

NVIDIA 「RTX Spark」 発表の真実と、ローカルAI PC時代が変える知財実務

The Truth Behind NVIDIA's "RTX Spark" and How Local AI PCs Will Transform IP Practices

Claude Opus 4.8

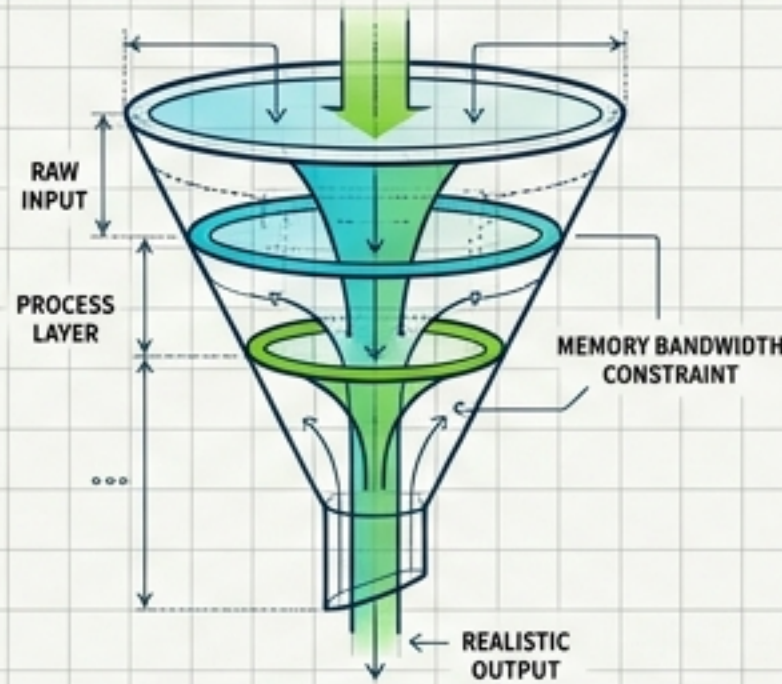
2026年6月1日

エグゼクティブ・サマリー：AI PCの実像と知財戦略の転換点



事実確認

「RTX Spark」は実在する。2026年秋発売予定のWindows向けAI PCチップであり、120B級LLMのローカル駆動を可能にする。



現実の性能

1ペタフロップスや100万トークンは「ピーク値」。実効性能はメモリ帯域 (273GB/s) がボトルネックとなり、生成速度や有効文脈長には一定の制約が生じる。



構造的変革

最大の価値は「機密情報のローカル処理」。これにより知財部門は「運用AI層」と「戦略人間層」の二層構造へシフトする。

