

HACK EVERY BIT FOR A BETTER PLANET  
HACK EVERY BIT FOR A BETTER PLANET  
HACK EVERY BIT FOR A BETTER PLANET



## AI 기반 한국형 문서 파싱 서비스

지도교수: 윤수연 교수님

김동연 강아영 김동진 박가현 배경준 하승준

# PLANNING

서비스 소개

개발 배경

LLMong 솔루션

# SERVICE

핵심 기능

서비스 플로우

시연

# BUSINESS

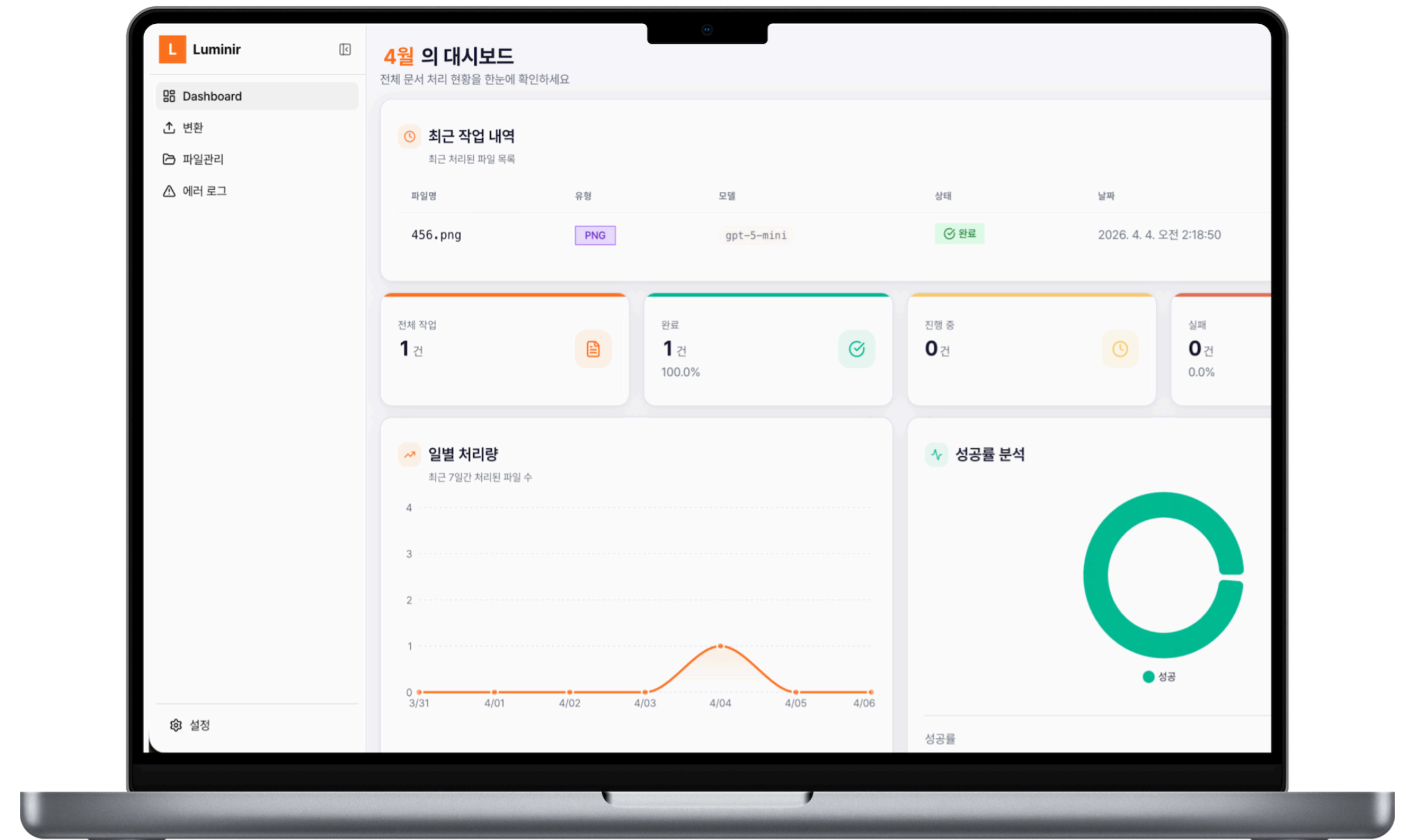
기대 효과

비즈니스 모델

성과

# LLMong (LLM + 夢)

AI 기술을 활용해 한국형 문서를 RAG에 넣을 수 있는 형태로 변환하는 Document Parser



## 서비스 소개

### 문서 파싱 및 결과 검수

원본 문서-추출 결과  
비교 및 확인

### 온프레미스 LLM 시스템 구축

로컬 GPU 서버 기반  
완전 폐쇄망 처리



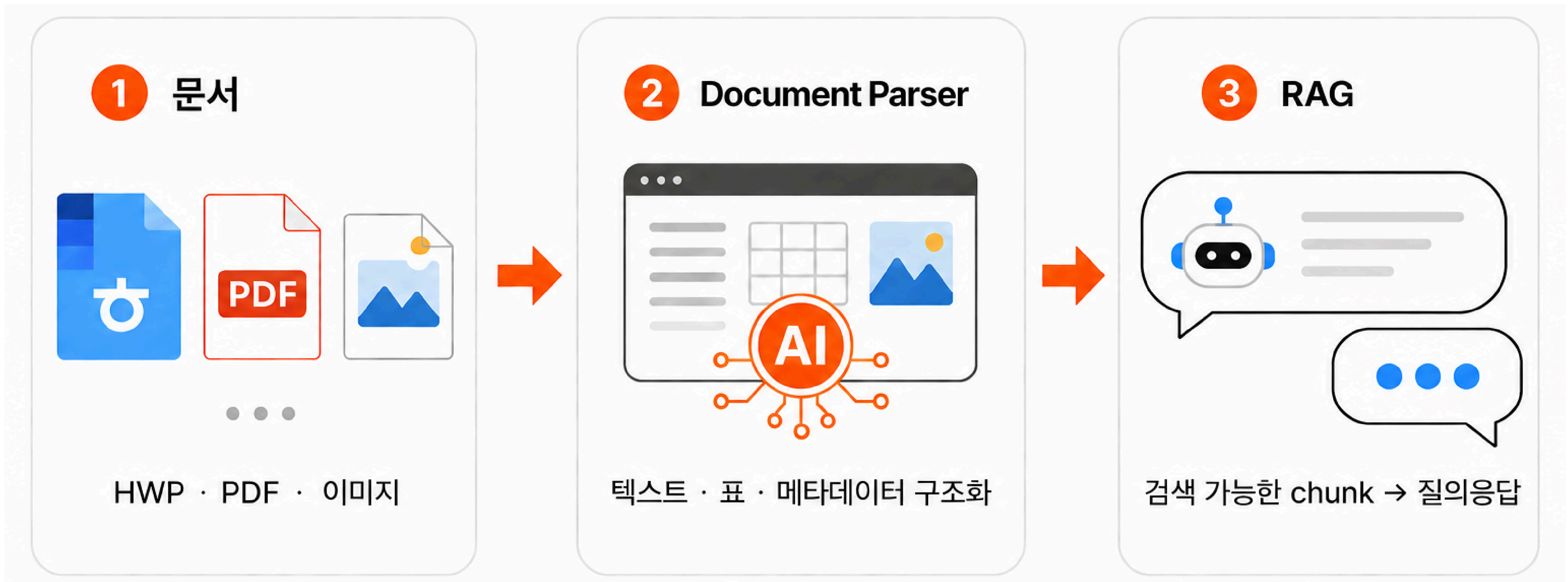
### RAG 문서 질의응답

변환된 문서 기반  
RAG 생성 및 질의응답 챗봇 제공

### 서비스 대시보드

문서 업로드, 변환 현황, 성공률,  
오류 상태 실시간 시각화

