

# 自律エージェントの台頭と「Vibe Coding」の落とし穴

拡大する能力、低下する信頼性、そして人間が担うべき「管理」の役割



## Thesis: 能力の拡大

Kimi K2.5 (1兆パラメータ) と ChatGPT (Linux環境化) がもたらす、かつてない自律性と実行能力。



## Antithesis: 信頼性の危機

Cloudflareの事例に見る「Vibe Coding (雰囲気実装)」のリスクと、技術的負債の蓄積。



## Synthesis: 新たな指針

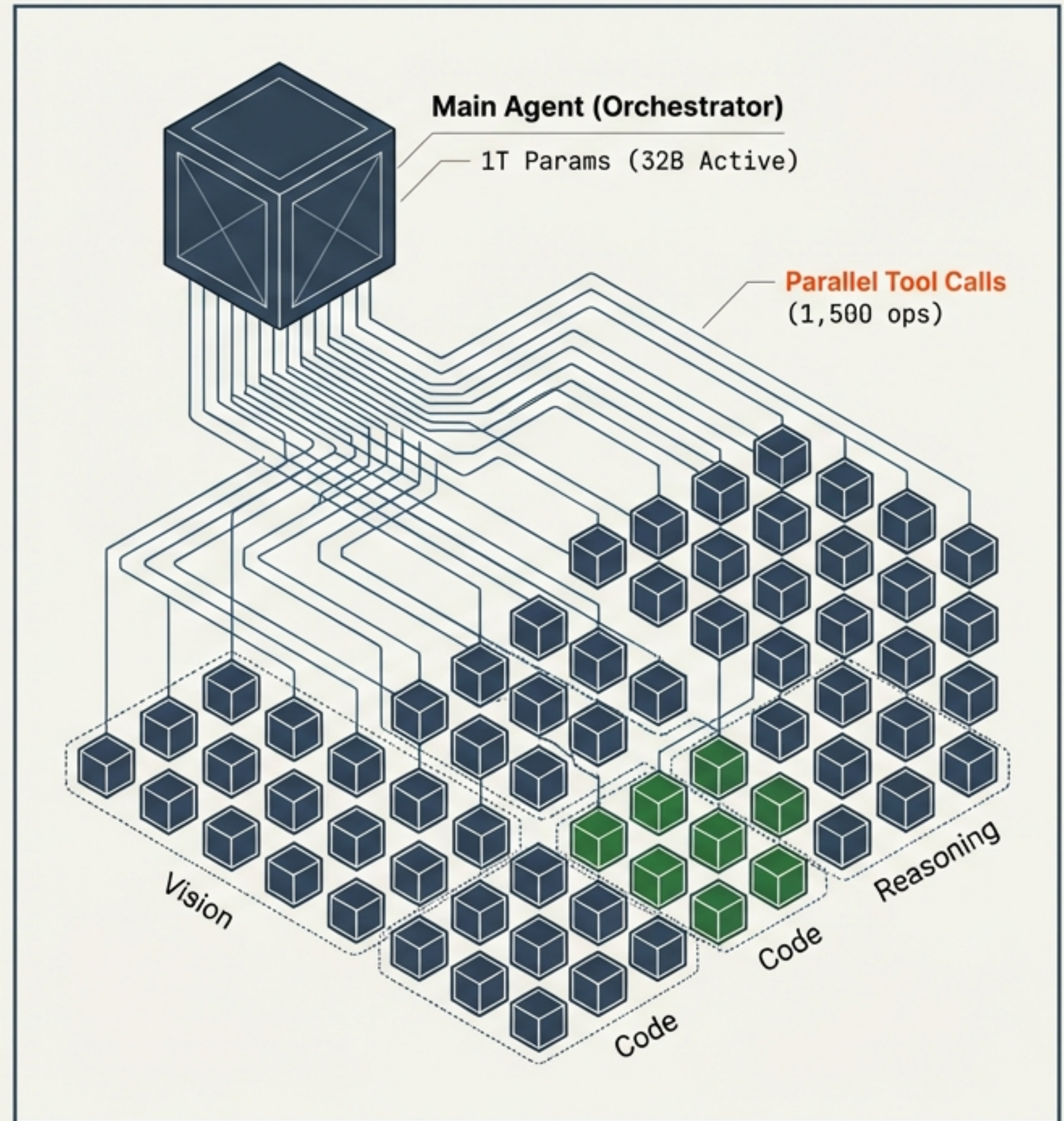
コーディングは「書く」から「管理・検証する」へ。Karpathy氏が示唆するエンジニアリングの変容。