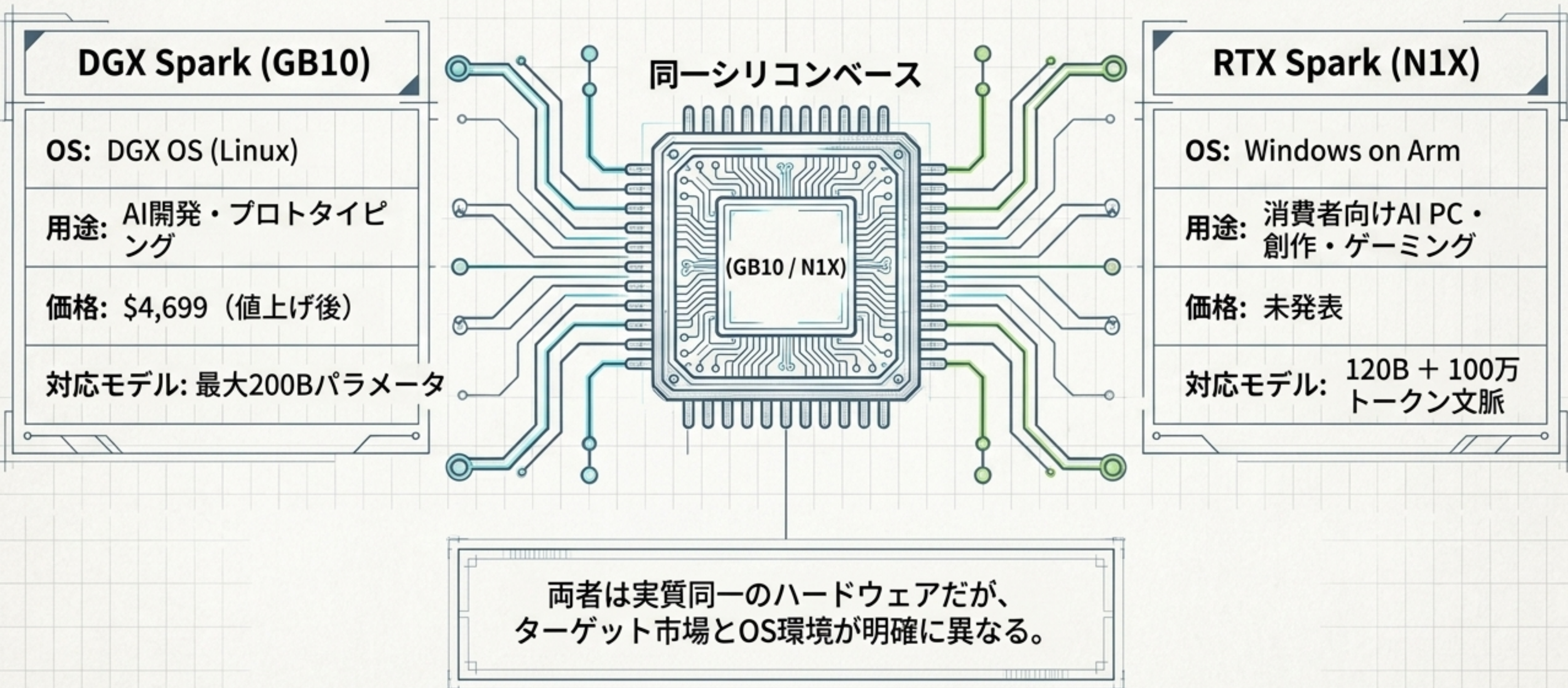
発表内容の事実確認：GTC Taipei 2026 基調講演

Verification Ledger

発表・噂の仕様	検証結果	公式ソース
20コア Arm Grace CPU	一致	MediaTek協業も公式記載
Blackwell世代 6144 CUDAコア	一致	公式発表と符合
最大128GB ユニファイドメモリ	一致	下位SKUは16GB～
最大1ペタフロップス AI性能	一致	NVFP4+スパース性利用時
1200億(120B)パラメータ・100万トークン文脈	一致	公式プレスリリースに明記
発売時期とOS	一致	同年秋発売、Windows on Arm、Microsoft共同開発

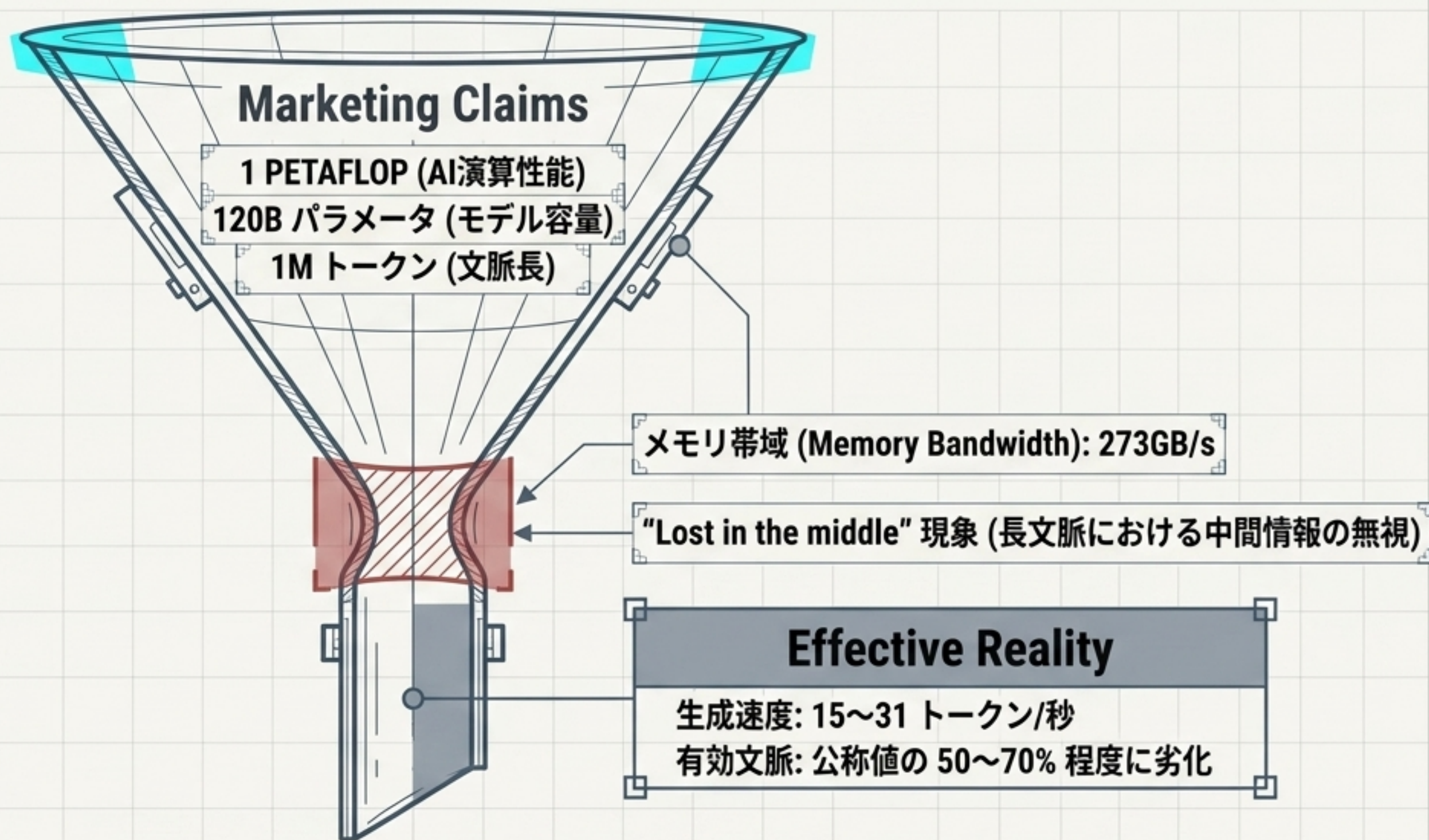
結論：提示された仕様は誤記や仮想シナリオではなく、100%実在の確認できる事実である。

