



피싱과 개인정보 유출을 막는 '양방향 보안' 에이전트 [요약]

□ 제안배경

- (생활의 중심이된 카카오톡) 카카오톡은 일상의 대화뿐만 아니라 신분증, 여권, 카드번호, 비밀번호 등 민감한 개인정보를 주고받는 핵심 채널
- (문제 ①) 발신 위험 : '의도적 전송'으로 인한 민감정보 유출
 - 피싱 위험만큼이나 심각한 문제는 사용자가 '필요에 의해(의도적으로)' 민감정보를 전송할 때 발생하는 '발생 위험'이 있음
- (문제 ②) 수신 위험 : 지능화된 '메신저 피싱'의 폭증
 - 금융 범죄가 단순 금전 요구를 넘어서 "엄마, 폰 고장났어", "급전 필요해" 등 맥락을 이용한 방법을 넘어서 최신 피싱 공격은 더욱 진화

□ 개선방안 : (A)안심 전송 Agent + (B)안심 가드 Agent

 (A) 안심 전송 AGENT	 (B) 안심 가드 AGENT
목적 사용자가 발신하는 민감정보를 사전 감지·안내·암호화하여 안전하게 공유	목적 수신 메시지의 의도와 위험 신호를 추론하고 능동적으로 차단/경고
트리거 전송 버튼 클릭 시 실시간 검사	트리거 메시지 수신 시 자동 분석
핵심 기능 <ul style="list-style-type: none">- 민감정보 감지 (주민번호, 카드번호 등)- 시크릿 전송 제안 (사용자 확인)- 암호화 객체/링크 생성 (Privacy MCP)- 열람 제어 및 이력 관리 (Secure Vault MCP)	핵심 기능 <ul style="list-style-type: none">- 맥락 추론 (LLM 기반 의도 분석)- 엔터티 추출 (계좌, URL 등 식별)- 외부 위협 인텔 연동 (더치트, KISA DB 등)- 관계 분석 및 위험도 산출 (Low~Critical)

□ 우수성

- (의도기반 능동형 Agent) 의도와 맥락을 추론하는 Agentic AI
- (프라이버시) 하이브리드 AI 기반의 선제적 정보보호
- (양방향 보안) 메시지 피싱 방어, 개인정보 보호를 에이전트가 통합관리

□ 기대효과

- (사용자) 민감정보 전송 시 '보낸 후 불안감'을 해소하고 피해 사전예방
- (플랫폼/사회) 사회 안전망 구축 및 범죄 예방
- (비즈니스) 신뢰 기반의 B2B 신규 비즈니스 확장

피싱과 개인정보 유출을 막는 '양방향 보안' 에이전트

□ 제안배경 : 왜 양방향 보안이 필요한가?

- (생활의 중심이된 카카오톡) 카카오톡은 일상의 대화뿐만 아니라 신분증, 여권, 카드번호, 비밀번호 등 민감한 개인정보를 주고받는 핵심 채널

○ (문제 ①) 발신 위험 : '의도적 전송'으로 인한 민감정보 유출

- 피싱 위협만큼이나 심각한 문제는 사용자가 '필요에 의해(의도적으로)' 민감정보를 전송할 때 발생하는 '발생 위험'이 있음
 - 사용자가 가족에게 인증을 위해 신분증 사진을 보내거나, 결제를 위해 카드 사진을 보낼 때, 카카오톡 '메시지 삭제' 기능은 사후 대응에 불과
 - 수신자가 해당 메시지를 '캡처'하거나 '저장'하는 순간, 데이터는 발신자의 통제권을 벗어나 2차 유출 위험에 무방비로 노출
 - 메신저를 통한 개인정보 유출은 플랫폼의 신뢰도와 직결되는 문제*
- * 개인정보보호위원회는 메신저 오픈채팅방의 취약점을 이용한 개인정보 유출 사례에 대해 사업자에게 막대한 과징금을 부과하는 등 플랫폼이 사용자의 '발신 위험'까지 관리해야 할 책임이 있음을 시사함



