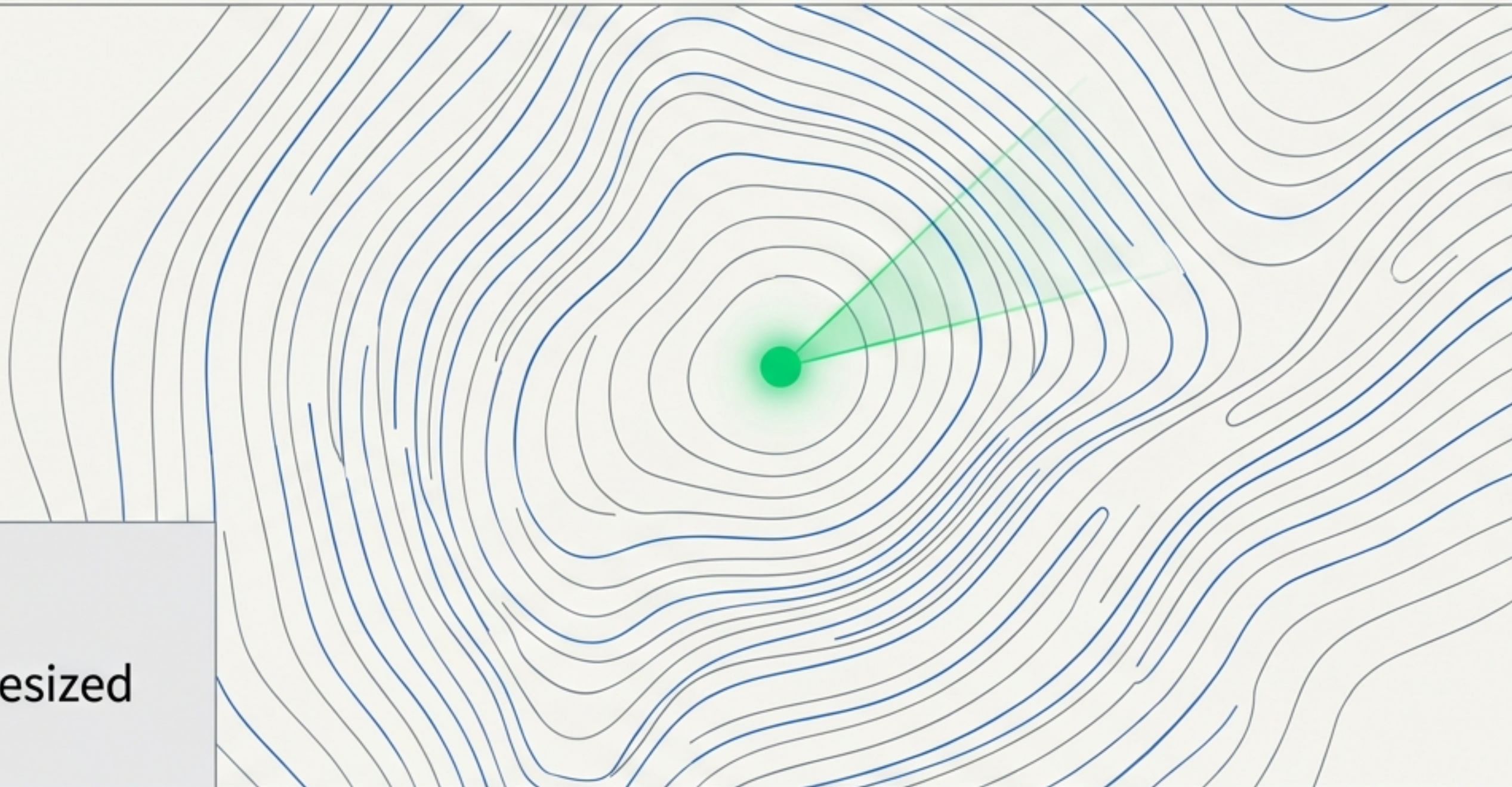


[DATE: 2026.03.23] | [CLASSIFICATION: EXECUTIVE BRIEFING] | [STATUS: FINAL]