シリコンの共有と用途の分化：DGX Spark vs. RTX Spark



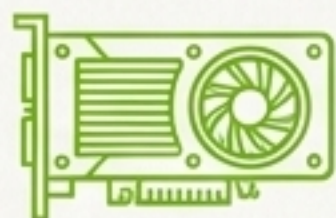
ピーク性能と実効性能の「現実ギャップ」

The Strategic Blueprint



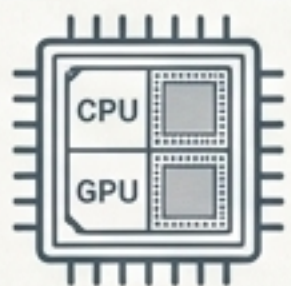
Strategic Action: 長大な特許明細書や契約書群の分析には、生コンテキストへの依存を避け、RAG（検索拡張生成）等との併用が必須。

ローカルAIハードウェアの技術動向（2025～2026年）



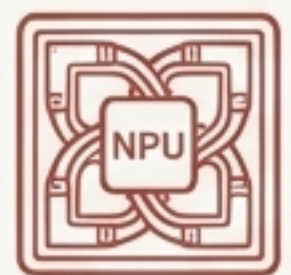
NVIDIA RTX / DGX Spark

128GB Unified / 273GB/s. 強み: CUDAエコシステムへの完全対応。AIツール互換性で圧倒的優位。



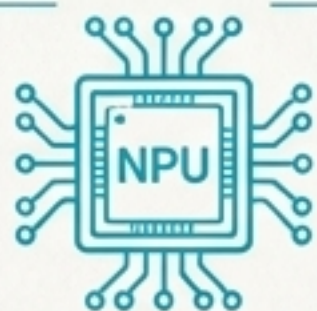
Apple Silicon - M4 Max/Ultra

帯域幅 546～800GB/s 超。強み: メモリ帯域によるLLM生成速度の優位性。
弱点: CUDA非対応。



AMD Ryzen AI Max+ 395 - Strix Halo

最大128GB / 約256GB/s。消費者向けでGPT-OSS 120Bが動作（最大30トークン/秒）。弱点: CUDA非対応。

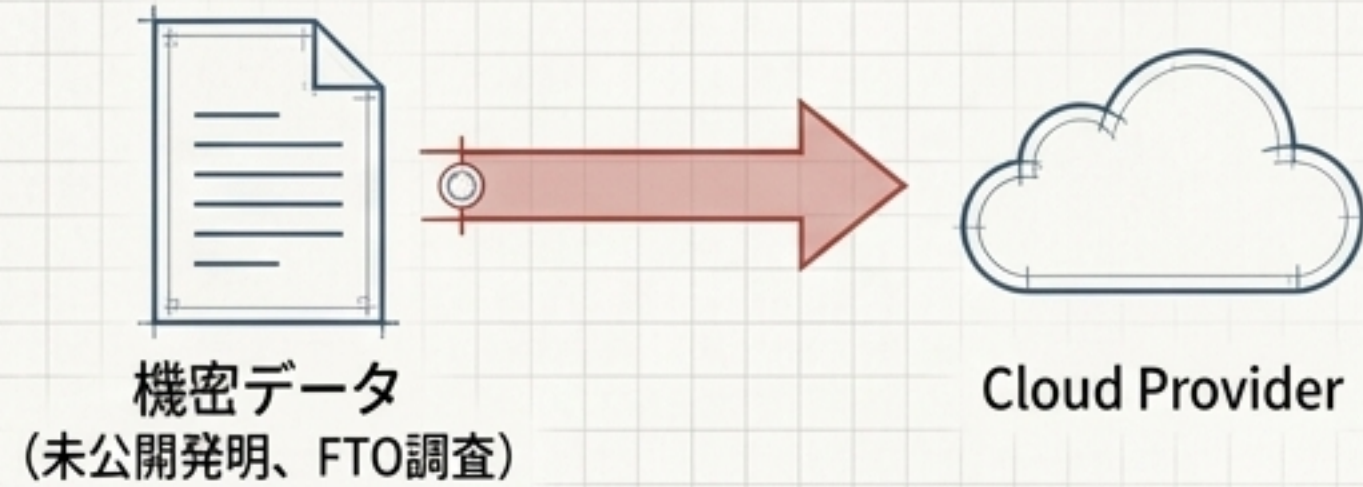


Copilot+ PC

NPU 40 TOPS以上 / 16GB RAM。軽量な日常業務 (Snapdragon X, Intel Core Ultra 200V) に特化。大規模な知財AIには不十分。

