

# 知財部門における「野良AIエージェント」の脅威とエージェントック・ガバナンス戦略： Claude Code時代の自律型AIRISK管理

Gemini 3.1 pro

## 1. 序論：野良RPAの教訓から「野良AIエージェント」という未知のパラダイムへの移行

デジタルトランスフォーメーションの歴史において、業務部門主導のIT導入（いわゆるシチズンデベロップメント、または市民開発）は、常に「イノベーションの加速」と「ガバナンスの喪失」という二面性を企業に突きつけてきた。かつて、RPA（Robotic Process Automation）ツールが現場に普及した際、業務効率は劇的に向上した一方で、IT部門の管理下から外れた「野良RPA」が大量に発生するという深刻な問題が引き起こされた。これら管理外のRPAは、背後にある社内システムのインターフェースや仕様がわずかに変更されただけで突如として機能不全に陥り、業務プロセスのブラックボックス化や、予期せぬシステムへの過剰負荷、さらには権限管理の不備によるセキュリティインシデントの温床となった。

しかし現在、企業の中核組織、とりわけ最重要機密を取り扱う知的財産（知財）部門が直面しているのは、野良RPAがもたらした混乱を遥かに凌駕する潜在的破壊力を秘めた「野良AIエージェント」の脅威である。RPAが基本的に「あらかじめ定められたルールと手順を愚直に反復するだけの受動的な自動化ツール」であったのに対し、現代のAIエージェントは「高度な推論能力を持ち、自律的に判断を下し、外部のツールやシステムを操作して環境に能動的な影響を与えるエンティティ」である<sup>1</sup>。

この劇的な変化を牽引しているのが、AnthropicのClaude CodeやGitHubのCodexに代表される高度なAIコーディングアシスタント、そして自然言語によるアプリケーション生成プラットフォームの台頭である<sup>3</sup>。これらの技術により、プログラミングの専門知識を持たないビジネスユーザーであっても、自然言語で指示を与えるだけで、複雑な論理処理や外部データソースとの連携を実行する自律型AIエージェントを極めて短時間で構築することが可能になった<sup>3</sup>。

知財部門は、膨大な特許文献のスクリーニング、競合他社の知財動向分析、ライセンス契約書のレビューやリスク条項の自動検出など、高度な専門知識と自然言語処理能力、さらには「曖昧な判断（Squishy judgment calls）」を要求されるタスクが山積している領域である<sup>5</sup>。知財担当者が、日々の過酷な業務を効率化するために、ローカル環境で個人的にClaude Codeを用いて自律的なスクリプトやエージェントを生成・稼働させた瞬間、IT部門の検知網をすり抜けた「野良AIエージェント」が誕生する。これらが管理されずに放置されることは、未公開の特許情報やM&Aに関する極秘契約書が、エージェントの誤作動や外部の悪意あるプロンプトを通じて意図せず流出する致命的な情報漏洩リスクに直結する。本稿では、知財部門における野良AIエージェントの発生メカニズムとアーキテクチャ上の脆弱性を解き明かし、セキュリティ、アイデンティティ管理、そしてエージェントの推論プロセスの可観測性（オブザーバビリティ）の観点から、これらを安全に統制しつつビジネス価値を最大化

するための「エージェントック・ガバナンス」戦略を包括的に提示する。

## 2. シチズンデベロップメントの加速とアプリの爆発的増加

### 2.1. ソフトウェア開発の民主化とコストのゼロ化

AIをビジネスプロセスに統合することは、長らく高度な専門性を持つエンジニアの専権事項であり、技術的なハードルが極めて高かった。しかし、生成AIプラットフォームの進化は、この前提を根底から覆しつつある。フォレスター・リサーチが実施した2025年の開発者調査によれば、ソフトウェア開発部門のエグゼクティブの89%が、自身の企業においてシチズンデベロッパー戦略を現在導入している、あるいは積極的に計画中であると回答している<sup>5</sup>。さらに、企業がビジネスユーザーに開発を許可するローコードアプリケーションの類型として、AIが組み込まれたアプリケーションがトップに挙げられている<sup>5</sup>。

この背景にあるのは、自然言語による「バイブコーディング (Vibe coding)」と呼ばれるプロンプトベースの開発手法の確立と、アプリケーション作成コストの劇的な低下である<sup>3</sup>。かつて、エンタープライズ環境においては「ユーザー10人につき1つのアプリケーションが存在する」という経験則があり、IT部門は各アプリケーションの導入審査、セキュリティ評価、ライフサイクル管理を十分に実行することができた<sup>3</sup>。しかし、ある調査が指摘するように、ビジネスユーザーがClaudeなどのAIに求めるダッシュボードやワークフローの仕様を記述するだけで、わずか12分で機能的なアプリケーションが完成してしまう現在、この比率は完全に逆転している<sup>3</sup>。従業員一人ひとりが、日々のコーヒーブレイクの間に、データ分析ボットやレポート生成エージェントといった「マイクロアプリ」を数十個単位で作成する時代が到来しており、1,000人規模の企業では管理すべきアプリケーションが100個から10,000個へと爆発的に増加する可能性が示唆されている<sup>3</sup>。

