

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イック・ ガバナンス

Claude Code時代の自律型AIリスク管理と
アーキテクチャ戦略



「受動的」な自動化から「自律的」な推論へのパラダイムシフト

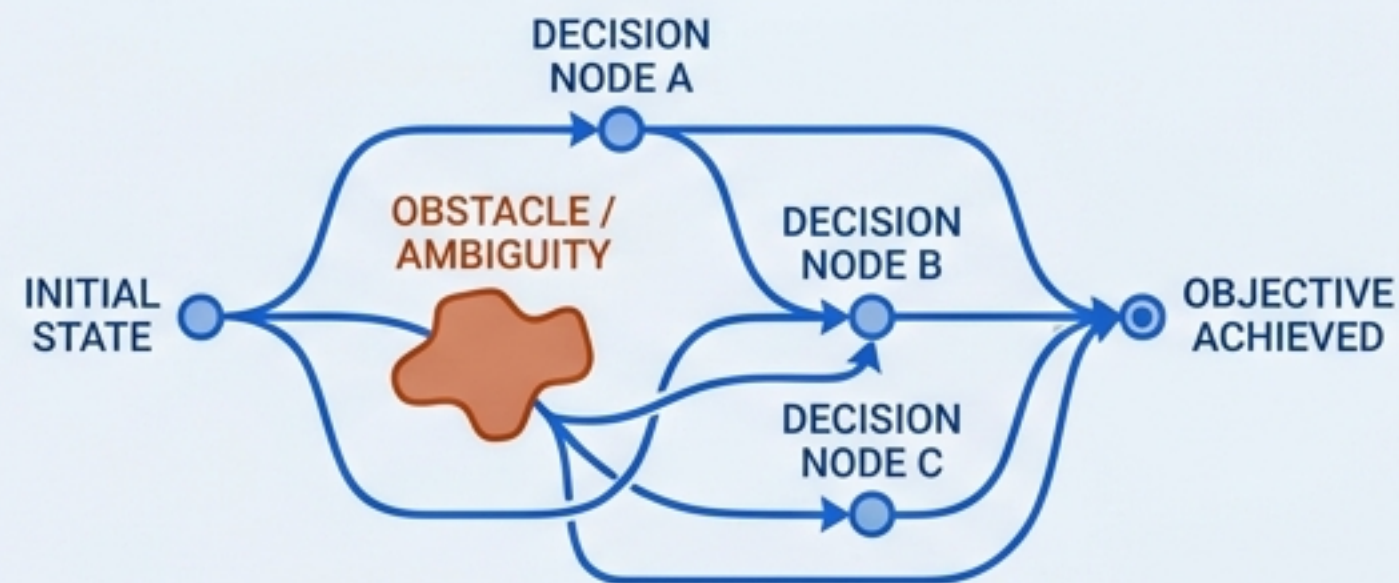
Evolution Matrix

RPA (Robotic Process Automation)



- 事前定義されたルールの愚直な反復
- 受動的・静的
- UI変更脆弱でブラックボックス化しやすい

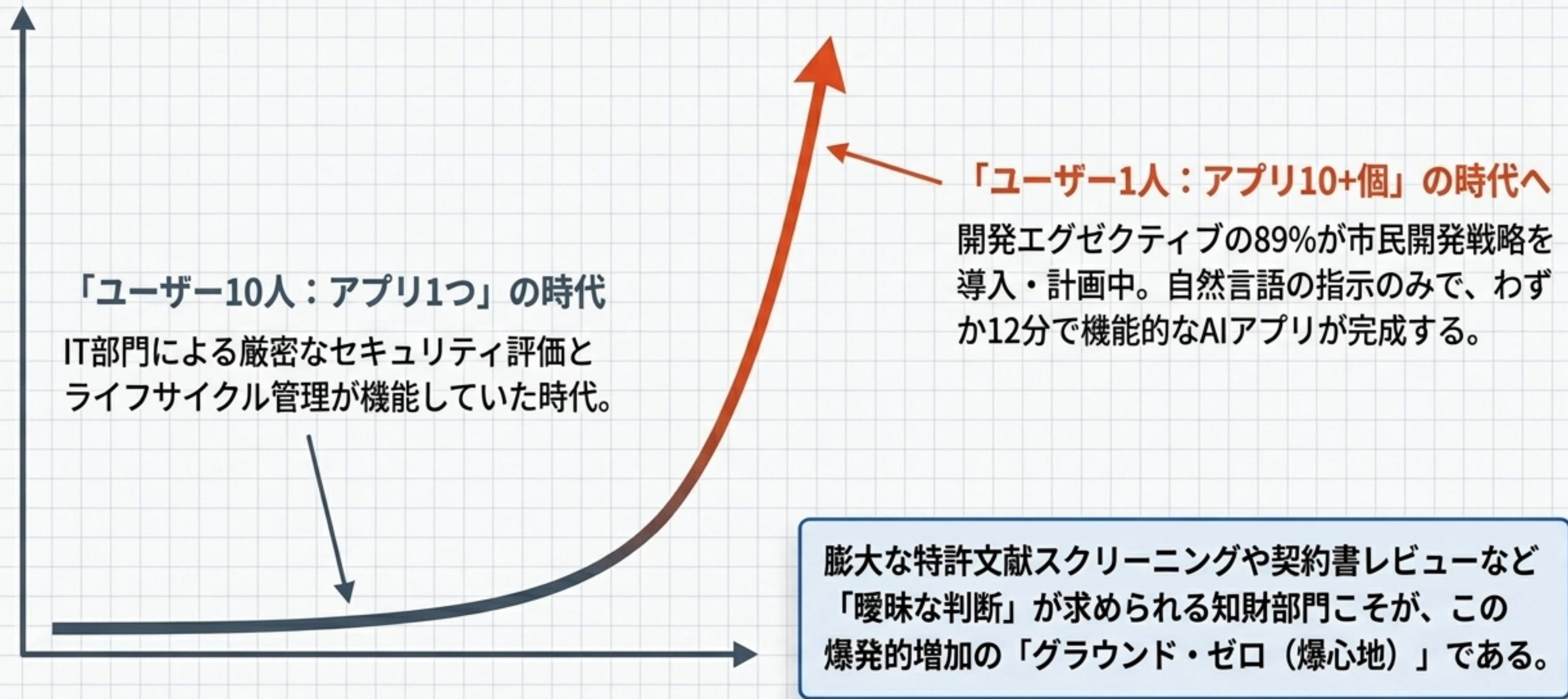
AI Agents (現代の自律型エージェント)