# PARSER / RAG, 어떤 연관이 있는가



Document Parser는 한국형 문서를 AI가 읽을 수 있는 **구조화 데이터로 변환**하는 **전처리 단계**입니다. 구조화된 문서는 RAG의 **검색 정확도와 답변 품질을 높이는 기반**이 됩니다.

# 서비스 개발 배경 : 공공·기업의 자체 LLM 도입 확산

성장하는 국내 전자문서 시장 규모

## 'AI 3대 강국' 내년 본격 시동...정부, 로드맵·예산 확정 발표

도시혁 기자

입력 2025.12.14 22:42 수정 2025.12.15 18:40

AI 기업

### 와이즈넷, 한국도로공사에 RAG 기반 챗봇 구축

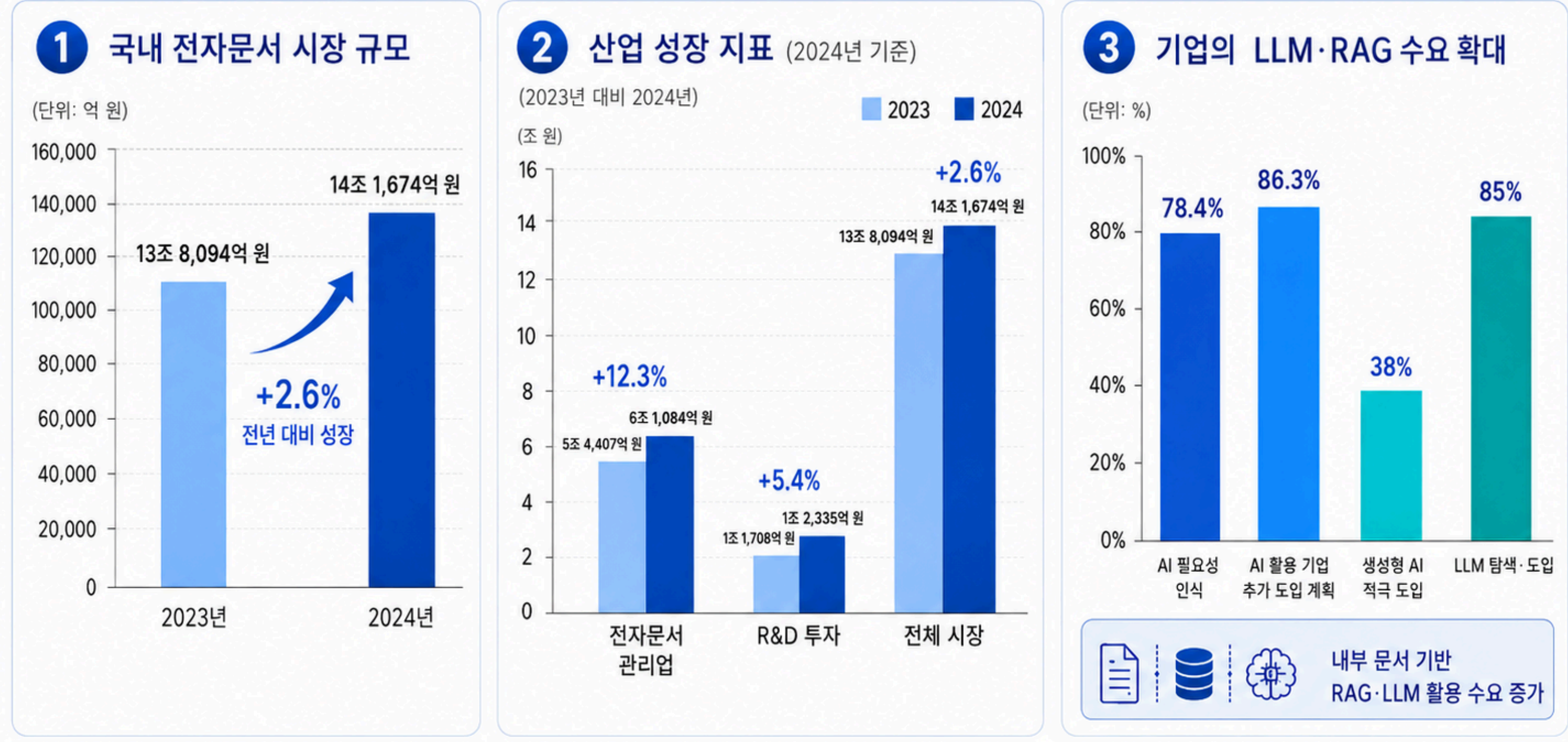
장세민 기자 입력 2025.01.15 15:57 댓글 0

공유 북마크 댓글

#### 행정 분야 자체 LLM 도입... `글로벌 AI 행정도시` 시동 건다

담당부서 | 디지털도시국·디지털정책과 문의 | 02-2133-2990 수정일 | 2025-08-11

- 행정업무, 생성형AI 적용하는 '챗봇 2.0' 시작... 내부망에 자체 LLM 단계적 도입
- 최적화된 생성형AI 행정지원 체계 구축, 반복 행정업무 지원... 하반기 시범 가동
- 시민용 챗봇 '서울톡'에도 생성형 AI 시범 적용... 응답률 및 시민 만족도 개선 기대
- 시 "공공행정 전반 AI 활용, 국내 넘어 AI행정 선도하는 글로벌 도시 거듭날 것"



