

KRA-AiGov 솔루션 개발 계획 업데이트에 대한 피드백 및 권고 사항

1. 개발 환경 및 도구 검토

크라세 대표님께서 제공해주신 개발 환경 및 도구 정보는 다음과 같습니다.

- **하드웨어:** GPU 탑재 게임용 노트북
- **운영체제:** Ubuntu 22.04
- **개발 환경:** Docker 컨테이너
- **버전 관리:** GitHub
- **편집기:** VS Code (보조: Cursor AI)

이러한 개발 환경은 AI 서비스 관련 솔루션 개발에 여러 가지 이점을 제공합니다.

1.1 하드웨어 및 운영체제

GPU 탑재 게임용 노트북과 Ubuntu 22.04 조합은 AI/ML 개발에 매우 적합한 환경입니다. GPU는 딥러닝 모델 학습 및 추론에 필수적인 연산 능력을 제공하며, Ubuntu는 개발자 친화적인 환경과 다양한 AI/ML 라이브러리 및 도구 지원을 통해 개발 효율성을 높일 수 있습니다.

- **장점:**
 - **고성능 GPU:** AI 모델 개발 및 테스트 시 빠른 연산 처리 속도를 제공하여 개발 시간을 단축할 수 있습니다.
 - **Ubuntu 22.04:** 안정적이고 광범위한 커뮤니티 지원을 통해 개발에 필요한 다양한 패키지 및 드라이버 설치가 용이합니다. 또한, Docker와 같은 컨테이너 기술과의 호환성이 뛰어납니다.
- **고려 사항:**
 - **드라이버 설치 및 관리:** GPU 드라이버(NVIDIA CUDA 등) 및 관련 라이브러리(cuDNN)의 올바른 설치와 버전 관리가 중요합니다. 이는 딥러닝 프레임워크(TensorFlow, PyTorch 등)와의 호환성에 직접적인 영향을 미칩니다.
 - **자원 관리:** 게임용 노트북이라 하더라도, 복잡한 AI 모델 학습 시에는 시스템 자원(RAM, CPU)이 부족할 수 있습니다. Docker 컨테이너 환경에서 자원 할당을 효율적으로 관리하는 것이 필요합니다.

1.2 개발 환경 (Docker 컨테이너)

Docker 컨테이너를 활용한 개발은 현대 소프트웨어 개발에서 권장되는 방식이며, KRA-AiGov 솔루션 개발에도 큰 이점을 제공합니다.

- **장점:**

- **환경 일관성:** 개발, 테스트, 운영 환경 간의 일관성을 보장하여 "내 컴퓨터에서는 되는 데..."와 같은 문제를 방지합니다. 이는 향후 운영 환경으로의 이관을 용이하게 합니다.
- **의존성 관리:** 프로젝트별로 필요한 라이브러리 및 종속성을 격리하여 관리할 수 있어, 개발 환경 충돌을 최소화합니다.
- **확장성 및 배포 용이성:** 컨테이너화된 애플리케이션은 확장 및 배포가 용이하며, Kubernetes와 같은 컨테이너 오케스트레이션 도구와 연계하여 효율적인 운영이 가능합니다.

- **고려 사항:**

- **GPU 패스스루:** Docker 컨테이너 내에서 GPU를 활용하기 위해서는 NVIDIA Container Toolkit과 같은 도구를 사용하여 GPU를 컨테이너에 올바르게 패스스루(pass-through)해야 합니다. 초기 설정에 대한 이해가 필요합니다.
- **볼륨 마운트:** 개발 중인 코드나 데이터를 컨테이너 내부와 공유하기 위해 볼륨 마운트 설정을 적절히 해야 합니다. 이는 개발 편의성과 데이터 영속성을 보장합니다.

1.3 버전 관리 (GitHub)

GitHub를 통한 버전 관리는 협업 및 코드 관리의 표준적인 방식입니다.

