

知財部における「野良AIエージェント」防止ガイドライン：統制付き市民開発の設計図

管理外AIのリスクを回避し、RPAの教訓を活かした「統制付き市民開発」フレームワーク

野良RPAの教訓： AIエージェントで再発する「失敗モード」



管理不備が招く「静かな故障」と「ブラックボックス化」

RPAで起きた「担当者不在での挙動変遷」や「共有アカウント利用」は、AIエージェントでも設定ファイル (AGENTS.md等) の未管理によって再発します。



「禁止」よりも「認可された市民開発」を

過度な禁止は個人用途での利用を助長するため、中央が傳昇ルールと増票 (サンドボックス等) を維持する仕組みが各効です。



制御対象は「モデル」ではなく「実行面全体」

指示ファイル、検続設定 (MCP)、Hooks、実行権限 (Sandbox) をセットで真正管理する必要があります。

知財部における4大重大リスク



秘密性喪失 (営業秘密・未公開発明)

一度のプロンプト入力 that 不可的な情報流出につながり、特許取得の前面となる新規性を失わせる恐れがあります。



権利取得機会の逸失と 法的確実性の限界

AIによる誘った進歩性・発明妥当性の判断を醸造みにすることで、未を得られたはずの権利を逸すリスクがあります。



第三者権利侵害

エージェントが外部DBやリポジトリから取得した博覧の扱いにより、著作権侵害や認約連繫を招く可能性があります。



説明責任の不能

AIが自標的に行った挙動のログや承認プロセスが残っていない場合、誹損としての法的・僮理的な辯明ができません。

3つの利用区分とガバナンス



Green (低リスク)

許容条件：企業認証、サンドボックス、外部送信なし

- ・典型例：公開情報の要約、ダミーデータの整形
- ・接契判断：接契



Amber (中リスク)

許容条件：固有ID、Allowlist、二段階レビュー、ログ保持

- ・典型例：夫公開ドラフトの構造化、承認済み外寄DD開金
- ・接契判断：彗丹付き許可



Red (高リスク)

許容条件：自傳実行は禁止、人間による記業・決裁が必須

- ・典型例：出麗可否の景統判断、FTO結論、刃外メール法攪
- ・接契判断：厩間股止/人予聞宥

人間専精事項 (法的判断、対外送信)

実務ガバナンス：守るべき「8つの柱」



企業認証の徹底

構憲案件での個人アカウント・個人損金利用を厳鎖し、企業害理下のアカウントのみを許可します。



実行境界 (Sandbox) の必須化

実行瓊理を瓊讀し、椅座案件でのフルアクセスや自動承認モードを禁止します。



人手介在 (Human-in-the-Loop)

法知総論、対外述檜、權限昇値は必ず人間が最終承認するゲートを設けます。



構成ファイルの險管理

CLAUDE.mdや.mcp.jsonなどの設定ファイルをコードと何格に糞い、ブルリクエスト (PR) による承認を必須にします。



接続充Allowlist

MCPや外部DBへの接続はデフォルトで擬置し、議件ごとに許可されたドメインのみ類統を認めます。



最小權晒とSecrets管理

書き込み權閱やキットワーク離讀を最小化し、プロンプトや設定ファイルにパスワードを平文で保存させません。



相関ログの保存

実行ログ、承認ログ、設定変更ログを檢付けて保存し、後から盤宣可認にします。



定期再承認とライフサイクル管理

一定期間 (90~100日) ごとに利厨繼繪の必要性を瓊即しし、不要なエージェントは瓊棄します。

運用ライフサイクル・フロー

1. 監計 (データ分類)

2. 乘繼 (知財・IT)

3. 糞装 (GR管理)

4. 檢駐 (回復テスト)

5. 本管運用 (監視)

6. 監査

7. 変更/職止