

이름 한지혁

지도교수 조현석

연구 배경

23년도 3월 삼성전자가 디바이스솔루션 부문 사업장 내 chatGPT 사용을 허가하였을 때, 기업 정보가 유출되는 사건이 발생했다. chatGPT와 같은 LLM은 유용하지만 이처럼 기업 정보가 유출될 수 있다.

삼성전자 기밀 유출 사고 당시, 회의 내용이 유출되었는데 만약 민감한 단어를 다른 키워드로 대체하고 문맥만을 유지한다면, GPT의 회의 요약 기능은 그대로 사용하되, 기밀 내용은 지킬 수 있지 않을까라는 생각에서 본 프로젝트를 진행한다.

본 연구를 통해, 여러 사용자가 On-premise LLM을 구축하지 않더라도 기밀 유출 없이 마음 편히 사용할 수 있는 것이 연구의 최종 목표이며, 현재는 회의록 데이터를 난독화하여 GPT에 추론을 요청한 후, 추론 데이터를 복호화하여 사용자가만 추론 데이터의 본 의미를 파악할 수 있는 텍스트로 변환하는 것이 목표다.

기업명
"기밀 유출 막아라"... 챗GPT 경계 나선 삼성-LG-SK

삼성-LG, 기밀 유출 우려에 사내 챗GPT 사용 제한
국내에서 챗GPT 비즈니스 적용 사례도 속속 나와
"챗GPT로 업무 생산성 향상, 운영비 절감 효과" 분석도

조선경제+ 경제일반
"챗GPT에 물다가 기밀 샌다" 기업마다 정보보호 골머리[NOW]

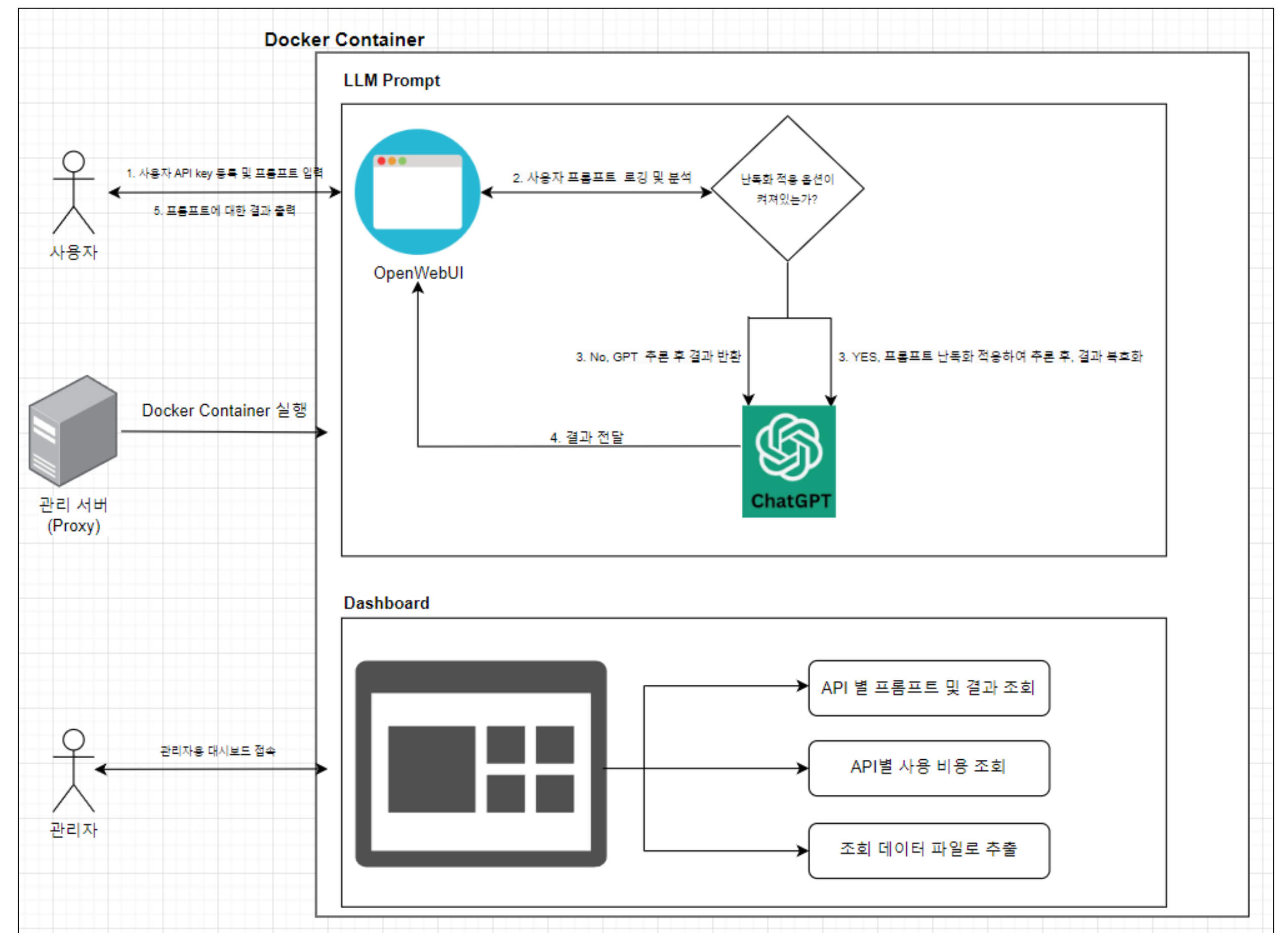
삼성 반도체 "프로그래밍 오류 해결 중 해라"
"기밀 유출 막아라" 챗GPT 사용 제한...
회사 기밀 유출과 생산성 향상 사이서 갈등 심화