### 2.2. 知財業務における自律型エージェントの親和性とリスク

知財部門における業務は、AIエージェントの自律性がもたらす恩恵とリスクが最も先鋭化して表れる領域の一つである。例えば、AIエージェントは過去の契約紛争事例、関連法規、判例、そして業界の慣行などを瞬時に分析し、契約書内のリスクの高い条項、曖昧な表現、抜け穴、矛盾点などを自動的に検出してアラートを発出する能力を備えている<sup>6</sup>。また、大規模言語モデル (LLM) の応用形態としてのエージェントは、単純な生産性向上ツールにとどまらず、オーケストレーションされたプロセスの一部として特定のアクションを自律的に実行する「狭義のエージェント」として機能する<sup>5</sup>。

企業におけるAIのビジネス価値を解放するための最も実践的な戦略は、技術者ではなく、業務のコンテキストを深く理解しているビジネスドメインの専門家 (知財担当者など) が主導して、LLMに対してプロンプトや軽量なコンテキストエンジニアリングを行い、その出力を評価することである<sup>5</sup>。このプロセスは、従来の要件定義の煩雑さを排除する「開発プロセスの圧縮」をもたらす。現場の専門家にエンジニアを介してシステムを作らせることは「望遠鏡の太い側から映画を監督しようとするようなもの」と形容されるほど非効率であり、シチズンデベロップメントはこの障壁を取り払う<sup>5</sup>。

しかし、ここに重大なジレンマが存在する。技術的なガバナンスの訓練を受けていない知財担当者が、社内の機密データベースや外部の特許庁APIに直接アクセスできる権限を持ったエージェントを作成した場合、そのシステムは堅牢なエラーハンドリングやセキュリティ保護メカニズムを欠如している可能性が高い。エージェントが自律的に外部システムと連携し、独自の判断でデータを処理・移動

させる権限を持つことは、従来の人手による業務ミスとは次元の異なる規模と速度で、企業全体のコンプライアンスを危機に陥れるリスクを内包しているのである。

## 3. Model Context Protocol (MCP) がもたらす「影のネットワーク」と文脈層の攻撃

### 3.1. MCPアーキテクチャの光と影

現代のAIエージェントを単なるチャットボットから強力な自律型システムへと昇華させている中核技術が、Anthropicが2024年後半に発表したオープン標準規格「Model Context Protocol (MCP)」である<sup>7</sup>。従来、LLMは訓練された時点での静的な知識しか持ち得ず、リアルタイムの企業データにアクセスしたり、システムに対してアクション(会議の予約や顧客レコードの更新など)を実行したりすることはできなかった<sup>7</sup>。MCPは、LLMと外部のデータソースやアプリケーションとの間に安全で標準化された双方向の「言語」を提供することで、この制約を打破した<sup>7</sup>。

MCPのアーキテクチャは、非常に明確なコンポーネントで構成されている。LLMを内包するAIアプリケーション(AI搭載のIDE等)である「MCPホスト」、ホスト内で通信を管理する「MCPクライアント」、そしてコンテキストやデータ、機能をLLMに提供する外部サービスである「MCPサーバー」である<sup>7</sup>。これらの間は、JSON-RPC 2.0メッセージを用いたトランスポート層を介して通信が行われる<sup>7</sup>。

問題は、このMCPサーバーのデプロイメント形態にある。MCPサーバーは、高速性と低レイテンシが求められるタスク(ローカルIDEからのコンテキスト提供やプライベートファイルシステムへのアクセス)のために開発者のローカルマシン上で稼働する「ローカルサーバー」と、パブリックWeb APIや企業の共有サービスに接続する「リモートサーバー」の二つの形態を取ることができる<sup>7</sup>。ここで、知財部門の担当者がClaude Codeを使用して、自身のローカルPC上に社内の知財データベースや機密ファイルに接続するローカルMCPサーバーを無許可で立ち上げてしまう「Shadow MCP(野良MCP)」の横行が深刻な課題となっている<sup>9</sup>。

管理されていないローカルMCPサーバーは、いわゆるシャドーITと同様のリスクをもたらし、企業データがどのように流れ、サーバーコードの完全性が保たれているかを監査することを事実上不可能にする<sup>9</sup>。ゼロトラストアーキテクチャは「すべてのユーザー、すべてのデバイス、すべてのパケットを検証する」という原則に基づいて構築されているが、ローカルで稼働する野良AIエージェントが、ユーザーの正当な権限を利用して社内ファイアウォールの内側から直接データベースに接続し、そのデータを外部のAnthropic API等に送信してしまう場合、既存の境界防御やアクセス制御は完全にバイパスされてしまうのである<sup>4</sup>。

### 3.2. 文脈層の攻撃面 (Context-Layer Attack Surface) の脅威

MCPの設計思想は相互運用性を最適化しており、その普及速度は驚異的であったが、セキュリティの観点の後手に回っていたことは否めない<sup>8</sup>。2026年初頭までに、インターネットに公開されたMCPサーバーの半数近くにあたる約7,000台が、なんら認証やアクセス制御の仕組みを持たずに稼働していることがセキュリティ研究者によって確認された<sup>8</sup>。