# 서비스 개발 배경 : 한국형 포맷 및 보안 문제

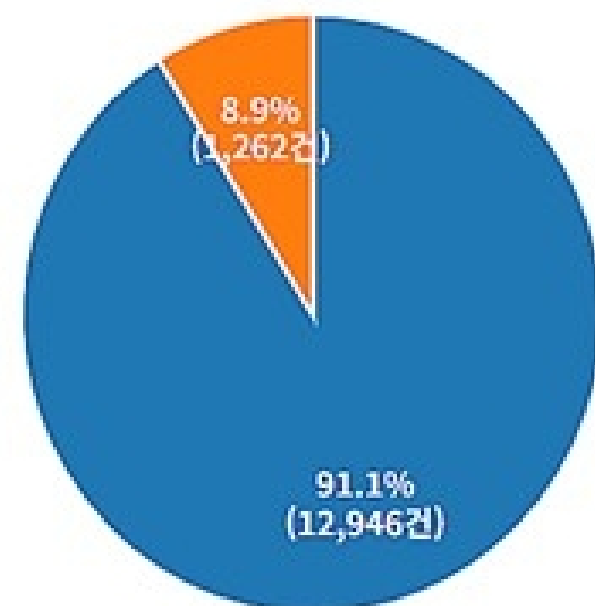
## 내부 문서의 AI 활용을 가로막는 기존 포맷

### 공무원 10명 중 7명 AI 사용...**행정문서 90% 'AI가 못 읽는 포맷'**

송고 2025-10-13 11:28 日本語



Q4. 귀 기관의 보고서·보도자료·사업계획 등 행정문서는 주로 어떤 형태입니까?



- LLM이 읽기 어려운 포맷 위주 (예: HWP(한글), 이미지/스캔 PDF[OCR 미적용] 등): 12,946 (91.1%)
- LLM이 읽을 수 있는 포맷 위주 (예: TXT, HTML, docx 등): 1,262 (8.9%)

## 외부로 유출되는 한국의 공공 문서

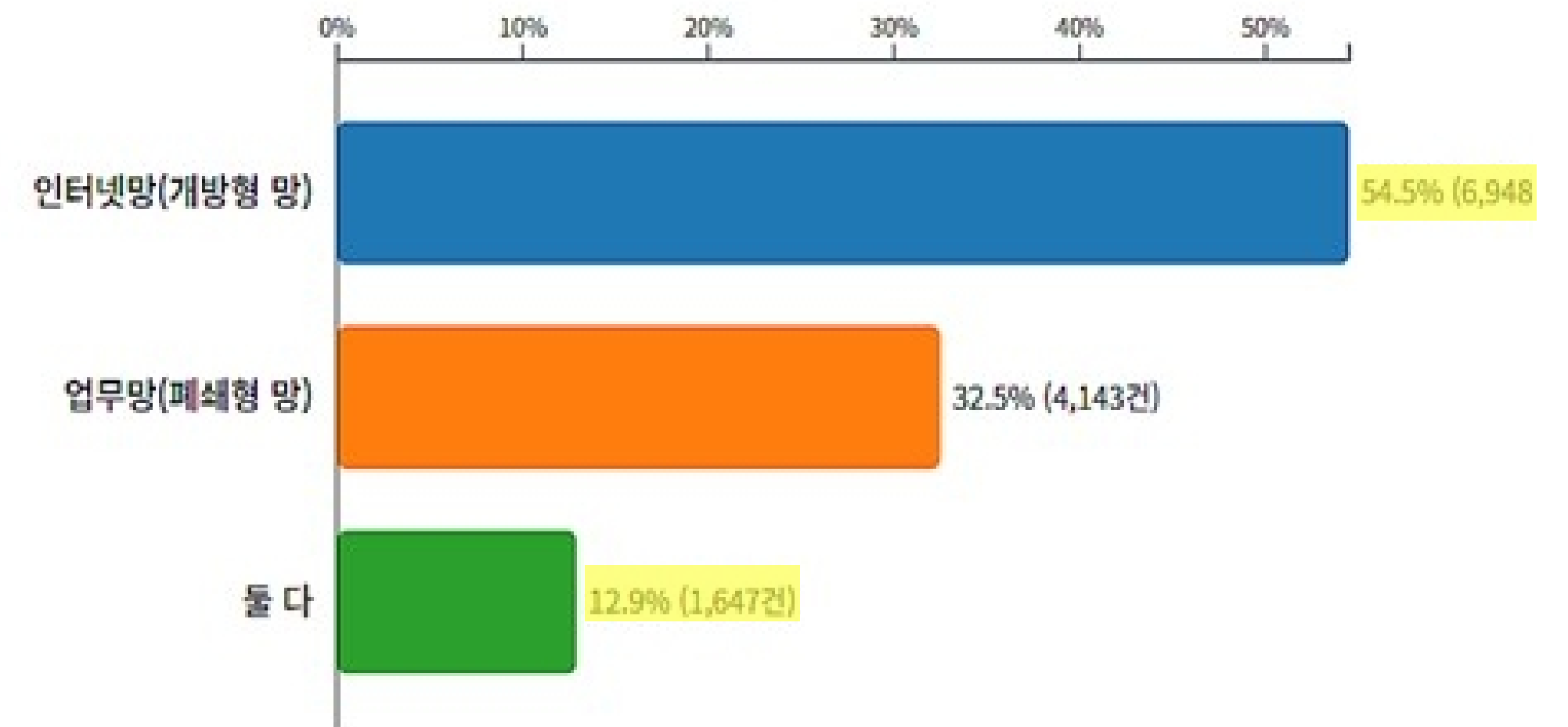
### [단독] "챗GPT로 일하는 서울시 공무원들"...보안 공백 우려

