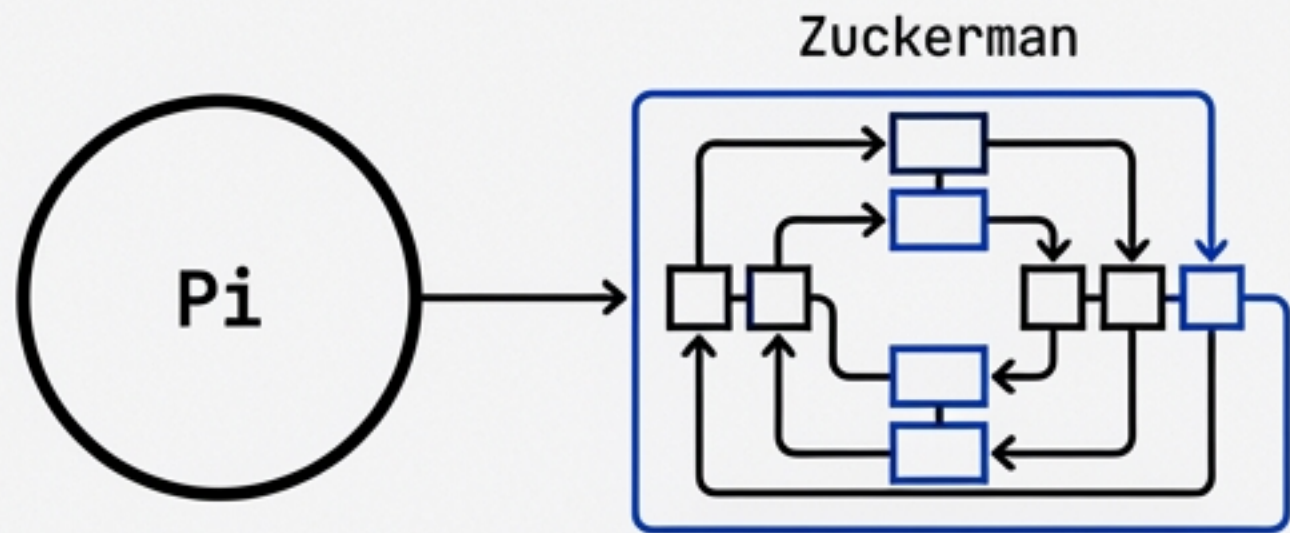
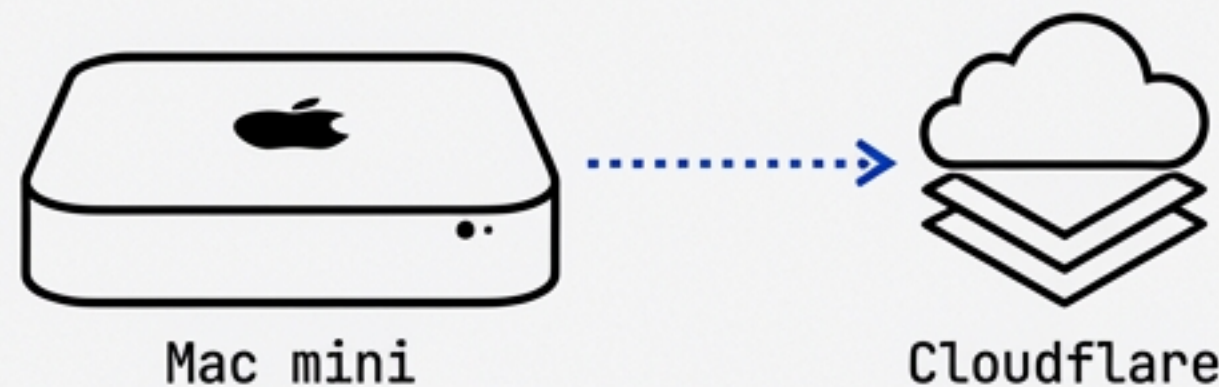