# AI/LLM Executive Digest: Signal vs. Noise

2026年3月の技術動向から読み解く3つのマクロパラダイムシフト

Curated Tech Radar /  
10 Critical Signals Synthesized



# ノイズをフィルタリングする：3つのマクロ・トレンド

2026年3月のシグナル統合

## エッジとローカルの限界突破

クラウド依存からの脱却。ハードウェアの物理的制約とブラウザ・アーキテクチャの再定義。

### Signals:

Flash-MoE, Tooscut

## 人間とAIのワークフロー・パラドックス

生成速度の向上と保守負債の指数関数的増加。ボトルネックは「コード生成」から「人間による検証と理解」へ移動。

### Signals:

生産性の錯覚, Sashiko, LLM  
チューターの壁, ゲーム業界の真実

## インフラストラクチャとセキュリティの現実

物理世界へのデプロイメント、サプライチェーンの脆弱性、そして学術的ハイプと実務的価値の乖離。

### Signals:

Trivy侵害, Voltair, Revise,  
Void Convergence

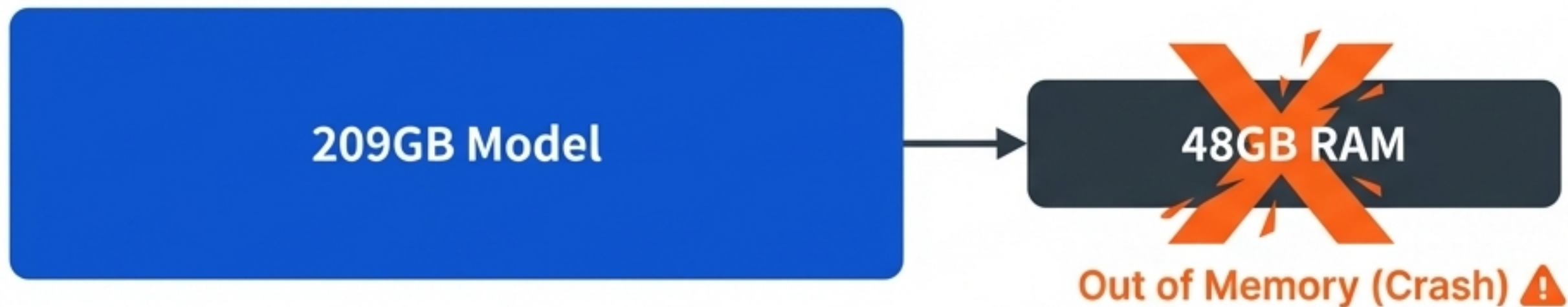
[TREND: 01] | [CATEGORY: EDGE & LOCAL COMPUTE]