관련이슈 이슈플러스

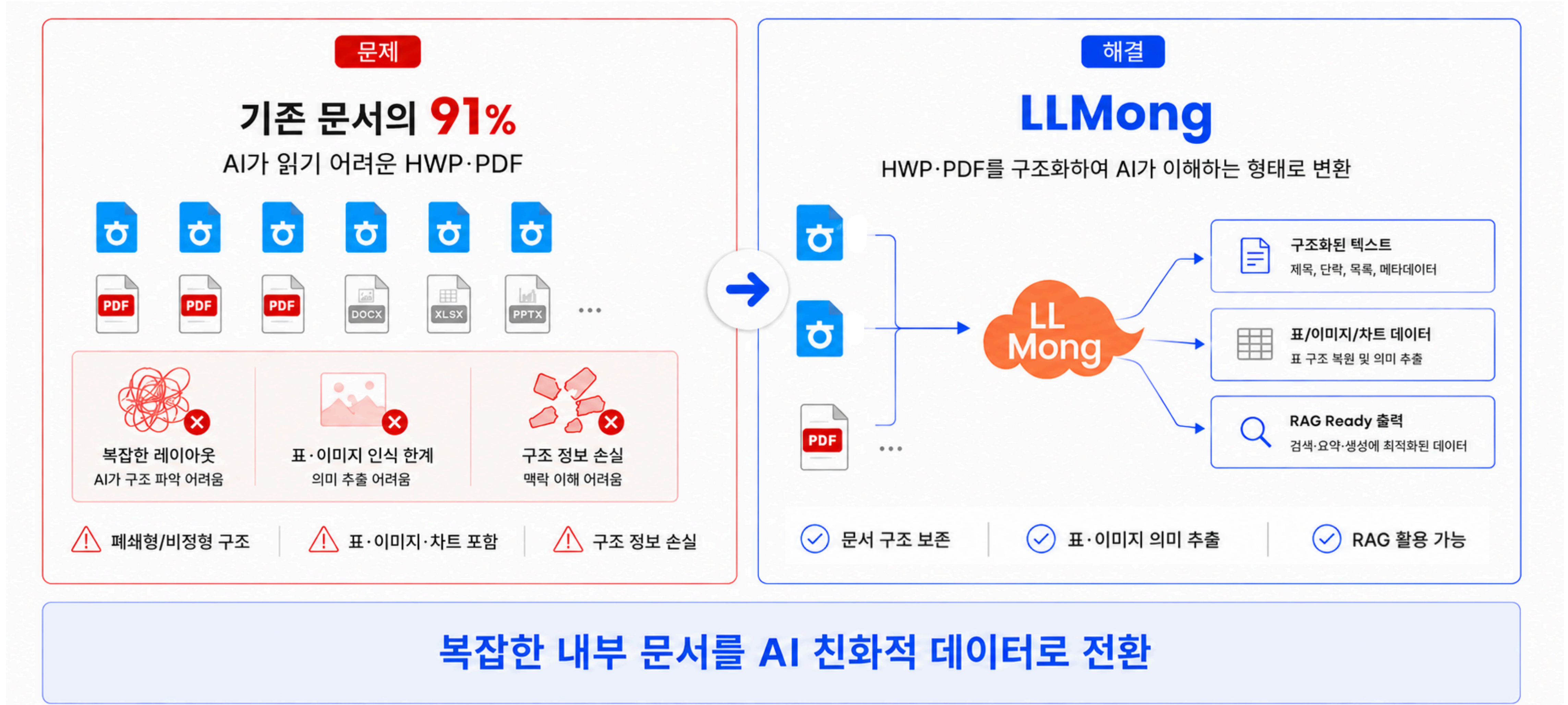
입력 : 2025-10-21 06:00 | 수정 : 2025-10-21 11:14

윤성연·김수연 기자

AI/LLM을 사용할 때 주로 어떤 환경에서 활용했습니까?



# LLMong Solution



# PLANNING

서비스 소개

개발 배경

LLMong 솔루션

# SERVICE

핵심 기능

서비스 플로우

시연

# BUSINESS

기대 효과

비즈니스 모델

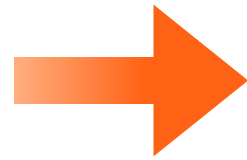
성과

# 서비스 소개 : 주요 핵심 기능

## 각종 문서를 HTML, JSON 등 다양한 LLM 학습용 포맷으로 변환

## 변환된 결과를 활용하여 RAG 생성 및 질의응답 챗봇 제공

구분	임무·역할
징후	<ul style="list-style-type: none"> <li>&lt;위기로 발전할 수 있는 경향이 나타나는 상태&gt;</li> <li>유해물질 등의 일부가 하천 등으로 유입되었을 때</li> <li>오염물질 유입으로 수질자동측정망의 경계경보가 지속 될 때</li> <li>국가하천에서 수질오염으로 추정되는 어류폐사가 상당량 발생하였을 때</li> </ul>
국가안보실 (위기관리센터)	<ul style="list-style-type: none"> <li>국가위기관리회의 운영 등(p14 참조)</li> </ul>
중앙안전관리 위원회 (국무조정실)	<ul style="list-style-type: none"> <li>재난사태 및 특별재난지역 선포 등 건의사항 심의 등(p14 참조)</li> </ul>
대통령비서실	<ul style="list-style-type: none"> <li>재난상황 종합평가, 초기 전략 대응반 운영 등(p14 참조)</li> </ul>
행정안전부	<ul style="list-style-type: none"> <li>위기경보 발령시 대응대책 및 단계별 조치사항 이행</li> <li>유관기관 협조체계 및 대비상황 점검 확인</li> <li>인명구조, 오염지역 방제활동 등 대응활동 지원</li> </ul>
환경부	<ul style="list-style-type: none"> <li>자체위기관리회의 개최</li> <li>종합상황실 설치,운영</li> <li>초동방제 등 사고수습상황 모니터링 및 대내외 보고,전파</li> <li>오염물질 이동, 확산 예측결과 제공,자문(국립환경과학원)</li> <li>사고 영향지역 광역 수질분석 지시(국립환경과학원)</li> <li>화학물질 관련정보 제공 및 기술지원(화학물질안전원)</li> <li>초동방제 등 사고수습상황 언론 보도,브리핑</li> </ul>
환경청	<ul style="list-style-type: none"> <li>지역 상황실 설치,운영</li> <li>현장의 사고수습상황 파악,보고</li> <li>지자체의 사고수습 및 방제활동 지원</li> <li>수질오염경보시스템 모니터링 강화</li> <li>사고영향지역 수질분석 실시</li> </ul>
지자체	<ul style="list-style-type: none"> <li>재난현장통합지원본부 설치,운영</li> <li>지역재난안전대책본부 설치,운영(필요시)</li> <li>오염물질 확산방지 및 방제활동 실시</li> <li>사고유발 사업장 추적조사 및 추가유출 방제조치</li> <li>사고지역 인근 배출입소 지도점검 강화</li> <li>어류폐사의 경우 폐사원인 조사 및 폐사체 수거,처리</li> <li>긴급구조통제단, 경찰서 등 유관기관 지원요청(필요시)</li> </ul>
국토교통부	<ul style="list-style-type: none"> <li>해당수계 유량조절을 위한 댐,보 운영계획 수립</li> <li>위기상황 모니터링 및 유관기관 지원인력 공조체계 유지</li> </ul>
국방부 경찰청	<ul style="list-style-type: none"> <li>위기상황 모니터링</li> <li>사고수습을 위한 교통통제 지원</li> <li>사고유발 업체에 대한 조사 협조</li> </ul>