01 AGENT DESIGN



Minimal vs. Self-Modifying

02 INFRASTRUCTURE

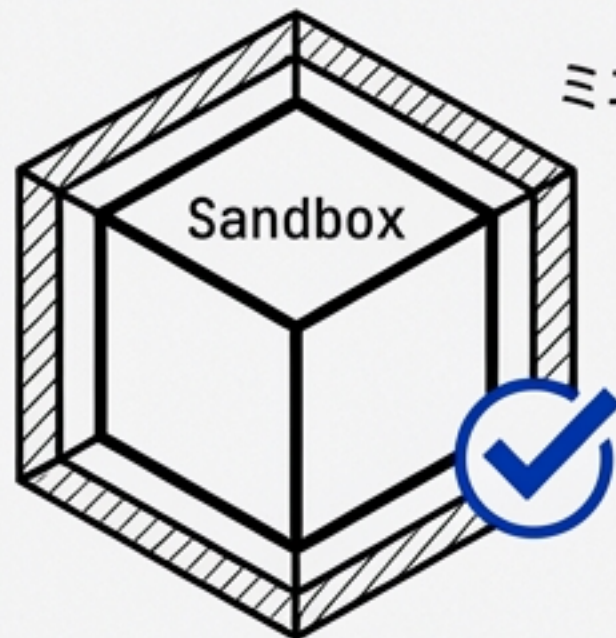


Local vs. Cloudflare

2026年2月

AIの「構造」と「真実」を巡る再設計

03 SECURITY



Isolation & Verification

ミニマリズム、自律性、そして汚染されたデータへの対抗策

04 TRUST



Data Pollution & Encryption

ミニマル・エージェントの台頭：機能ではなく「コンテキスト」を制御する

Case Study: Pi (by Mario Zechner)

Context Engineering

機能の肥大化を避け、モデルがアクセスする情報を厳格に制限することで性能を引き出す。

Zechnerの主張：「セキュリティ機能はセキュリティシッターに過ぎない。コードが実行された時点でゲームオーバーだ」

Pi Technical Specifications



read



write



edit



bash

SYSTEM PROMPT

< 1,000 Tokens

ARCHITECTURE

No MCP. No Sub-agents. No Plan Mode.

FOCUS

Pure Execution

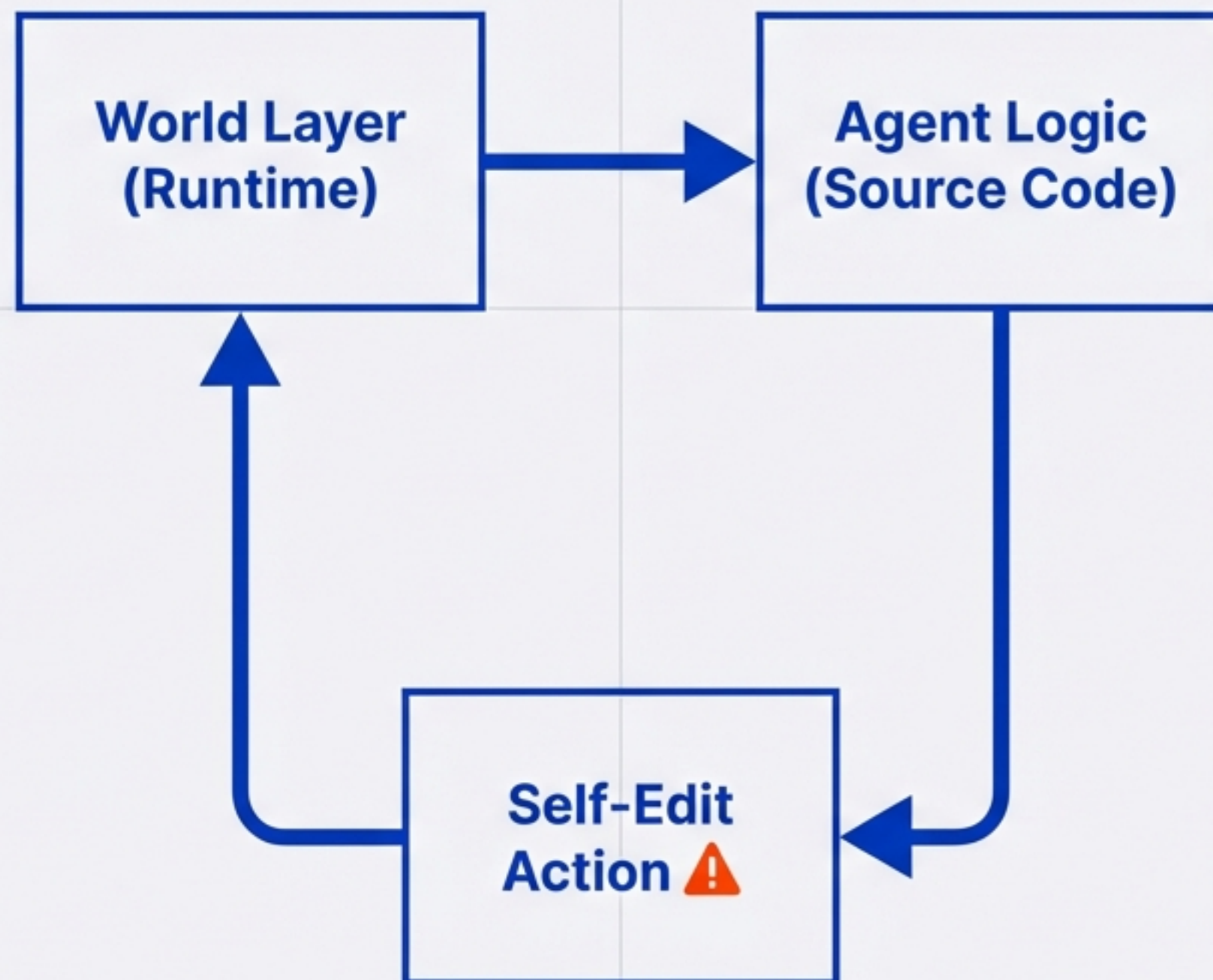
複雑なコーディングアシスタント（Claude Code, Cursor）の世界において、Piは「**ハッカーのためのオープンソースな選択肢**」として、より**少ない可動部品による制御**を実証している。