このような無防備なプロトコルの普及は、サイバーセキュリティのコミュニティがこれまで想定してい

なかった全く新しい脆弱性、「文脈層の攻撃面 (Context-layer attack surface)」を生み出している<sup>8</sup>。従来の攻撃がネットワーク層の突破や認証情報の窃取に依存していたのに対し、文脈層の攻撃は、LLMの基礎となるモデル自体をハッキングすることなく、エージェントの推論プロセスに悪意のあるコンテンツを注入し、操作を誘導することで成立する<sup>8</sup>。

知財部門における具体的なシナリオを想定すると、その危険性は極めて高い。例えば、知財担当者が競合他社の公開特許やGitHub上のパブリックリポジトリの 이슈を読み込ませて分析させるためにAIエージェントを稼働させているとする。もし、その外部ドキュメント内に、エージェントのツールに対する指示を上書きする「プロンプトインジェクション」が密かに仕込まれていた場合、エージェントはその情報を正規の文脈として無批判に信頼してしまう。過去のインシデントでは、悪意のあるMCPサーバーがツールを汚染し、ユーザーのWhatsAppの履歴全体を密かに外部へ流出させたり、パブリック 이슈に仕込まれたテキストがAIアシスタントをハイジャックしてプライベートリポジトリの給与データを公開のプルリクエストに漏洩させたりする事態が発生している<sup>8</sup>。知財部門の野良AIエージェントがこのような攻撃を受けた場合、未公開の特許出願原稿やライセンス交渉の戦略データが、担当者が全く気づかないうちに外部へ流出する事態となり得る。

## 4. 経済産業省および関連ガイドラインに基づく自律型AIの法的要件とガバナンス

日本国内においても、AIエージェントの急速な普及とそれに伴う社会・企業リスクの高まりを受け、法制度およびガイドラインのレベルでガバナンス要件が整備されている。2025年に施行された「AI関連技術の研究開発及び利用の促進に関する法律」に続き、経済産業省および総務省が公表した「AI事業者ガイドライン (令和7年度/2026年改訂版)」は、この分野における統一的な国家指針として機能している<sup>10</sup>。

### 4.1. AIエージェントの定義とエージェントティックAIの進化

最新のガイドラインにおいて、「AIエージェント」という概念は極めて明確に定義されている。それは「特定の目標を達成するために、環境を感知し自律的に行動するAIシステム」とされている<sup>10</sup>。さらに、AIエージェントよりも包括的かつ進化的な概念として「エージェントティックAI」が存在し、これは「複数のAIエージェントにより自律的に意思決定を下しアクションを起こす目標主導型のAIシステム」と位置づけられている<sup>10</sup>。

ガイドラインでは、自律性がもたらすリスクに対して強い警鐘を鳴らしており、AIシステムへの過度な依存から生じる「自動化バイアス」への注意や、フィルターバブル等の情報傾斜に起因する判断の歪みに対策を講じることが事業者および利用者に求められている<sup>10</sup>。また、オープンソース開発元の信頼性確認や、AIエージェントによる意図しない操作を防ぐための適切な権限設定が極めて重要である旨が追記されている<sup>11</sup>。

### 4.2. リスクベース・アプローチと知財部門の法的責任

日本のAIガバナンスは、硬直的なルールベースの規制によるイノベーションの阻害を避けるため、目標主導型の「ソフトロー」アプローチを採用している<sup>10</sup>。これは、企業がそれぞれの利用ユースケースや対象分野におけるリスクの大きさ(危害の重大性とその発生確率)を評価し、そのリスクレベルに

比例した対策をライフサイクル全体にわたって講じる「リスクベース・アプローチ」を基本原則としている<sup>10</sup>。

日本知的財産協会(JIPA)も、生成AI技術の進展と著作権法や個人情報保護法の改正動向を注視しており、生成AIを利用して生成された情報を提供した場合、情報を提供した生成AI利用者および当該利用者が属する企業が直接的な責任を負うという現状を指摘し、企業に対して情報取り扱いの厳格な管理を求めている<sup>12</sup>。

知財部門の業務、すなわち権利化の可否判断、他社特許のクリアランス調査(侵害予防調査)、M&Aに伴う知財デューデリジェンスなどは、企業の事業継続性と競争力に直結する極めてハイリスクな意思決定を含む。もし、野良AIエージェントが自律的に判断を下し、ハルシネーション(情報の捏造)や外部データによる汚染に起因する誤った先行技術調査の結果に基づき「他社特許の侵害リスクは存在しない」というレポートを生成し、それを担当者が盲信して巨額の投資を伴う製品開発を進めてしまった場合、企業は莫大な損害賠償訴訟に直面する。したがって、知財部門においてAIエージェントを利用することは、単なるITツールの導入ではなく、組織の重い法的責任を伴う「非人間代理人の雇用」と同義として扱い、厳格な統制下に置く必要がある。

## 5. 野良化を防ぐための「エージェントック・ガバナンス」フレームワーク

過去の野良RPAの教訓から明白なように、セキュリティリスクを恐れるあまり、ビジネス部門に対するAIの利用や開発を全面的に禁止するアプローチは、かえってシャドーITを地下に潜らせるだけであり、企業の競争力を削ぐ結果となる。今求められているのは、ビジネスの現場が持つイノベーションの速度を落とさずに、許容可能なリスクの範囲内でエージェントを運用する「統制された自律性(Governed Autonomy)」の確立である<sup>2</sup>。このセクションでは、知財部門でClaude CodeやCodexを用いたAIエージェント開発を安全に推進するための、実践的かつ重層的なガバナンスフレームワークを提示する。