작성일자
2023.03.23. 10:20

국내 최대 기업 삼성전자는 세계 최대 인공지능 기업인 OpenAI(이하 '오픈AI') 챗GPT 서비스를 도입하고 나서다. 삼성전자가 도입한 챗GPT는 '기밀 유출 막아라'라는 키워드로 운영된다. 삼성전자는 '기밀 유출 막아라'라는 키워드로 운영된다. 삼성전자는 '기밀 유출 막아라'라는 키워드로 운영된다.

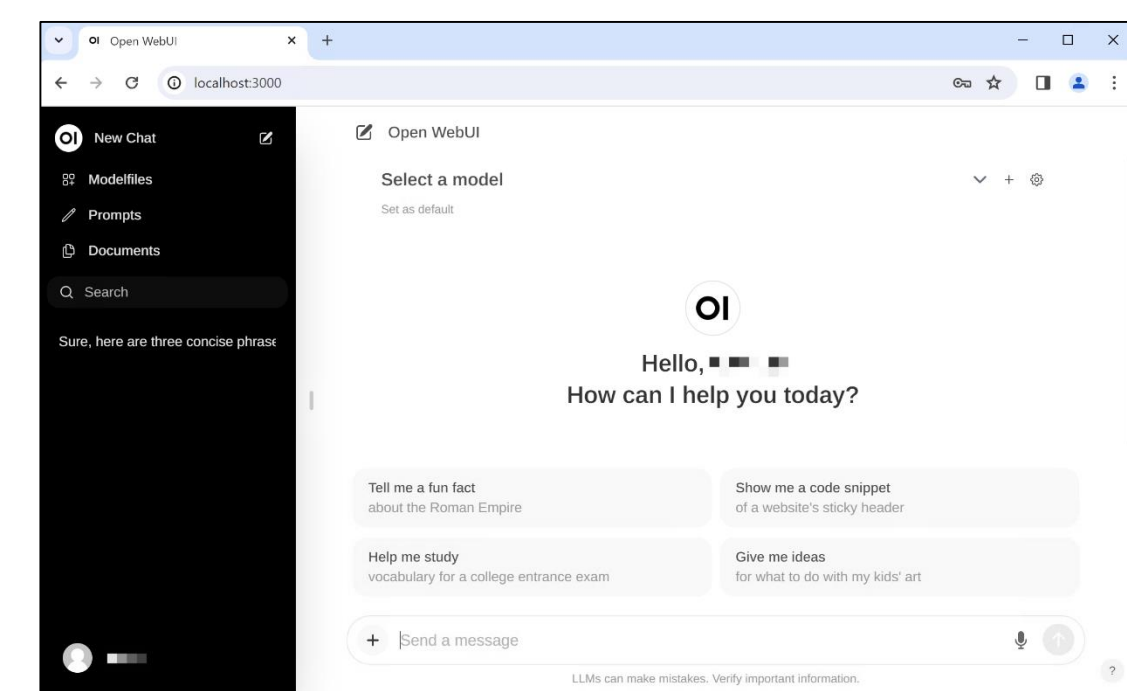
연구 진행 과정

프로젝트 구조도



진행 과정

1. OpenWebUI 필터 개발
 - 본 프로젝트 기능을 사용하기 위해서는 웹이 아닌 GPT API 토큰을 사용하여 추론을 진행해야 함.
 - API 질의를 사용하기 위해서는 Python과 같은 프로그래밍 언어를 사용해야 하며, 이는 비전문가에게 어려울 수 있음
- OpenWebUI는 GPT Web과 유사한 웹 UI를 제공하여, 웹에서만 ChatGPT를 사용한 사람들의 진입 장벽을 낮춰줌



- 필터를 개발하여, OpenWebUI 프레임워크 내 파이썬 라이브러리의 데이터를 컨트롤
- 이를 사용하여 사용자의 원본 회의록이 있는 프롬프트를 가로채 난독화 진행 후 GPT에게 전달 및 추론된 난독화 회의록의 요약본을 사용자에게 제공하기 전에 복호화하여 사용자의 화면에서는 원문 회의록 입력과 요약 결과만 확인할 수 있도록 개발