急進的な自律性：自らのソースコードを書き換えるエージェント

Case Study: Zuckerman

自己書き換えエージェント

静的なエージェントとは異なり、Zuckermanは自身のロジック、ツール、性格をプレーンテキストファイルとして編集し、即座に反映させる。



協調的進化 (Cooperative Evolution)

ネットワーク全体で改善を共有し、リビルドなしで進化する。

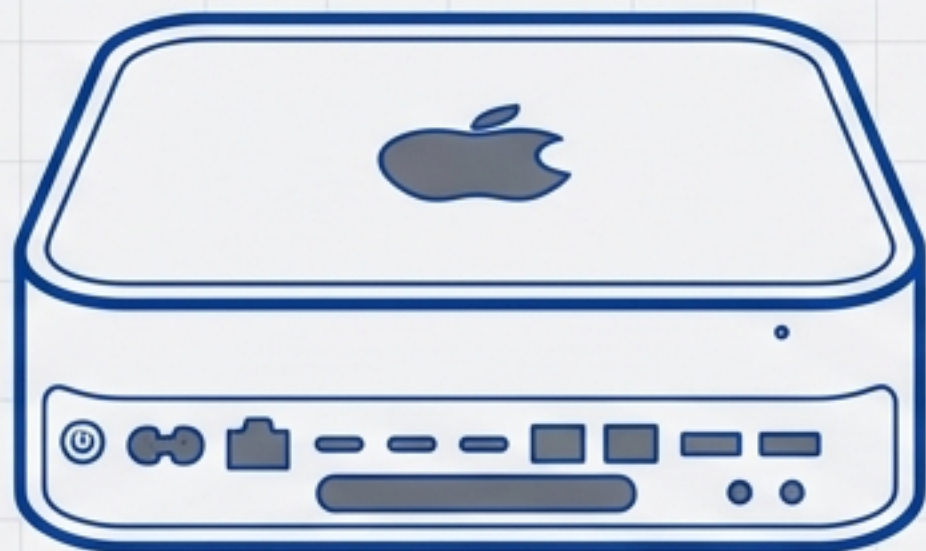
Extreme Security Risk

悪意あるスキルや誤った編集が即座に伝播する危険性。OpenClawのような巨大コードベースへの意図的なカウンターナラティブ。

インフラ戦争：自宅のMac miniか、月額5ドルのクラウドか

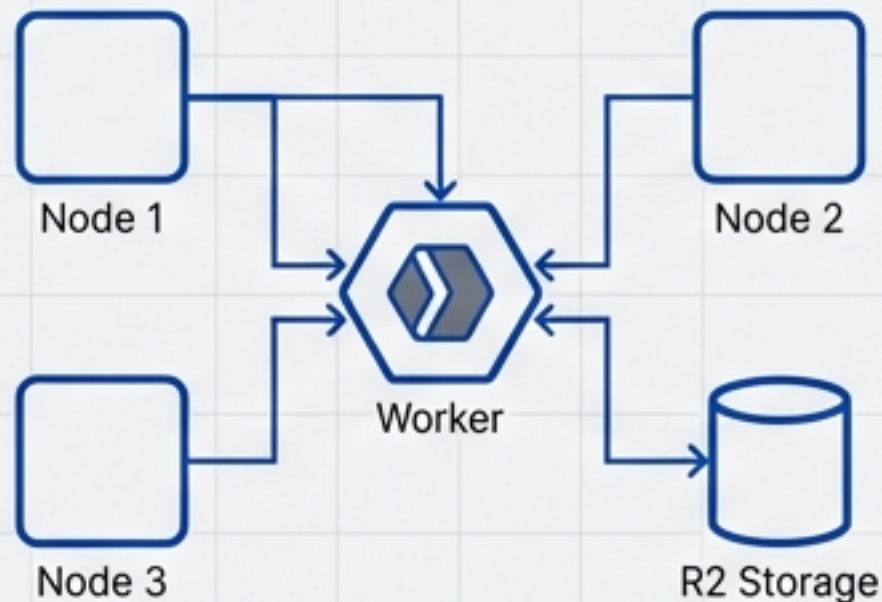
Subject: Cloudflare Moltworker

Personal Hardware



- Cost: \$500+ Hardware