- **장점:**

- **협업:** 여러 개발자가 동시에 작업하고 코드를 병합하는 데 효율적입니다.
- **버전 이력 관리:** 코드 변경 이력을 추적하고, 필요한 경우 이전 버전으로 쉽게 되돌릴 수 있습니다.
- **코드 리뷰:** Pull Request 기능을 통해 코드 리뷰 프로세스를 도입하여 코드 품질을 향상시킬 수 있습니다.

- **고려 사항:**

- **민감 정보 관리:** AI 서비스 관련 솔루션의 특성상 민감한 코드나 설정 정보(API 키 등)가 포함될 수 있습니다. GitHub에 직접 업로드하기보다는 `.gitignore` 파일을 통해 민감 정보를 제외하거나, 환경 변수, Secret Management 서비스 등을 활용하여 안전하게 관리해야 합니다.
- **운영 환경 이관:** GitHub에 있는 코드를 운영 환경으로 이관할 때, CI/CD(Continuous Integration/Continuous Deployment) 파이프라인을 구축하여 자동화된 배포 프로세스를 고려하는 것이 좋습니다.

1.4 편집기 (VS Code 및 Cursor AI)

VS Code와 Cursor AI의 조합은 개발 생산성을 크게 향상시킬 수 있는 좋은 선택입니다.

- **장점:**

- **VS Code:** 강력한 기능, 풍부한 확장 프로그램, 사용자 친화적인 인터페이스를 제공하여 개발 편의성을 높입니다.
- **Cursor AI:** AI 기반 코드 자동 완성, 코드 생성, 디버깅 지원 등 개발 생산성을 극대화하는데 도움을 줍니다. 특히 AI 관련 코드 작성 시 유용할 수 있습니다.

- **고려 사항:**

- **AI 보조 도구의 한계:** Cursor AI와 같은 AI 보조 도구는 개발 생산성을 높여주지만, 생성된 코드의 정확성이나 보안 취약점 여부를 항상 검토해야 합니다. AI가 생성한 코드에 대한 의존도를 맹신하기보다는, 개발자의 검증과 이해가 필수적입니다.
- **성능:** AI 보조 도구가 시스템 자원을 많이 사용할 수 있으므로, 노트북의 성능과 VS Code의 확장 프로그램 사용량을 적절히 조절하는 것이 좋습니다.

전반적으로, 제시해주신 개발 환경은 KRA-AiGov 솔루션 개발에 매우 적합하며 효율적인 개발을 위한 기반을 잘 갖추고 있습니다. 각 요소별 고려 사항을 염두에 두시면 더욱 안정적이고 생산적인 개발이 가능할 것입니다.

2. 모듈명 검토

이전에 잠정 결정했던 모듈명(PromptShield, Shadoweye, DashAILog, FortLLM)에서 PromptGate, ShadowScan, DashIQ, TrustLLM으로 변경하신 것을 확인했습니다. 변경된 모듈명에 대한 검토 의견은 다음과 같습니다.

2.1 변경된 모듈명 분석

- **PromptGate:** PromptShield 에서 PromptGate 로 변경되었습니다. Gate 는 Shield 보다 통제 및 관문의 역할을 더 명확하게 나타내므로, 프롬프트 필터링 프록시의 기능적 역할을 직관적으로 전달합니다. 상표권 측면에서는 PromptGate 가 PromptShield 보다 유사 상표 존재 가능성이 더 높았으나, Gate 라는 단어가 일반적이므로 특정 분야에서 고유성을 확보한다면 문제가 없을 수 있습니다. 다만, 상표권 재검토가 필요합니다.
- **ShadowScan:** Shadoweye 에서 ShadowScan 으로 변경되었습니다. Scan 은 탐지 및 분석 행위를 명확하게 나타내므로, Shadow AI 탐지 시스템의 역할을 잘 표현합니다. Shadoweye 는 유사 상표 존재 가능성이 있었으나, ShadowScan 역시 기존에 SHADOWSCAN 이라는 상표가 등록된 이력이 있으므로, 이 부분에 대한 재검토가 필요합니다. (이전 검토에서 DarkWatch 대신 Shadoweye 를 제안했으나, ShadowScan 으로 변경되었으므로 다시 확인이 필요합니다.)
- **DashIQ:** DashAILog 에서 DashIQ 로 변경되었습니다. Dash 는 대시보드를 의미하고, IQ 는 지능 또는 정보의 질을 의미하는 것으로 보입니다. DashAILog 가 AI 사용 로그라는 직

관적인 의미를 가졌다면, DashIQ 는 좀 더 추상적이지만, AI 사용 현황에 대한 지능적인 분석 및 통찰력을 제공한다는 의미를 내포할 수 있습니다. 상표권 측면에서는 DashAILog 와 마찬가지로 현재까지 등록된 상표는 확인되지 않았습니다.