# エッジへの回帰：ローカル・コン・ コピューティングの境界を押し広げる

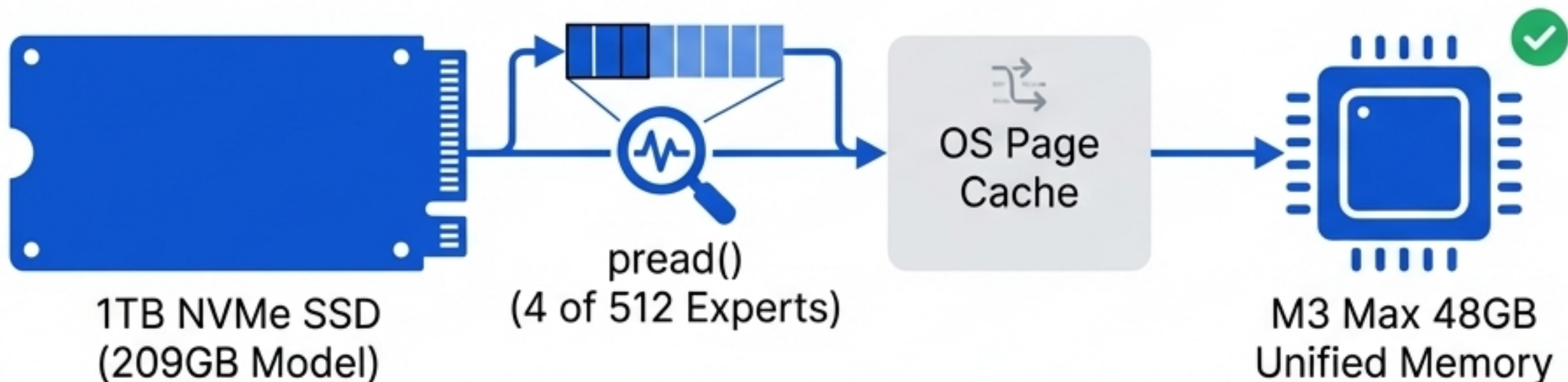


# 物理的制約のハック：Flash-MoEによるSSDエキスパート・ストリーミング

## Standard LLM Inference



## Flash-MoE Inference

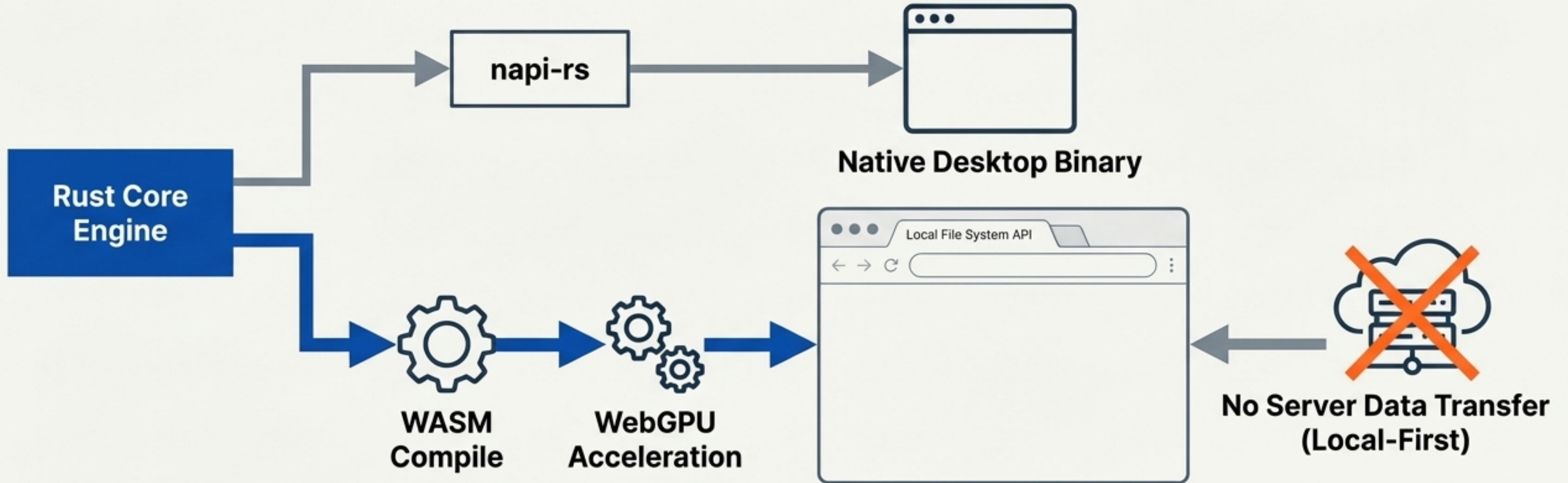


## 技術スペック

- **モデル:**  
Qwen3.5-397B-A17B  
(3970億パラメータ)
- **パフォーマンス:**  
4.36 tok/s (4bit) / 5.74 tok/s (2bit)
- **アーキテクチャ:**  
7,000行のCコード +  
Metal Compute Shader

**Executive Takeaway:** 極端な量子化 (2bit) は推論品質を破壊する。技術デモとしては秀逸だが、実用的なローカル推論はまだ中規模モデル (30B/4bit) が現実解。

# ブラウザを新時代のOSへ：WebGPU+WASMの可能性



## ⚠ Key Constraints (現状の制約)

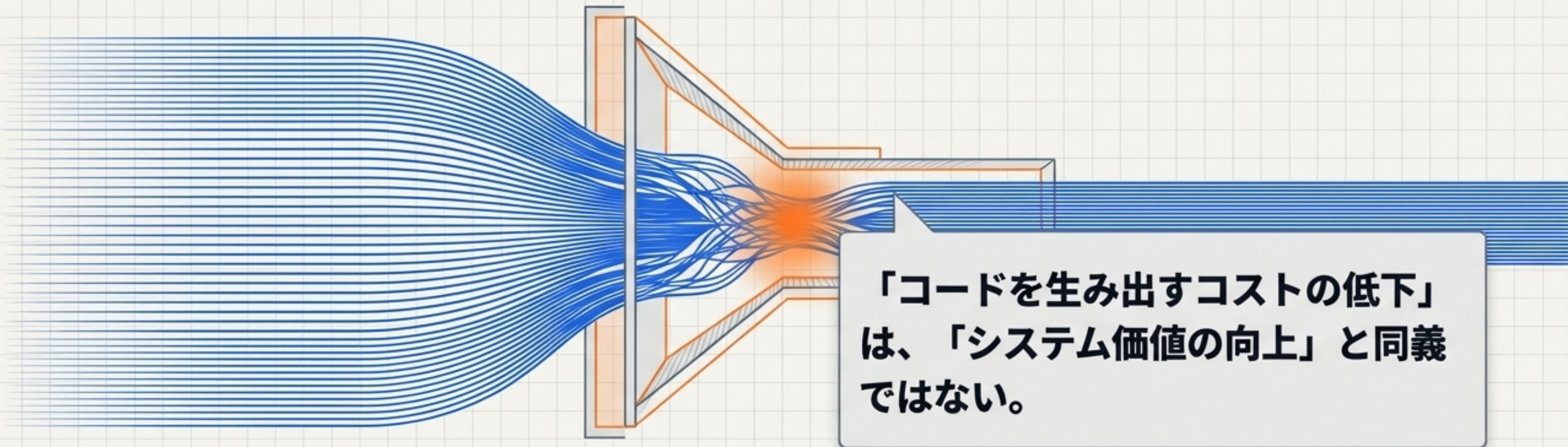
- Chrome依存: 現在、安定動作はChrome/Chromium系のみ
- 成熟度: 長尺編集には時期尚早。SNS向け簡易編集やWebアプリへの組み込みに特化