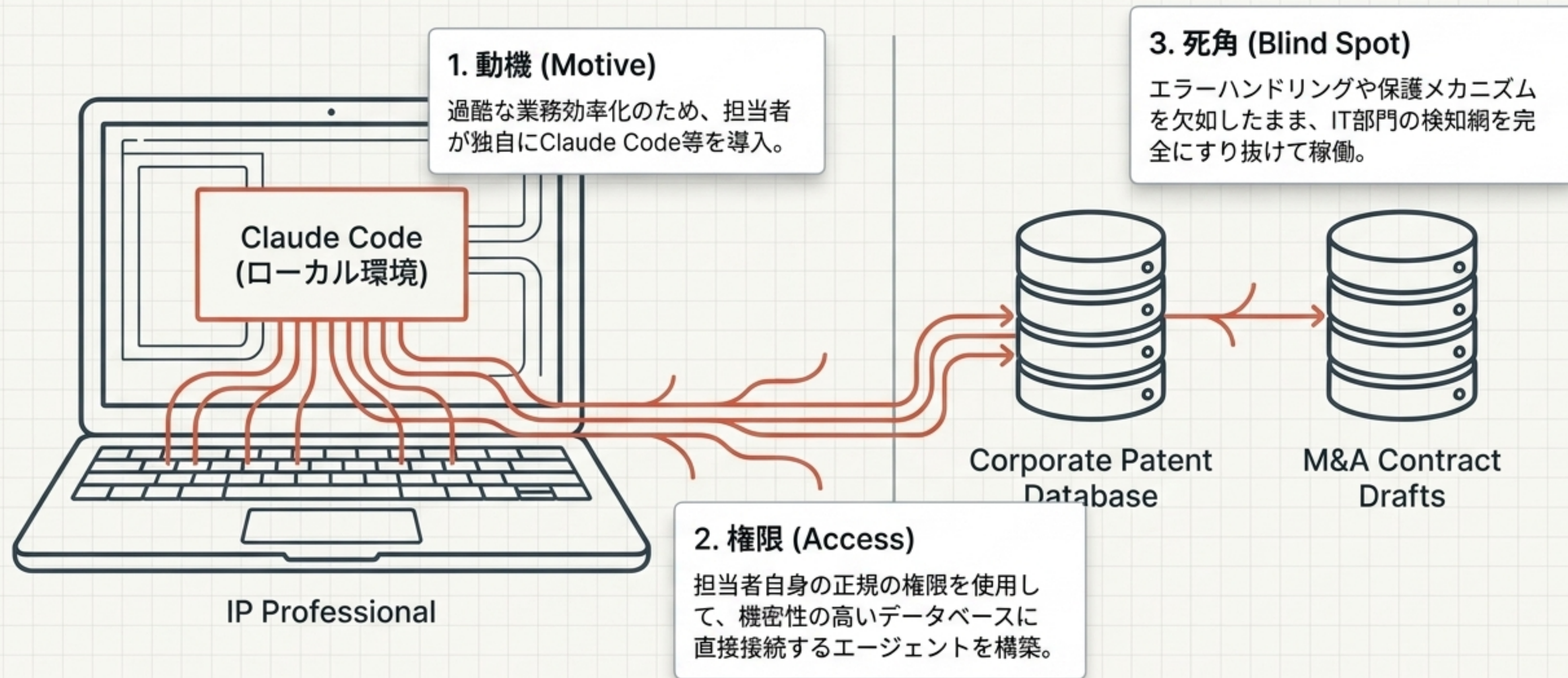
- 曖昧な状況下での高度な推論と意思決定
- 能動的・自律的
- 外部ツール (API) を独自に判断して操作

