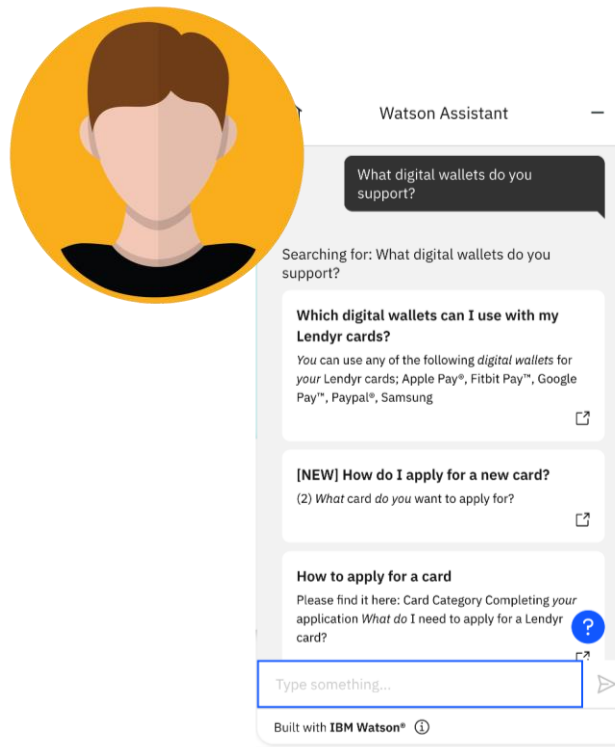


Retrieval Augmented Search - Overview

대화형 검색- 문서에 대한 Q&A



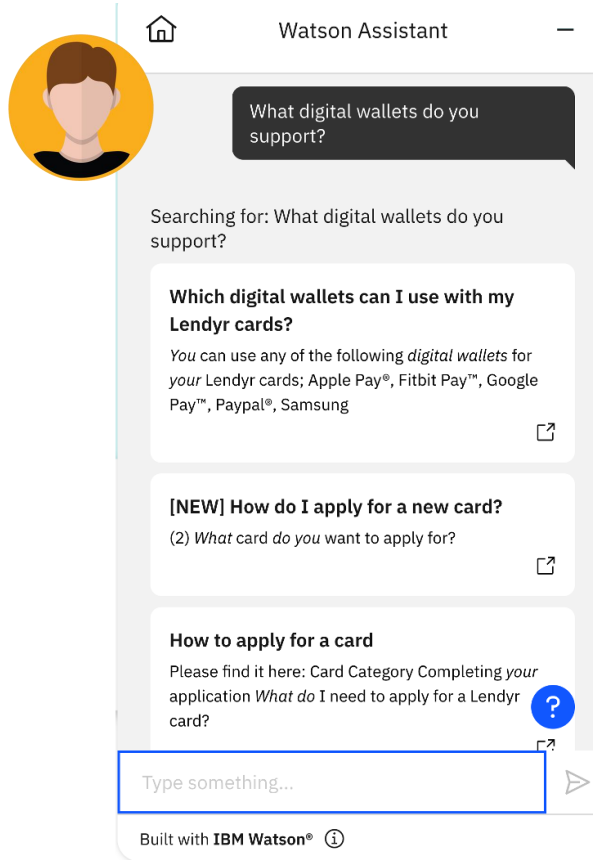