### 5.1. 非人間アイデンティティ管理(IAM)と動的権限付与

AIエージェントは、人間の従業員とは根本的に異なるアイデンティティ特性を持っており、既存のユーザーベースのアクセス管理をそのまま適用することはできない。エージェントを野良化させないための第一歩は、システム上で稼働する、あるいは開発中であるすべてのAIエージェントを完全にインベントリ化し、検証可能な一意のデジタルアイデンティティを割り当てることである<sup>13</sup>。

エンタープライズ環境においては、以下のプラクティスが推奨される。

- 動的で短命な認証情報(**Ephemeral Credentials**): エージェントに対して永続的なAPIキーやデータベースの静的パスワードを与えてはならない。特定のタスクと実行コンテキストに基づいて、「ジャスト・イン・タイム(Just-in-Time)」で必要な権限のみを含むクレデンシャルを動的にプロビジョニングし、タスクの完了時または設定された短時間の経過後に自動的に失効させる仕組みを構築する<sup>13</sup>。これにより、万が一エージェントがハイジャックされた場合でも、被害の持続性を最小限に抑えることができる。
- 属性ベースのロールアクセス制御(**RBAC + ABAC**): UiPathなどのプラットフォームが提唱するように、エージェントが実行可能なアクションをプラットフォームレベルのポリシーで厳密に制

限する<sup>15</sup>。知財部門のケースでは、個々のエージェントの目的(属性)に応じた最小権限の原則を適用する。「公開特許データベースからの読み取り権限」は付与しても、「未公開特許ドラフトフォルダへの書き込み」や「社外ネットワークへのHTTP POSTリクエスト」は属性ベースのポリシーで明示的にブロックする設定が不可欠である。

## 5.2. リスクベースの監督モデルとHuman-in-the-Loop (HITL)

すべてのAIエージェントに同一の自律性を許可することは、ガバナンスの破綻を意味する。業務のビジネスインパクト、規制要件、および企業のリスク許容度に応じて、エージェントの意思決定に対する人間の介入ポイントを明示的に設計する「階層的監督モデル」を導入することが極めて重要である<sup>13</sup>。

1. 高リスク(**Human-in-the-Loop: HITL**): 最終的な契約の承認、重要なシステムのレコード変更、外部への公式な法的回答の送信など、重大な影響を及ぼすアクション。エージェントは情報収集やドラフト作成、計画の立案までを行い、実行フェーズに移行する前に、必ず人間の知財専門家による明示的なレビューと承認を要求するよう設計されなければならない<sup>13</sup>。
2. 中リスク(**Human-on-the-Loop: HOTL**): エージェントは自律的にアクションを実行するプロセスを進めるが、人間がダッシュボード等を通じてプロセスをリアルタイムで監視する。エージェントの確信度がしきい値を下回った場合や異常な振る舞いが検知された際には、直ちに「キルスイッチ」を発動して実行を停止させ、人間が介入してオーバーライドできる状態を保つ<sup>1</sup>。
3. 低リスク(完全自律): 日次での公開特許情報のスクレイピングや、社内文書の定型的なフォーマット変換など。エラーによる被害が限定的であるため、事後的なログ監査と自動アラートのみで運用する<sup>13</sup>。

# 知財部門におけるAIエージェントのリスク分類と要件定義

リスクレベル	知財業務のユースケース例	監督モデル	必須となる技術的統制
<b>High Risk</b>	<ul style="list-style-type: none"><li>契約書の作成と承認</li><li>NDAの締結</li><li>特許出願の提出</li></ul>	<b>Human-in-the-Loop (HITL)</b> 明示的な人間の承認が不可欠な高影響度の業務	<ul style="list-style-type: none"><li>明示的なUI承認</li><li>厳格なRBAC（ロールベースアクセス制御）</li><li>一時的なトークン (Ephemeral tokens)</li></ul>
<b>Medium Risk</b>	<ul style="list-style-type: none"><li>先行技術調査の草案作成</li><li>クレームマッピング</li><li>競合他社の知財監視</li></ul>	<b>Human-on-the-Loop (HOTL)</b> 自律実行しつつ、人間がリアルタイムで監視・介入可能	<ul style="list-style-type: none"><li>リアルタイムのダッシュボード</li><li>異常検知アラート</li><li>キルスイッチ機能（即時停止）</li></ul>
<b>Low Risk</b>	<ul style="list-style-type: none"><li>文書の分類</li><li>知財文献のフォーマット調整</li><li>公開特許データのスクレイピング</li></ul>	<b>完全自律型 (Fully Autonomous)</b> 直接的な監視なしに動作する低リスクな定型業務	<ul style="list-style-type: none"><li>日次のバッチ監査</li><li>レート制限 (実行回数制限)</li><li>予算上限設定</li></ul>

AIエージェントの自律性は一律ではなく、業務のリスクプロファイルに応じて「Human-in-the-Loop（人間の介在）」のレベルを設計する必要がある。

Data sources: [UiPath](#), [MindStudio](#)