- **TrustLLM:** FortLLM 에서 TrustLLM 으로 변경되었습니다. Fort 가 요새, 보안을 강조했다면, Trust 는 신뢰를 강조합니다. 내부 LLM의 경우 기업 내부에서 신뢰할 수 있는 AI 모델을 제공한다는 의미를 강조하는 데 적합합니다. iTrustLLM 은 유럽 내 유사 상표 존재 가능성이 있었으나, TrustLLM 은 i 가 제거되어 좀 더 일반적인 명칭이 되었습니다. 하지만 TrustLLM 이라는 용어 자체가 AI 신뢰성 분야에서 널리 사용되므로, 상표권 재검토가 필요합니다.

2.2 모듈명 변경에 대한 권고 사항

전반적으로 변경된 모듈명들은 각 기능의 역할을 더 명확하게 나타내거나, 솔루션의 지향점을 잘 반영하고 있습니다. 다만, 상표권 측면에서 재검토가 필요한 부분이 있습니다.

- **상표권 재검토:** PromptGate , ShadowScan , TrustLLM 에 대해서는 다시 한번 상표권 및 유사 상표 존재 여부를 면밀히 검토하는 것을 권장합니다. 특히 ShadowScan 은 기존에 등록 이력이 있는 명칭이므로, 법률 전문가와 상의하여 사용 가능 여부를 판단하거나 대체 명칭을 고려하는 것이 안전합니다.
- **일관성 유지:** 변경된 모듈명들이 KRA-AiGov라는 메인 브랜드명과 잘 어울리는지, 그리고 솔루션 전체의 아이덴티티를 일관되게 유지하는지 확인하는 것이 중요합니다.
- **내부 커뮤니케이션:** 변경된 모듈명에 대해 개발팀 및 관련 이해관계자들에게 충분히 설명하고, 명칭 변경의 배경과 의미를 공유하여 혼란을 최소화해야 합니다.

새로운 모듈명은 솔루션의 특징을 잘 나타내고 있지만, 법적 리스크를 최소화하기 위한 추가적인 검토가 필수적입니다. 특히 상표권 문제는 추후 큰 문제로 이어질 수 있으므로 신중하게 접근해야 합니다.

3. PromptGate 구현 방안 검토

PromptGate를 Rebuff에서 제공하는 Python-SDK로 구현하고자 하는 계획에 대해 검토하고 의견을 드립니다.

3.1 Rebuff (Python-SDK) 분석

Rebuff는 LLM(Large Language Model) 기반 애플리케이션을 프롬프트 인젝션(Prompt Injection) 공격으로부터 보호하기 위해 설계된 오픈소스 프레임워크입니다. 다층 방어 메커니즘을 통해 AI 애플리케이션의 보안을 강화하며, 특히 다음과 같은 특징을 가집니다.

- **자기 강화(Self-hardening) 기능:** 공격 시도에 노출될 때마다 스스로를 강화하여 유사한 공격을 미래에 인식하고 방지할 수 있도록 벡터 데이터베이스에 이전 공격의 임베딩을 저장합니다.

- **다층 방어:** 휴리스틱 규칙, 벡터 유사성 검색, LLM 기반 탐지 등 여러 계층의 탐지 방법을 사용합니다.
- **오픈소스 및 API 연동:** Python-SDK 형태로 제공되어 기존 LLM 서비스에 쉽게 통합할 수 있습니다.