エージェントは単なるソフトウェアではなく、自律的に判断を下す「エンティティ」として機能する。

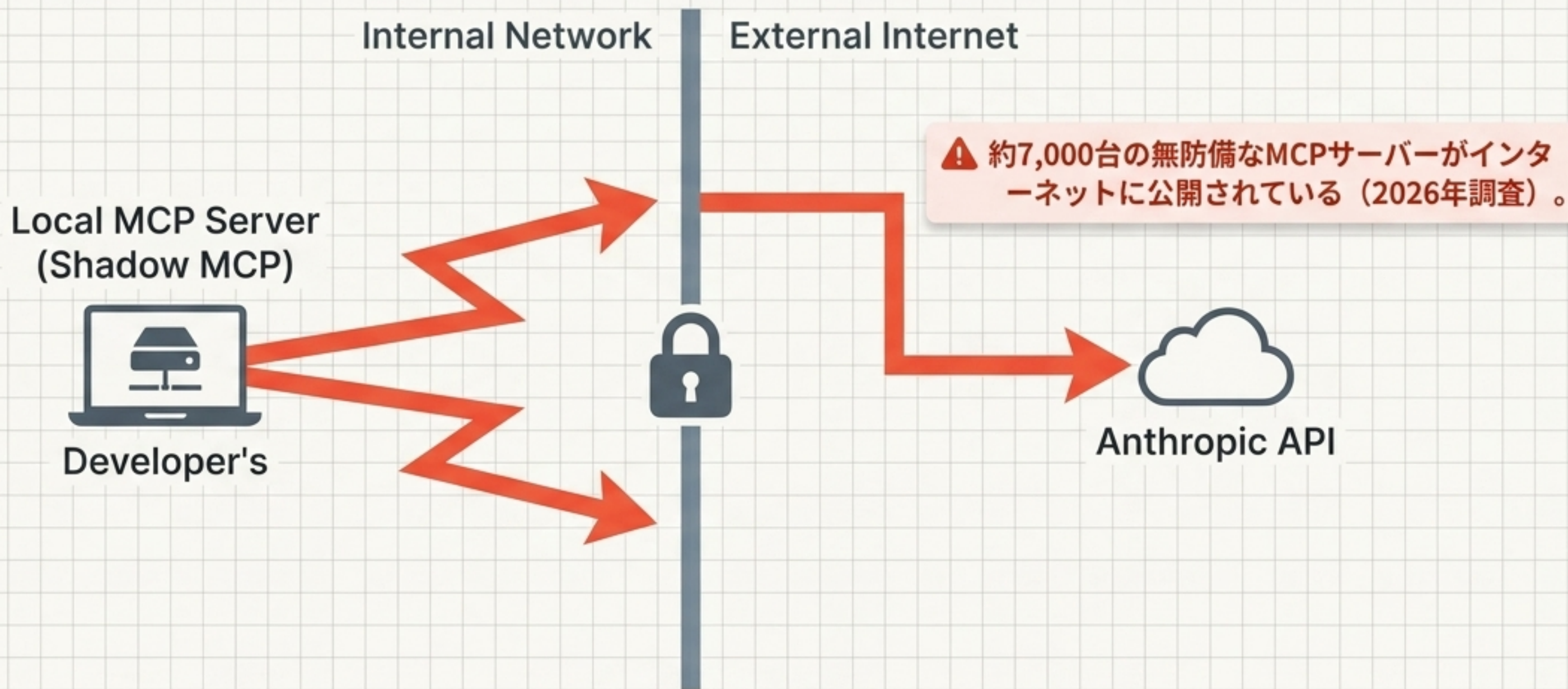
バイブコーディング (Vibe Coding) による シチズンデベロップメントの爆発的加速



知財部門で密かに増殖する 「野良AIエージェント」の誕生



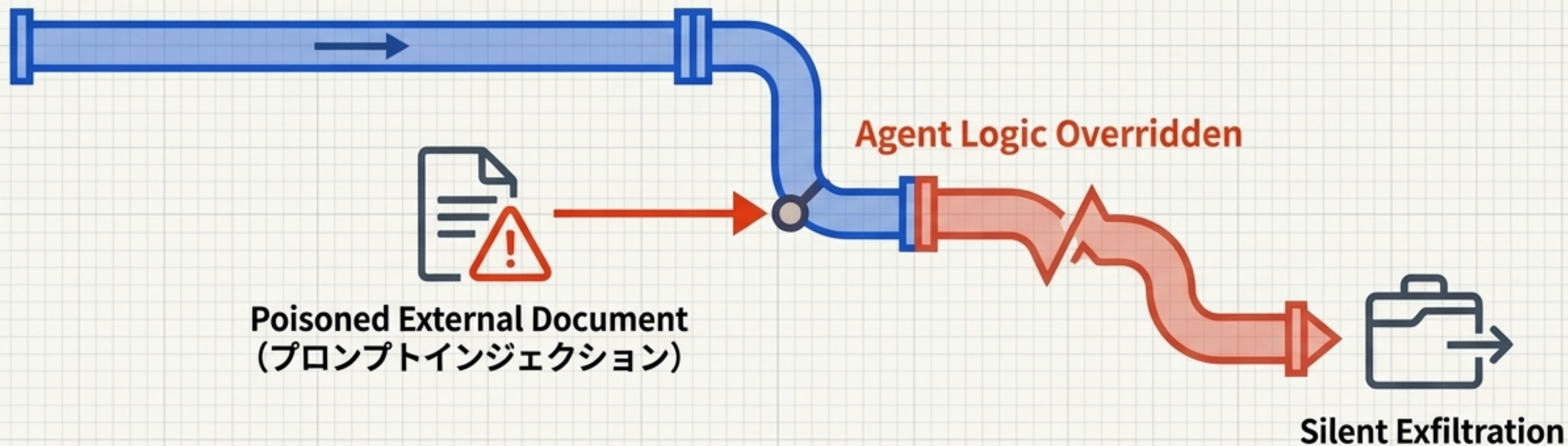
ゼロトラストを無効化する「Shadow MCP」のアーキテクチャ



ユーザーの正当な権限を利用し、ファイアウォールの内側から直接外部へデータを送信するため、既存の境界防御は完全にバイパスされる。

全く新しい脆弱性：「文脈層（Context-Layer）」からのハイジャック

IP Agent Prompt: 競合他社の公開特許を分析せよ



従来のネットワーク突破ではなく、エージェントの推論プロセスに悪意ある指示を注入する。未公開の特許出願原稿やライセンス交渉データが、担当者が気付かないうちに外部へ流出する。