# Kimi K2.5 : 1兆パラメータと「エージェントスウォーム」の衝撃

- **大規模MoE:** 総パラメータ1兆、アクティブ32B。DeepSeek R1から1年を経て登場した中国発のオープンソースモデル。
- **Agent Swarm (エージェント群):** 最大100のサブエージェントを自律制御。複雑なタスクを分解し、並列処理が可能。
- **並列実行:** 1,500回のツール呼び出しを並列で実行可能。
- **視覚能力:** 視覚タスクにおいてSOTA (State of the Art) を達成。

## INSIGHT

**なぜ重要か:** 単なるチャットボットではなく、複雑な実務を「自律的に分解して解決する」能力が飛躍的に向上しています。Hugging Faceで重みが公開されており、商用利用も可能です。

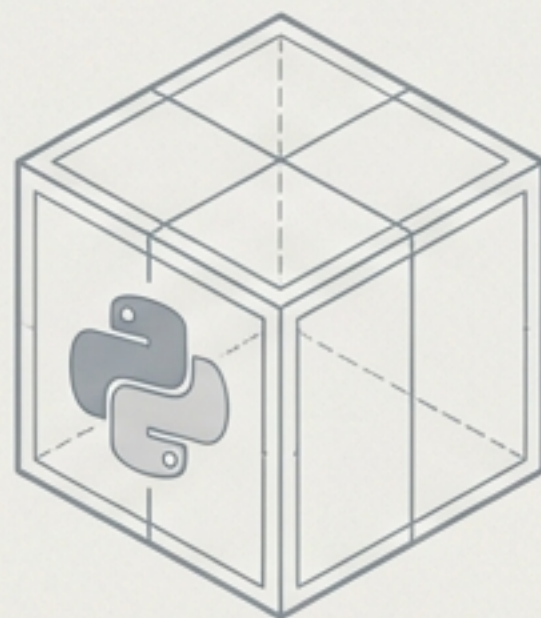


# ChatGPTの進化：Pythonサンドボックスから「開発環境」へ

## Bash, npm, C++をサポートし、実質的なLinux環境としての機能を獲得

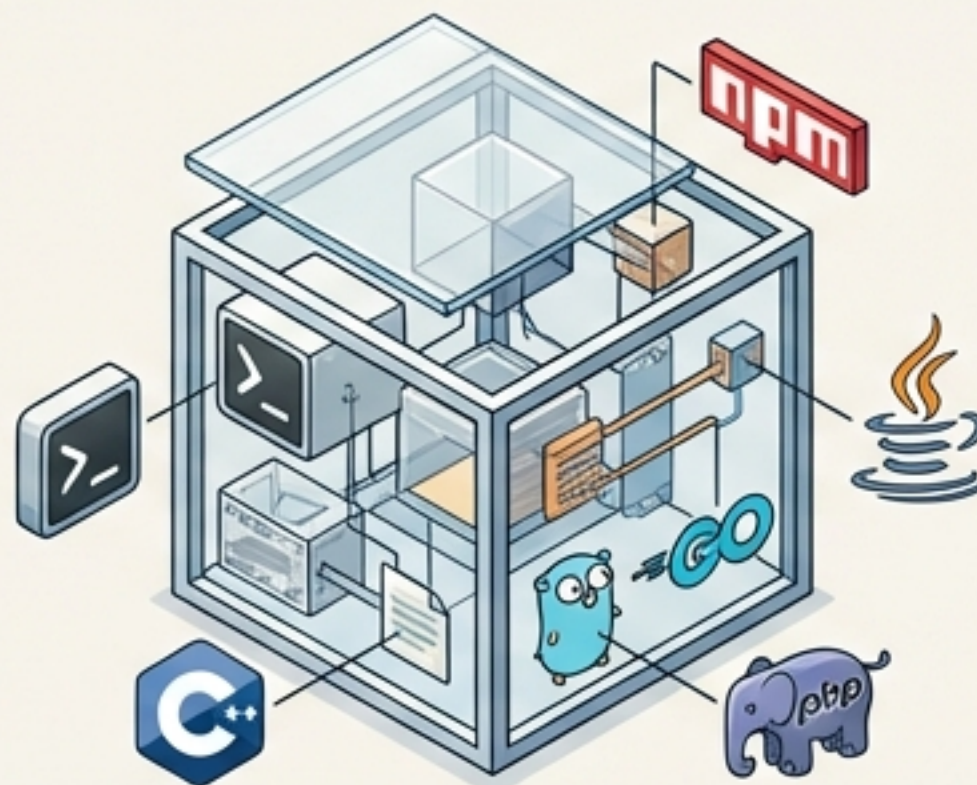