Phase 1

데이터 준비



Phase 2

데이터 쿼리



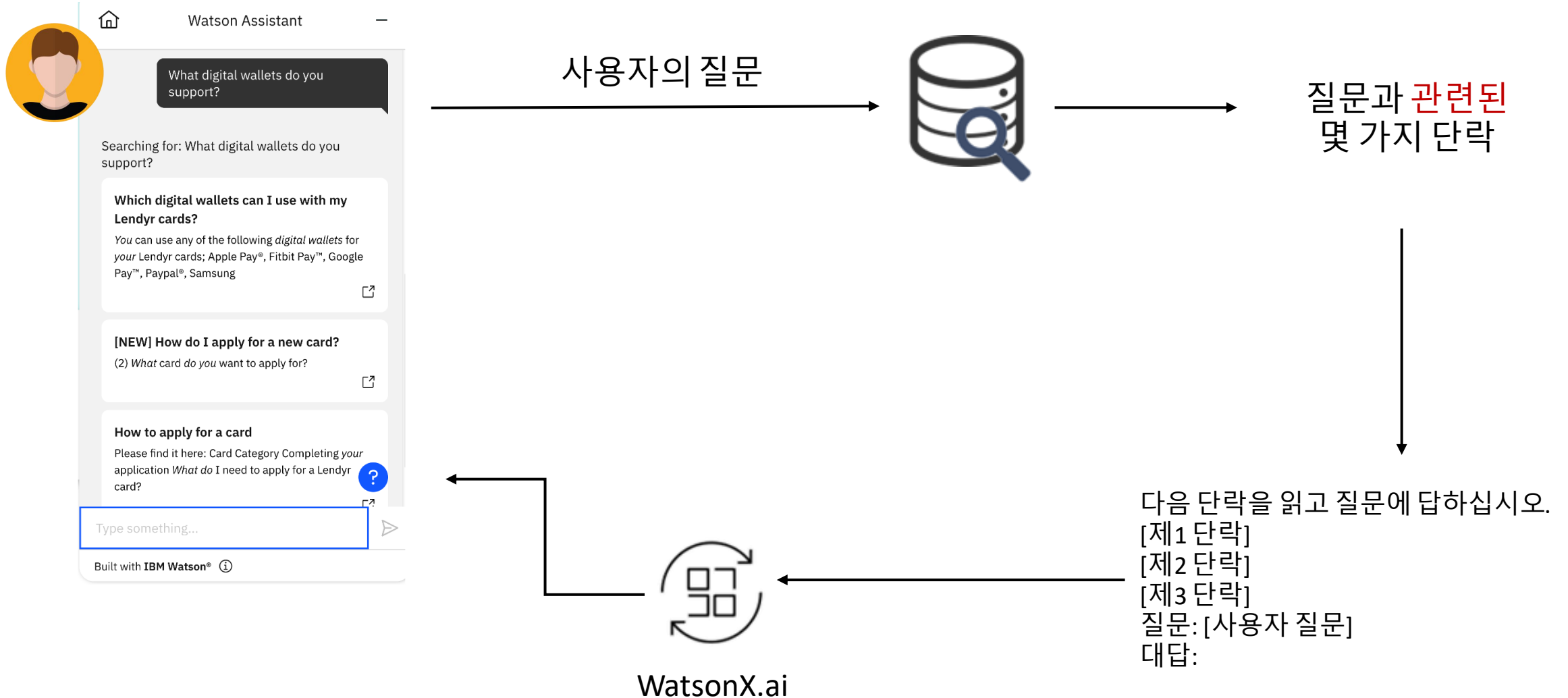
사용자의 질문



질문과 **관련된**
몇 가지 단락

Phase 2 (LLM을 기반으로 하는 새로운 단계)