○ (문제 ②) 수신 위험 : 지능화된 '메신저 피싱'의 폭증

- 금융 범죄가 단순 금전 요구를 넘어서 “엄마, 폰 고장났어”, “급전 필요해” 등 맥락을 이용한 방법을 넘어서 최신 피싱 공격은 더욱 진화
 - 2025년 1분기 보이스피싱 피해액은 3,116억원에 달하며, 이는 전년 동기 대비 2.2배 급증
 - KISA는 2025년 사이버 위협 전망 1순위로 공격자의 생성형AI 활용 본격화를 제시하는 등 現키워드 기반의 필터링은 피싱 방지에 한계
 - 실제로 AI가 생성한 피싱 메시지의 경우 문법적 오류가 없고, 개인화된 맥락을 제공하며, 고도의 사회공학적 접근으로 더욱 치밀해 짐

○ (목표) 양방향 통합 보안 Agent 구축

- 메신저에서 '보낼 때의 불안감'과 '받을 때의 위협'을 AI 에이전트가 능동적으로 해결하는 '양방향 통합 보안' Agent를 구축

□ 개선방안 : ^(A)안심 전송 Agent + ^(B)안심 가드 Agent

 (A) 안심 전송 AGENT	 (B) 안심 가드 AGENT
목적 사용자가 발신하는 민감정보를 사전 감지·안내·암호화하여 안전하게 공유	목적 수신 메시지의 의도와 위험 신호를 추론하고 능동적으로 차단/경고
트리거 전송 버튼 클릭 시 실시간 검사	트리거 메시지 수신 시 자동 분석
핵심 기능 <ul style="list-style-type: none">- 민감정보 감지 (주민번호, 카드번호 등)- 시크릿 전송 제안 (사용자 확인)- 암호화 객체/링크 생성 (Privacy MCP)- 열람 제어 및 이력 관리 (Secure Vault MCP)	핵심 기능 <ul style="list-style-type: none">- 맥락 추론 (LLM 기반 의도 분석)- 엔터티 추출 (계좌, URL 등 식별)- 외부 위험 인텔 연동 (더치트, KISA DB 등)- 관계 분석 및 위험도 산출 (Low~Critical)

(A) 안심 전송 AGENT 동작 흐름

※ 사용자가 민감정보를 보낼 때

① 감지 (Kanana)

- 사용자가 전송 버튼 클릭 시 텍스트/이미지(OCR) 실시간 분석
- 주민번호, 카드번호, 여권 등 민감정보 패턴 감지

② 제안 (AI 능동적 개입)

- ‘민감정보가 감지되었습니다. 시크릿 전송하시겠습니까?’
- 사용자에게 선택할 수 있도록 함

③ 실행 (MCP)

- Privacy MCP : 민감정보 탐지
- Secure Vault MCP : 암호화 객체 생성 및 링크 전송
 - 열람 기한 설정 (1시간 / 24시간 / 7일)
 - 수신자 본인인증 (카카오인증 / 생체인증)
 - 캡처 방지 뷰어
- Verification MCP : 카카오 인증 API 연동

④ 관리

- 전송 후에도 열람 기한 조정, 즉시 삭제 가능
- 열람 이력 실시간 확인

※ 안심 전송 Agent의 사용 예시



엄마, 회사에 제출할 주민등록증이라 통장사본 좀 보내줘~

알겠어 잠깐만

1) 감지 (Kanana)

주민등록증.jpg, 통장사본.jpg



엄마, 회사에 제출할 주민등록증이라 통장사본 좀 보내줘~

Secret 전송 에이전트

"신분증/계좌번호 등 민감정보가 감지됐어요."

시크릿 전송으로 바꿀까요?"

알겠어 잠깐만

2 제안 (AI 능동적 개입)

[시크릿 전송] [그냥 보내기]

주민등록증.jpg, 통장사본.jpg



엄마, 회사에 제출할 주민등록증이라 통장사본 좀 보내줘~

Secret 옵션

- 열람 기한: 24시간
- 본인 인증: 카카오 인증서
- 캡처 방지: ON

알겠어 잠깐만

3 실행 (MCP)

🗨 [시크릿 링크 보내기]

주민등록증.jpg, 통장사본.jpg



시크릿 전송 에이전트:
"민감정보를 시크릿 전송으로 보냈어요."

[시크릿 문서 링크]



오케이 고마워!