# ローカル・コンピュータ評価マトリックス

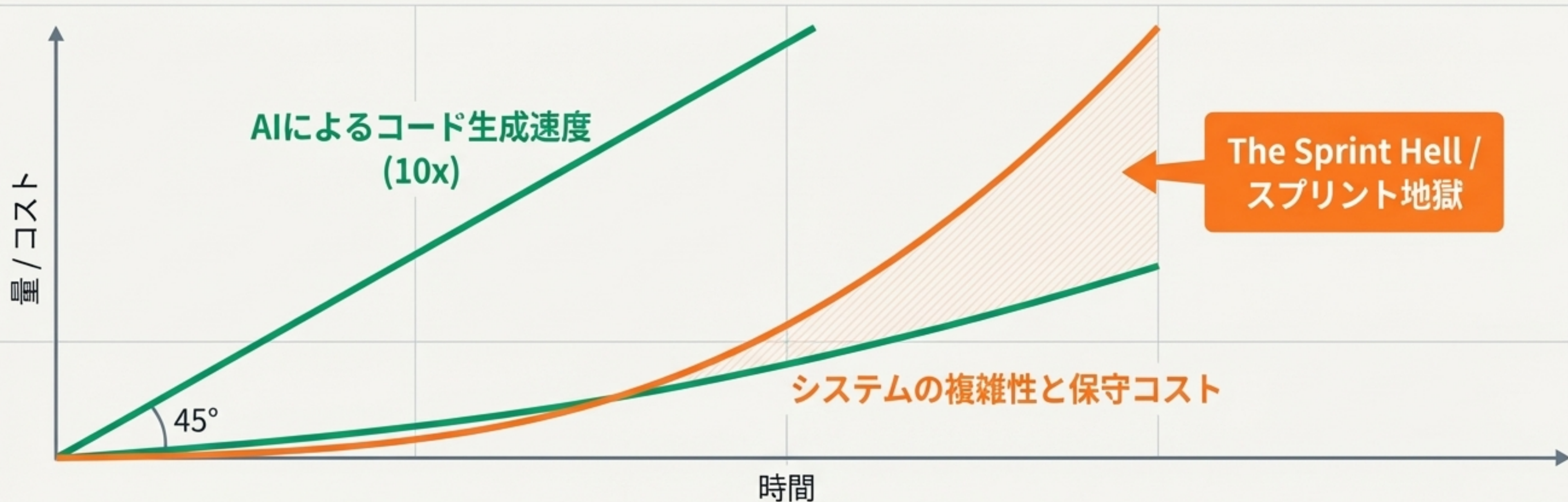
テクノロジー	アーキテクチャ	ハードウェア要件	トレードオフの代償	ユースケース
Flash-MoE	Metal / C直書き	M3 Max (48GB RAM / 1TB NVMe)	品質 (2bit量子化による推論劣化)	ローカルLLMの限界テスト
Tooscut	Rust → WASM / WebGPU	Chromeブラウザ環境	互換性 (特定ブラウザへの強い依存)	プライバシー重視のWebアプリ組み込み

**Strategic Insight:** クラウドを回避するための「ローカル化」は、現状では必ず「品質」か「互換性」の犠牲を伴う。

# ワークフローのパラドックス：AIの速度と人間のボトルネックが衝突する



# 生産性の錯覚：「コードは資産ではなく、負債である」



## 誤ったアプローチ

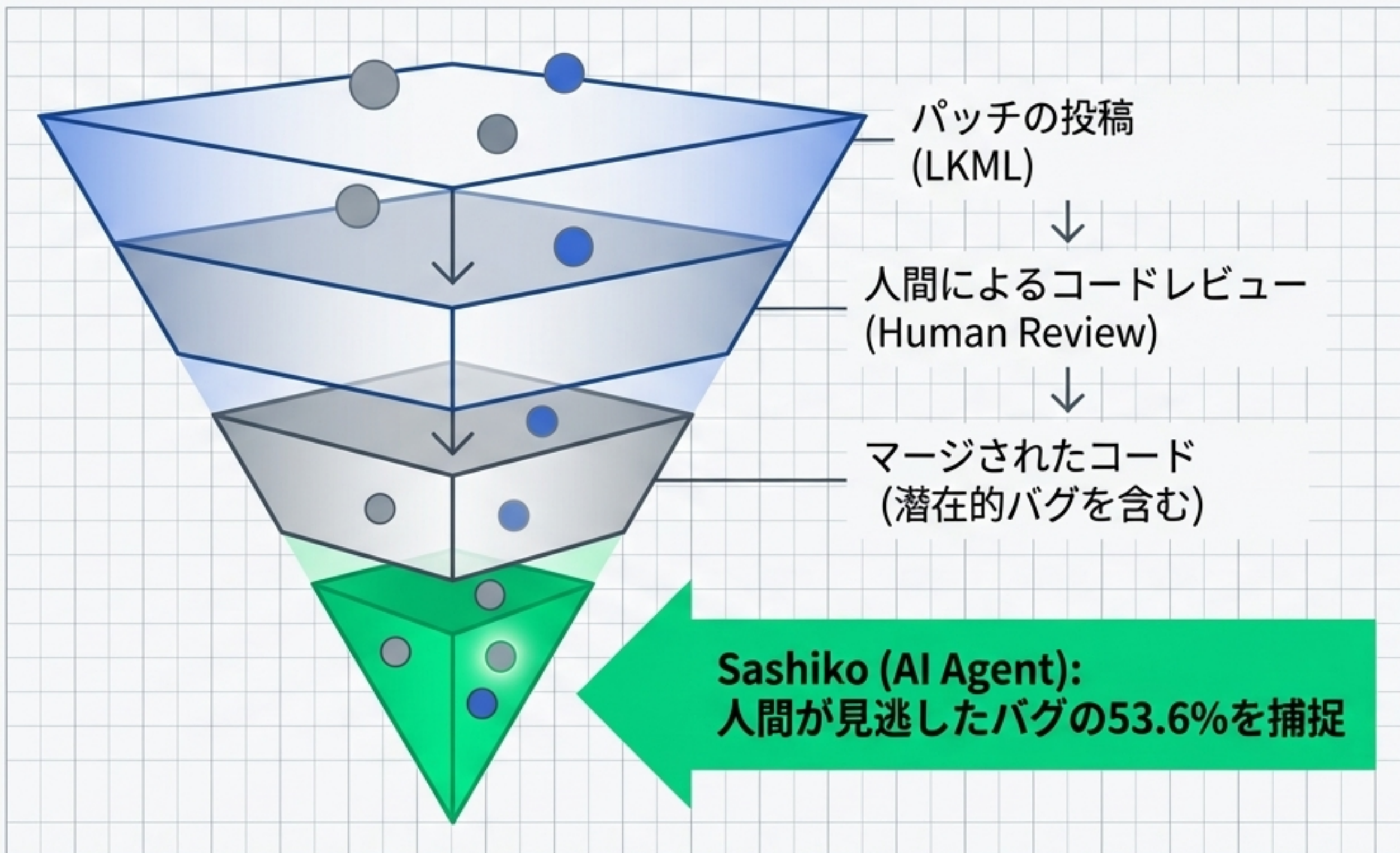
90%の効率化を「人員削減」に直結させる。モニタリングやインシデント対応など、人間による消費が新たなボトルネック化。