구분	임무·역할
징후	<ul style="list-style-type: none"> <li>&lt;위기로 발전할 수 있는 경향이 나타나는 상태&gt;</li> <li>유해물질 등의 일부가 하천 등으로 유입되었을 때</li> <li>오염물질 유입으로 수질자동측정망의 경계경보가 지속 될 때</li> <li>국가하천에서 수질오염으로 추정되는 어류폐사가 상당량 발생하였을 때</li> </ul>
국가안보실 (위기관리센터)	<ul style="list-style-type: none"> <li>국가위기관리회의 운영 등(p14 참조)</li> </ul>
중앙안전관리 위원회 (국무조정실)	<ul style="list-style-type: none"> <li>재난사태 및 특별재난지역 선포 등 건의사항 심의 등(p14 참조)</li> </ul>
대통령비서실	<ul style="list-style-type: none"> <li>재난상황 종합평가, 초기 전략 대응반 운영 등(p14 참조)</li> </ul>
행정안전부	<ul style="list-style-type: none"> <li>위기경보 발령시 대응대책 및 단계별 조치사항 이행</li> <li>유관기관 협조체계 및 대비상황 점검 확인</li> <li>인명구조, 오염지역 방제활동 등 대응활동 지원</li> </ul>
환경부	<ul style="list-style-type: none"> <li>자체위기관리회의 개최</li> <li>종합상황실 설치,운영</li> <li>초동방제 등 사고수습상황 모니터링 및 대내외 보고,전파</li> <li>오염물질 이동, 확산 예측결과 제공,자문(국립환경과학원)</li> <li>사고 영향지역 광역 수질분석 지시(국립환경과학원)</li> <li>화학물질 관련정보 제공 및 기술지원(화학물질안전원)</li> <li>초동방제 등 사고수습상황 언론 보도,브리핑</li> </ul>
환경청	<ul style="list-style-type: none"> <li>지역 상황실 설치,운영</li> <li>현장의 사고수습상황 파악,보고</li> <li>지자체의 사고수습 및 방제활동 지원</li> <li>수질오염경보시스템 모니터링 강화</li> <li>사고영향지역 수질분석 실시</li> </ul>
지자체	<ul style="list-style-type: none"> <li>재난현장통합지원본부 설치,운영</li> <li>지역재난안전대책본부 설치,운영(필요시)</li> <li>오염물질 확산방지 및 방제활동 실시</li> <li>사고유발 사업장 추적조사 및 추가유출 방제조치</li> <li>사고지역 인근 배출입소 지도점검 강화</li> <li>어류폐사의 경우 폐사원인 조사 및 폐사체 수거,처리</li> <li>긴급구조통제단, 경찰서 등 유관기관 지원요청(필요시)</li> </ul>
국토교통부	<ul style="list-style-type: none"> <li>해당수계 유량조절을 위한 댐,보 운영계획 수립</li> <li>위기상황 모니터링 및 유관기관 지원인력 공조체계 유지</li> </ul>
국방부 경찰청	<ul style="list-style-type: none"> <li>위기상황 모니터링</li> <li>사고수습을 위한 교통통제 지원</li> <li>사고유발 업체에 대한 조사 협조</li> </ul>

Document Chat

B1-ICT-질자-0003\_00\_중북본부 SCADA\_운영하면,표준화,질자서,전문.pdf.pdf

Loaded "B1-ICT-질자-0003\_00\_중북본부 SCADA\_운영하면,표준화,질자서,전문.pdf.pdf". Ask a question about this document.

SCADA DB 포인트 선정 요청은 연동시험 예정일 60일 전까지 선정하고, SCADA DB 작업예정일 7일 전까지 공문과 SCADA 운영관리시스템으로 요청하여야 합니다.

154kV 변전소에 설비 용량이 예정되어 있는데, SCADA 포인트 선정부터 DB 작업, 시험까지 어떤 순서로 진행해야 할까요?

154kV 변전소에 설비 용량이 예정되어 있을 경우, SCADA 포인트 선정부터 DB 작업, 시험까지 다음과 같은 순서로 진행해야 합니다.

- \*\*포인트 선정 순서\*\***: 송변전설비 감독부서에서 설계단계부터 선,중설 내역을 SCADA DB 관리시스템에 등록하고 관련부서에 통보합니다. 여기에는 공사일정, 선,중설 설비내역, 공사 내용 등이 포함됩니다.
- \*\*포인트 선정\*\***: 변전자동화부서에서 송변전설비 감독부서에서 등록된 공사정보와 SCADA DB 포인트 선정기준을 기반으로 SCADA에 수용할 포인트를 선정합니다.
- \*\*DB 작업\*\***:
  - 변전소 생성 및 RTU 정의
  - 포인트 DB 입력 (감시, 제어, 측정 포인트)
  - 실시간 계산 포인트 입력 (현재부하, 예비율 등)
  - DISPLAY 입력 (단상경선도, 계통도, 조류도, 최대부하표 등)
  - 보고서 DB 입력 (일보, 월보, 속보 등 보고서)
  - 연계DB 입력 (PIS, 고정관전DB, 일괄제어DB, 조작지원시스템 등)
  - 계통운영센터와 EMS 간 연계DB 입력
  - 신규 입력 DB는 변전설비 연동시험 전까지 반드시 제어금지 TAG를 설정
- \*\*포인트 시험\*\***:
  - SCADA 서버 자체 시험을 실시합니다. 워킹스텝치 오류 시 중설 워킹스텝치 DR이러 지명계통운영센터 크레딧스 DR이러이 완료된 후 시 중설 워킹스텝치에 대해 시험을 진행합니다.