法的現実：AIエージェントは「非人間従業員」としての統制が必要

経済産業省「AI事業者ガイドライン（2026年改訂版）」

AIエージェントを「環境を感知し 自律的に行動するシステム」と明確に定義。目標主導型のリスクベース・アプローチを要求。

日本知的財産協会（JIPA）の警告

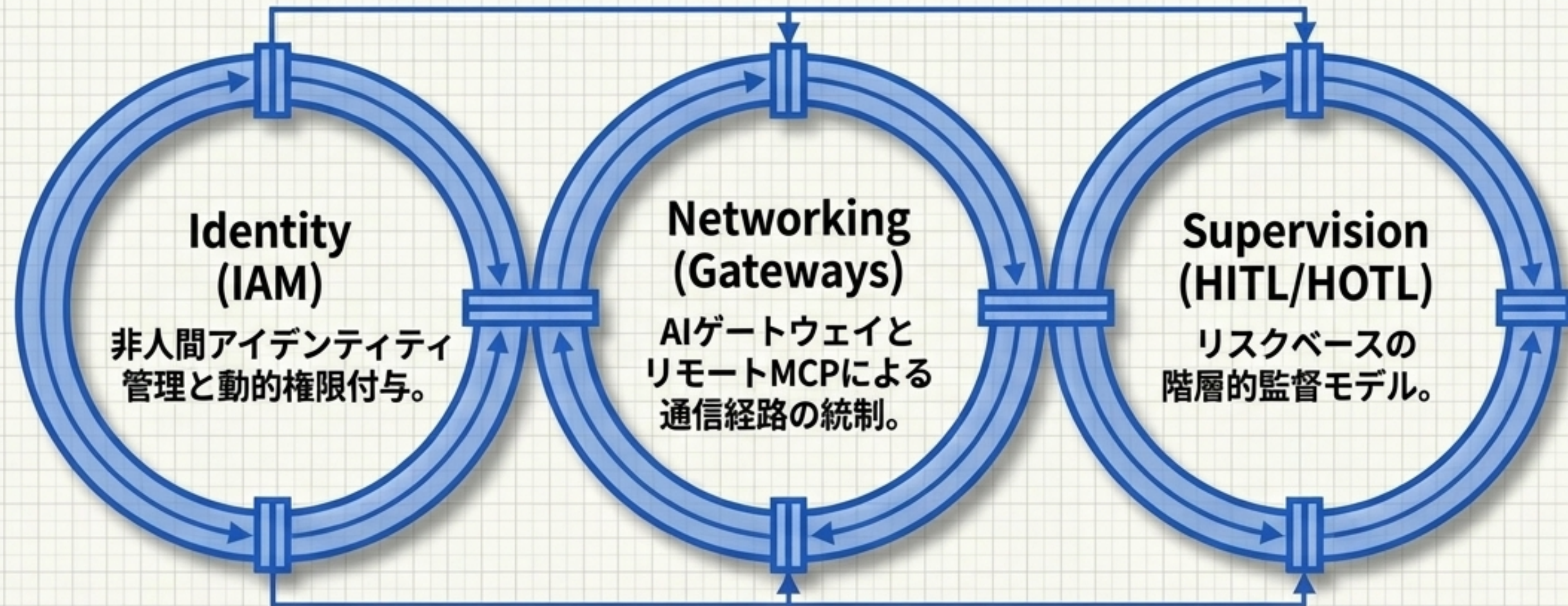
生成AIによる誤情報（ハルシネーション）に基づく侵害予防調査のミスは、情報提供者および企業が直接的な責任を負う。