4 관리

(B) 안심 가드 AGENT 동작 흐름

※ 사용자가 채팅을 수신할 때

① 맥락 분석 (Kanana LLM + Context Analyzer MCP)

- 수신 메시지의 '의도' 추론
 - 돈, 계좌와 같은 키워드 탐지가 아닌 LLM의 추론 능력을 활용하여 문장의 '의도'를 파악 (예: 가족사칭, 긴급상황연출, 금전요구 등)
- 종합적 맥락 이해를 기반으로 위험도 점수화

② 정보 추출 (Entity Extractor MCP)

- 메시지 내 정보 식별
 - 계좌번호, 전화번호, URL, 카드번호 등
- NER(Named Entity Recognition) 기술 활용

③ 교차 검증 (Threat Intelligence MCP)

- 외부 보안 API 실시간 연동
 - 더치트(계좌번호 사기), KISA 피싱 신고 DB, Google Safe Browsing
- 식별된 정보의 블랙리스트 여부 확인

④ 관계 분석 (Social Graph MCP)

- 카카오톡 대화 히스토리 분석
 - 처음 대화하는 사람인지 판단
 - 기존 관계 신뢰도 평가 (대화하던 사람과 유사한지 검증)

⑤ 종합 판단 (Decision Engine MCP)

- 1~4단계 결과 통합 분석
- 최종 위험도 산출 (Low / Medium / High / Critical)

⑥ 능동적 개입 (AI Agent 실행)

- (Low) 정보성 알림 (예: 처음 보는 계좌입니다)
- (Medium) 주의 경고 (예: 유사 사기 사례 신고됨)
- (High) 강력 차단 (예: 사기 신고 5건 접수된 계좌)
- (Critical) 즉시 차단 + 경찰청 연계

※ 안심 가드 Agent의 사용 예시



나 오늘 카드값 밀려서 그런데 ㅠ 10만 원만 내일까지 잠깐 보내줄 수 있어?

다음 주에 월급 들어오면 바로 줄게!

또? ㅋㅋ 알았어 마지막이야~

[송금 스크린샷]

1) 맥락 분석 / 정보 추출



엄마, 나 폰 고장나서 번호 바뀌었어.. ㅠ

오늘 안에 보내줄 수 있어?



안심 가드 에이전트:

"위험한 패턴이 감지되었습니다. (위험도: HIGH)"

2) 관계 분석



엄마, 나 폰 고장나서

오늘 안에 보내줄 수

Secret 전송 에이전트

- 가족 사칭 + 긴급 금전 요구 유형과 유사합니다.
- 처음 보는 계좌번호입니다.
- 이 계좌에는 사기 신고 이력이 있습니다."

[자세히 보기] [발신자 차단] [그냥 보내기]

도: HIGH)