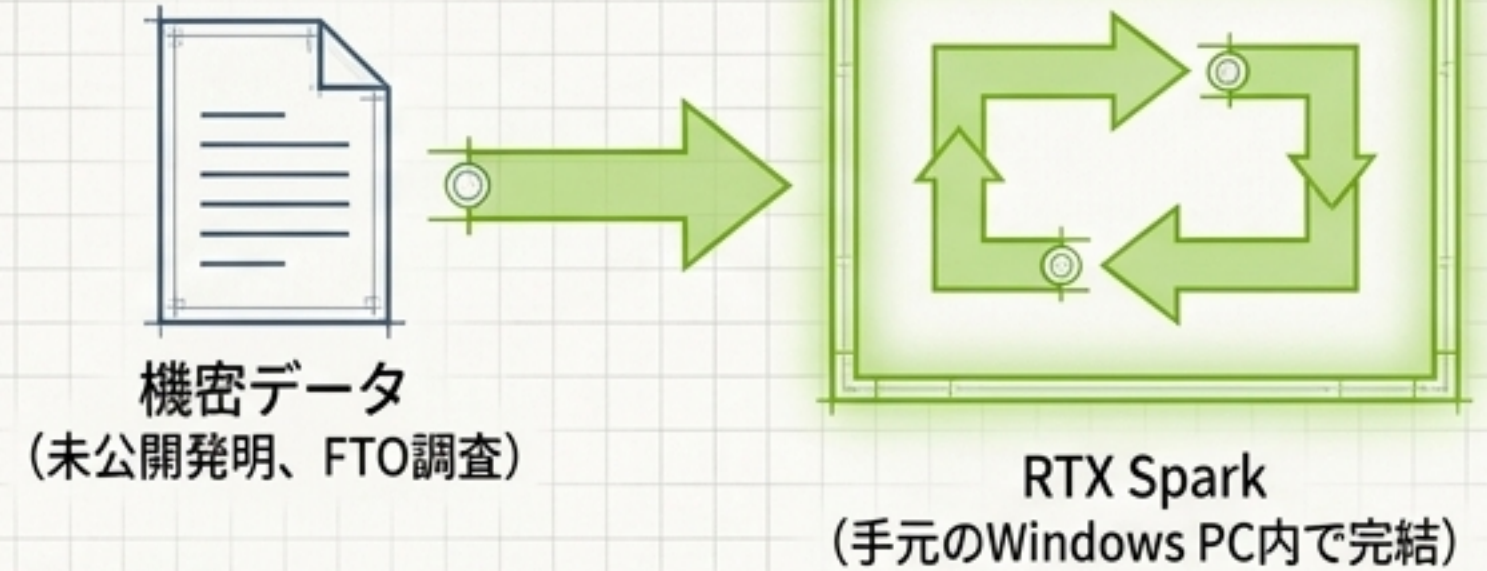
知財セキュリティのジレンマ：クラウド vs. ローカル

Cloud Flow



法的リスク: 経済産業省「営業秘密管理指針」における秘密管理性喪失の懸念 / 新規性・進歩性の破壊リスク。

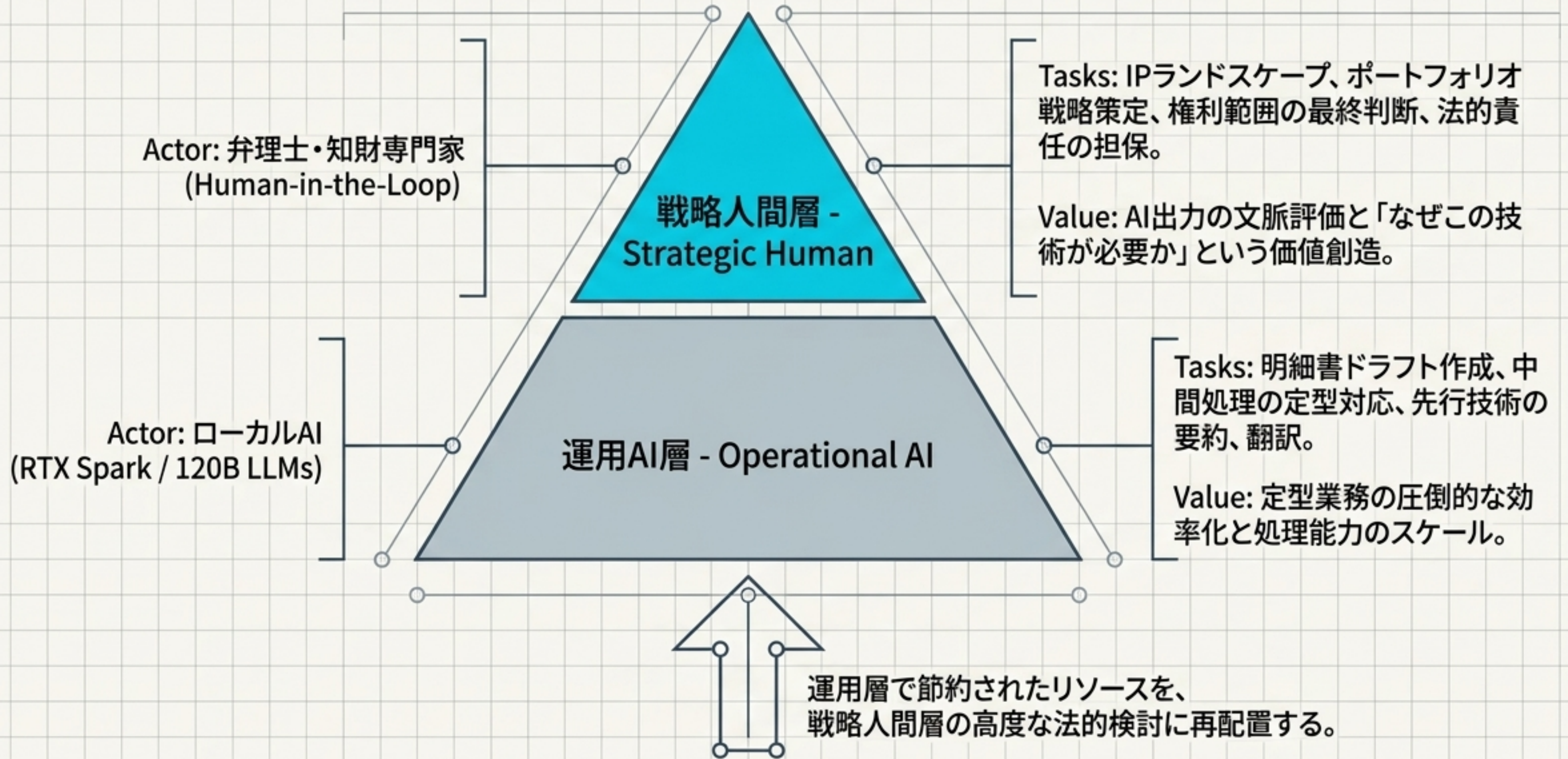
Local Flow