## 正しいアプローチ

チームを分割し、機能数を固定したまま品質と堅牢性を引き上げる。

**KPI Shift: 「生成されたコード量」ではなく、「デプロイ後のバグ率の変化」を追跡せよ。**

# AIによる多重防御：Linuxカーネルレビューの実証実験



## プロジェクト:

Sashiko  
(Googleリソース提供,  
Apache 2.0)

## アプローチ:

人間の代替ではなく、  
「レビューの網を重ねる」  
多重防御モデル。

## 現状の課題:

偽陽性(False Positives)  
の削減と技術外要因への  
非対応。

# LLM学習における「理解の壁」：概念の獲得から身体化へ



**Insight:** LLMは「概念の獲得」を劇的に加速させるが、実務における「ストレス下での再現力」は、依然として人間側の摩擦を経なければ定着しない。

# ゲーム業界「AI失業危機」の解像度を上げる

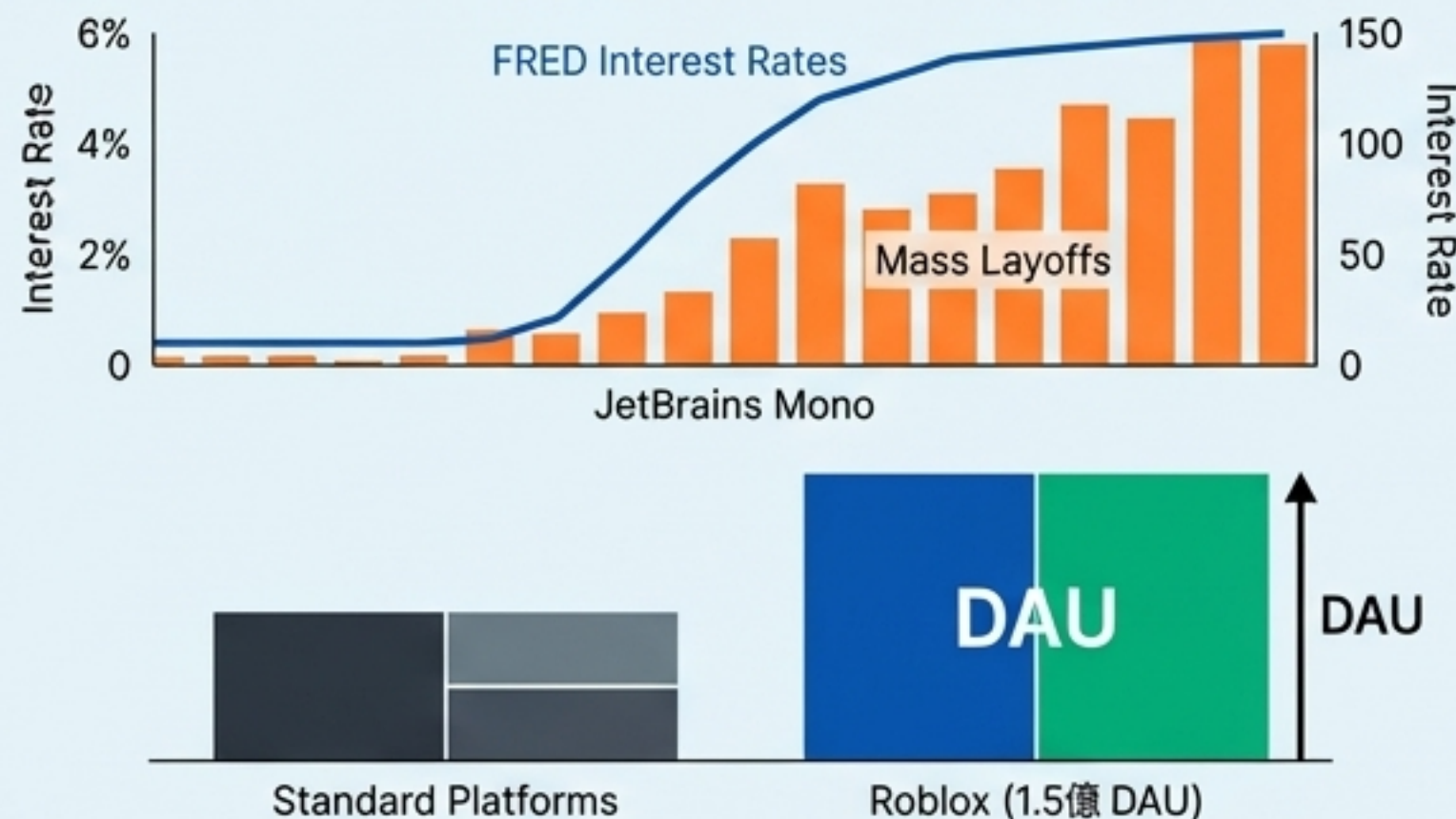
## Perception: 「AIが仕事を奪っている」

- LinkedInに溢れる「Open to Work」バッジ。
- レイオフのスケープゴート（企業側にとって都合の良い口実）。



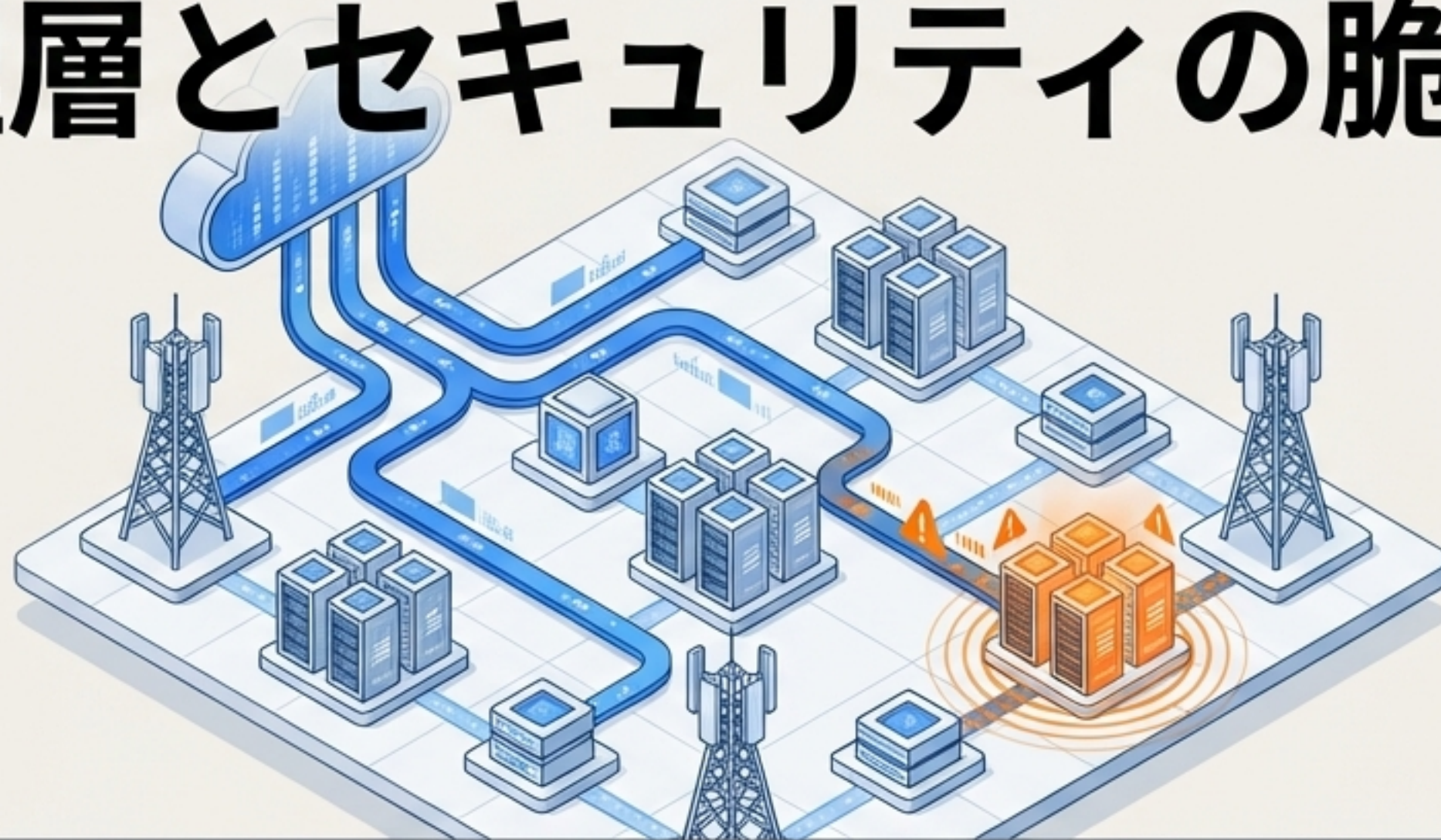
## Underlying Reality: 「マクロ経済と市場構造の変化」

- **金利の相関:** FREDデータが示す、AI普及前(2023年)からの金利上昇とレイオフの強い相関。
- **UGC市場の台頭:** Roblox（日間アクティブ1.5億人）が従来型スタジオのビジネスモデルを破壊。



