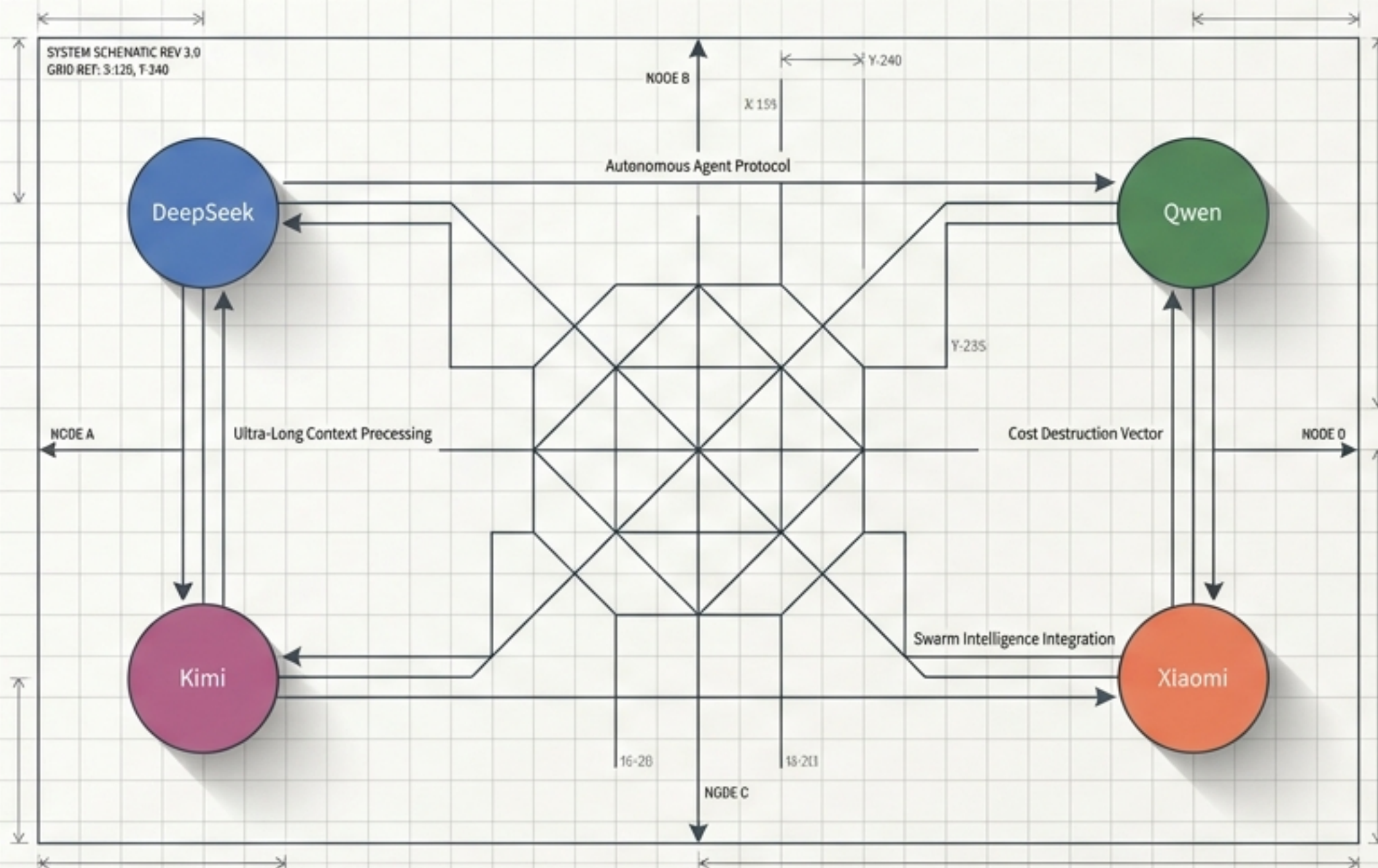


2026年 AIパラダイムの転換点

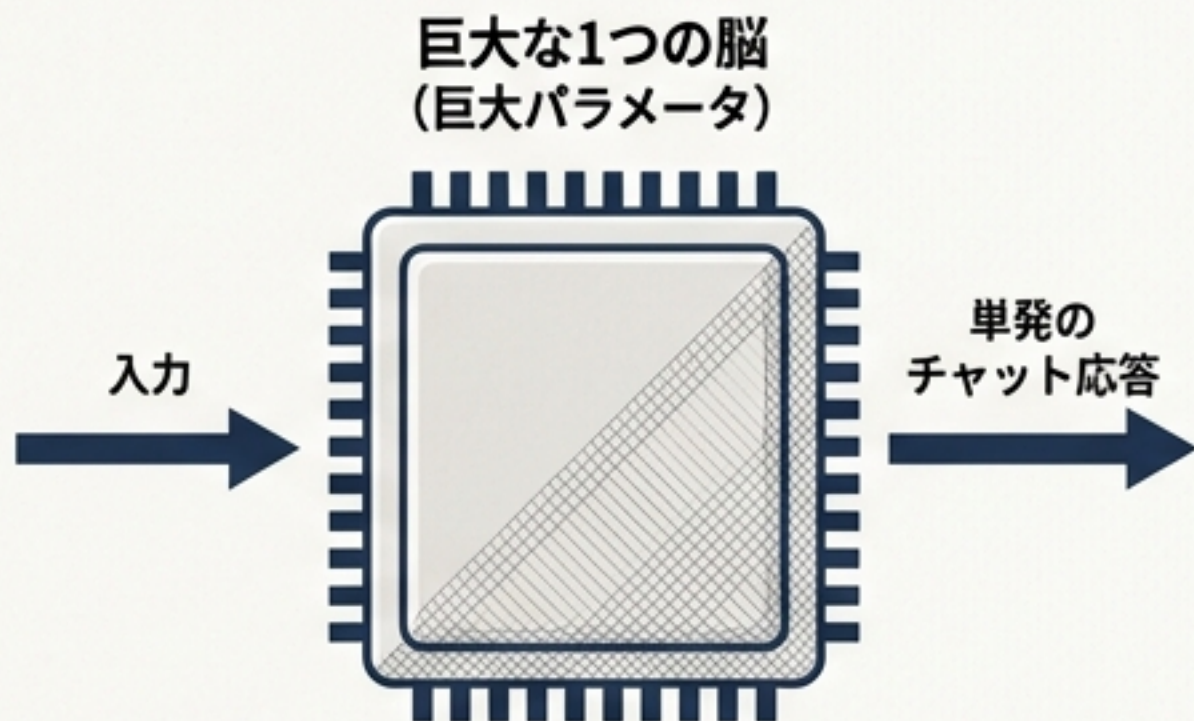
自律型エージェントと超長文脈の最前線

中国発の最先端LLMエコシステムが牽引する「限界費用の破壊」と「群知能」の実装



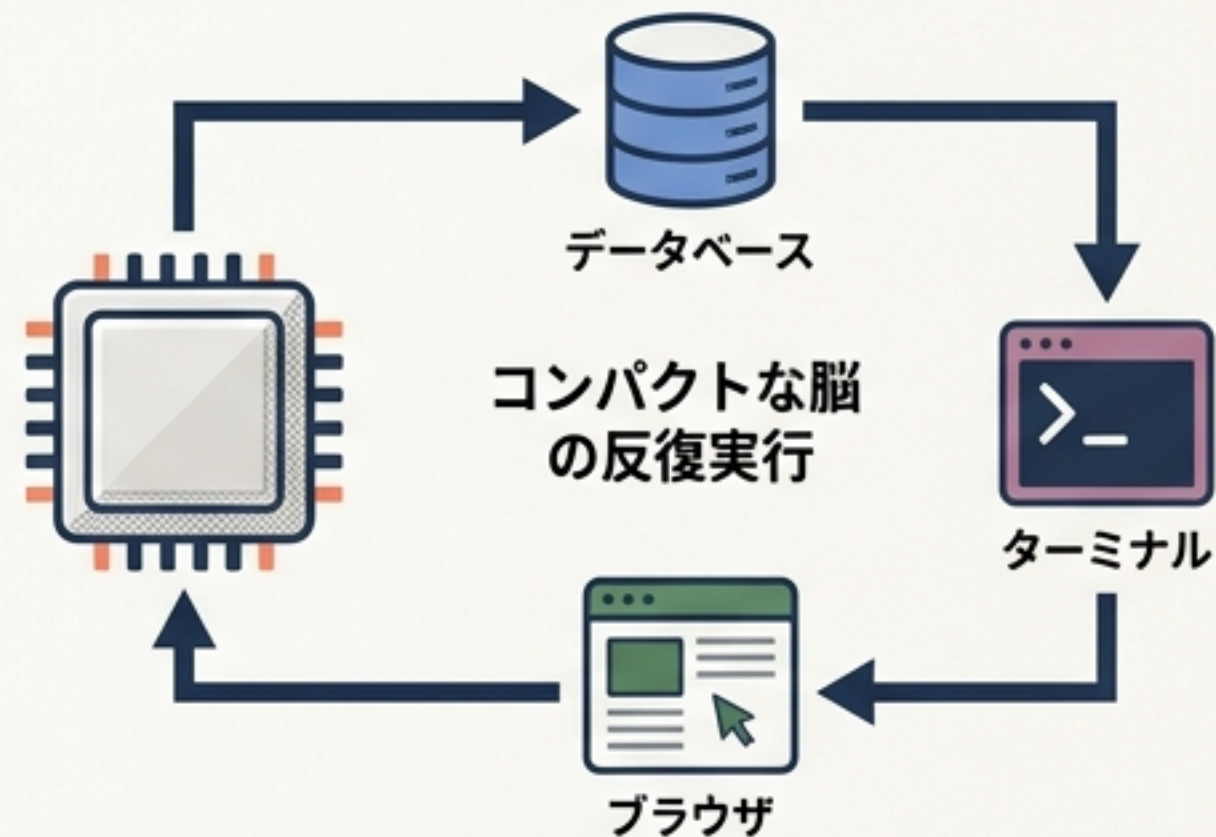
競争の焦点は「パラメータ規模」から「推論コンピュート効率と長期自律性」へ完全に移行した

2024-2025: 過去 (Past)



- スケーリング則への依存
- 単発のチャット応答
- パラメータのインフレーション

2026: 現在 (Present)



- 推論時のコンピュート効率
- 自律型エージェントの長期オーケストレーション
- 米国主導クローズドモデルとの差の逆転

米国の主要ラボ（GPT-5.5やClaude 4.7 Opus等）が一部で優位を保つ一方、中国の研究機関はアーキテクチャの根源的革新により、同等以上の性能を圧倒的な低コストで実現し、コーディングや自律推論タスクにおいて逆転現象を引き起こしている。