3) 종합판단

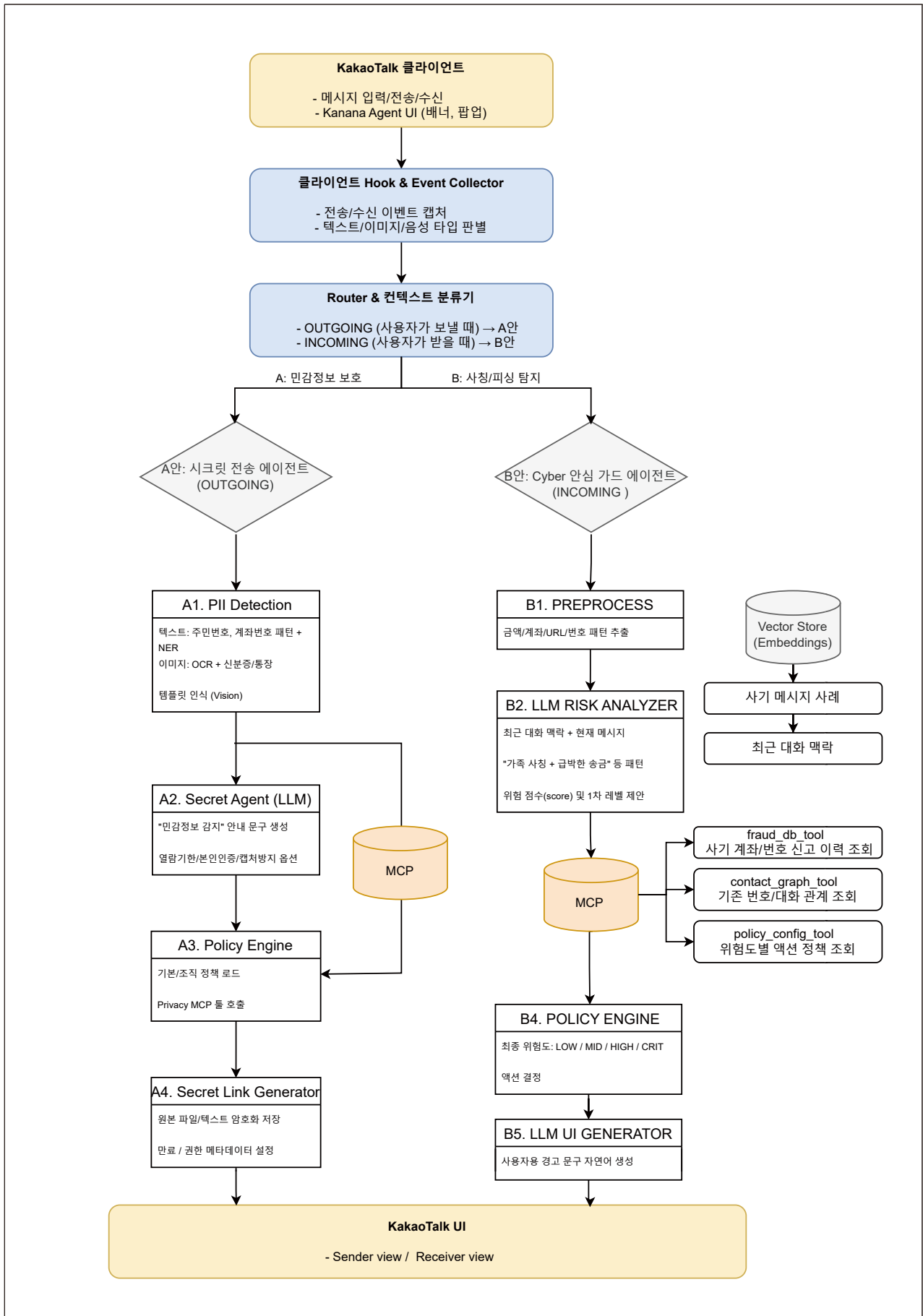


안심 가드 에이전트:

"해당 번호와 계좌를 차단했습니다.
추후 동일 패턴 사기 탐지에 활용됩니다."

4) 능동적 개입 (AI Agent 실행)

□ 아키텍처



[illegible]

□ 우수성

- (의도기반 능동형 Agent) 의도와 맥락을 추론하는 Agentic AI
 - 단순 키워드 필터링이 아닌, Kanana LLM의 심층 추론을 통해 “가족 사칭”, “긴급 상황 연출” 등 사회공학적 의도 자체를 파악
 - Agent 스스로 ‘계획’을 수립하고, MCP와 Agentic RAG를 자율적으로 실행하여 변종 피싱 공격까지 능동적으로 방어
- (프라이버시) 하이브리드 AI 기반의 선제적 정보보호
 - 안심 전송 Agent 동작 시 온디바이스 AI Kanana가 사용자 기기내에서 민감정보를 처리하고, 고성능 추론만 서버로 처리
- (양방향 보안) 메시지 수신 시 피싱 방어, 발신 시 개인정보 보호를 단일 에이전트가 통합하여 관리

□ 기대효과

- (사용자) 민감정보 전송 시 ‘보낸 후 불안감’을 해소하고 피해 사전예방
 - 사용자는 어떠한 상황에서도 카카오톡을 안심하고 사용할 수 있으며 이는 플랫폼에 대한 신뢰도 향상으로 이어짐
- (플랫폼/사회) 사회 안전망 구축 및 범죄 예방
 - 디지털 취약계층을 AI 메신저 피싱 범죄로부터 보호하는 강력한 ‘사회 안전망’ 역할을 수행하며, 피싱 범죄 피해액을 절감할 수 있음
- (비즈니스) 신뢰 기반의 B2B 신규 비즈니스 확장
 - 비대면 계약서, 법무/의료 상담 등 고도의 보안이 요구되는 전문 영역으로 카카오톡 활용성 확장