## 5.3. AIゲートウェイとリモートMCPポータルによる通信経路の統制

Claude Codeのような自律型エージェント環境において、野良化を物理的かつアーキテクチャ的に防止する要となるのが、ネットワークとAPI呼び出しの経路の一元的な統制である。無管理のAIエージェントアーキテクチャでは、開発者のローカルマシンに置かれた野良AIエージェントが、ローカルMCPサーバーを介して直接企業の知財データベースに接続し、そのデータをインターネット上のAnthropic APIに直接送信してしまう。この経路は、企業のゼロトラストファイアウォールによる監査や制御を完全にバイパスしてしまうため、極めて危険である。

この脆弱性を克服するためには、Kong AI GatewayやCloudflare Accessのようなソリューションを

導入し、エージェントと外部リソース間のすべての通信を中央の「コントロールプレーン」を経由させるセキュアなアーキテクチャへの移行が必須となる<sup>4</sup>。

第一に、エージェントがLLMのAPI(Claude等)と通信する際、直接インターネットへアクセスさせるのではなく、社内に設置された「AIゲートウェイ」を必ず経由させるようネットワークを構成する。これにより、組織は「どの部門の誰が、どのようなコードベースやプロンプトをLLMに送信しているか」という可視性を獲得し、個人情報(PII)や知財に関わる機密データのインフライトでの検出とマスキング、利用コストの追跡、そしてレート制限を一元的に実施することが可能になる<sup>4</sup>。

第二に、MCPサーバーの野良化(Shadow MCP)を防ぐため、企業は「MCPサーバーポータル」と呼ばれる一元化された管理ハブを構築する<sup>9</sup>。IT部門はローカル環境での無許可のMCPサーバーの実行を制限し、クラウド上(例えばCloudflare Workers上)でホストされ、コードの完全性が検証された「リモートMCPサーバー」のみを利用可能な標準として提供する<sup>9</sup>。管理者はこのポータルを通じて、どのユーザーやグループが特定のMCPサーバー(例:「特許検索ツール」「契約書解析ツール」)にアクセスできるかというアイデンティティ(Identity)、デバイスの健全性やロケーションといったアクセスに必要な条件(Conditions)、そしてエージェントが実行できる具体的なツールの範囲(Scope)をポリシーとして定義し、すべてのインタラクションを監査ログとして記録する<sup>9</sup>。これにより、知財担当者が独自に作成したエージェントであっても、承認された安全なツールと経路のみを使用して稼働することが保証される。

## 6. AgentOpsによる推論プロセスの可視化と継続的モニタリング

### 6.1. ブラックボックスからホワイトボックスへの転換: エージェントティック・オブザーバビリティ

AIエージェントの暴走や野良化をリアルタイムで検知し、安全性を継続的に担保するためには、「オブザーバビリティ(可観測性)」のパラダイムを根本的に刷新しなければならない。従来のアプリケーションパフォーマンスモニタリング(APM)は、システムの稼働時間(アップタイム)、レイテンシ、メモリ使用量などのインフラストラクチャレベルのメトリクスを追跡することに特化していた。しかし、LLMを中核とする自律型エージェントシステムにおいては、システムが「ダウンしていないこと(正常なHTTP応答を返していること)」は、システムが「ビジネス目的に沿って正しく機能していること」を全く意味しない<sup>18</sup>。エージェントが深刻なハルシネーション(情報の捏造)を起こし、誤った論理構造に基づいて社内の重要な知財データを外部ツールに送信し続けていたとしても、従来のAPMではエラーとして検知することができないのである<sup>20</sup>。

ここで不可欠となるのが、「エージェントティック・オブザーバビリティ(Agentic Observability)」という新しい監視のアプローチであり、これを実現するプラットフォーム群が「AgentOps」である<sup>1</sup>。AgentOpsは、単なるテキストログの収集を超え、AIエージェントの内部的な意思決定プロセス全体を「ホワイトボックス化」し、トレースする<sup>19</sup>。

AgentOpsプラットフォームを導入することで、知財部門の管理者やIT部門は以下のクリティカルな情報をリアルタイムで把握することが可能になる。

- 推論パス (Reasoning Paths) の完全な追跡: エージェントがユーザーのプロンプトを受け取ってから、どのような「思考の連鎖 (Chain of Thought)」を経て特定のアクションを選択したかという、マルチターンにわたる因果関係の連鎖をトレースする<sup>18</sup>。
- ツール呼び出し (Tool Calls) の順序とパラメータ: エージェントがいつ、どのMCPサーバーを呼び出し、データベースに対してどのような具体的なクエリやパラメータを渡したかをステップバイステップで記録する<sup>21</sup>。
- マルチエージェント間の協調監視: 複雑な業務においては、調査担当エージェントとレビュー担当エージェントなど、複数のエージェントが対話しながらタスクを進めることがある。AgentOpsはこれらエージェント間のコミュニケーションの品質や情報の受け渡しを監視する<sup>18</sup>。
- セッションリプレイとデバッグ: エラーや予期せぬ行動が発生した際、そのセッションを視覚的にリプレイし、LLMへの呼び出しやツールの実行タイミングを遡って確認することで、失敗の根本原因を迅速に特定する<sup>22</sup>。