3.2 PromptGate에 Rebuff 적용 검토

PromptGate의 주요 기능은 사용자의 AI 서비스 요청(프롬프트)을 가로채어 사전 정의된 정책에 따라 필터링 및 마스킹 처리하고, 민감 정보 유출을 방지하는 것입니다. Rebuff는 특히 프롬프트 인젝션 공격 방어에 특화되어 있으므로, PromptGate의 핵심 기능 중 하나인 '입력 필터링'에 매우 효과적으로 활용될 수 있습니다.

• 장점:

- **프롬프트 인젝션 방어 강화:** Rebuff의 자기 강화 및 다층 방어 메커니즘을 통해 PromptGate의 프롬프트 인젝션 방어 능력을 크게 향상시킬 수 있습니다. 이는 KRA-AiGov 솔루션의 보안 신뢰도를 높이는 데 기여합니다.
- **개발 시간 단축:** 프롬프트 인젝션 탐지 로직을 직접 개발하는 대신, 검증된 오픈소스 라이브러리인 Rebuff를 활용함으로써 개발 시간을 단축하고 안정성을 확보할 수 있습니다.
- **유연한 통합:** Python-SDK 형태로 제공되므로, 기존 Python 기반의 PromptGate 구현 계획과 쉽게 통합될 수 있습니다.

• 고려 사항:

- **민감 정보 마스킹:** Rebuff는 주로 프롬프트 인젝션 탐지에 초점을 맞추고 있습니다. PromptGate의 또 다른 핵심 기능인 '민감 정보 탐지 및 자동 마스킹' 기능은 Rebuff만으로는 완벽하게 구현하기 어려울 수 있습니다. 따라서 민감 정보 마스킹을 위해서는 별도의 NLP 기반 민감 정보 탐지 및 처리 모듈을 Rebuff와 함께 연동하여 사용해야 합니다.
- **성능 오버헤드:** Rebuff의 다층 방어 메커니즘은 프롬프트 처리 과정에서 추가적인 연산 오버헤드를 발생시킬 수 있습니다. 특히 대규모 트래픽을 처리해야 하는 경우, 성능 최적화에 대한 고려가 필요합니다.
- **커스터마이징 필요성:** Rebuff는 일반적인 프롬프트 인젝션 공격에 대한 방어를 제공하지만, 기업 특유의 보안 정책이나 특정 유형의 민감 정보에 대한 맞춤형 필터링/마스킹 로직을 추가해야 할 수 있습니다. 이 경우 Rebuff의 확장성을 고려하여 커스터마이징 방안을 마련해야 합니다.

3.3 권고 사항

Rebuff를 PromptGate 구현에 활용하는 것은 매우 긍정적인 접근 방식입니다. 다음 사항들을 고려하여 구현 계획을 구체화하는 것을 권고합니다.

1. **Rebuff와 민감 정보 처리 모듈의 통합:** Rebuff를 프롬프트 인젝션 방어 핵심으로 사용하되, 민감 정보 탐지 및 마스킹을 위한 별도의 모듈(예: 정규표현식 기반 필터링, NLP 라이브러리 활용)을 개발하여 PromptGate 내에서 함께 작동하도록 설계합니다.
2. **성능 테스트 및 최적화:** MVP 개발 단계에서 Rebuff 적용 시 발생할 수 있는 성능 오버헤드를 측정하고, 필요한 경우 캐싱 전략, 비동기 처리 등을 통해 성능을 최적화하는 방안을 모색합니다.
3. **정책 엔진과의 연동:** PromptGate의 정책 엔진이 Rebuff의 탐지 결과를 활용하여 특정 프롬프트에 대한 처리 방식(차단, 경고, 마스킹 등)을 결정할 수 있도록 연동 방안을 설계합니다.
4. **지속적인 업데이트 및 관리:** Rebuff는 오픈소스 프로젝트이므로, 최신 보안 위협에 대응하기 위해 정기적인 업데이트 및 취약점 관리가 필요합니다.

Rebuff는 PromptGate의 핵심적인 보안 기능을 강화하는 데 큰 도움이 될 것입니다. 다만, PromptGate의 모든 기능을 Rebuff 하나로 해결하기보다는, Rebuff의 강점을 활용하면서 부족한 부분을 보완하는 통합적인 접근 방식이 필요합니다.