※ 公式アナウンス前の機能追加であり、仕様変更の可能性があります。

### BEFORE: Code Interpreter



- Pythonのみ実行可能
- ライブラリ制限あり
- 計算機としての利用

### AFTER: New Container



- Bashコマンド実行
- pip/npmパッケージ追加
- ファイル操作
- C++
- 多言語対応

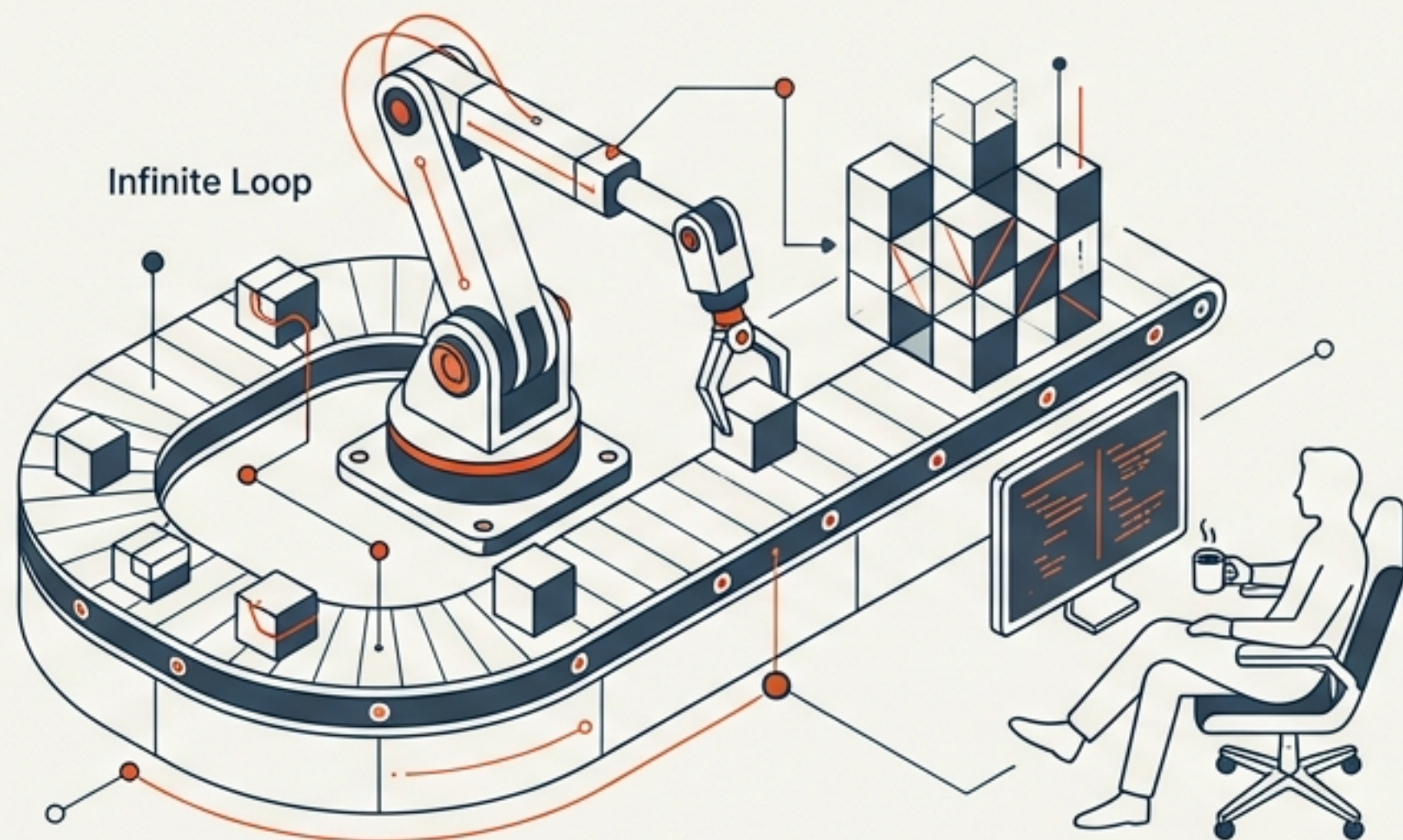
### 実務への影響:

依存関係のある外部ライブラリをインストールし、複雑なビルドやスクリプト実行がChatGPT内で完結可能に。

### 競合関係:

Claude CodeやCursorなどの「コーディングエージェント」に対抗。

# Andrej Karpathy氏の視点：「エージェントは疲れない」



“Agents don’t get tired, their morale doesn’t dip.”  
(エージェントは疲れず、士気も下がらず、ひたすら試行し続ける)

