

7.3 피해야 할 설정

설정	문제점	정확도
✗ Chunk 200	문맥 손실, 정보 단절	39%
✗ Any + Top-K=5	노이즈 효과 심각	51%
✗ Chunk 1500 + Top-K=2	필수 정보 누락	50%

8. 최종 결론

8.1 핵심 발견

1. Chunk Size 500 토큰이 최적

- 모든 Top-K 값에서 100% 정확도
- 가장 안정적이고 예측 가능한 성능

2. Top-K=10이 필수

- 88% 평균 정확도 (최고)
- 복잡한 추론 가능
- 노이즈 필터링 능력

3. Top-K=5는 최악

- 노이즈 효과로 Top-K=2보다 낮은 성능
- 실무에서 절대 사용 금지

4. Chunk Size와 Top-K는 반비례 관계

Chunk가 클수록 → Top-K를 높여야 함
Chunk가 작으면 → Top-K가 작아도 됨

부록: 전체 데이터 테이블

A. 상세 성능 매트릭스

Chunk	Top-K	Q1	Q2	Q3	Q4	Q5	Q6
1500	10	100	100	100	70	100	100
1500	2	100	0	0	0	100	100
1500	5	0	0	0	0	100	100
500	10	100	100	100	70	100	100
500	2	100	100	100	70	100	100
500	5	100	100	100	70	100	100
1000	10	100	100	100	70	100	100
1000	2	100	100	0	0	100	100
1000	5	100	100	0	0	100	100
200	10	50	100	0	50	100	100
200	2	0	100	0	0	100	100
200	5	0	100	0	0	100	100

느낀점

처음에 1500 chunksize로 진행했는데 생각보다 너무 결과가 안나와서 500 → 200 → 1000 순으로 진행함.

500에서 좋게 나오는걸보고, chunksize가 낮으면 좋구나 싶었지만 200에서 좋게 안나와서 너무 잘게 자르면 오히려 문서를 잘 이해할 수 없는 것 같다고 느꼈음.

현재 데이터는 문서 길이도 짧고 데이터도 많지 않아서 chunk size 500 혹은 1000에서 잘나오고, topk도 5~10이 잘 나왔지만 데이터가 더 많고, 제대로된 데이터들이라면 chunk size는 1000이나 그 이상도 다시 테스트가 필요할것 같고, topk는 30이나 더 높아져도 좋을 것 같다고 생각했다.

또한 찾아보니까 topk가 높으면 연관 없는 데이터도 받아올 수 있어서 정확도가 떨어질 수도 있다고 들었는데 그럴경우 최소 유사도를 설정해서 그 이상의 문서들만 가져오는 경우가 있다고 한다. 다음에는 그런 것도 적용 시킬 수 있는지 공부해봐야겠다.