데이터 쿼리

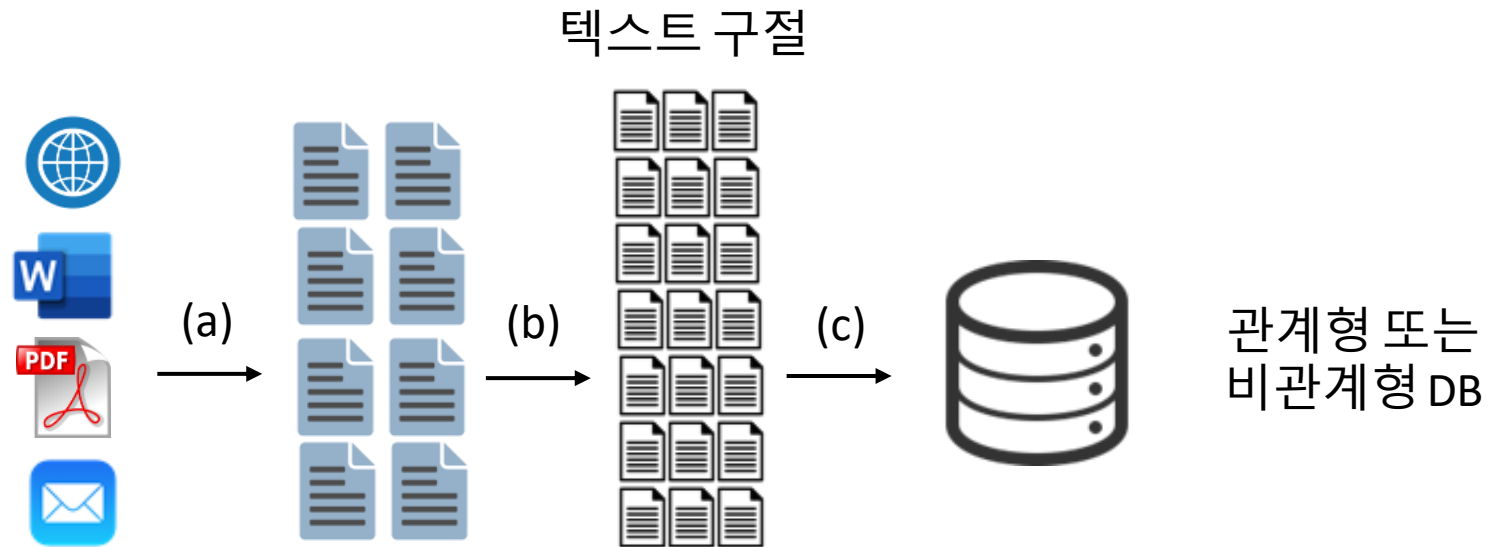


Phase 1: "전통적인" 방식

Phase 1

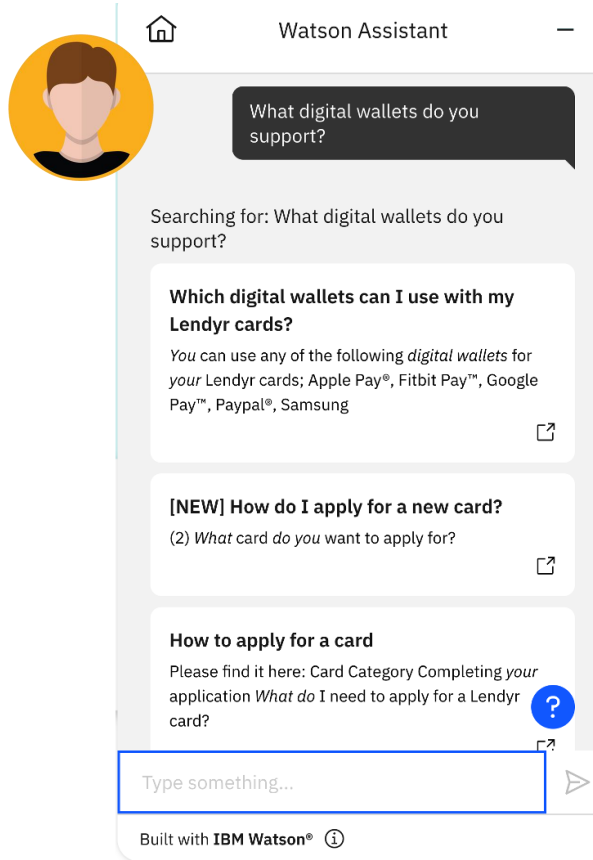
데이터 수집

- (a) 원본 문서 파일의 문서화
- (b) 문서를 청크(chunk)로 분할
- (c) 분할된 청크를 DB로 저장



Phase 2: syntactic(구문) 검색

데이터 쿼리



카드 신청은 어떻게 하나요?



질문과 관련된
몇 가지 단락

청크에 “신청” 또는 “카드”가
포함되도록 모든 행/문서를
검색합니다.