DeepSeek-V4：実務環境における「エージェント経済学」の再定義



コストの破壊

キャッシュミス時でもGPT-5.5の約1/7のコスト。100万トークンあたり合計\$5.22。キャッシュヒット時は\$3.625まで劇的に低下。

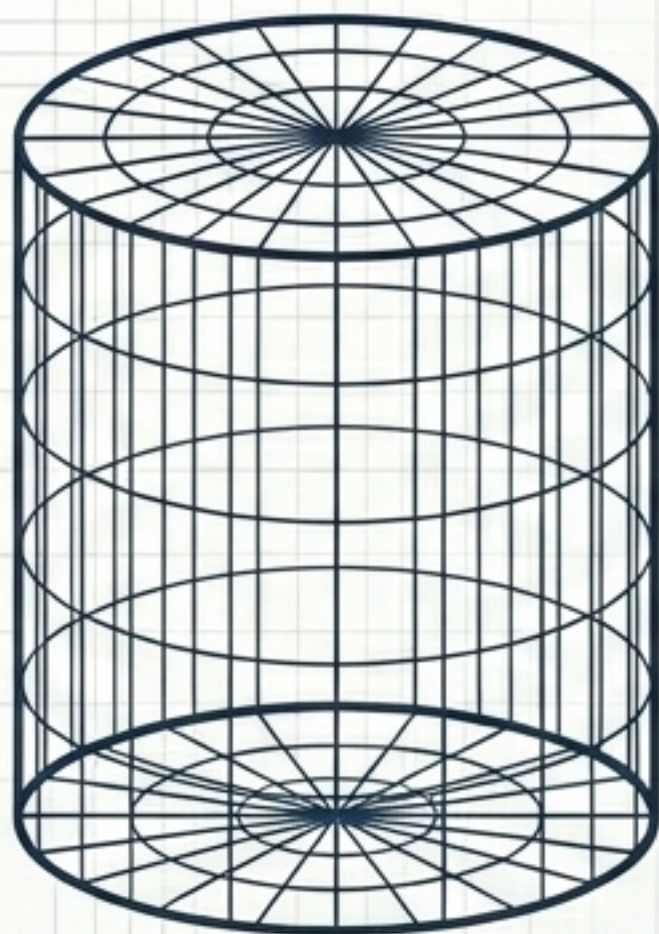
$$\text{自律ループ回数 (何千回)} \times \text{超低APIコスト (入力\$1.74/出力\$3.48)} = \text{累積限界費用の圧倒的優位性}$$

Flashモデルの逆転現象

CLIのコード編集テストにおいて、V4-Proを抑えて軽量のV4-Flashが1位を獲得。エージェント環境では「知能のわずかな優位」より「低レイテンシと累積コストの低さ」が実用上の勝者となる。

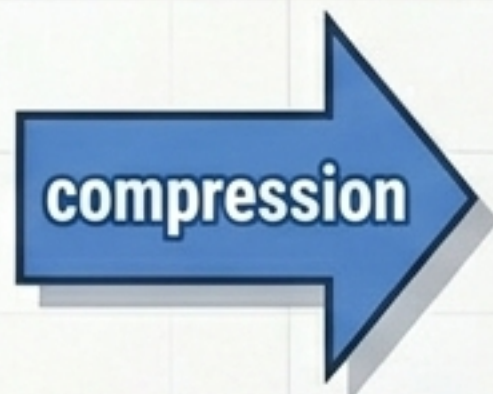
100万トークンをインフラ化する、KVキャッシュの「2%」への極限圧縮

従来の標準モデル
(100% Volume)



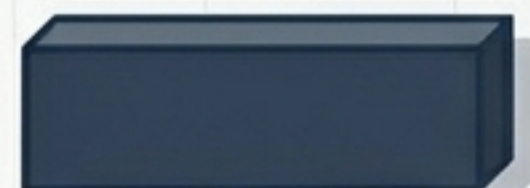
従来の標準モデル
(100% Volume)

- BF16保存
- $O(n^2)$ の計算負荷



DeepSeek-V4
(わずか2%に圧縮)

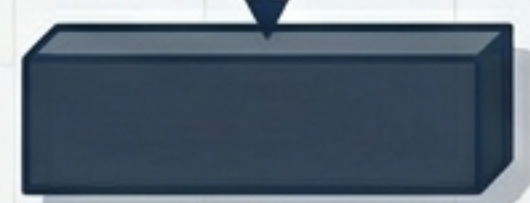
- FP8ストレージ
- FP4 Lightning Indexer



層 0-1: HCA
(Heavily Compressed Attention)