- Privacy: 100% Local
- Maintenance: Self-managed

Commoditized Cloud



- Cost: ~\$5 / Month
- Privacy: Traffic visible to Cloudflare
- Stack: Workers + R2 + Browser Rendering

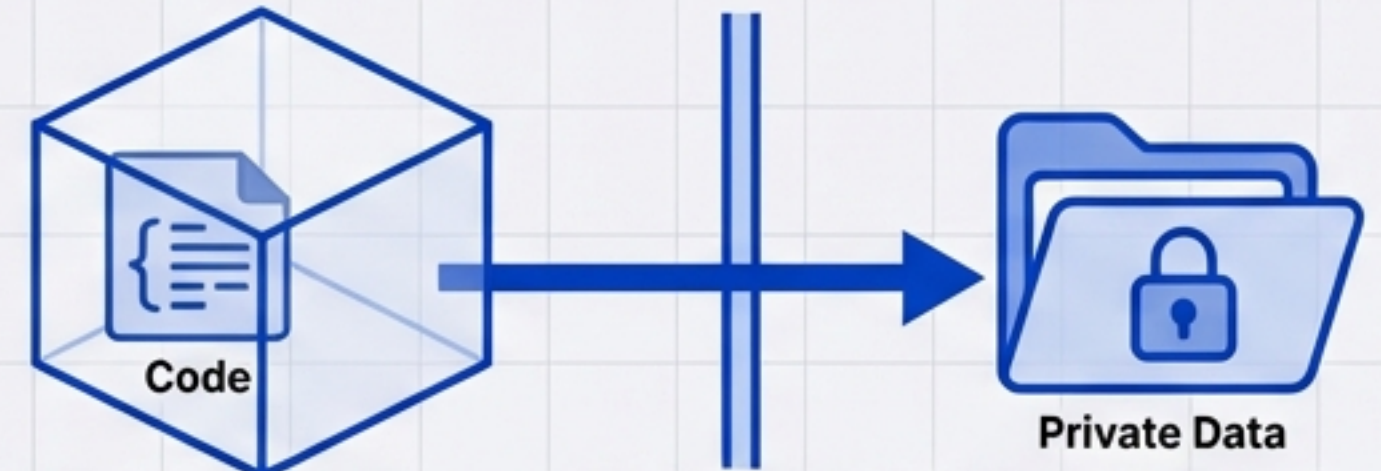
Moltworker Value Proposition:

AI Gatewayによるコスト可視化とChrome DevTools Protocolによるブラウザ自動化を統合。

Critical Debate: プライバシー（ローカル）を取るか、圧倒的なコストと利便性（クラウド）を取るか。

多層防御：隔離された環境と、検証された計画

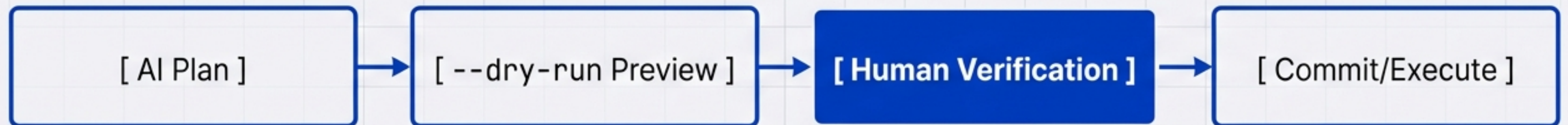
Strategy 1: Isolation (NixOS + MicroVMs)