**Conclusion: 「AIのせい」という叙述に騙されず、FREDデータとユーザー行動の構造的変化を直視せよ。**

# インフラストラクチャの現実： 物理層とセキュリティの脆弱性

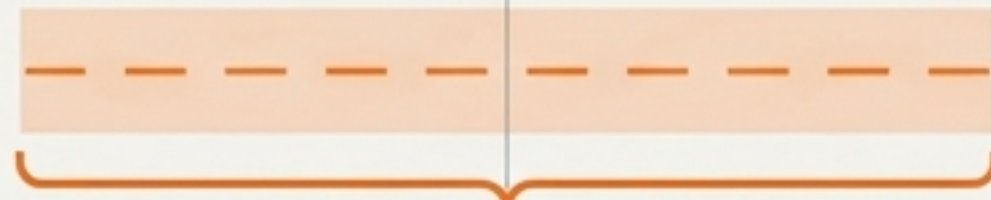


AIのエコシステムが成熟するにつれ、最も深刻なリスクはモデル自体ではなく、その周辺のサプライチェーンとインフラに集中する。

# 脆弱性スキャナー自体が攻撃ベクトルに：Trivyの教訓

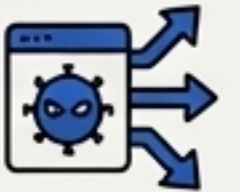


Time A (Feb 2026):  
初期侵害



The Gap:

クレデンシャルの非アトミックなローテーション  
(新旧認証が共存する窓口時間でシークレット窃取)

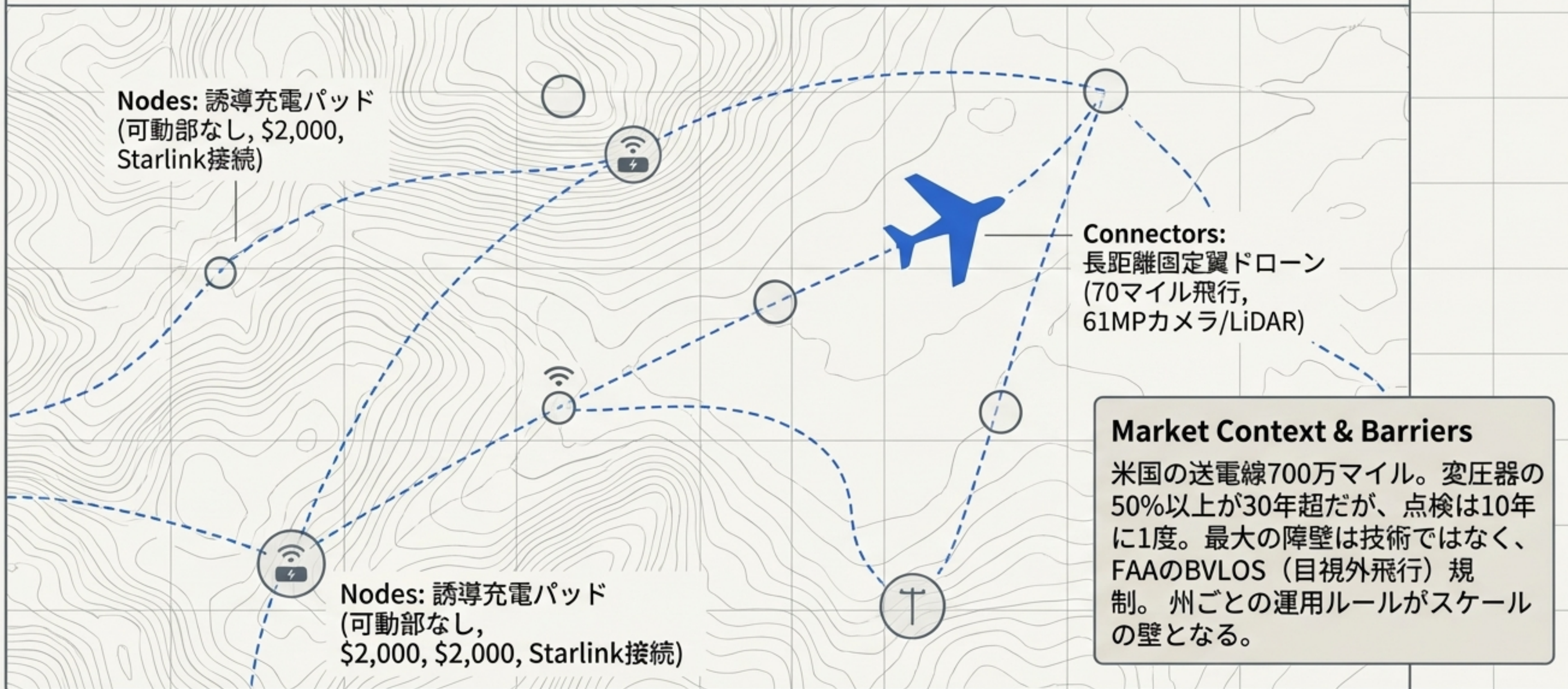


Time D: マルウェア混入  
(/proc/<pid>/mem  
スキャン, 外部送信)