層 2-60: CSA + HCA の
交互配置による計算分散

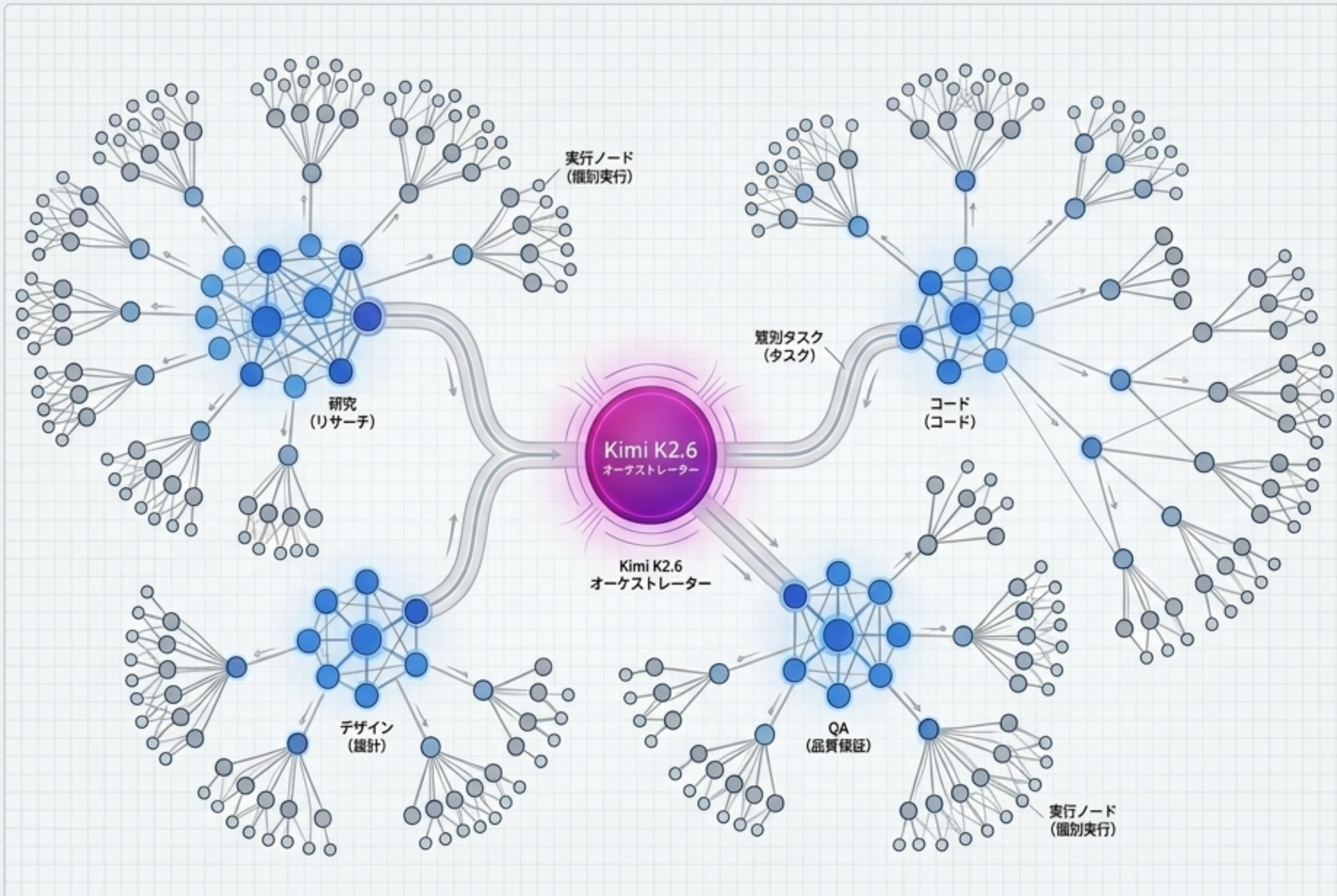


終端: MTP
(Multi-Token Prediction) ブロック

推論モード構成

- | | |
|----------------|-----------------------------|
| 1. Non-think: | 日常タスク用の直感的な高速応答 |
| 2. Think High: | 論理分析用 (<think>プロセス出力) |
| 3. Think Max: | 限界推論 (GPQA Diamond 90.1%達成) |

Moonshot Kimi K2.6：スケールアップから「スケールアウト（群知能）」への思想転換



同時並列オーケストレーション

複雑なタスクを動的に分解。最大300のサブエージェントが4,000のステップを同時実行。

驚異的な持続力

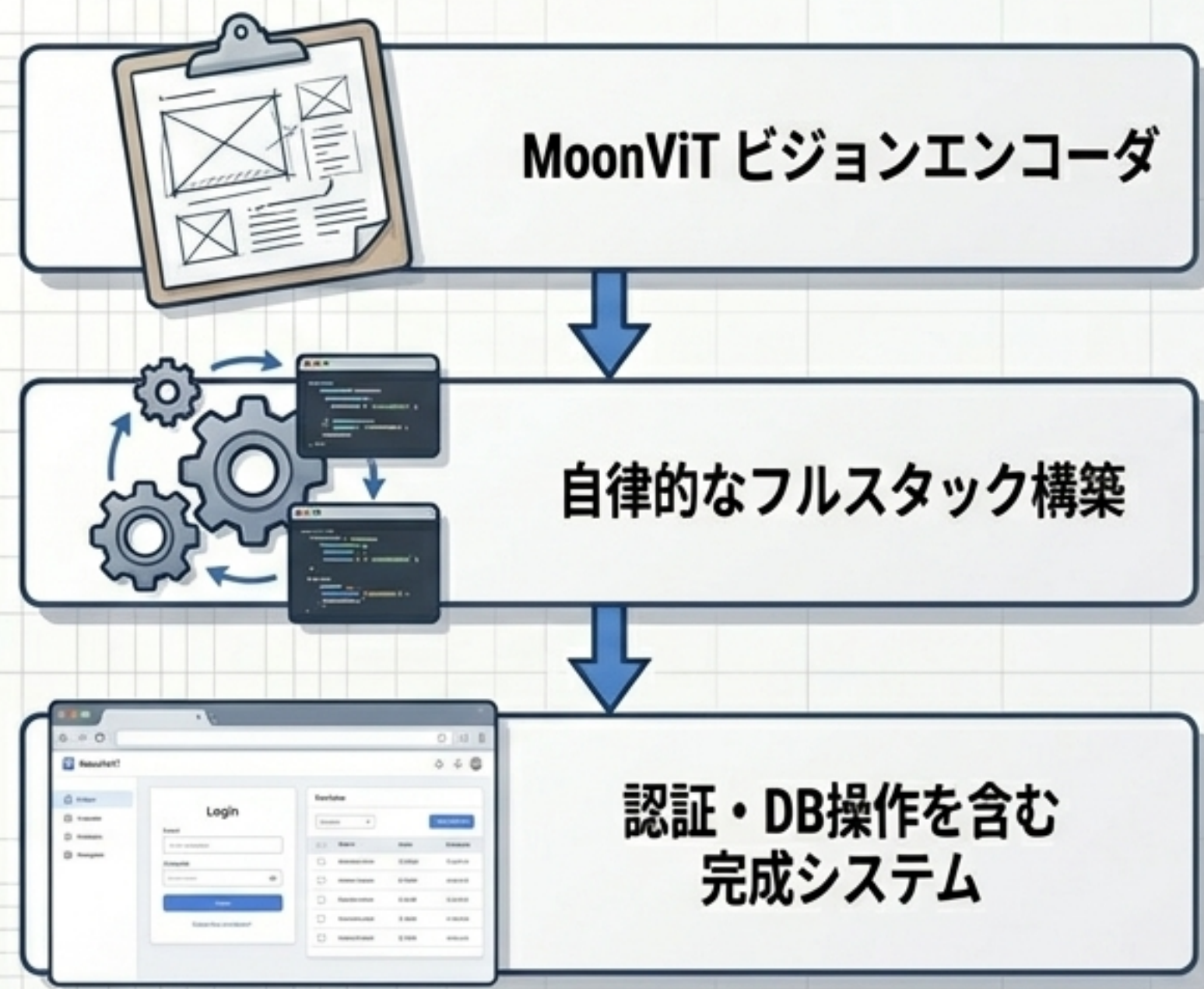
バックグラウンドで人間を介さず5日間にわたり完全自律稼働。監視からインシデント解決まで完遂。

コード改修の自動化

Zig言語の推論最適化において12時間で14イテレーション（4000ツールコール）を完遂。スループットを15から193 tokens/secへ劇的向上。

視覚駆動開発の実現と、エコシステムを支配する「戦略的ライセンス」

コーディング主導設計 (Coding-Driven Design)



MMMU-Proで80.1を記録。単なるテキストからアプリをゼロベース構築。

高度な知財戦略 (Modified MIT License)