4. 종합 의견 및 권고 사항 제시

크라세 대표님께서 제시해주신 KRA-AiGov 솔루션의 개발 환경, 모듈명, 그리고 PromptGate 구현 계획에 대한 종합적인 의견과 권고 사항은 다음과 같습니다.

4.1 종합 의견

전반적으로 KRA-AiGov 솔루션의 개발 계획은 매우 구체적이고 현실적입니다. 특히 GPU 탑재 노트북과 Ubuntu, Docker, GitHub, VS Code/Cursor AI 조합은 AI 서비스 관련 솔루션 개발에 최적화된 환경을 제공하며, 효율적인 개발을 위한 기반을 잘 갖추고 있습니다. 모듈명 변경 또한 각 기능의 역할을 명확히 하고 솔루션의 지향점을 잘 반영하고 있습니다. PromptGate에 Rebuff를 활용하려는 계획은 프롬프트 인젝션 방어라는 핵심 보안 기능을 강화하는 데 매우 효과적인 전략입니다.

4.2 권고 사항

성공적인 KRA-AiGov 솔루션 개발 및 출시를 위해 다음 권고 사항들을 고려해주시기 바랍니다.

1. **개발 환경 설정의 안정화:** GPU 드라이버(NVIDIA CUDA) 및 관련 라이브러리(cuDNN)의 정확한 설치와 버전 관리에 특히 주의를 기울여야 합니다. Docker 컨테이너 내에서 GPU를 효율적으로 활용하기 위한 NVIDIA Container Toolkit 설정 및 볼륨 마운트 전략을 초기에 명확히 수립하는 것이 중요합니다.

2. **모델명 상표권 재검토 및 법적 자문:** 변경된 모델명(PromptGate, ShadowScan, DashIQ, TrustLLM) 중 PromptGate, ShadowScan, TrustLLM은 유사 상표 존재 가능성이 있거나 과거 등록 이력이 있는 명칭이므로, 법률 전문가(변리사)와 상의하여 상표권 침해 여부를 면밀히 검토하고 필요한 경우 대체 명칭을 고려하거나 정식 상표 등록을 진행하는 것이 필수적입니다. 특히 ShadowScan은 기존 등록 이력이 명확하므로 신중한 접근이 필요합니다.
3. **PromptGate 구현 시 Rebuff 활용 전략 구체화:** Rebuff는 프롬프트 인젝션 방어에 강력한 도구이지만, PromptGate의 모든 기능을 대체할 수는 없습니다. 민감 정보 탐지 및 자동 마스킹 기능은 별도의 NLP 기반 모듈을 개발하여 Rebuff와 연동하는 통합적인 접근이 필요합니다. 또한, 대규모 트래픽 처리 시 발생할 수 있는 성능 오버헤드에 대비하여 초기 단계부터 성능 테스트 및 최적화 방안을 고려해야 합니다.
4. **보안 취약점 관리 및 지속적인 업데이트:** AI 보조 도구(Cursor AI)가 생성하는 코드의 보안 취약점 여부를 항상 검토하고, GitHub에 민감 정보가 노출되지 않도록 .gitignore 설정 및 Secret Management 전략을 철저히 수립해야 합니다. Rebuff와 같은 오픈소스 라이브러리 또한 최신 보안 위협에 대응하기 위해 정기적인 업데이트 및 취약점 관리가 필요합니다.
5. **MVP 개발 로드맵의 유연성 유지:** 제시된 MVP 개발 로드맵은 현실적이지만, 실제 개발 과정에서 예상치 못한 문제나 새로운 요구사항이 발생할 수 있습니다. 애자일(Agile) 방법론을 적용하여 유연하게 로드맵을 조정하고, 사용자 피드백을 적극적으로 반영하는 것이 중요합니다.

KRA-AiGov 솔루션이 성공적으로 개발되어 기업 내 AI 서비스 사용의 보안과 효율성을 동시에 확보하는 데 기여하기를 기대합니다. 추가적인 문의사항이나 지원이 필요하시면 언제든지 말씀해주십시오.