다음 단락을 읽고 질문에 답하십시오.
[제1 단락]
[제2 단락]
[제3 단락]
질문: [사용자 질문]
대답:



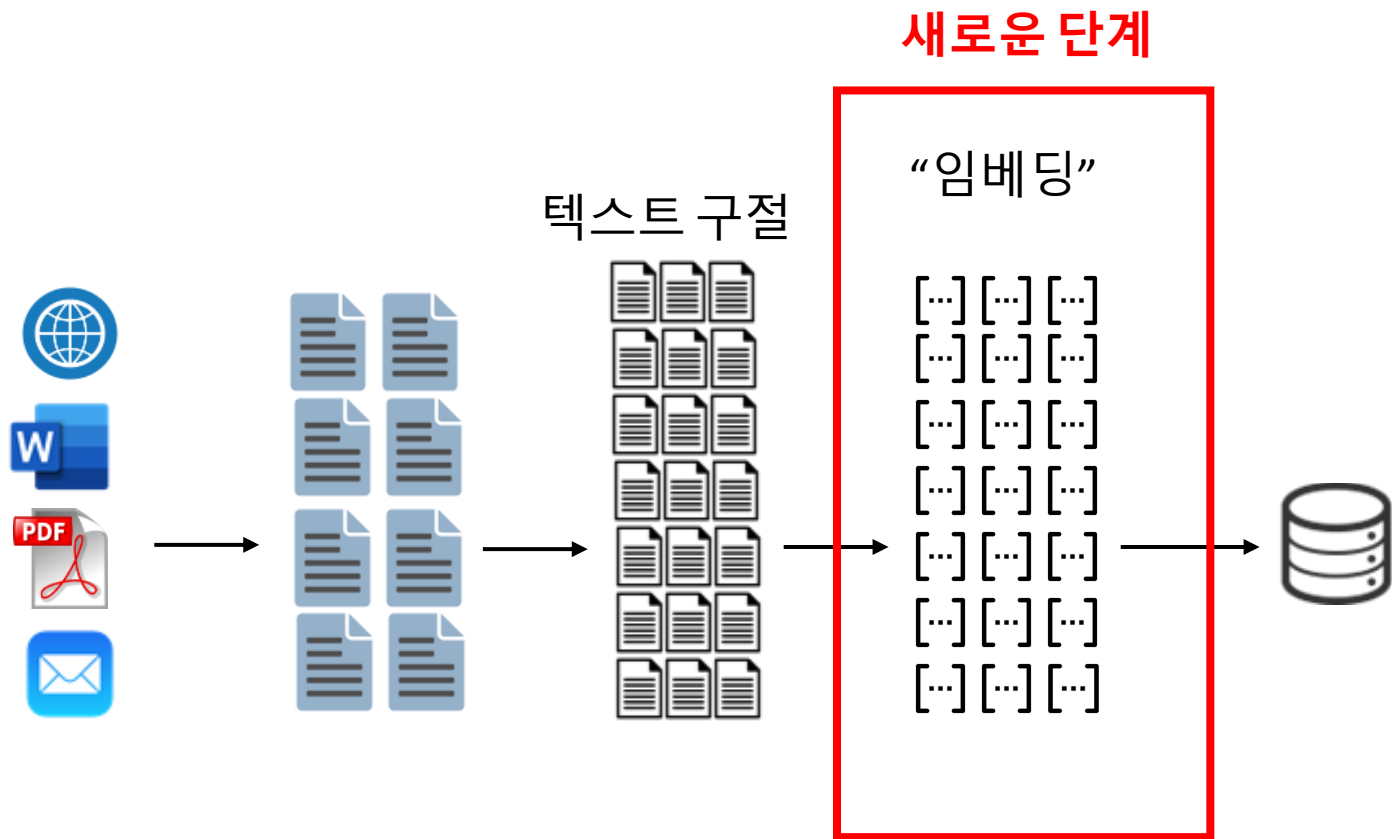
watsonx.ai

Phase 1: "임베딩" 방식

Phase 1

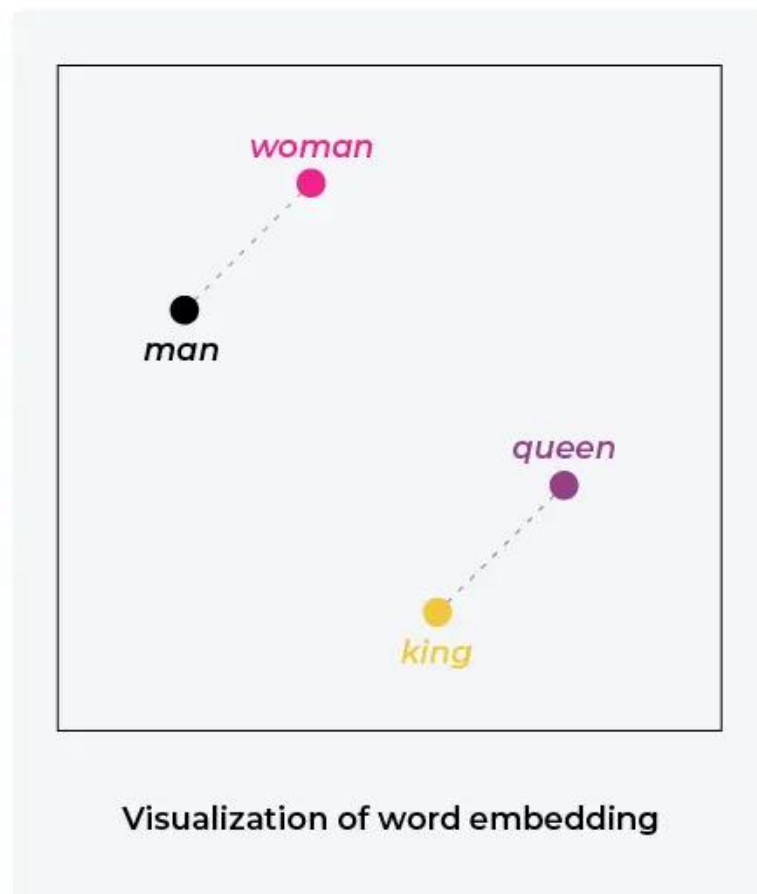
데이터 수집

- (a) 원본 파일을 문서로
- (b) 문서를 청크로 분할
- (c) 청크를 임베딩으로 변환
- (d) 임베딩을 벡터 스토어로 전달



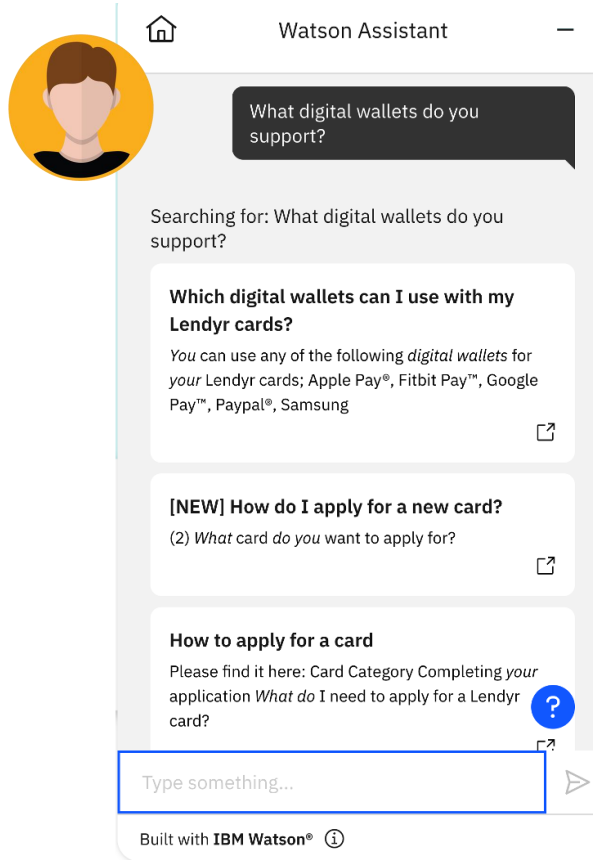
Phase 1: "임베딩" 방식 (예시)

		living being	feline	human	gender	royalty	verb	plural
<i>man</i> →		0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →		0.7	0.3	0.8	-0.7	0.1	-0.5	-0.4
<i>king</i> →		0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →		0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9
word		Word embedding						



Phase 2: 시멘틱 검색

데이터 쿼리



카드 신청은
어떻게 하나요?