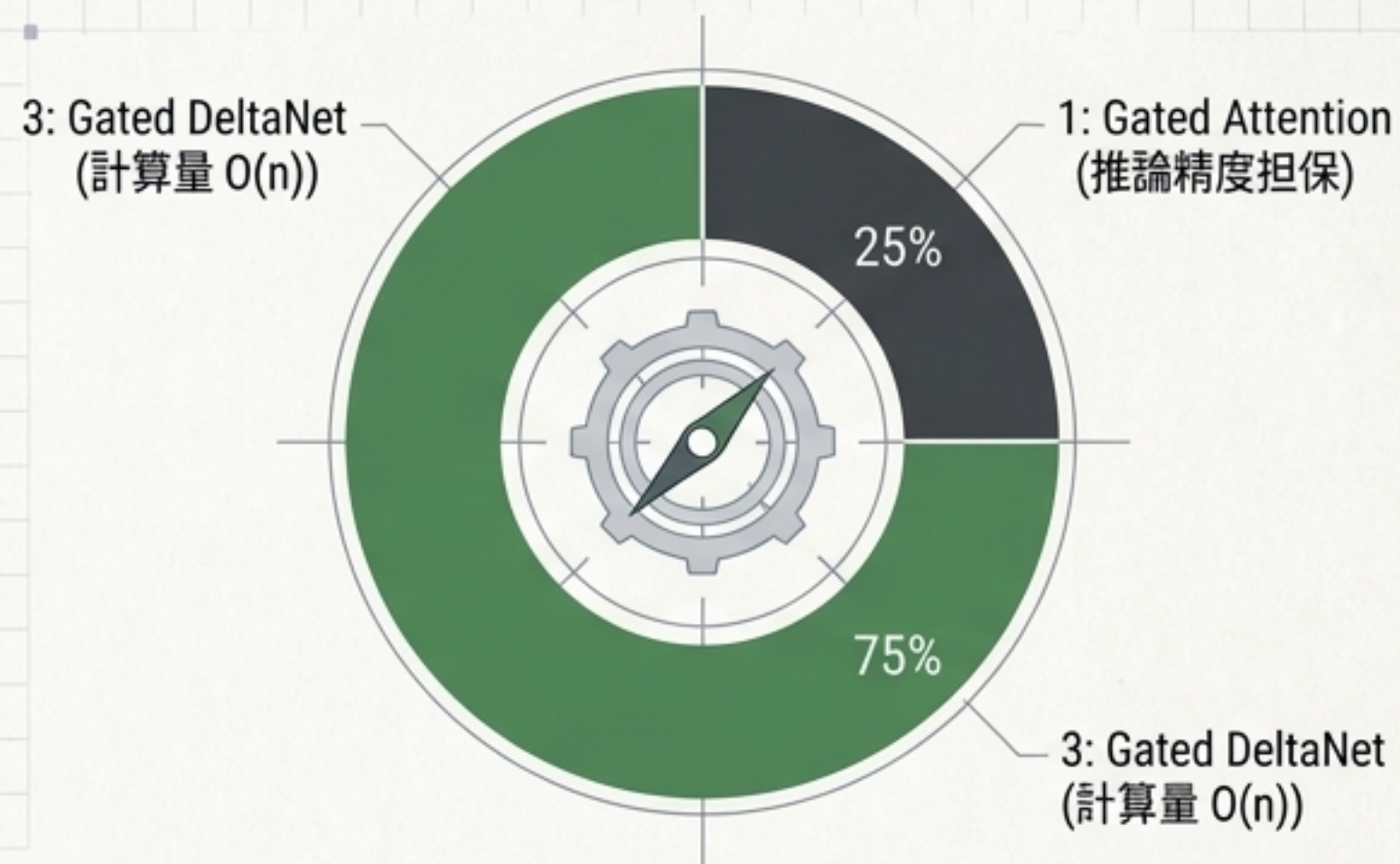
基本無料・オープンソースだが、以下の条件に該当する商用製品への組み込み時は、UI上に「Kimi K2」の表示を義務付ける。

条件 1: 月間アクティブユーザー数 1億人超

条件 2: 月間収益 2000万ドル超

戦略的意義: 巨大テックのフリーライドを防ぎつつ、自社ブランドをAIエコシステムの最上位層へ強制的に浸透させる。

Alibaba Qwen3.6-27B：長文脈の計算爆発を回避するハイブリッド線形アテンション



3:1 の黄金比による計算量圧縮

- **Denseモデルの勝利:** 27Bという軽量のDenseアーキテクチャでありながら、397Bクラスの巨大MoEモデルをコーディング領域で凌駕。
- **Gated DeltaNetの導入:** Transformer最大の弱点である $O(n^2)$ の計算爆発を解決。長大なシーケンスの安価な処理と、精密な論理展開を見事に両立。

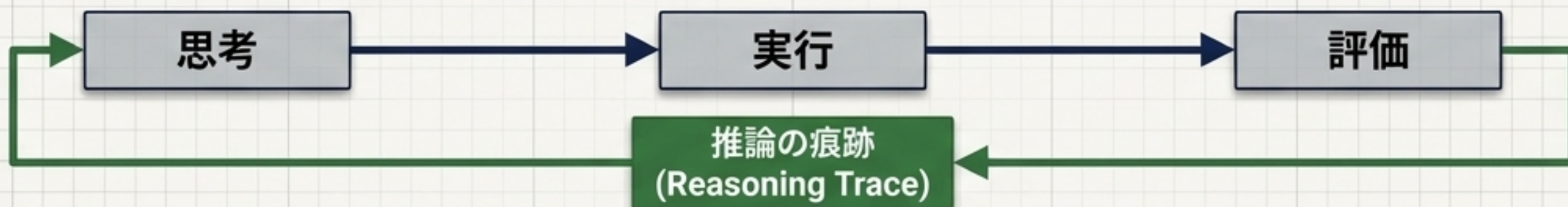
「金魚の記憶」からの脱却：Thinking Preservationによる自律ループの安定化

従来 (Goldfish-brained)