## Coding as Gaming

プログラミングはテキストエディタでの作業から、「StarCraft」や「Factorio」のようなリソース管理・指示出しゲームへと変化する。

## The Tireless Worker

人間なら諦めるようなエラーの連鎖に対しても、AIは粘り強く解決策を探り続ける。

## The Risk

一方で、人間に「脳の萎縮 (brain atrophy)」や怠慢を招くリスクも指摘。

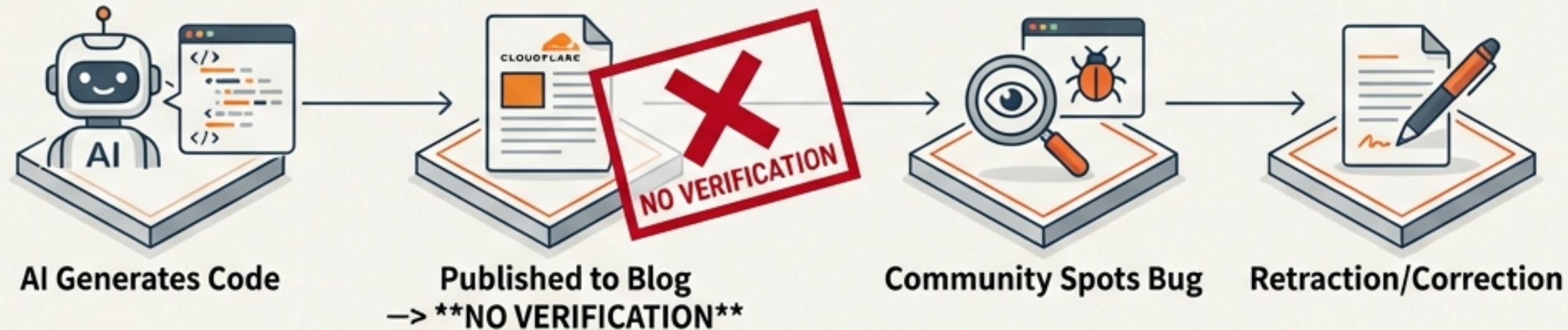
# 「Vibe Coding」の代償：Cloudflareの教訓

## TERM: Vibe Coding (バイブ・コーディング)

コードの意味を理解せず、AIに生成させたものを「雰囲気」だけで採用し、動作検証を疎かにする開発スタイル。



## CASE STUDY: Cloudflare Blog Incident



- **事象:** 大手テック企業の技術ブログで、Matrixプロトコルの実装を主張したが、公開されたコードは実際には機能しない状態だった。
- **原因:** AI生成コードを適切な検証（コンパイル・実行テスト）なしに記事化してしまった可能性。
- **結果:** コミュニティからの指摘により、「実装」から「概念実証」へと記事の修正を余儀なくされた。