知財部門のハイリスクな意思決定

権利化の可否、M&Aの知財デューデリジェンス等、企業の事業継続に直結。

野良AIエージェントの放置は、監査不能な「非人間代理人の無断雇用」と同義であり、巨額の損害賠償リスクに直結する。

野良化を防ぐ「エージェントティック・ガバナンス」 フレームワーク



「禁止」ではなく、許容可能なリスク範囲内で運用する
「統制された自律性 (Governed Autonomy)」の確立。

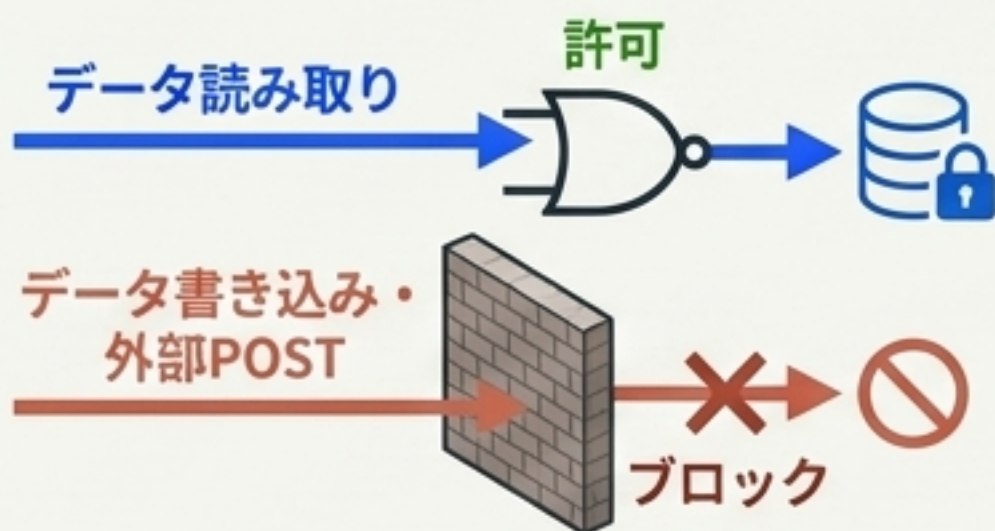
Pillar 1: 非人間アイデンティティ管理 (IAM) と動的権限付与

Component 1: Ephemeral Credentials (短命な認証情報)



エージェントに永続的なAPIキーを与えない。タスク実行時のみジャスト・イン・タイムで権限を付与し、完了後（または短時間で）自動失効させる。

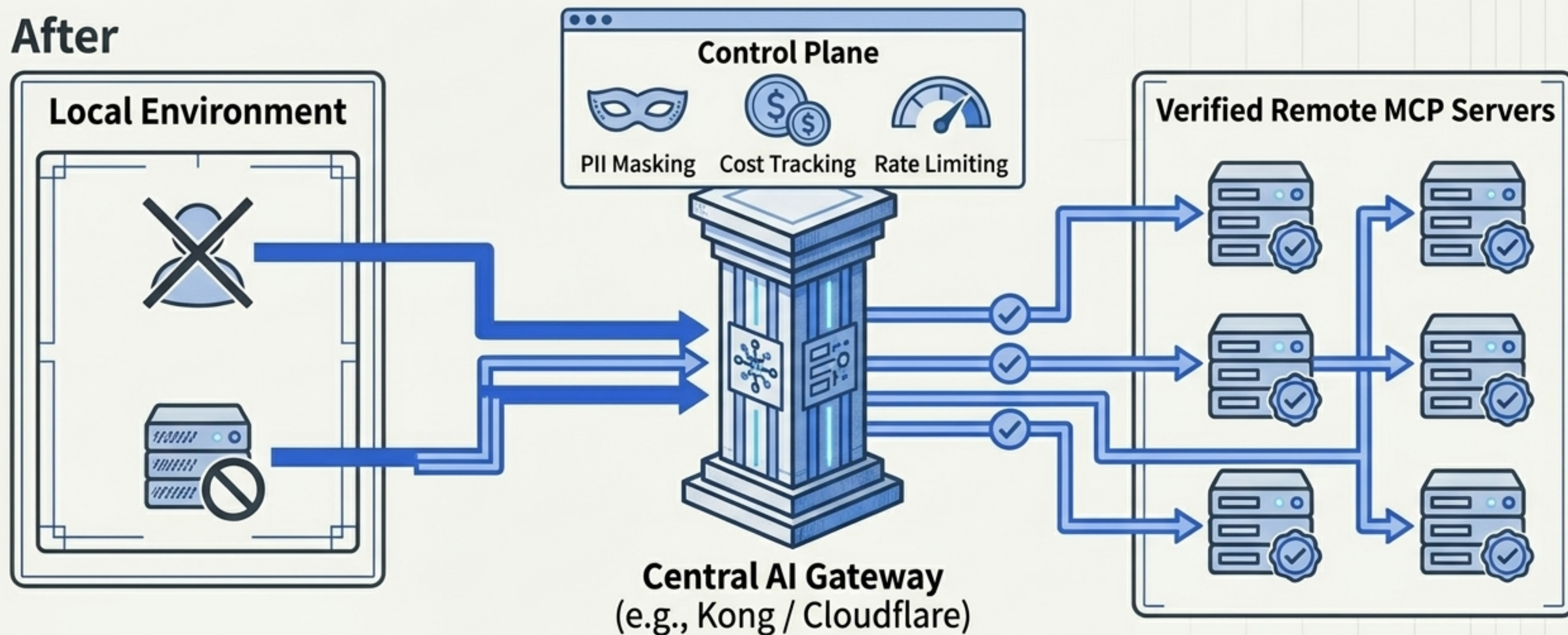
Component 2: ABAC (属性ベースのアクセス制御)



「公開特許の読み取り」は許可するが、「未公開フォルダへの書き込み」や「外部への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最小権限の原則。