ターン終了ごとに推論コンテキストが断絶。次ターンでゼロベースから再計算（高いオーバーヘッドとループエラーの頻発）。

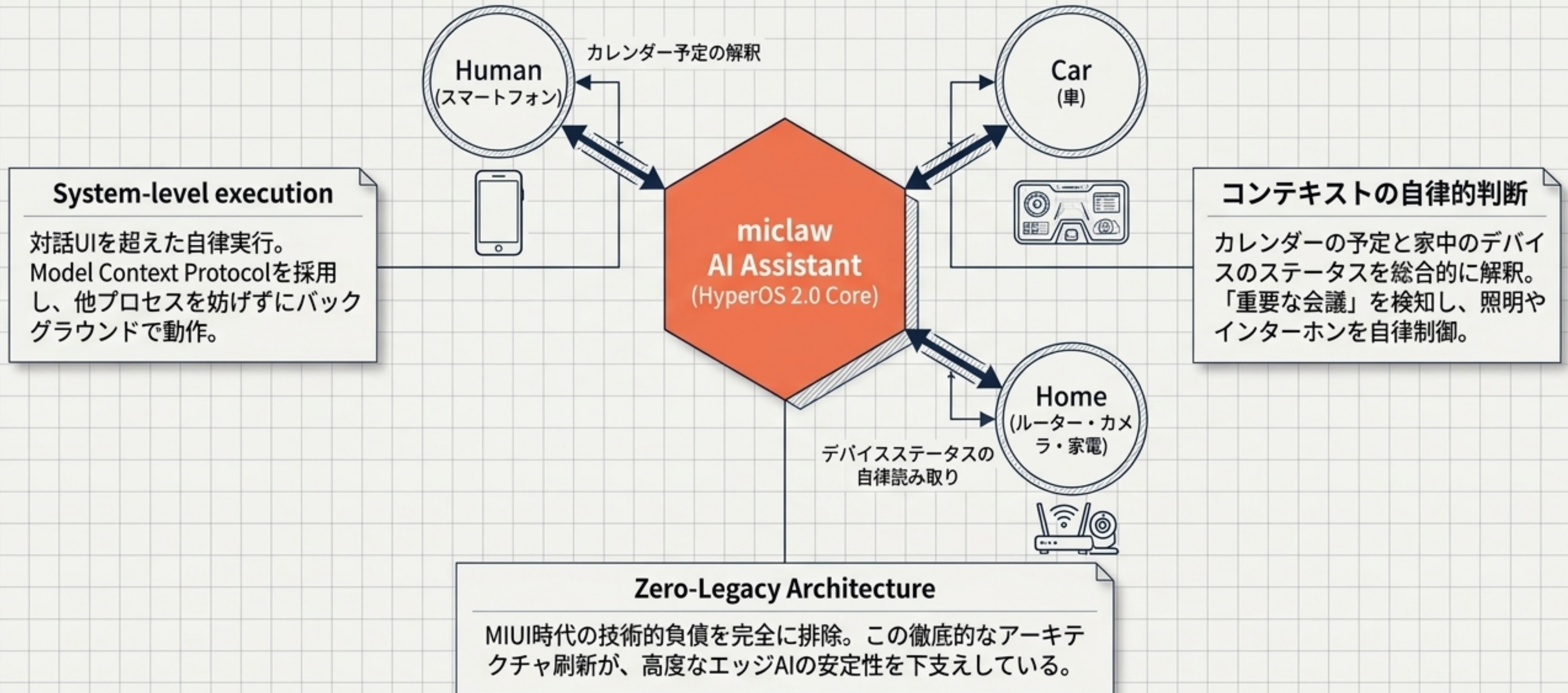
Qwen3.6 (Thinking Preservation)



ループ間で「推論の痕跡」がシームレスに継承される。無駄な再計算を省き、同じエラーを繰り返す失敗を排除。

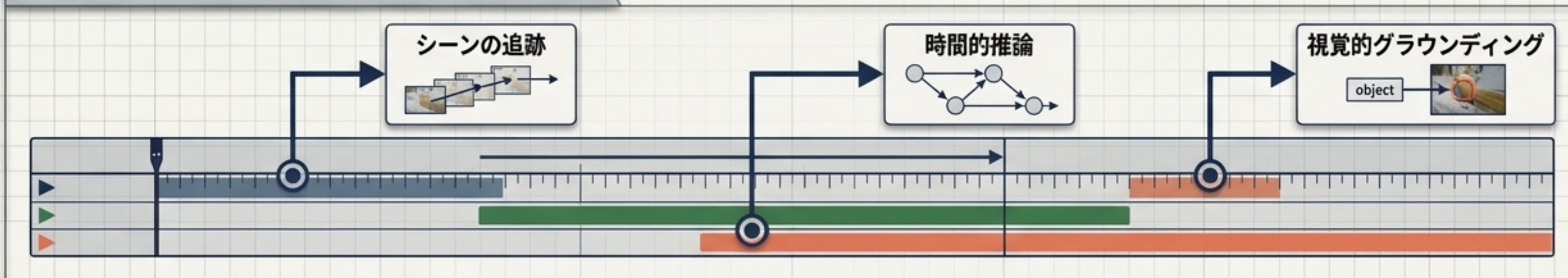
実用性: OpenClawやClaude Code等のターミナル環境アシスタントとAPI経由でシームレスに統合可能。

Xiaomi MiMo-V2.5 : ハードウェア層へ溶解する知能と「miclaw」エコシステム



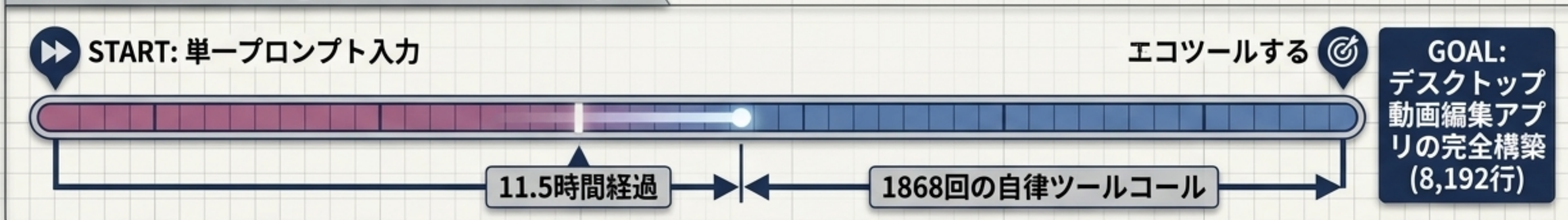
「長尺動画の時空間推論と、11.5時間に及ぶ完全自律コーディング」

最高峰の映像理解 (Video-MME: 87.7)