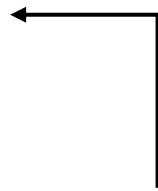
[...]



질문과 관련된 몇
가지 단락을 검색
결과 기준 선별



사용자 질문의 임베딩에 **가장
가까운** 모든 임베딩을 검색합니다.



watsonx.ai

다음 단락을 읽고 질문에 답하십시오.
[제1 단락]
[제2 단락]
[제3 단락]
질문: [사용자 질문]
대답:

구문 기반 검색 vs. 시맨틱(의미 기반) 검색

시맨틱 검색이 정보를 검색하는 "더 나은" 방법인 이유

사용자는 자신의 방식으로 표현하는 반면 문서는 주로 "전문적인" 용어를 사용합니다.

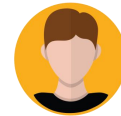
예제:



유급 휴가
(IBM HR 문서)

기업 자산
(은행 윤리강령)

매출, 이익, 편익
재무 보고서



연차

회사의 노트북

수익

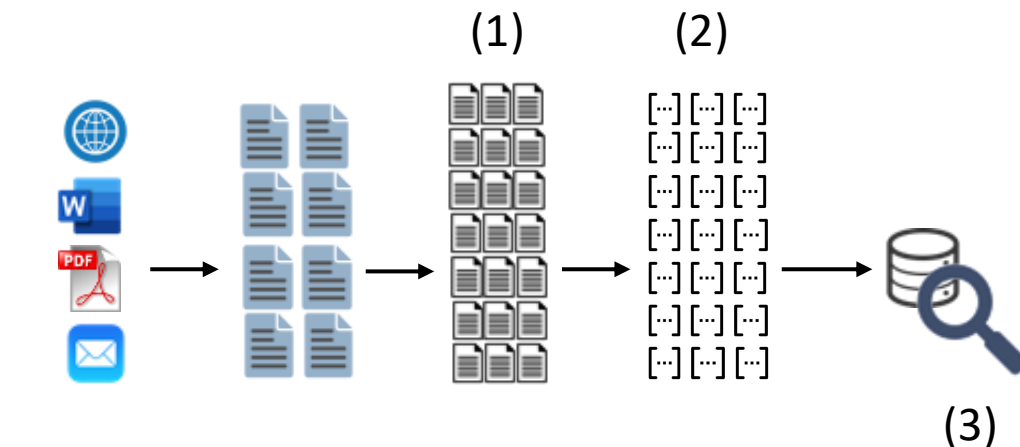
정확도를 향상시키는 방법은 무엇일까요?

예시문서 : 시몬스 은행 - 윤리 강령

각 단계에서 config 구성요소를 최적화합니다.

1. 텍스트 청크의 길이
2. 임베딩 라이브러리 선택
3. 임베딩 간 거리 함수
4. 데이터베이스에서 검색되는 청크 수
5. 프롬프트
6. LLM 파라미터(temperature, topK, top 등)
7. LLM 선택
8. 기타