Ask about this document

문서 분석 결과는 HTML, JSON, Markdown의 다양한 포맷으로 제공  
LLM 입력 형식에 맞춰 학습에 사용할 수 있도록 구조와 포맷 정제

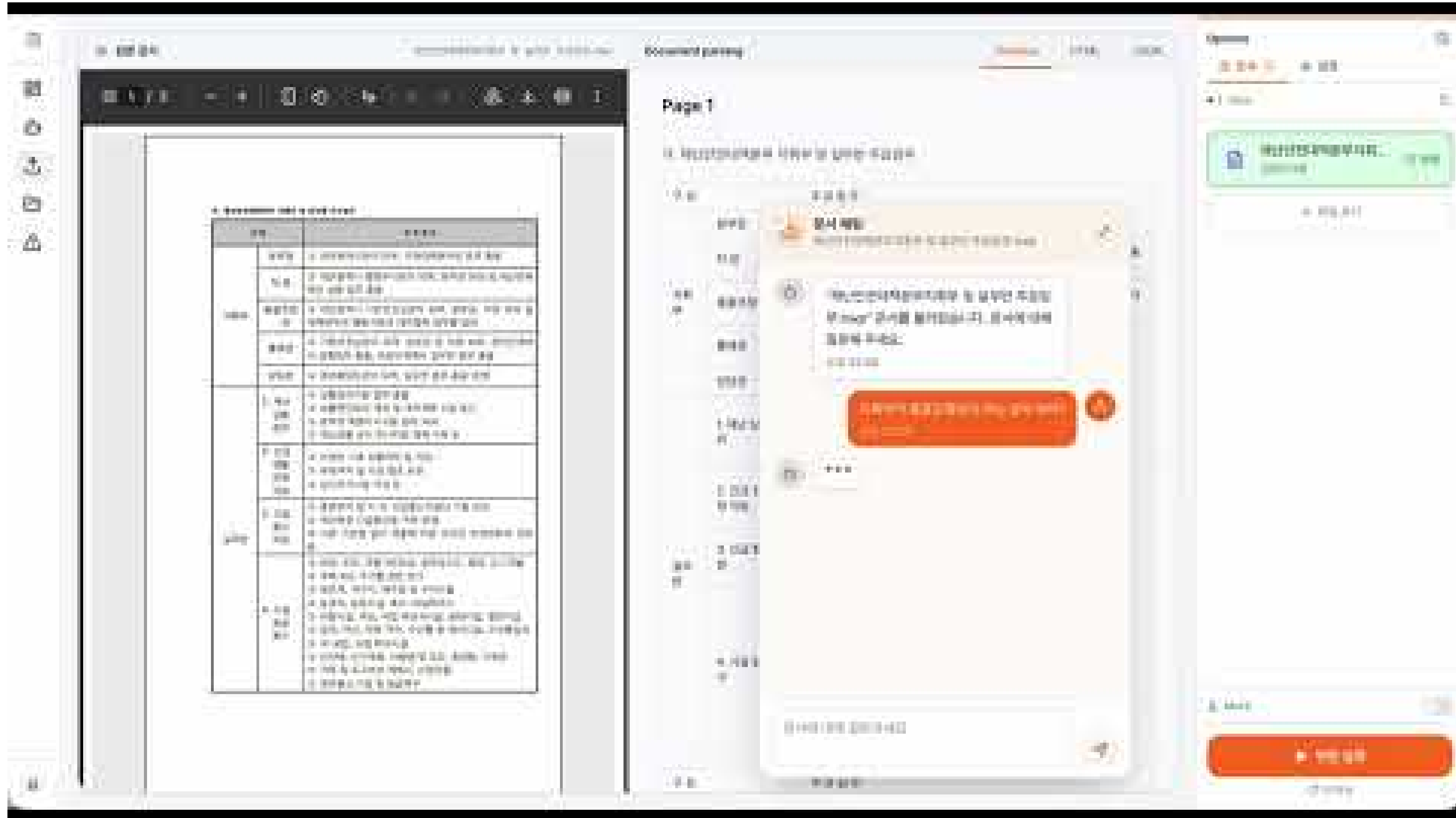
문서 분석 결과를 통해 RAG를 생성  
파싱이 잘 되었는지 확인 가능한 질의응답 챗봇 제공

# LLMong (LLM + 夢)

AI 기술을 활용해 한국형 문서를 시가 읽을 수 있는 형태로 변환하는 Document Parser

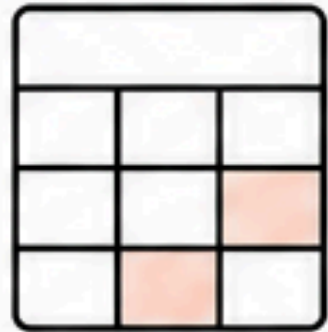


# 시연



## 서비스 소개 : 모든 Document Parser의 핵심

### 핵심 포인트



#### 표 구조 복원

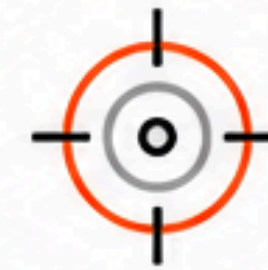
복잡한 표도 행/열 구조를  
정확히 인식하여 데이터 손실 최소화



#### 정확한 텍스트 추출

다양한 문서에서도 깨끗하고  
정확한 텍스트를 안정적으로 추출

### 중요한 이유



#### 데이터 신뢰성의 기반

정확한 표 구조 복원과 텍스트 추출은  
데이터의 신뢰성과 품질을 보장하는 핵심 요소입니다.



#### 업무 자동화의 필수 조건

정확한 데이터 추출 없이는 자동화된 처리와  
의사결정이 불가능합니다.



#### 시간과 비용 손실 방지

수작업 검수와 재작업을 줄여  
운영 효율과 비용 절감 효과를 제공합니다.