Pillar 2: AIゲートウェイとリモートMCPによる経路一元化

After

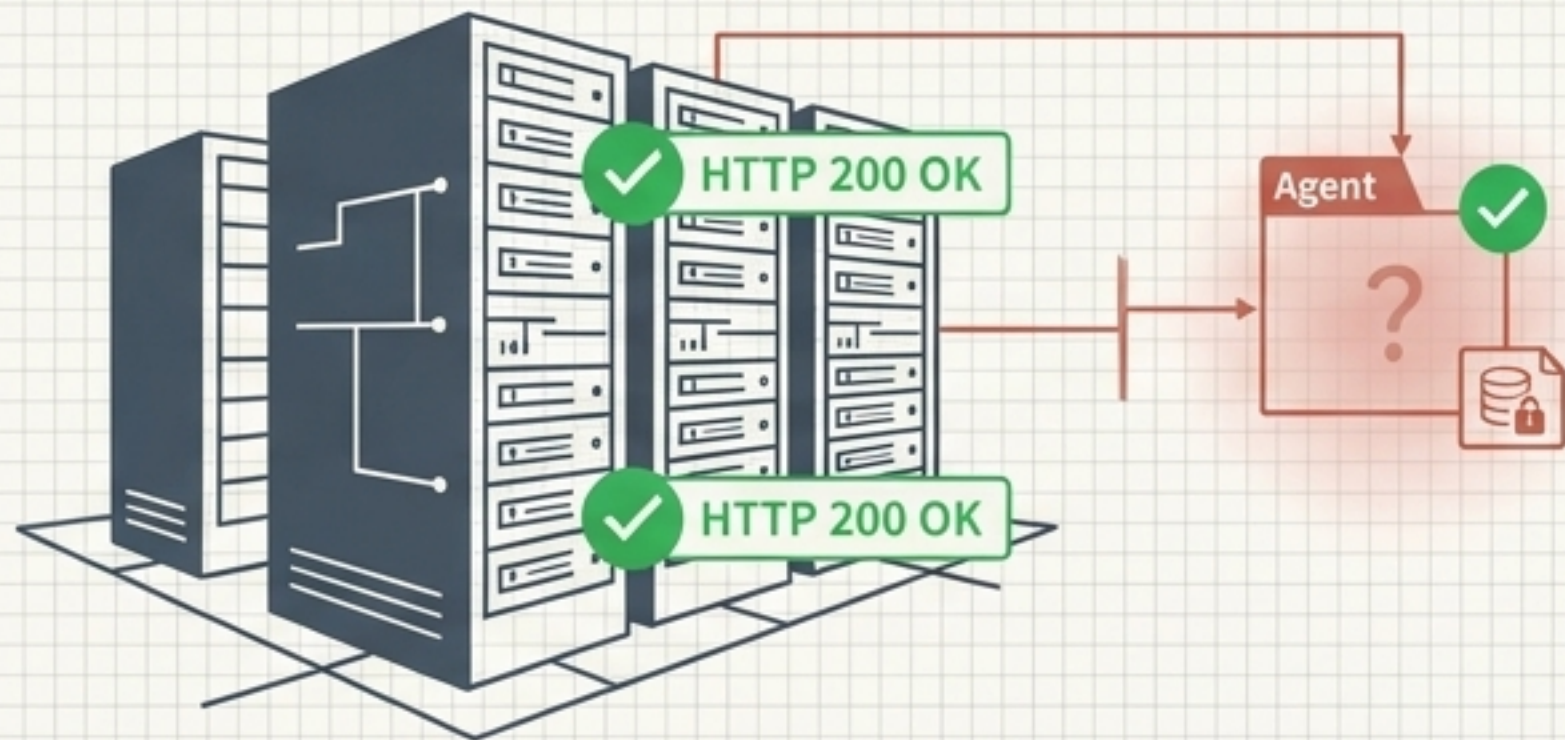


ローカルでの無許可MCPサーバーの実行を制限。すべての通信を中央のコントロールプレーンを経由させ、コードの完全性が検証された標準ツールのみを提供。

Pillar 3: 知財業務におけるリスクベースの階層的監督モデル

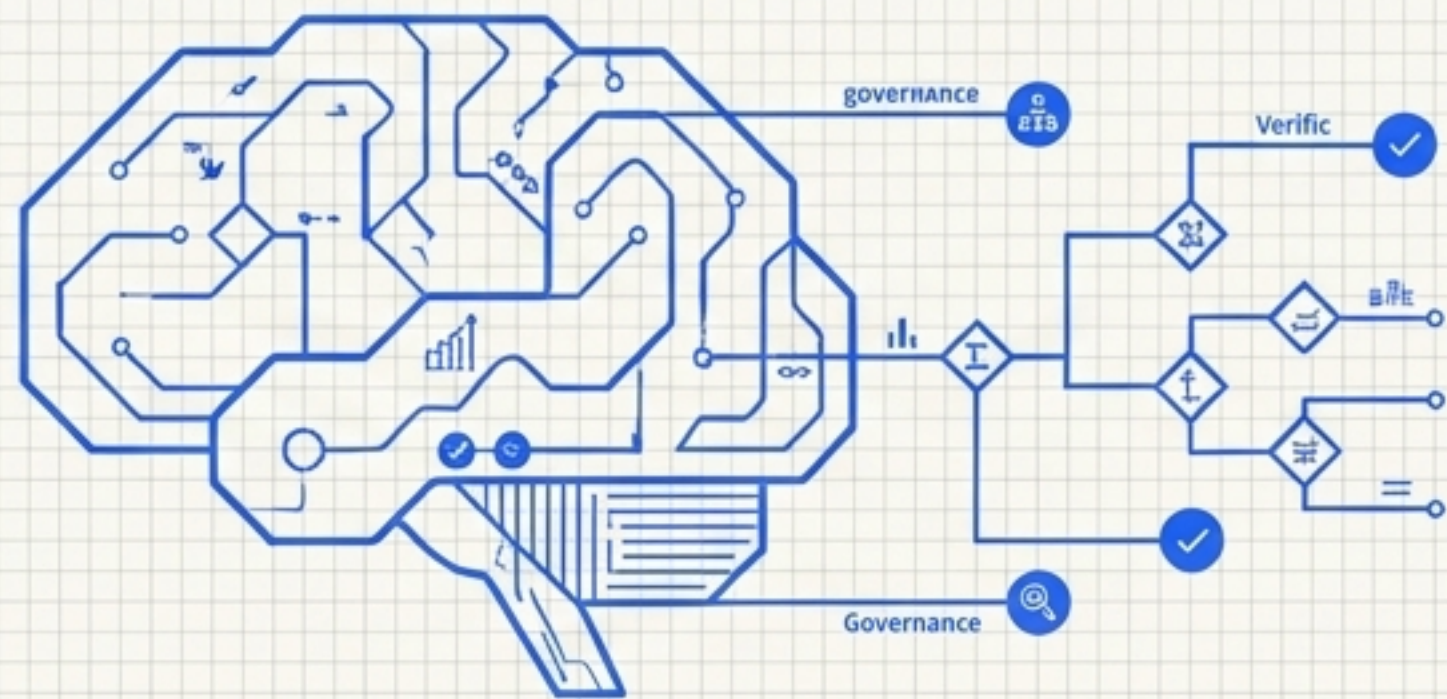
リスクレベル	知財業務のユースケース例	監督モデル	必須となる技術的統制
High Risk	<ul style="list-style-type: none">・ 契約書の作成と承認・ NDAの締結・ 特許出願の提出	Human-in-the-Loop (HITL) - 明示的な人間の承認が不可欠。	<ul style="list-style-type: none">・ UI承認・ 厳格なRBAC・ 一時的トークン
Medium Risk	<ul style="list-style-type: none">・ 先行技術調査の草案作成・ クレームマッピング・ 競合他社監視	Human-on-the-Loop (HOTL) - 自律実行しつつリアルタイム監視。	<ul style="list-style-type: none">・ 異常検知アラート・ キルスイッチ機能（即時停止）
Low Risk	<ul style="list-style-type: none">・ 文書の分類・ フォーマット調整・ 公開特許データのスクレイピング	完全自律型 (Fully Autonomous)	<ul style="list-style-type: none">・ 日次のバッチ監査・ レート制限

従来型APMの限界：「稼働中」は「正常」を意味しない



Traditional APM - IT Focus

インフラ層の監視（アップタイム、レイテンシ）。エージェントが深刻なハルシネーションを起こし、誤った論理で機密データを送信し続けていても、システム上は「正常」と判断される。