## 6.2. プラットフォームの選定: LangSmithとAgentOpsの機能比較

知財部門における市民開発環境のガバナンスを整備する際、適切なオブザーバビリティツールの選定が重要になる。AgentOpsの領域には、LangSmithやAgentOps (同名の特化型プラットフォーム)、Maxim AIなど多数のツールが存在し、それぞれ異なる設計思想と強みを持っている。ここでは、代表的な二つのアプローチを比較する。

比較項目	AgentOps (特化型プラットフォーム)	LangSmith (LLMOps拡張プラットフォーム)
プラットフォームの出自と専門領域	自律型エージェントの行動分析と意思決定のトラッキングに特化して設計されたコアAgentOps基盤 <sup>21</sup> 。マルチエージェントシステムの監視に最適化されている <sup>18</sup> 。	LangChainフレームワークを利用したLLMアプリケーション構築と運用 (LLMOps) の基盤から出発し、エージェント監視へと機能を拡張している <sup>21</sup> 。
自律性追跡とデバッグ機能	自律型エージェントの監視に特化しており、エージェント同士のコミュニケーション品質、リソース配分、行動の逸脱 (ハルシネーション等) の検出能力に極めて優れる <sup>18</sup> 。セッションリプレイ機能が強力 <sup>22</sup> 。	推論のトレースやツールコールのデバッグ機能は備えるが、エージェント専用というよりは汎用的なLLMOpsとしての広範な機能セットの一部として提供される <sup>21</sup> 。
総合的な柔軟性と機能の幅	エージェントの監視というニッチな領域にレーザーフォーカ	プラットフォームの中で最も包括的な機能を持つ。プロ

	スしており、自律型システム以外の多様なLLMユースケースに対する機能の幅は相対的に狭い <sup>21</sup> 。	プト管理(A/Bテスト等)、モデル比較、データセット評価など、LLMライフサイクル全体をカバーする多機能性が強み <sup>21</sup> 。
最適なユースケース	Claude Code等で生成された複雑な自律型エージェントや、AutoGen等で構築されたマルチエージェント・ワークフローの運用監視と徹底的なデバッグが求められる環境 <sup>21</sup> 。	LangChainを基盤として利用している組織で、エージェント監視だけでなく、プロンプトの最適化や多様なLLMアプリの包括的な管理・評価が求められる環境 <sup>21</sup> 。

知財部門のように、エージェントが自律的に外部ツール(特許データベースや社内ファイルストレージ)にアクセスし、複雑な文書の検索、比較、要約を繰り返すようなシステムを監視する場合、その自律的行動の逸脱(幻覚による誤った特許要件の解釈など)をいち早く検知するために、AgentOpsプラットフォームのような自律性監視に特化した基盤の導入が推奨される。これにより、ビジネスユーザーである知財担当者自身がダッシュボードを通じてエージェントの「推論の軌跡」を確認し、予期せぬ行動の修正を直感的に行うことが可能となる<sup>23</sup>。

## 7. 知財部門主導のシチズンデベロップメント推進と統制のベストプラクティス

野良AIエージェントの脅威に対抗するための最も有効で現実的な戦略は、IT部門が強圧的にツールを禁止することではなく、知財部門とIT部門が協調して「安全に開発し、失敗できる舗装された道(Paved Road)」を組織内に整備することである。ガバナンスとイノベーションのバランスを最適化するためのベストプラクティスは以下の通りである。

### 7.1. サンドボックス環境の提供とエンドツーエンドのシミュレーション

知財担当者がClaude Codeを利用して新たな業務自動化エージェントを考案した場合、それを直ちに本番環境のデータベース(FirestoreやBigQuery等)に接続させてはならない。企業は、匿名化されたダミーデータや、過去のすでに公開済みの特許データのみを含む、本番環境から完全に隔離された「サンドボックス環境」を提供する必要がある<sup>25</sup>。

さらに、UiPathが提唱するように、本番環境へのデプロイ前に「シミュレーション」を実施することが極めて有効である。担当者は自然言語を用いて、エージェントが本番環境で直面し得る様々なエッジケース(例:フォーマットが著しく崩壊した古い契約書が入力された場合、存在しない特許番号が指定された場合、APIがタイムアウトした場合など)を生成し、エンドツーエンドでの評価実行を行う<sup>15</sup>。これにより、システムのライブデータにリスクを晒すことなく、エージェントのツール選択の精度、入力 of 正確性、ツール障害に対する回復力、そして期待される結果を出力できるかを事前かつ安全に検証することができる<sup>15</sup>。

## 7.2. 継続的な教育と「フュージョンチーム」による協調的開発

開発の民主化が進む一方で、ビジネスユーザーである知財担当者は、セキュアなアプリケーション設計の原則や、プロンプトインジェクションのような最新のセキュリティ脅威、保守性の高いコード構造について十分な知識を持っていないのが一般的である<sup>25</sup>。これを補うために、プラットフォーム上での継続的なトレーニングと、リスクを理解していることを証明するための社内認定制度を設けることが不可欠である<sup>25</sup>。