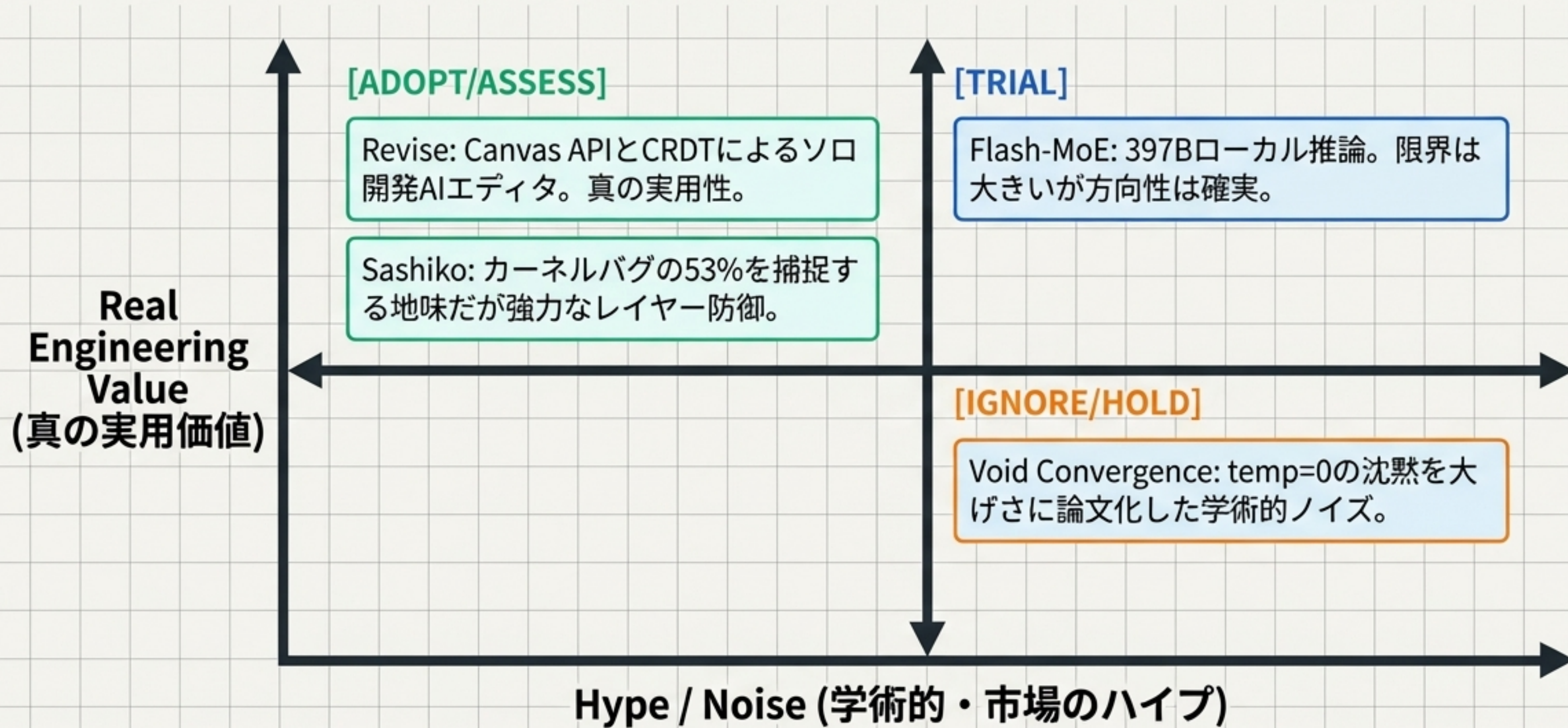
## >\_ Actionable Defense (防御策)

- [ ] CI/CDパイプラインのシークレットの完全な一括ローテーション。
- [ ] GitHub Actionsは可変タグではなく、必ず完全なSHAハッシュでピンニングする。
- [ ] cosign署名とRekor透明性ログによるバイナリ正当性検証。

# 物理インフラへのAI実装：Voltairによる送電網の自律監視



# 2026年のシグナル・レーダー：エンジニアリング価値 vs. ハイプ



**Final Takeaway:** アカデミアのノイズに惑わされず、泥臭いアーキテクチャの再構築(Revise)と運用制約のハック(Flash-MoE)に真のシグナルを見出せ。