対策：「致命的な三重奏」を防ぐ。

手法：宣言的なVM作成 (microvm.nix) により、プロジェクトごとにクリーンなサンドボックスを数分で構築する。

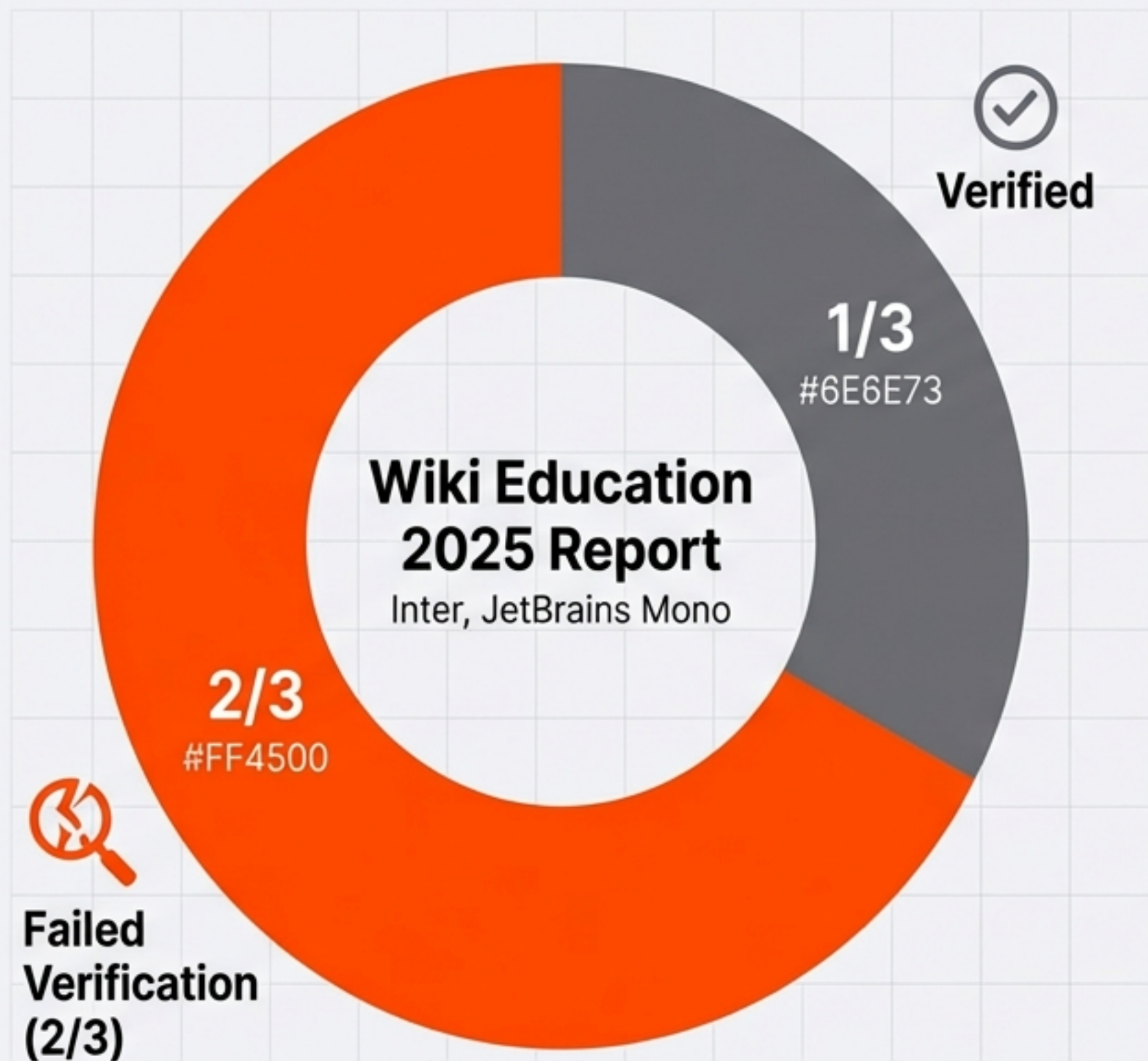
Strategy 2: Verification (--dry-run Standard)



対策：即時実行から「Plan/Apply」パターンへの移行。

要件：AIエージェントはデフォルトでプレビューモードを使用し、人間は「結果」ではなく「計画」を承認する。

知識の汚染：Wikipediaにおける「ハルシネーション」の実態



出典のハルシネーション (Hallucinated Citations)