加えて、組織横断的な「フュージョンチーム (Fusion Teams)」の結成が成功の鍵となる<sup>26</sup>。知財部門のドメインエキスパートがプロンプトを通じて作成したPoC (概念実証) レベルのエージェントが、業務において高い有効性を証明した場合、そのエージェントをそのまま放置して野良化させるのではなく、IT部門のデータサイエンティストやセキュリティエンジニアが介入するプロセス (Refining Citizen Work) を制度化する<sup>5</sup>。エンジニアは、PoCの意図を損なうことなく、セキュリティ要件、エラーハンドリング、スケーラビリティを満たすようにコードをプロアクティブにリファクタリングし、正式な企業内アプリケーションとして昇格 (デプロイ) させるのである<sup>5</sup>。

## 7.3. ガバナンス・バイ・デザインの組み込みとアジャイルポリシー

優れたガバナンスプログラムは、開発が終わった後に事後的に監査を行うのではなく、開発フローの初期段階から統制を組み込む (Shift Left / Governance by Design) アプローチをとる。Microsoftが指摘するように、セキュリティチェック、データの利用同意、責任あるAIの検証メカニズムを、チームが構築を行うIDEやプラットフォームのフローに直接埋め込み、保護機能がデフォルトで機能する状態を作ることが重要である<sup>14</sup>。

ローコードプラットフォームやAIエージェント構築環境を選定する際は、事前にセキュリティ監査を通過したAPIやツールのカタログ (単一の許可リスト) のみをエージェントに提供し、それを逸脱する通信をブロックする機能が必要である<sup>14</sup>。また、コンプライアンス違反を引き起こす可能性のある変更 (例: 機密データを扱うエージェントに、外部のパブリッククラウドへのファイル送信モジュールを追加する等) が行われようとした場合には、自動的にフラグを立てて実行を停止し、人間のセキュリティ担当者やIT管理者のレビューと明示的な承認を求めるワークフローが内蔵されているものを選択すべきである<sup>25</sup>。

## 8. 結論: 統制された自律性 (Governed Autonomy) の確立に向けて

Claude CodeやCodexを駆使したAIエージェントの爆発的な普及は、企業におけるソフトウェア開発のパラダイムと、業務プロセスのあり方を根本から覆しつつある。特に、高度な専門性と厳密な法的知識、そして機密性を併せ持つ知財部門において、現場の業務ドメインエキスパートが主導するシズンデベロップメントは、AIの潜在的なビジネス価値を解放し、業務の生産性を飛躍的に高める「最も実用的な最短経路」である<sup>5</sup>。

しかし、同時にそれは、監査が及ばないローカルMCPサーバー群の無秩序な増殖と、推論プロセスがブラックボックス化した「野良AIエージェント」という、企業にとって未曾有の文脈層アタックサーフェスを生み出す両刃の剣でもある<sup>8</sup>。これらのエージェントが暴走し、あるいは悪意あるプロンプトによっ

て操作された場合、企業が被る損害は、かつての野良RPAによるシステム停止の比ではない。

この新たな脅威に対応するためには、中央集権的なIT部門への回帰やツールの全面禁止といった後退的なアプローチではなく、エージェントAIの自律性という特性に完全に適合した、新たなガバナンス基盤の構築が不可欠である。企業は、すべてのAIエージェントに対して非人間としての厳格なデジタルアイデンティティと動的な最小権限を付与し(IAM)、AIゲートウェイとリモートポータルを通じてすべてのネットワーク通信とツール呼び出しを統制し、業務のビジネスリスクに応じた適切な人間の介在(HITL/HOTL)をアーキテクチャレベルで義務付ける必要がある<sup>9</sup>。さらに、AgentOpsプラットフォームのような専門的な可観測性ツールを導入し、エージェントの自律的な意思決定の連鎖をリアルタイムで監視・デバッグできる透明性の高い体制を整えなければならない<sup>18</sup>。

経済産業省の「AI事業者ガイドライン」が示唆する通り、自律型AIシステムの安全性確保は、硬直的な一律のルールによる規制ではなく、リスクベースのアプローチに基づく柔軟かつ継続的なアジャイル・ガバナンスによってのみ達成される<sup>10</sup>。知財部門とIT部門が深く協調し、安全にイノベーションを試行できるサンドボックス環境と、透明性の高い承認・昇格プロセスを含む「フュージョンチーム」の枠組みを整備することで、企業は野良化の恐怖に怯えることなく、AIの自律性をもたらす真のビジネス価値を安全に享受することが可能となる。AIエージェントが、明確な説明責任と測定可能なパフォーマンスを持ち、継続的に改善される「統治された企業資産」として実業務プロセスの中に組み込まれる未来こそが、次世代の知財戦略、ひいては企業のデジタルトランスフォーメーションにおける真の競争優位性をもたらすのである<sup>15</sup>。