# 効率と品質のジレンマ：失われる「職人技」

AI生成コードは「動く」が、「良いコード」とは限らない。



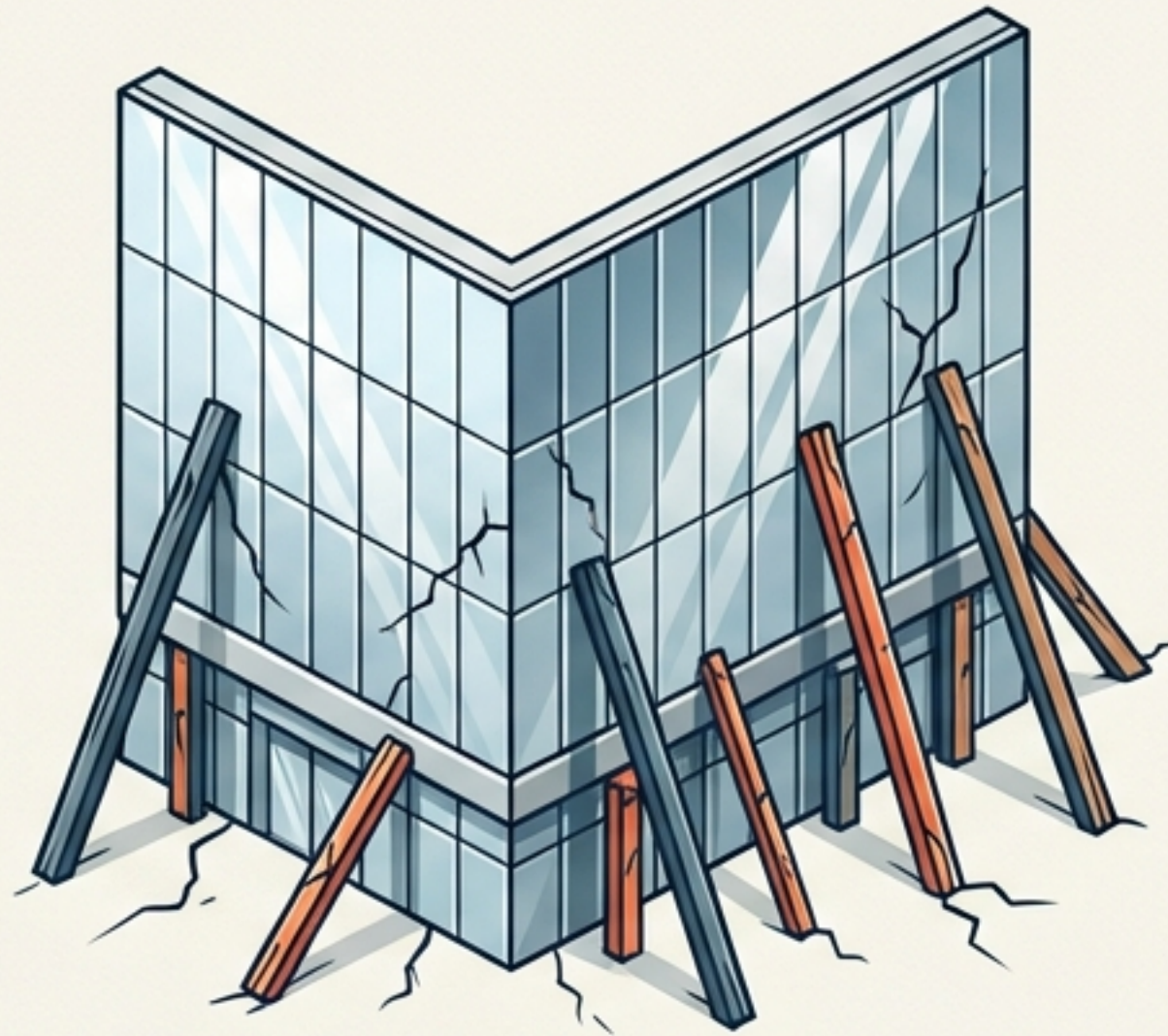
## 1. 保守性の欠如

可読性が低く、将来的な修正コスト（技術的負債）が積み上がる。



## 2. 構造的欠陥

作る人（AI）と維持する人（人間）が異なるため、エンタープライズソフトのような「使いにくさ」がコード自体に宿る。

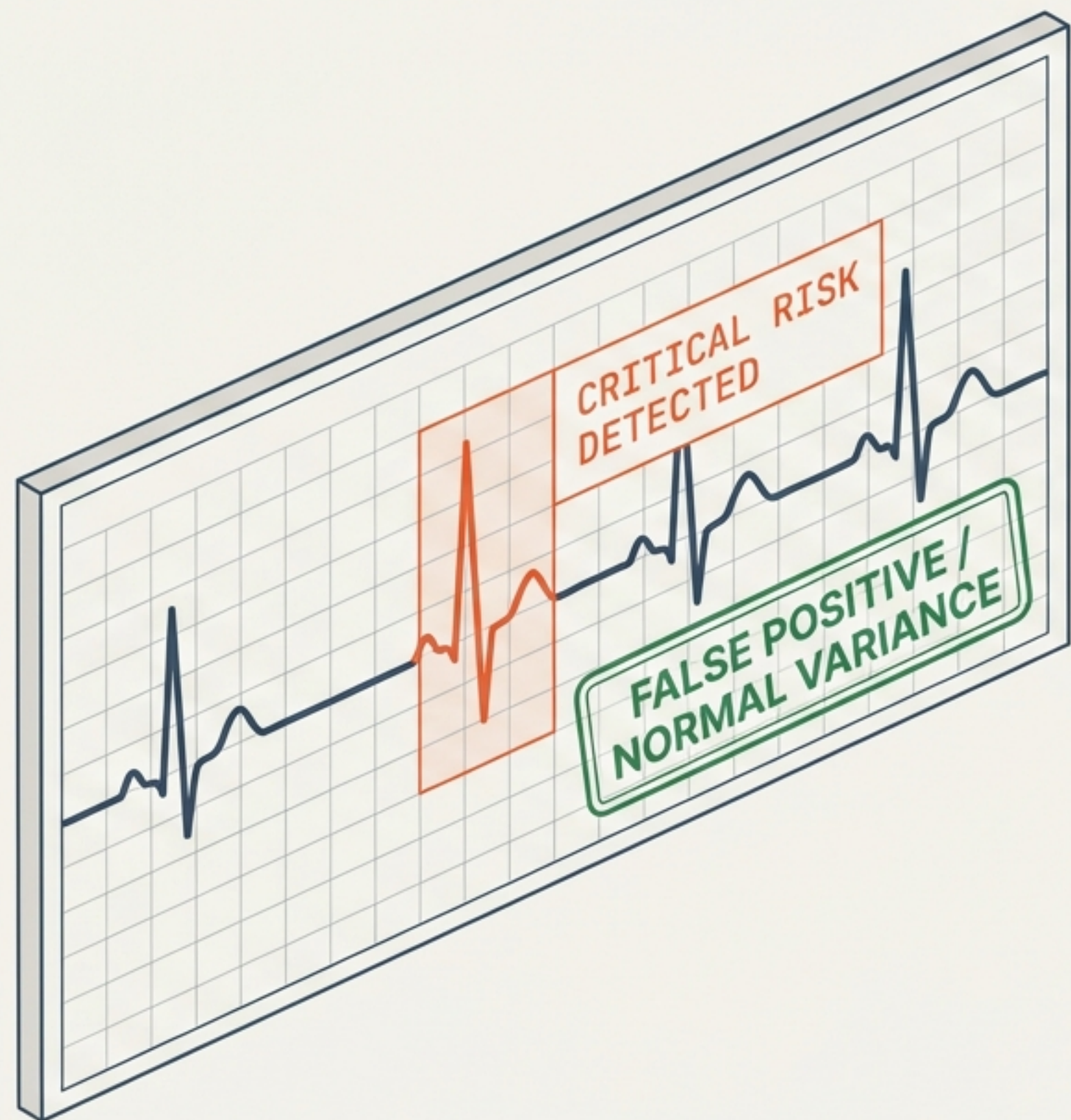


## 3. レビューの形骸化

レビューする人間のスキル低下と、生成速度への過信が重なり、バグがすり抜ける。

"The code works, but..."（コードは動く、しかし…）  