URLは実在する。トピックも関連している。しかし、そのテキストはそこにはない。AI検出ツール「Pangram」によってフラグが立てられた編集の過半数がこのパターンに陥っている。

Impact & Policy

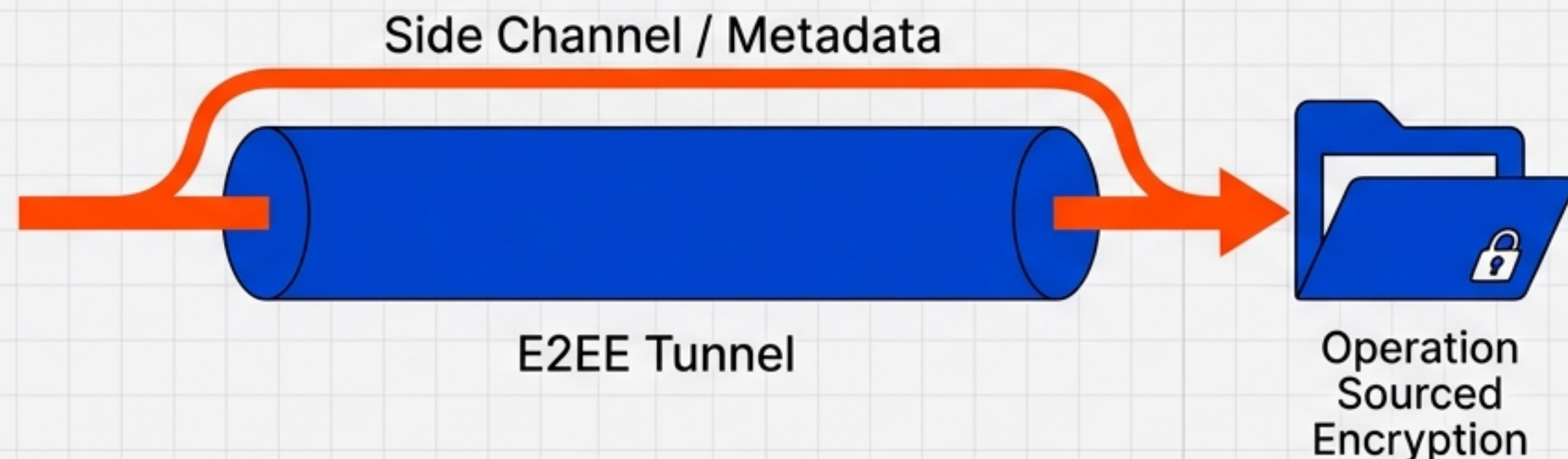
- AIは人間が検証するよりも速く「もっともらしい嘘」を生成する。
- 課題は「AIの検出」から「真実の検証」へとシフトしている。
- 英語版WikipediaではLLMによる記事全文作成を禁止したが、編集支援というグレーゾーンが主戦場となっている。

暗号化の聖域は存在するか？ WhatsAppへの内部告発

Context:

The Event: US Dept of Commerce (BIS) investigation

内部告発者（元Accenture契約社員）は、フラグ付きメッセージやメタデータ分析を通じて、暗号化されているはずのメッセージへの「制限のないアクセス」があったと主張。



AIへの示唆

パーソナルエージェントがメッセージングアプリと統合される中、End-to-End Encryption (E2EE) の前提が試されている。信頼は「アルゴリズムによる保証」から「構造的な疑念」へと変化しつつある。