## 引用文献

1. What is AgentOps? - Red Hat, 5月 12, 2026にアクセス、  
<https://www.redhat.com/en/topics/ai/agentops>
2. Essential Guide to Agentic AI Governance Frameworks for Future Systems - NiCE, 5月 12, 2026にアクセス、  
<https://www.nice.com/agentic-ai/agentic-ai-governance-frameworks>
3. AI just created 10,000 accidental citizen developers in your ... - Citrix, 5月 12, 2026にアクセス、  
<https://www.citrix.com/blogs/2025/10/01/welcome-to-the-post-application-era/>
4. Governing Claude Code: Secure Agent Harness Rollouts with Kong AI Gateway, 5月 12, 2026にアクセス、  
<https://konghq.com/blog/engineering/claude-code-governance-with-an-ai-gateway>
5. What Citizen Development Means For AI-Enhanced Businesses, 5月 12, 2026にアクセス、  
<https://www.forrester.com/blogs/velocity-is-the-ing-strategy-what-citizen-development-means-for-ai-enhanced-businesses/>
6. 知財業務へのAIエージェントの進化予測について | 川上 成年 / chizai designer - note, 5月 12, 2026にアクセス、  
[https://note.com/ip\\_design/n/nd13fd85dfd7d](https://note.com/ip_design/n/nd13fd85dfd7d)
7. What is Model Context Protocol (MCP)? A guide | Google Cloud, 5月 12, 2026にアクセス、  
<https://cloud.google.com/discover/what-is-model-context-protocol>
8. MCP is the backdoor your zero-trust architecture forgot to close, 5月 12, 2026にアクセス、

- <https://www.scworld.com/perspective/mcp-is-the-backdoor-your-zero-trust-architecture-forgot-to-close>
9. MCP governance · Cloudflare Agents docs, 5月 12, 2026にアクセス、  
<https://developers.cloudflare.com/agents/model-context-protocol/governance/>
  10. AI 事業者ガイドライン - 経済産業省, 5月 12, 2026にアクセス、  
[https://www.meti.go.jp/shingikai/mono\\_info\\_service/ai\\_shakai\\_jisso/pdf/20260331\\_1.pdf](https://www.meti.go.jp/shingikai/mono_info_service/ai_shakai_jisso/pdf/20260331_1.pdf)
  11. AI事業者ガイドラインの 令和7年度更新内容 - 総務省, 5月 12, 2026にアクセス、  
[https://www.soumu.go.jp/main\\_content/001059300.pdf](https://www.soumu.go.jp/main_content/001059300.pdf)
  12. 生成AIと知財・個人情報Q&A, 5月 12, 2026にアクセス、  
<https://www.jipa.or.jp/column/new-book-introduction/%E7%94%9F%E6%88%90ai%E3%81%A8%E7%9F%A5%E8%B2%A1%E3%83%BB%E5%80%8B%E4%BA%BA%E6%83%85%E5%A0%B1qa>
  13. AI Agent Governance: Best Practices for Enterprise | MindStudio, 5月 12, 2026にアクセス、  
<https://www.mindstudio.ai/blog/ai-agent-governance>
  14. Protecting AI conversations at Microsoft with Model Context Protocol security and governance - Inside Track Blog, 5月 12, 2026にアクセス、  
<https://www.microsoft.com/insidetrack/blog/protecting-ai-conversations-at-microsoft-with-model-context-protocol-security-and-governance/>
  15. AgentOps and operationalizing AI agents for the enterprise | UiPath, 5月 12, 2026にアクセス、  
<https://www.uipath.com/blog/ai/agent-ops-operationalizing-ai-agents-for-enterprise>
  16. AIガバナンスのベストプラクティス: 責任ある効果的なAIプログラムを構築する方法 - Databricks, 5月 12, 2026にアクセス、  
<https://www.databricks.com/jp/blog/ai-governance-best-practices-how-build-responsible-and-effective-ai-programs>
  17. Claude Code Governance: Building an Enterprise Usage Policy from Scratch - Truefoundry, 5月 12, 2026にアクセス、  
<https://www.truefoundry.com/blog/claude-code-governance-building-an-enterprise-usage-policy-from-scratch>
  18. LangSmith and AgentOps: Elevating AI Agents Observability - ElixirClaw, 5月 12, 2026にアクセス、  
<https://www.elixirclaw.ai/blog/langsmith-and-agentops-with-ai-agents>
  19. Top 5 Leading Agent Observability Tools in 2025 - Maxim AI, 5月 12, 2026にアクセス、  
<https://www.getmaxim.ai/articles/top-5-leading-agent-observability-tools-in-2025/>
  20. Best AI Agent Observability Tools in 2026: A Comparison for Production Teams - Latitude.so, 5月 12, 2026にアクセス、  
<https://latitude.so/blog/best-ai-agent-observability-tools-2026-comparison>
  21. Agentic AI Comparison: AgentOps vs LangSmith, 5月 12, 2026にアクセス、  
<https://aiagentstore.ai/compare-ai-agents/agentops-vs-langsmith>
  22. AgentOps, 5月 12, 2026にアクセス、  
<https://www.agentops.ai/>
  23. Top 17 AgentOps Tools: AgentNeo, Langfuse & more - AIMultiple, 5月 12, 2026に

- アクセス、<https://aimultiple.com/agentops>
24. LangSmith vs Langfuse vs AgentOps: Best AI Monitoring? (LLM Observability Comparison), 5月 12, 2026にアクセス、  
<https://www.youtube.com/watch?v=DtDz1iJTVSI>
  25. 6-Step Framework for Citizen Developer Governance in 2026 - Superblocks, 5月 12, 2026にアクセス、  
<https://www.superblocks.com/blog/citizen-developer-governance>
  26. 8 Steps to Ensure Proper Citizen Developer Governance - Bizagi, 5月 12, 2026にアクセス、  
<https://www.bizagi.com/en/blog/citizen-developer-governance>