그러나 더 효율적인 방법이 있습니다



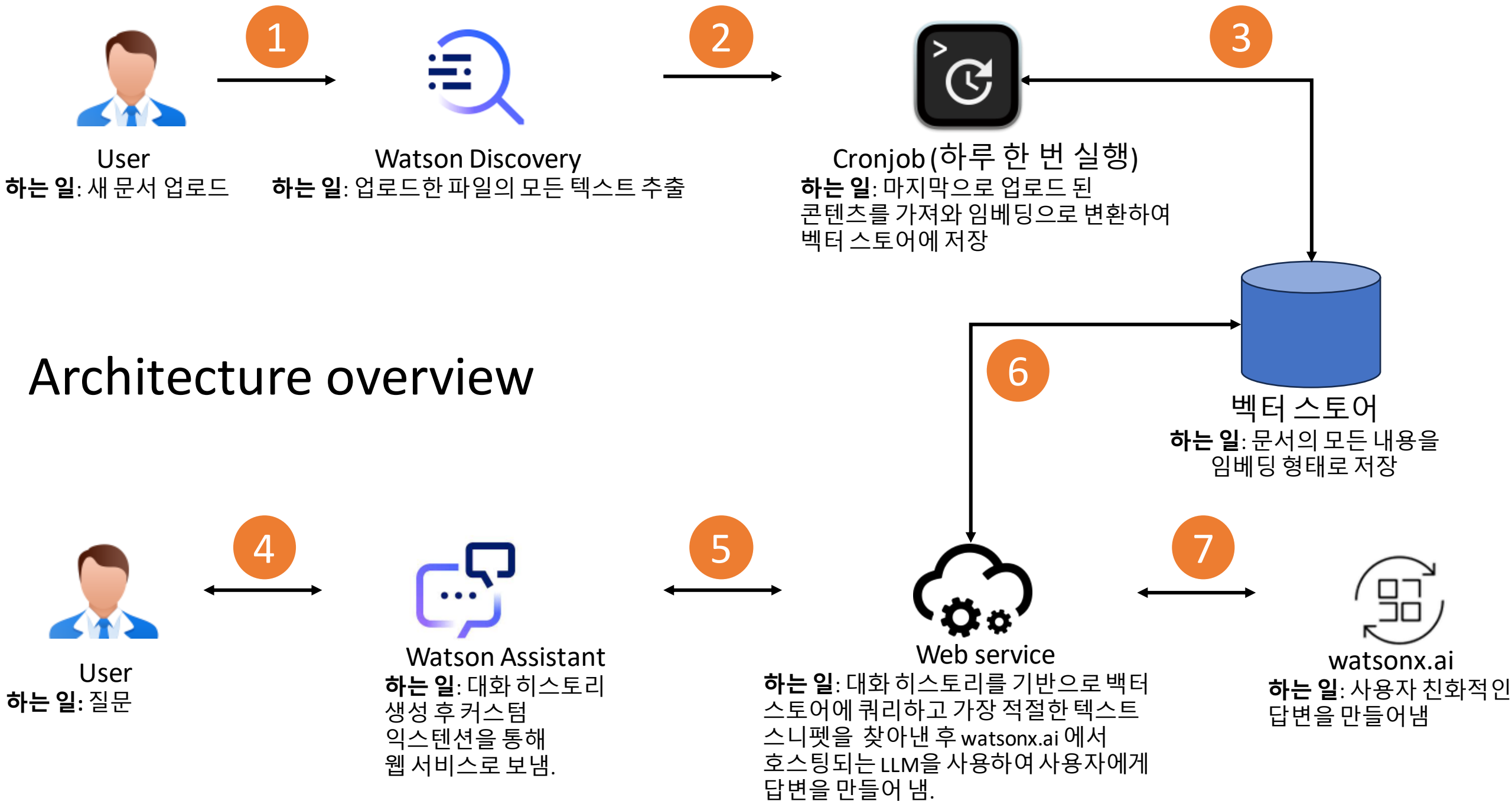
(5)

(4) { 다음 단락을 읽고 질문에 답하십시오.
[제1 단락]
[제2 단락]
[제3 단락]
질문: [사용자 질문]
대답:

(6), (7)



watsonx.ai



Retrieval Augment Generation - Applications



내용이 긴 문서 요약

Problem:

대규모 corpus, 잠재적인 정보 손실 및 비용이 많이 드는 프로세스

Solution:

질문시 RAG 활용 후 결과 수집 후 요약



API 검색 에이전트

Problem:

LLM이 사용할 정보를 전달하기 위해 YAML 파일의 대규모 corpus에 쿼리 및 처리

Solution:

RAG를 이용한 Q&A를 위해 처리된 YAML file에 쿼리합니다



Wisper Bot

Problem:

대화 기록에서 자동 RAG 검색 어려움

Solution:

사용자의 채팅 기록을 실시간으로 처리, RAG를 활용하여 다음 단계를 위한 관련 정보를 얻습니다