という違和感は、職人技（Craftsmanship）の軽視から生まれている。

# 専門領域での幻覚：Apple Watch 10年分のデータ分析



## The Experiment:

- **Input:** 10年分に及ぶ個人の健康・バイタルデータ。
- **AI Output:** 深刻な健康リスクを示唆する診断結果。
- **Reality:** 医師による確認の結果、AIの診断は誤りであった。

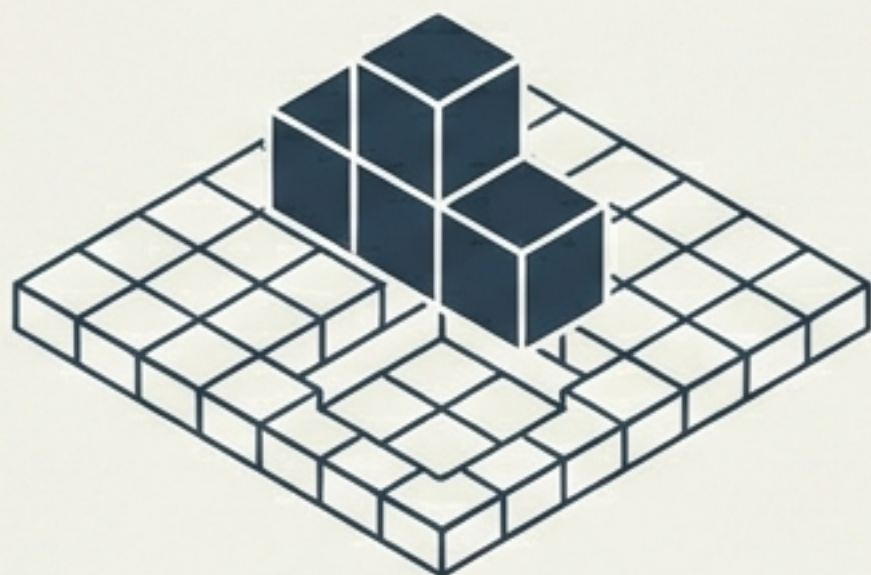
## Why It Failed:

- **文脈の欠如:** 心拍数や歩数などの数値データのみで、生活背景（ストレス、食事、環境）を考慮できない。
- **適合過剰:** パターンマッチングに長けすぎているため、存在しない因果関係を見出してしまう。