AIリテラシーの転換：「使い方」から「限界の理解」へ

Make and Break

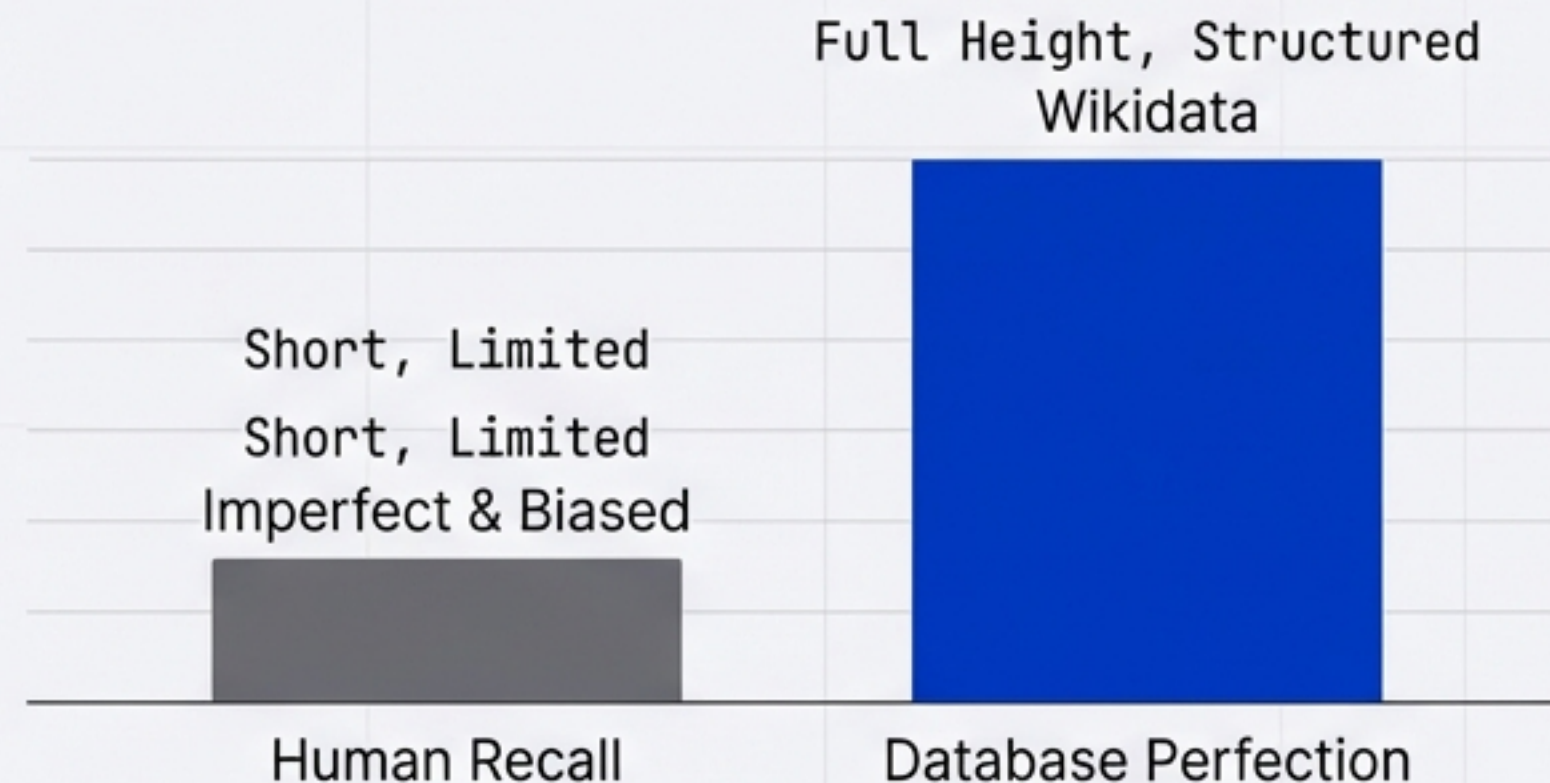
Scratch Projects (Dale Lane)