### 활용 분야



보고서 자동화



데이터 입력 자동화

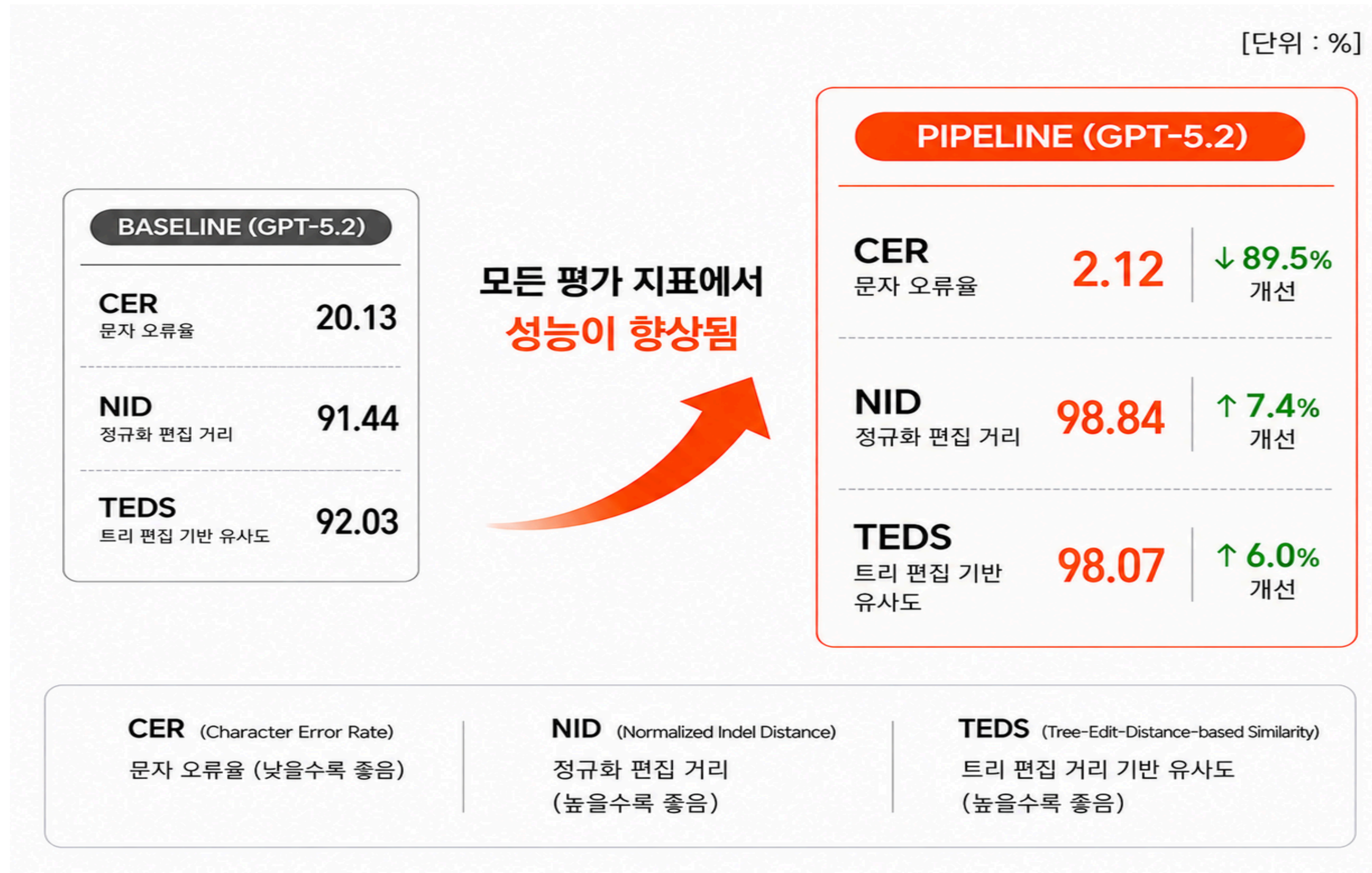


데이터 분석 전처리



문서 검색/관리 시스템

## 서비스 소개 : 문서 복원 성능 비교



기준 모델 : GPT-5.2 | Baseline 및 Pipeline 모두 동일한 기준 모델(GPT-5.2)을 기반으로 비교하였습니다.

# PLANNING

서비스 소개

개발 배경

LLMong 솔루션

# SERVICE

핵심 기능

서비스 플로우

시연

# BUSINESS

기대 효과

비즈니스 모델

성과

## 기대 효과 : 한국 문서의 AI 활용 기반 구축



**“궁극적 목표 : 공공, 기업 문서의 안전한 AI 활용 촉진”**

# 비즈니스 모델

고객 유형에 따라 가장 적합한 도입 방식을 제공합니다.



# 확장 가능성

HWP, HWPX 포맷 이외에도 향후 모든 형태의 문서 처리를 목표로 확장 중입니다.



현재

진행 중

목표

## HWP · PDF · 이미지 스캔 처리

다양한 한국형 문서를  
높은 정확도로 처리 중

## 손글씨 파싱

기능 구현 완료,  
현재 정확도 고도화 진행 중

## 모든 형태의 문서 처리

문서 유형 커버리지 100%를 목표로  
지속적인 확장 예정



**한국형 포맷만이 아닌, 모든 문서 처리 플랫폼으로 확장합니다.**

# 연구 프로젝트 성과 - KICS(한국 통신학회) 논문 발표

## PDF 문서 내 표 구조 복원을 위한 VLM 파이프라인의 성능 분석에 관한 연구

강아영, 김동연, 김동진, 배경준, 박가현, 하승준, \*윤수연

chexxish@kookmin.ac.kr, angrybird0@kookmin.ac.kr, kdj22250@kookmin.ac.kr, seungj03@kookmin.ac.kr, rudwns7g@kookmin.ac.kr, gahyeon1022@kookmin.ac.kr, 1104py@kookmin.ac.kr

Performance Analysis of VLM Pipelines for Table Structure Recovery in PDF Documents  
AYeong Kang, DongYeon Kim, DongJin Kim, GyeongJun Bae, GaHyeon Park, SeungJun Ha, SooYeon Yoon

Kookmin Univ., \*Kookmin Univ.  
요약

본 연구는 PDF 구조 파싱에서 VLM 파이프라인이 직접 호출 방식 대비 어떤 효과를 보이는지 평가한다. 특히 본문 텍스트·캡션·꼬리말과 표가 함께 배치된 PDF 문서에서, 표 영역 분리와 입력 정제가 구조 요소 변환 품질에 미치는 영향을 분석한다. 이를 위해 AI-Hub 표 이미지-텍스트 쌍 데이터 중 50개 샘플을 기반으로 합성 PDF 10개를 구성하고, Qwen3-VL-8B와 Qwen3-VL-32B에서 파이프라인과 직접 호출 베이스라인을 비교하였다. 표 이미지만 입력한 ablation에서는 pipeline-baseline TEDS 차이가 Qwen3-VL-8B +0.0063, Qwen3-VL-32B -0.0023으로 일관되지 않았고 두 비교 모두 통계적으로 유의하지 않았다( $p>0.05$ ). 반면 주변 텍스트가 포함된 합성 PDF에서는 두 모델 모두 파이프라인이 직접 호출 베이스라인보다 낮은 CER/CER-NS와 높은 NID/TEDS/TEDS-S를 보였다. 이는 PDF 구조 파싱에서 VLM 파이프라인의 이득이 단순한 VLM 호출보다 문서 내 구조 영역을 분리하고 입력을 정제하는 단계에서 주로 발생함을 시사한다.

## PDF 문서 내 표 구조 복원 연구

성능에서 가장 큰 부분을 차지하는 PDF 내 표 복원 성능 분석 연구  
Qwen, GPT 등 다양한 모델을 활용한 실험으로 파이프라인의 높은 성능 확인

## RAG 적용을 위한 VLM 기반 문서 파싱 파이프라인의 병렬처리 성능 분석에 관한 연구

배경준, 강아영, 김동연, 김동진, 박가현, 하승준, \*윤수연

rudwns7g@kookmin.ac.kr, chexxish@kookmin.ac.kr, angrybird0@kookmin.ac.kr, kdj22250@kookmin.ac.kr, gahyeon1022@kookmin.ac.kr, seungj03@kookmin.ac.kr, 1104py@kookmin.ac.kr

A Study on the Parallel Processing Performance Analysis of a VLM-Based Document Parsing Pipeline for RAG Applications

GyeongJun Bae, AYeong Kang, DongYeon Kim, DongJin Kim, GaHyeon Park, SeungJun Ha, \*SooYeon Yoon  
Kookmin Univ., \*Kookmin Univ.