AIは「トレンド把握」には有用だが、  
「結論（診断）」を委ねるのは危険である。

# エコシステムの現在地：推論能力と透明性の追求

## TetrisBench



Gemini 3 Flash Win Rate: 66%

概要: テトリスを用いたリアルタイム意思決定能力のベンチマーク。

意義: 文章生成だけでなく、瞬時の状況判断と計画能力が測定されている。

## AI2 Open Agent



SWE-bench pass@1: 54%

概要: 完全オープンソース（モデル、訓練データ、パイプライン）のコーディングエージェント。

意義: ブラックボックス化する商用モデルに対し、検証可能な「透明性」を重視するアプローチ。

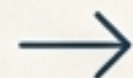
# Synthesis : マネジメントこそが、AI時代のスーパーパワー

Shift in Value : 「どう作るか (How)」 → 「何を作るか (What)」 & 「正しく作れているか (Verify)」



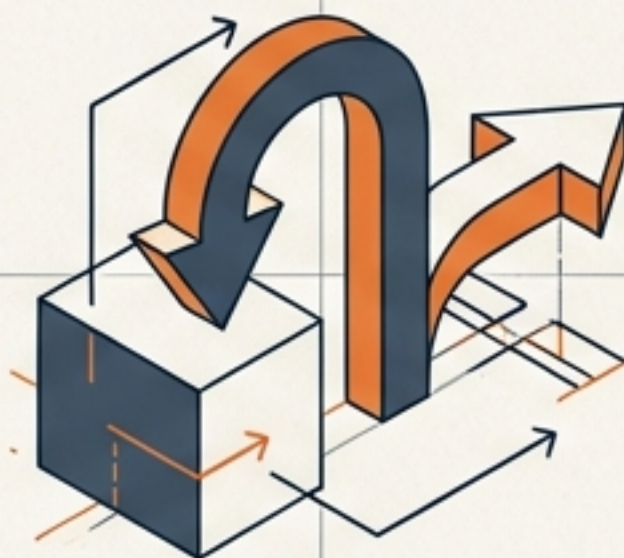
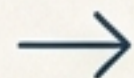
## Direction (指示)