2. 회의록 내, 중요 키워드 선정 모듈 개발
 - 입력으로 제공되는 회의록에서 중요한 키워드를 알아서 추출하는 것이 목표
 - 회의록으로부터 구(Phrase)를 추출하여 중복 항목 제거, 조사 제거를 통해 1차 가공
 - 1차 가공한 구에서 '와', '과' 와 같은 접속 조사가 포함된 구를 분리하는 2차 가공

선택된 키워드

주요 키워드 목록

1. 키워드 선정 모듈 개발
2. 키워드 선정 모듈 개발
3. 키워드 선정 모듈 개발
4. 키워드 선정 모듈 개발
5. 키워드 선정 모듈 개발
6. 키워드 선정 모듈 개발
7. 키워드 선정 모듈 개발
8. 키워드 선정 모듈 개발
9. 키워드 선정 모듈 개발
10. 키워드 선정 모듈 개발

난독화 요약 결과

키워드 선정 모듈 개발

1. 키워드 선정 모듈 개발
2. 키워드 선정 모듈 개발
3. 키워드 선정 모듈 개발
4. 키워드 선정 모듈 개발
5. 키워드 선정 모듈 개발
6. 키워드 선정 모듈 개발
7. 키워드 선정 모듈 개발
8. 키워드 선정 모듈 개발
9. 키워드 선정 모듈 개발
10. 키워드 선정 모듈 개발

원문 회의록 요약 결과와 비교

원문 회의록 요약 결과와 비교

1. 원문 회의록 요약 결과와 비교
2. 원문 회의록 요약 결과와 비교
3. 원문 회의록 요약 결과와 비교
4. 원문 회의록 요약 결과와 비교
5. 원문 회의록 요약 결과와 비교
6. 원문 회의록 요약 결과와 비교
7. 원문 회의록 요약 결과와 비교
8. 원문 회의록 요약 결과와 비교
9. 원문 회의록 요약 결과와 비교
10. 원문 회의록 요약 결과와 비교

- 선정된 키워드는 회의록에서 Key1 과 같이 의미가 없는 단어로 변환하여 회의록 난독화

키워드 1: [오류 발생]	- 오류 발생
키워드 2: [인원 증가]	- 인원 증가
키워드 3: [계정 관리]	- 계정 관리
키워드 4: [데이터 분석]	- 데이터 분석
키워드 5: [시스템 점검]	- 시스템 점검
키워드 6: [유지보수]	- 유지보수
키워드 7: [기밀 유출]	- 기밀 유출
키워드 8: [데이터 백업]	- 데이터 백업
키워드 9: [보안 강화]	- 보안 강화
키워드 10: [신규 직원]	- 신규 직원
키워드 11: [예산 관리]	- 예산 관리
키워드 12: [고객 서비스]	- 고객 서비스
키워드 13: [프로젝트 진행]	- 프로젝트 진행
키워드 14: [팀 회의]	- 팀 회의
키워드 15: [업무 보고]	- 업무 보고
키워드 16: [시스템 업데이트]	- 시스템 업데이트
키워드 17: [데이터 마이그레이션]	- 데이터 마이그레이션
키워드 18: [신규 소프트웨어]	- 신규 소프트웨어
키워드 19: [보안 감사]	- 보안 감사
키워드 20: [인사 관리]	- 인사 관리
키워드 21: [예산 승인]	- 예산 승인

결과 및 분석

OpenWebUI 필터 개발 후 데이터 조작



GPT의 난독화된 회의록 안에 있는 키워드 예측률

키워드	예측률
오류 발생	0.8225
인원 증가	0.8225
계정 관리	0.8225
데이터 분석	0.8225
시스템 점검	0.8225
유지보수	0.8225
기밀 유출	0.8225
데이터 백업	0.8225
보안 강화	0.8225
신규 직원	0.8225
예산 관리	0.8225
고객 서비스	0.8225
프로젝트 진행	0.8225
팀 회의	0.8225
업무 보고	0.8225
시스템 업데이트	0.8225
데이터 마이그레이션	0.8225
신규 소프트웨어	0.8225
보안 감사	0.8225
인사 관리	0.8225
예산 승인	0.8225

회의록 난독화 및 요약 후 복호화하여 원문 회의록 요약 결과와 비교

원문 회의록 요약 결과

원문 회의록 요약 결과

난독화 회의록 요약 결과

난독화 회의록 요약 결과

분석

- OpenWebUI 필터 개발을 통해 정상적으로 사용자 채팅 데이터 수집 가능 확인
- 난독화 된 회의록을 넣었을 때, 중요 키워드 리스트를 GPT가 예측하지 못하는 것을 통해, 난독화 의미 있음 확인
- 요약된 난독화 회의록 결과와 회의록 원문 요약 결과에 유사도를 측정했을 때, 평균 0.8225의 유사도를 보임

오픈소스 URL

OpenWebUI 프로젝트 저장소 - <https://github.com/open-webui/open-webui/>