요약

본 논문은 RAG 적용을 위한 문서 데이터 생성을 목적으로, 본 연구에서 구현한 시각언어모델(Vision-Language Model, VLM) 기반 문서 파싱의 파이프라인과 병렬처리 성능과 이를 통한 RAG 성능을 분석한다. 단순 텍스트 추출 방식은 문서 내 표, 차트, 흐름도 등의 구조를 보존하지 못한 채 추출하는 문제가 있고, 순차 처리 방식은 이미지 추출, 문서 변환, VLM 추론, 결과 병합의 과정이 차례대로 수행되어 CPU를 사용하는 처리에서 GPU가 유휴 상태가 되거나, VLM 추론 중 CPU 기반 작업이 대기하는 문제가 있다. 이를 완화하기 위해 본 연구에서는 문서 처리 과정을 전처리, VLM 추론, 후처리 단계로 분리, 각 작업을 큐에 등록한 뒤 비동기 워커가 처리하는 큐 기반 비동기 병렬 구조와 문서별 파이프라인을 적용하였다. 본 연구는 문서별 파이프라인과 큐 기반 비동기 병렬 구조 처리 대기시간이 감소와 정보 추출량 증가를 확인하였다.

## 파이프라인의 병렬처리 성능 분석 연구

효율적인 비동기 문서처리 구조, 변환 방식 연구  
순차 처리 방식 대비 속도 향상, 파이프라인을 통한 RAG 정확도 향상 확인

# 팀 소개



**김동연**  
PM & Full Stack

프로젝트 일정 및 기능 기획  
프론트엔드 개발  
서버 배포 자동화 및 품질관리  
전체 서비스 흐름 관리

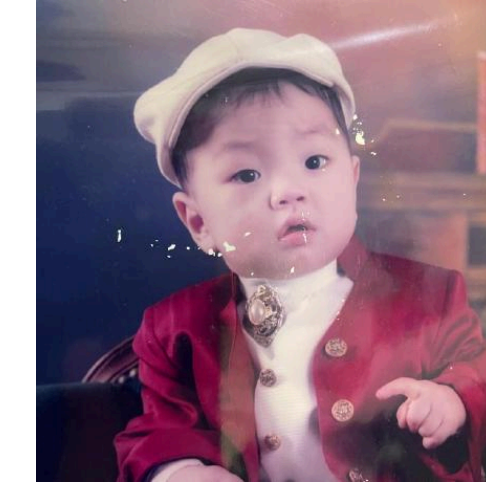
@Oyeonnnn0



**강아영**  
AI

LLM 파서 성능 평가 시스템 설계·구축  
HWP 다중 페이지 표 파싱 고도화  
정량 지표 기반 GT 데이터 구축·검수

@kaye0ng



**김동진**  
Frontend

프론트엔드 화면 개발  
반응형 인터페이스 개선  
공통 UI 컴포넌트 구조화

@K-Dongjin



**박가현**  
Backend

FastAPI 서버 개발  
작업 상태 API 구현  
데이터 처리 흐름 관리

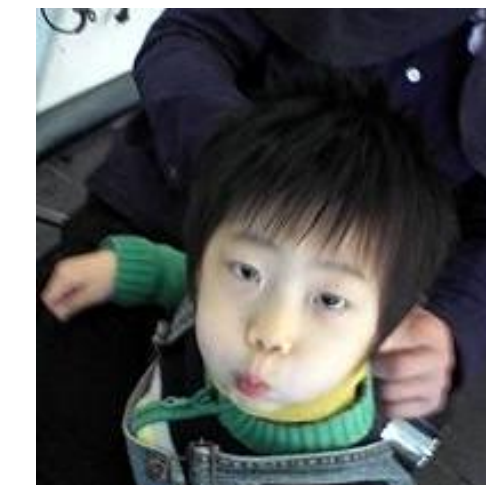
@gahyeon1022



**배경준**  
Backend

문서 처리 API 개발  
RAG 질의응답 구현  
비동기 작업 안정화

@jun-kookmin



**하승준**  
Backend & AI

VLM 기반 문서 추출 파이프라인 고도화  
Worker 처리 구조 설계  
복합 문서 요소 구조화 로직 개선

@seunG-Zzun

# 문서를 읽는 AI에서 문서를 이해하는 AI로

“한국 문서의 AI 활용, LLMong이 선도합니다.”

## 구조화 데이터 변환

복잡한 한국형 문서를 AI가 이해 가능한 구조화 데이터로 변환

## 강력한 보안성 확보

온프레미스 기반 처리로 공공·기업 문서의 기밀성 완벽 보호

## RAG 품질 최적화

RAG 연계를 통해 검색·질의응답 품질 향상

## 문서 AI 플랫폼 지향

HWP·PDF·이미지·손글씨까지 확장 가능한 문서 처리 플랫폼

Thank you



# APPENDIX

HACK EVERY BIT FOR A BETTER PLANET  
HACK EVERY BIT FOR A BETTER PLANET  
HACK EVERY BIT FOR A BETTER PLANET

## 표 테스트 진행 결과

Engine	CER↓	NID↑	TEDS↑	sec/page↓
pipeline: qwen3-vl-8b	0.0480	0.9749	0.9708	11.94
baseline: qwen3-vl-8b	0.0586	0.9654	0.9709	10.69
pipeline: qwen3-vl-32b	0.0295	0.9829	0.9807	9.80
baseline: qwen3-vl-32b	0.1965	0.9168	0.9433	11.24
pipeline: gpt-5.2	0.0212	0.9884	0.9818	7.72
baseline: gpt-5.2	0.2013	0.9144	0.9203	21.79

## RAG 성능 지표

평가조건 : 검색 결과 N위 chunk 안에 정답 근거가 있는지 평가

평가 조건	방식	질의 수	Hit	정확도
Top5	gt_flat	300	80.3%	89.3%
Top5	gt_structured	300	97.7%	97.7%
Top1	gt_flat	300	35.1%	55.9%
Top1	gt_structured	300	87.6%	87.6%

gt\_flat : GT에 포함된 표 내용을 일반 텍스트처럼 평탄화하여 chunk로 나눈 방식

gt\_structured : GT 표를 행·열 구조가 유지되도록 table/row 단위 chunk로 변환해 색인한 방식