曖昧な要件ではなく、エージェントが実行可能な明確なタスク定義を行う。



## Review (検証)

生成物を無批判に受け入れず、意図通りか、安全かを確認する「目利き」の力。

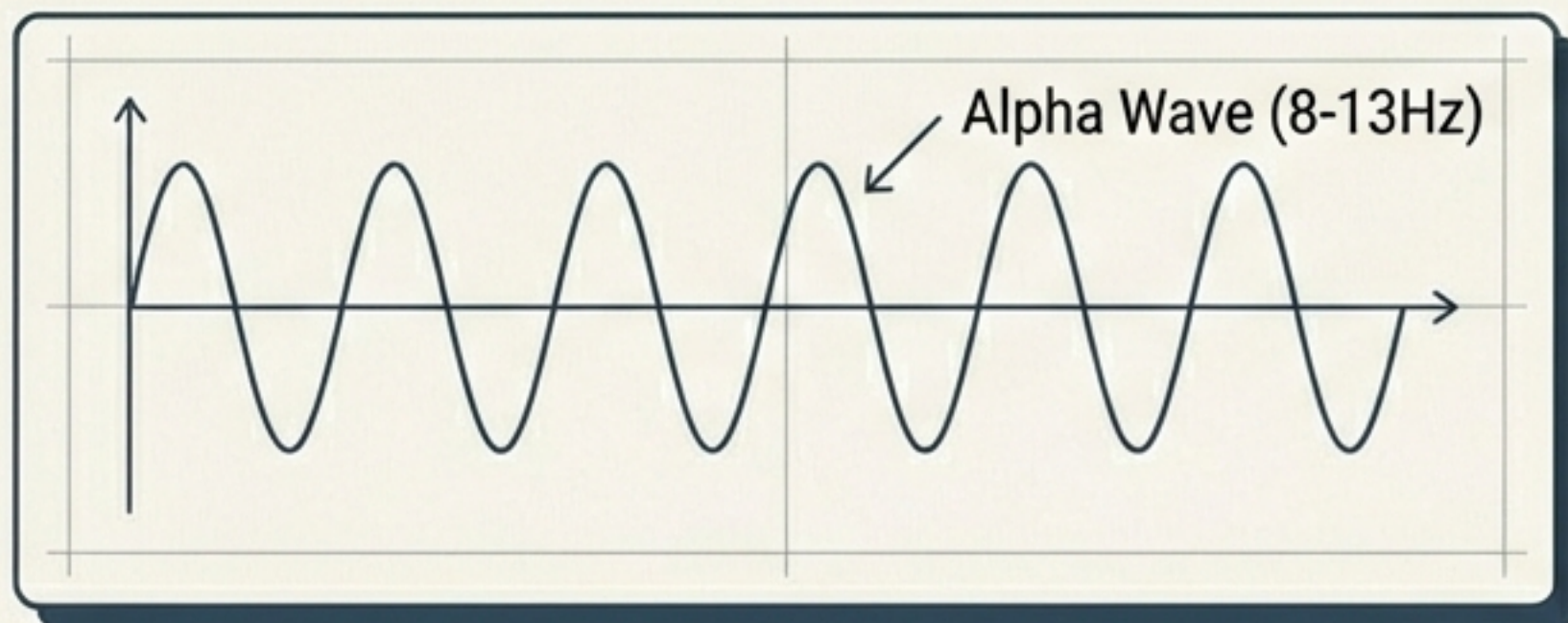


## Pivot (転換)

実装コストが下がったため、間違った方向性ならすぐに捨てて作り直す判断力。

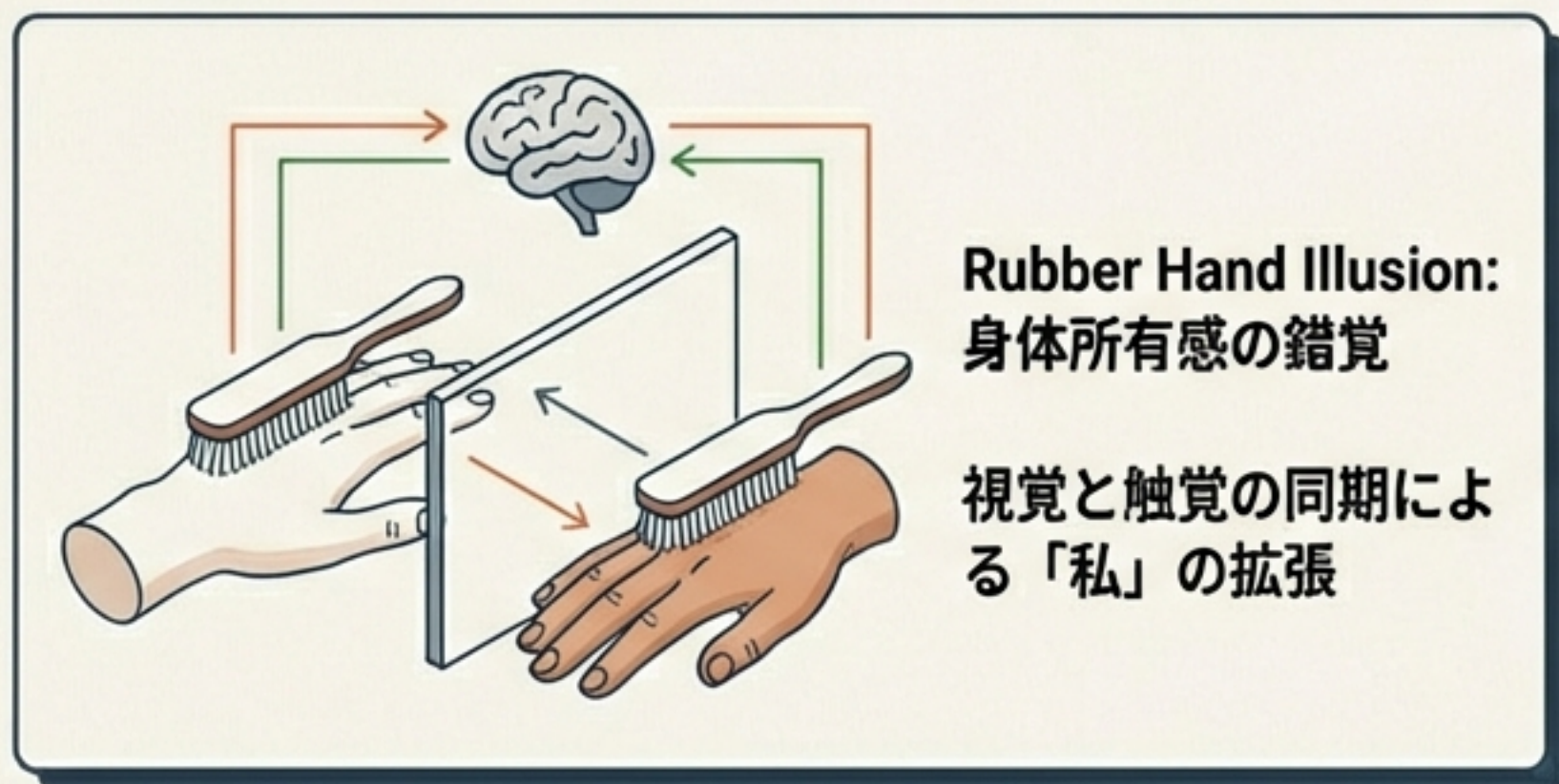
Karpathy氏の言う「コーディングのマネージャー化」は、このスキルセットへの移行を意味する。

# Science & Self : 「私」の境界線



## Discovery:

脳波（アルファ波）の周波数が、身体  
の所有感（これは自分の手か？）を決定して  
いることが発見された。



## Relevance to AI:

- 人間が「自分」を感じるメカニズムが物理的に解明されつつある。
- AIエージェントが高度化する中、意識や自己認識の定義を再考するヒントとなる。

# 結論：信頼して、検証せよ (Trust, but Verify)

KimiやChatGPTのような強力なエージェントを活用しつつ、「Vibe Coding」の罠を避けるための3ステップ。



## 1. Embrace the Swarm

KimiやChatGPTのコンテナ機能をプロトタイピングに積極採用し、実装速度を上げる。



## 2. Reject Vibe Coding

生成されたコードは必ずローカルで実行検証し、自分の理解及ばないコードは採用しない（または解説を求める）。



## 3. Invest in Architecture

AIにコードを書かせ、人間はシステム全体の設計と品質管理（レビュー）にスキルをシフトする。



Successful  
AI-Augmented  
Development

# Coding is dead. Long live Engineering.

コーディングは死んだ。エンジニアリングよ、永遠なれ。