法的保護: プロンプト・データが外部に送信されないため、漏洩リスクと秘密管理性喪失リスクを構造的に排除。

知財という典型的な高機密領域において、「機密を外に出さずに高性能AIを使う」ニーズにローカルAIが正面から応える。

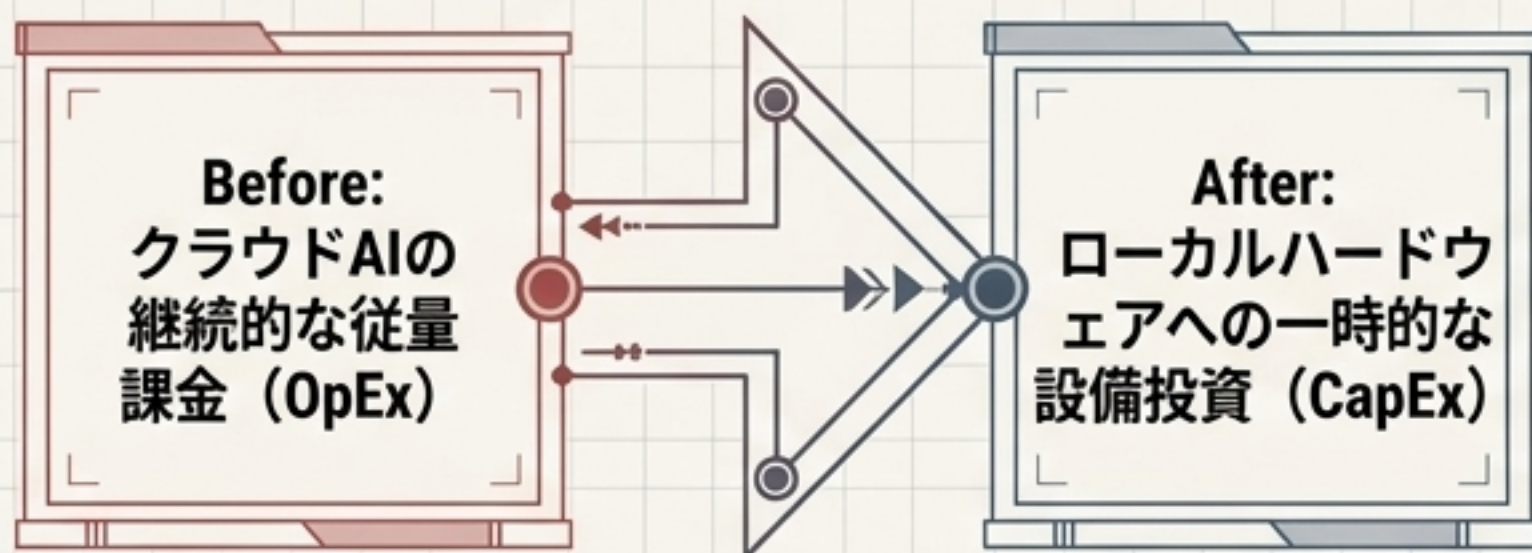
知財部門の再定義：「二層組織モデル」への移行



コスト構造の変化と出願の「民主化」

Shift 1: コストモデルの転換