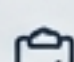
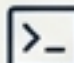
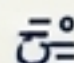
Gemini 3 Proと同等の映像理解力。単なるキャプション生成ではない真の時空間推論を実証。

驚異の持続力 (Long-horizon autonomy)



開発者シェアの奪取: クレジットカード登録不要の100%無料APIティアを提供。VS Code拡張を通じた社会実装を急加速。

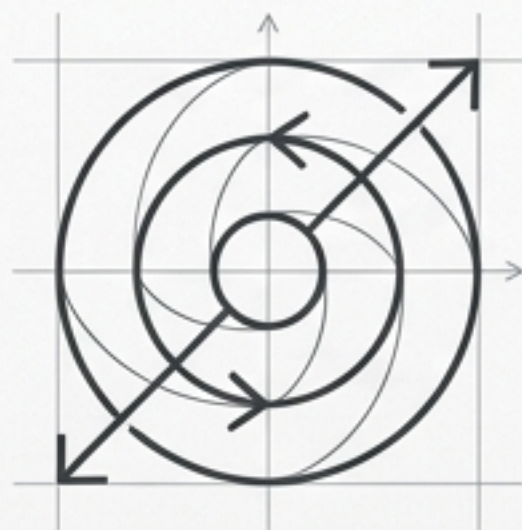
統合分析：2026年 最先端モデル「戦略的機能・性能マトリクス」

特徴・ベンチマーク	DeepSeek-V4-Pro	Kimi K2.6	Qwen3.6-27B	MiMo-V2.5-Pro	(参考) Claude 4.6 Opus
 総パラメータ	1.6兆 (アクティブ49B)	1兆 (アクティブ32B)	27B (Dense)	1兆超 (アクティブ42B)	-
 最大文脈長	100万トークン	262,144トークン	101万トークン	100万トークン	-
 コアアーキテクチャ	圧縮アテンション (CSA/HCA)	ネイティブ・マルチモーダルMoE	O(n) 線形アテンション	視覚・音声エンコーダ統合	-
 SWE-Bench Verified	80.6%	80.2%	-	-	80.8%
 Terminal-Bench 2.0	67.9%	66.7%	-	43.2%	65.4%
 特筆機能	3層の推論モード切替	300エージェント並列スウォーム	Thinking Preservation	IoT自律連動 (miclaw)	-



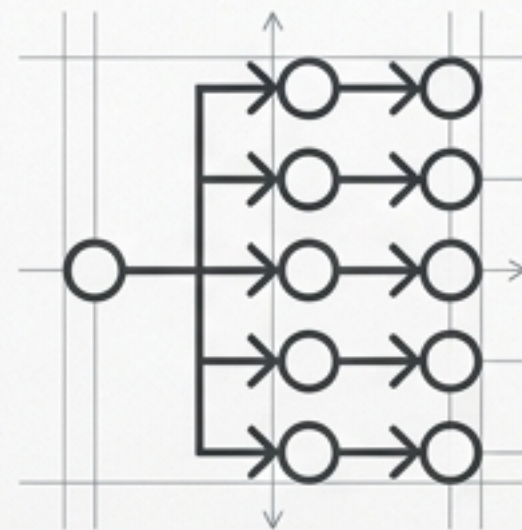
戦略的インサイト: 中国発のオープンソース/低価格API群が、実世界のバグ修正能力においてClaude 4.6 Opusと完全に同等の80%台へ到達。技術的優位性の独占は事実上崩壊している。

2026年以降のエンタープライズAI：戦略的メガトレンド (1/2)



Trend 1: 長期自律性 (Long-Horizon Autonomy) とループ経済の支配

- 限界稼働時間が「数時間」から「数日単位」へ劇的に延長。
- ソフトウェア開発の反復ループにおいては、モデル単体の知能の高さよりも、「十分な賢さ」と「圧倒的なAPIコスト/計算負荷の低さ」の掛け合わせが実務的価値を完全に決定づける。



Trend 2: スウォーム知能による「直列から並列」へのパラダイムシフト

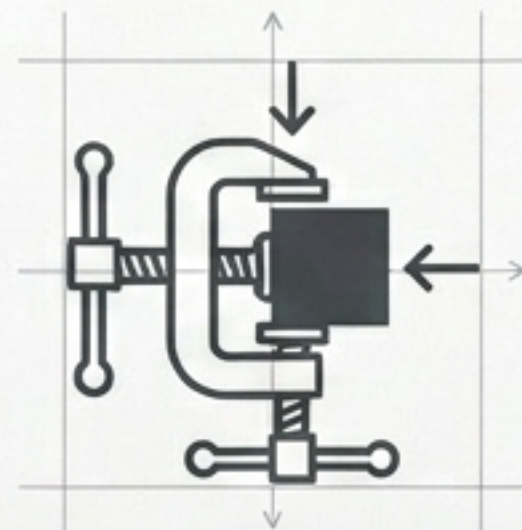
- 巨大な1つのモデルに長文を解釈させる「スケールアップ (直列処理)」の限界。
- タスクを動的分解し、数百のエージェントに同時処理させる「スケールアウト (並列処理)」へ移行。複数部門にまたがる複雑なワークフローを人間不在で一気通貫に生成する。

2026年以降のエンタープライズAI：戦略的メガトレンド (2/2)



Trend 3: ハードウェア／OSレベルへの「知能の溶解」

- LLMはブラウザのチャットUIから脱却し、OS (HyperOS 2.0等) やエッジデバイスへとネイティブ統合。
- カメラ映像や数百万のIoTデバイスからコンテキストをリアルタイムに直接読み取り、自律的なオーケストレーションを実行する段階へ突入。