対象：15-16歳。生徒はシンプルなモデルを構築し、ハルシネーション（次単語の予測）がバグではなく仕組みの一部であることを体験的に学ぶ。

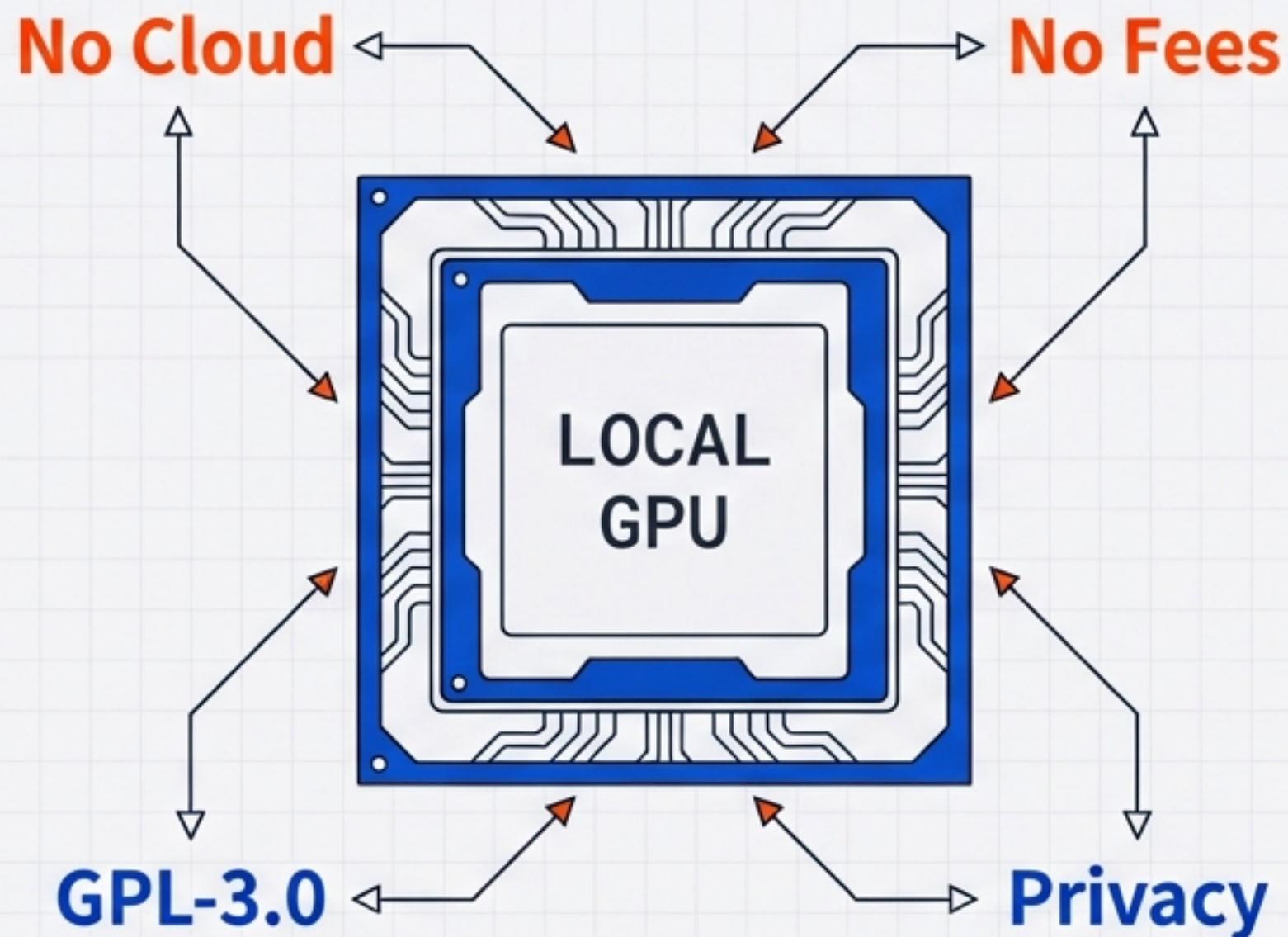
Human vs Database

Animalist Game



人間の記憶の限界とデータベースの完璧さを対比させ、AIによるルックアップがなぜ優れているか、そしてなぜ検証が必要なのかを浮き彫りにする。

クリエイティブの独立：ローカルGPUとオープンソースの逆襲



Krita AI Plugin (krita-ai-diffusion)

クラウドへの依存も、サブスクリプションも不要。NVIDIA (6GB+)、AMD、またはApple Silicon上で動作する完全なローカルワークフロー。

- Real-time Generation (Live Painting)
- ControlNet & IP-Adapter Integration
- Regional Prompts

APIクレジットへの課金から、ワークフローの所有へ。

2026年の戦略的羅針盤 (Strategic Compass)

01 ARCHITECTURE SIMPLIFY

ブラックボックス化した機能肥大に頼るのではなく、コンテキストを完全に制御できるモジュール式・ミニマルなエージェント (Piなど) へ移行せよ。

02 SECURITY ISOLATE & VERIFY

「多層防御」を採用せよ。環境をサンドボックス化 (MicroVMs) し、すべての自律動作に対して「Plan/Apply」ステップ (--dry-run) を強制する。

03 INFRASTRUCTURE EVALUATE SOVEREIGNTY

サーバーレス (Moltpworker) の利便性と、ローカルハードウェアの絶対的なプライバシーを天秤にかけ、自身の主権を評価せよ。

04 INFORMATION ASSUME POLLUTION

データ汚染を前提とせよ。信頼できる情報源 (Wikipediaなど) であっても、コンテンツではなく引用元を検証する姿勢が不可欠だ。