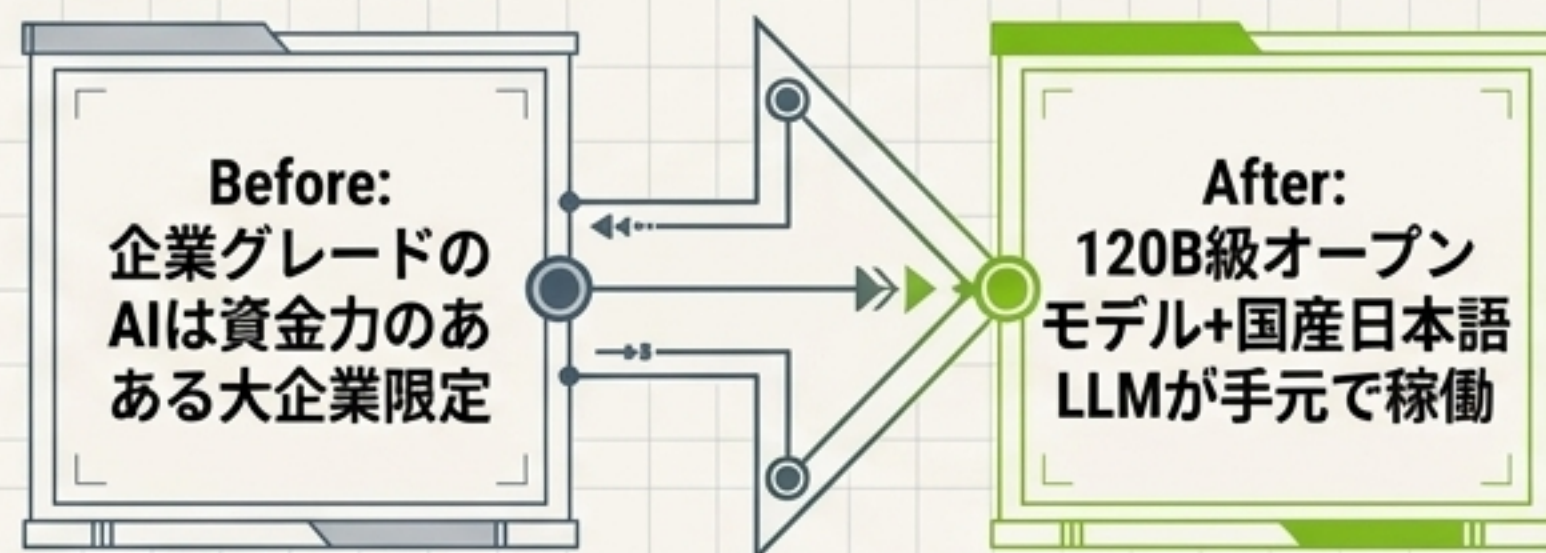
CapEx vs. OpEx



Impact: 大量利用する企業知財部では、数年でオンプレミス（ローカル）構成が圧倒的に割安になる。

Shift 2: 中小・個人事務所のエンパワーメント

規模の壁の崩壊



Impact: 高品質・低コストの文書作成環境が民主化される。
(※ただし弁理士のレビュー体制は必須)