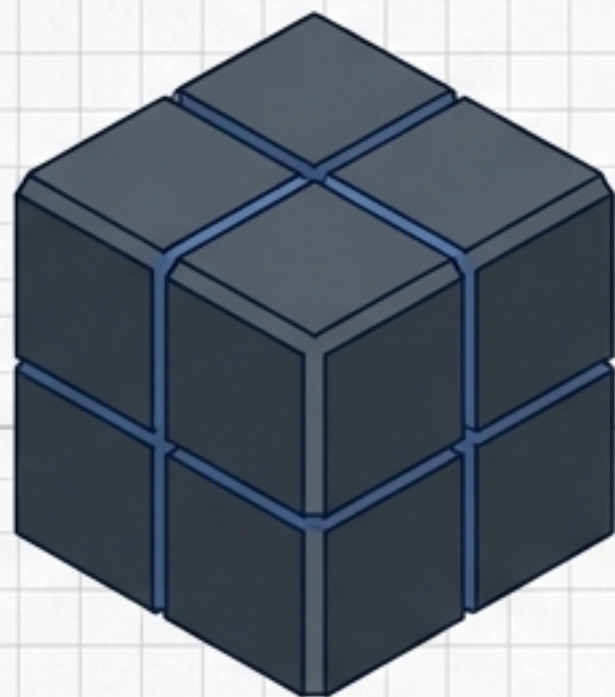


Trend 4: 計算資源の制約がもたらした「効率性のイノベーション」

- 米国の半導体輸出規制というハードウェア制約が、逆に中国の研究者たちに「アーキテクチャの極限最適化」を追求させる強いインセンティブとして機能。
- 資本の力技で肥大化した欧米モデルの価格決定権を根本から揺るがす「AI地政学の皮肉な現実」。

結論：「モデル単体の選択」から「自律コンポーネントの統合」へ

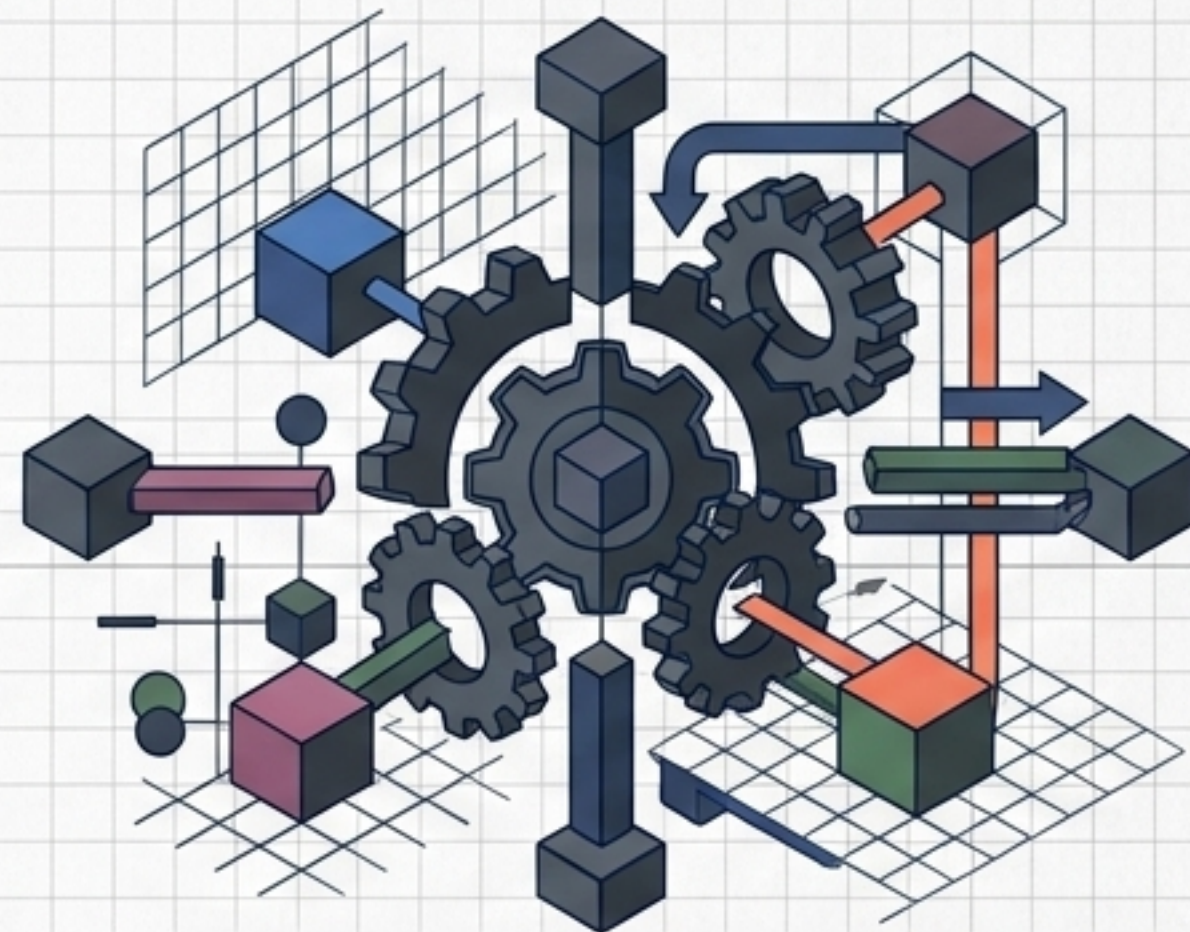
過去の競争軸



どのモデルを使うか
(単発の質問応答)



今後の競争軸



どうオーケストレーションするか
(24時間稼働のデジタル労働力)

AIの役割は不可逆的にシフトした。競争優位性は「単一モデルの性能」ではなく、「これら安価で強力な自律コンポーネントを、自社のビジネスロジック（スウォームやIoT）にいかに深く組み込み、統合できるか」というシステムエンジニアリング領域へ完全に移行している。