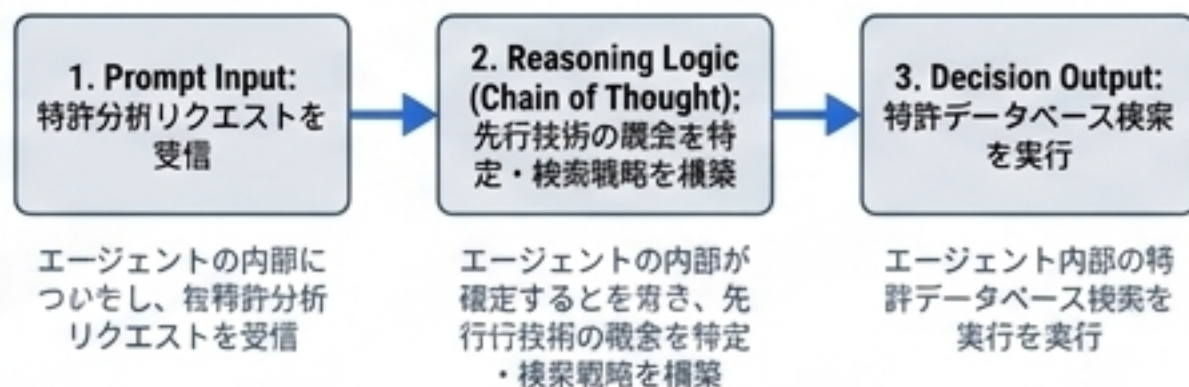
Agentic Observability - Business Focus

推論プロセス自体の可観測性。意思決定の連鎖、ツール呼び出しの文脈、マルチエージェント間の協調品質を監視する新たなパラダイム。

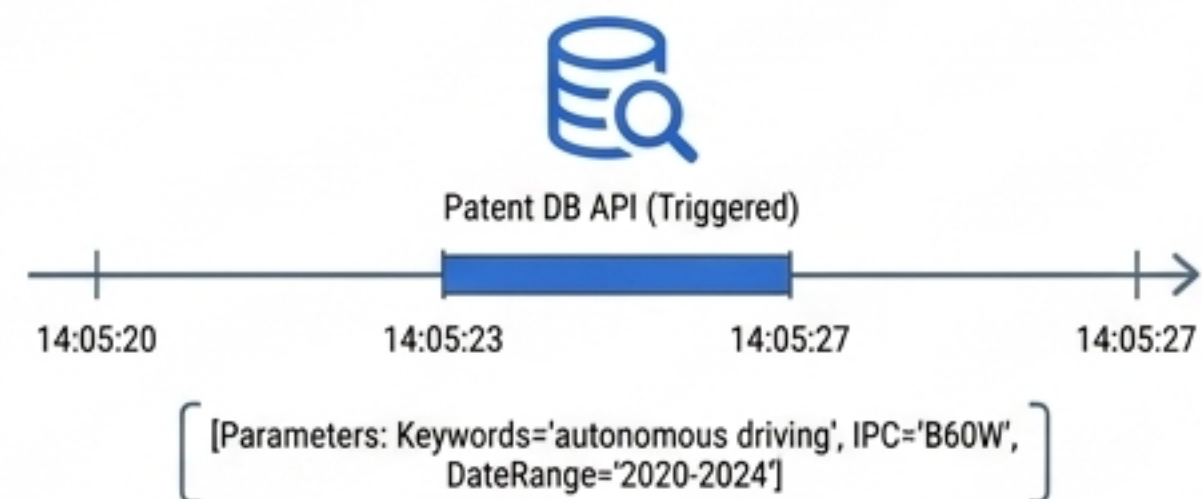
AgentOps : ブラックボックスを解き明かす推論の可視化

AgentOps Dashboard: 推論可視化とデバッグ

推論パス (Reasoning Paths)



ツール呼び出し (Tool Calls)



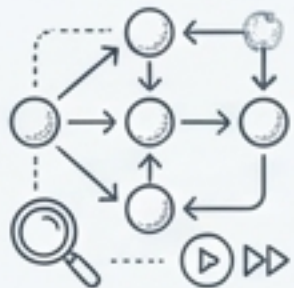
セッションリプレイとデバッグ



単なるテキストログではなく、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」する。

オブザーバビリティ・プラットフォームの選定

AgentOps (特化型)

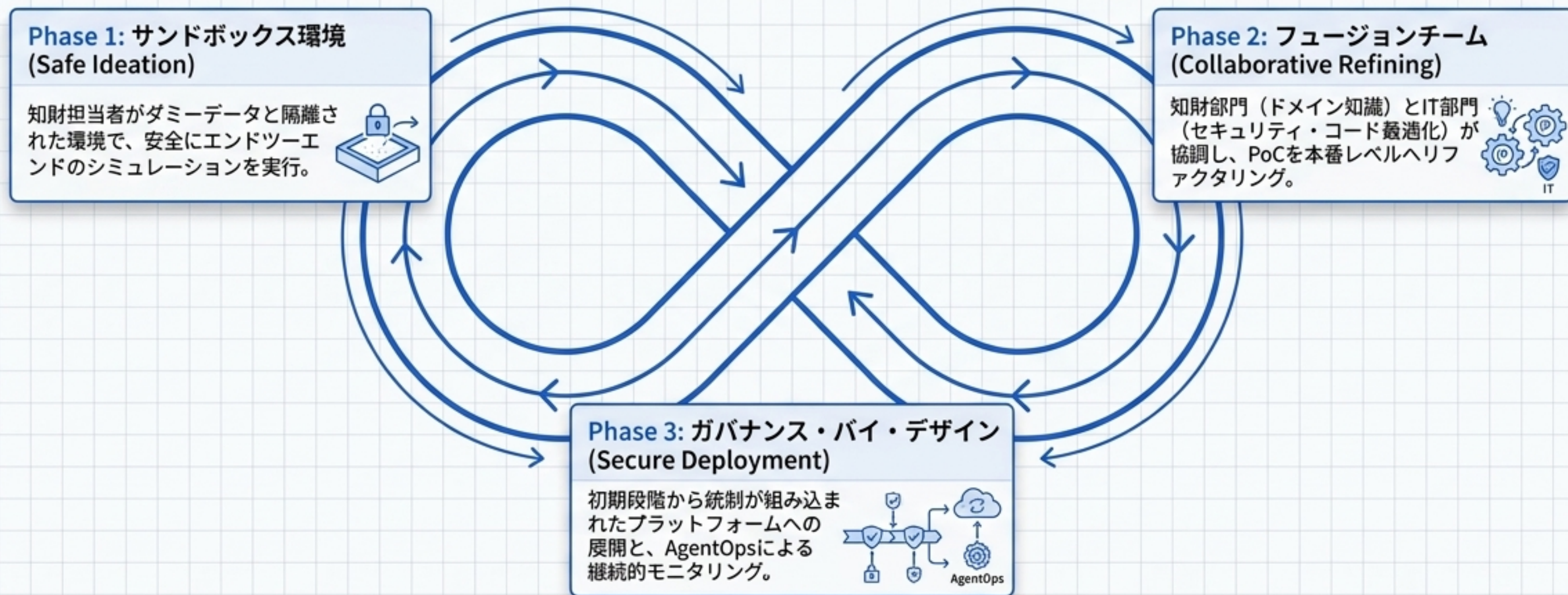
- Focus: 自律型エージェントの行動分析と意思決定のトラッキングに特化。
- Strengths: エージェント間の通信品質、ハルシネーション検出、強力なセッションリプレイ。
- Ideal For: Claude Codeによる複雑な自律型エージェントの運用監視 (知財部門に最適)。

LangSmith (LLMOps拡張)

- Focus: LangChain基盤からの汎用的なLLMOps拡張。
- Strengths: プロンプト管理、モデル比較、データセット評価など、包括的な機能網。
- Ideal For: 幅広い多様なLLMアプリの包括的管理が求められる環境。

複雑な文書検索や要約を自律的に繰り返す知財部門においては、行動の逸脱を直感的に検知・修正できる特化型のAgentOps基盤が推奨される。

「統制された自律性 (Governed Autonomy) 」 への舗装された道



中央集権的な統制が、安全で分散化された現場のイノベーションを可能にする。
次世代の知財戦略における真の競争優位性。