ローカル特有の新たなリスク： 「野良AIエージェント」

OWASP Top 10 for Agentic Apps 2026

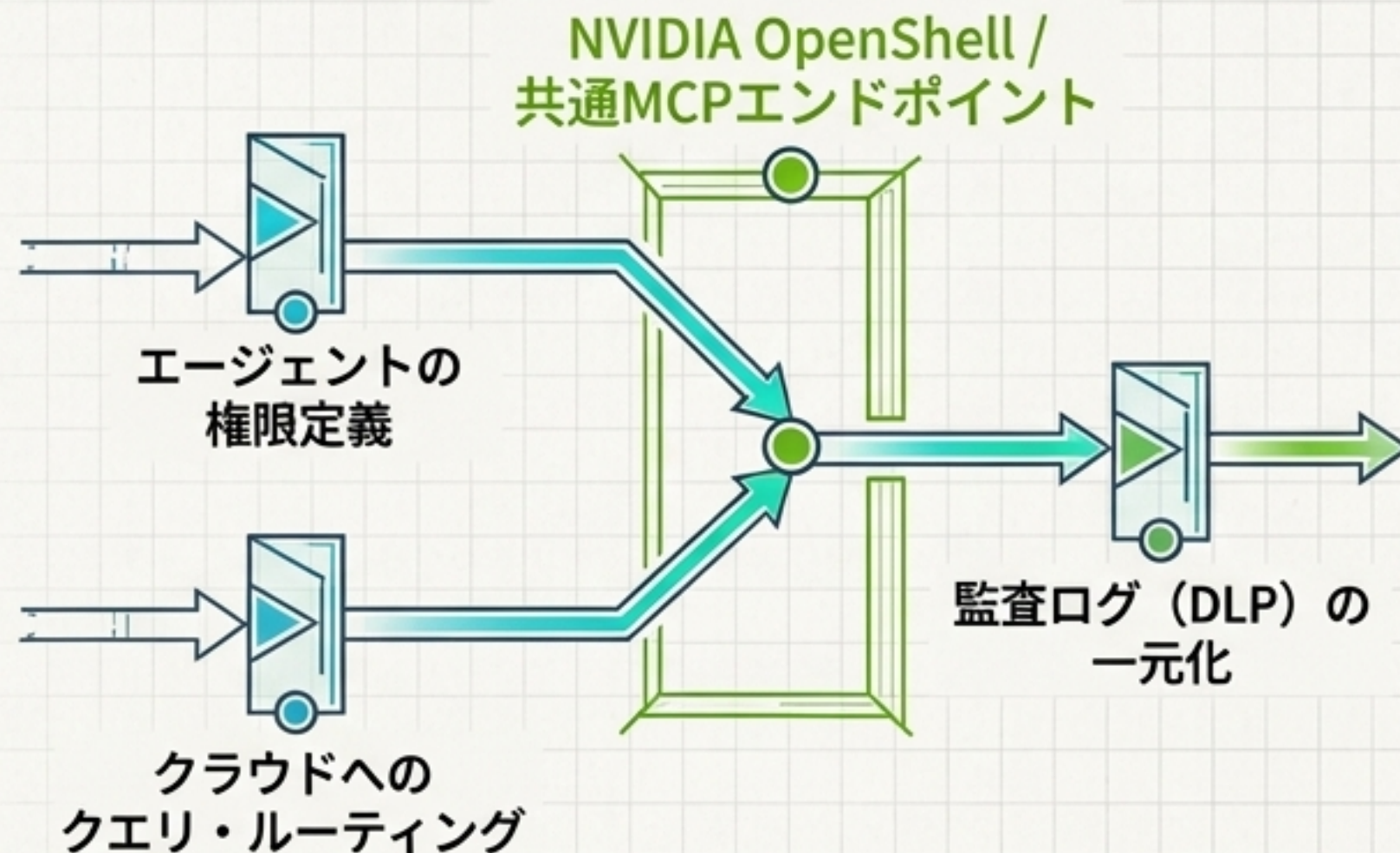
シャドーAIがローカル完結型エージェントとして再構成される。

プロンプトインジェクション、メモリ汚染、外部システムへの無断自律アクセス。

未承認SaaS以上に、ログ取得や権限統制が困難に。

統制された自由 (Controlled Freedom)

全面禁止ではなく、安全な経路の提供を。



ローカルAI導入への3フェーズ・ロードマップ

Phase 1: 即時



事実認識と期待値設定

- 社内資料の修正（DGX vs RTXの混同防止）。
- ピーク値と実効性能（15-31 t/s）の乖離を前提とした現実的なプロジェクト設計。

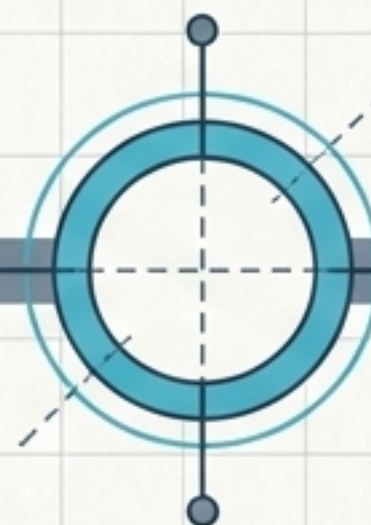
Phase 2: 3～6ヶ月



高機密業務のPoC

- 出願前発明開示、FTO、他社特許分析などクラウド禁止領域でのパイロット導入。
- 「AI事業者ガイドライン第12版」に準拠した知財部門向け利用ルールの策定。

Phase 3: 6～18ヶ月



二層体制の本格稼働

- 運用AI層と戦略人間層への役割再配置。
- 監査ログ・HITLを組み込んだエージェント・ガバナンス（OpenShell等）の構築。

本格展開のトリガー（アクションの閾値）

Threshold / Decision Gate

Critical Threshold Line

1 効率化の実績

ローカルLLMが生成した日本語明細書に対する、弁理士の「初稿修正工数」が従来比で半減した時点。

2 ハードウェア実測値

RTX Spark等の実機において、大型モデルが「30トークン/秒以上」で安定動作し、長文脈の有効活用率が実用レベルに達した時点。

3 国産モデルの成熟

120B級の国産LLMが海外モデルに匹敵する日本語性能を持ち、ローカル環境で機密性と性能が完全両立した時点。

ハードウェアの準備は整った。最大のボトルネックは、技術の制約を理解し組織を再設計する「人間の戦略的判断」にある。