

Google I/O 2026: Gemini 3.5 Flash と「エージェントAI」の衝撃

Passive Assistant



Gemini 3.5 Flash: 次世代エージェントの心臓部



「速度」と「知能」のトレードオフを解消
Mixture of Experts (MoE) ブレーキテクチャを採用。
アクティブパラメータを100種~160種に削減、フラグリップ等、秒間約280トークンの超高速出力を実現。

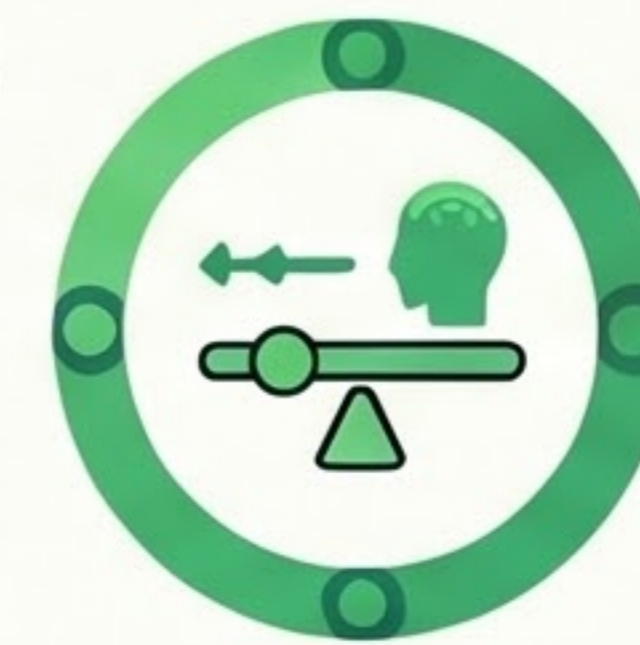
縮小済パラメータ数: 2508~300B



100万トークンの巨大なコンテキストウィンドウ
膨大なソースコードやAPIドキュメントを一度に理解。
外続ツールへのアクセスを最小限に抑えた長期的で
多様なタスク実行が可能。

最大入力: 1,048,576 トークン (1M)

最大出力: 68,935 トークン

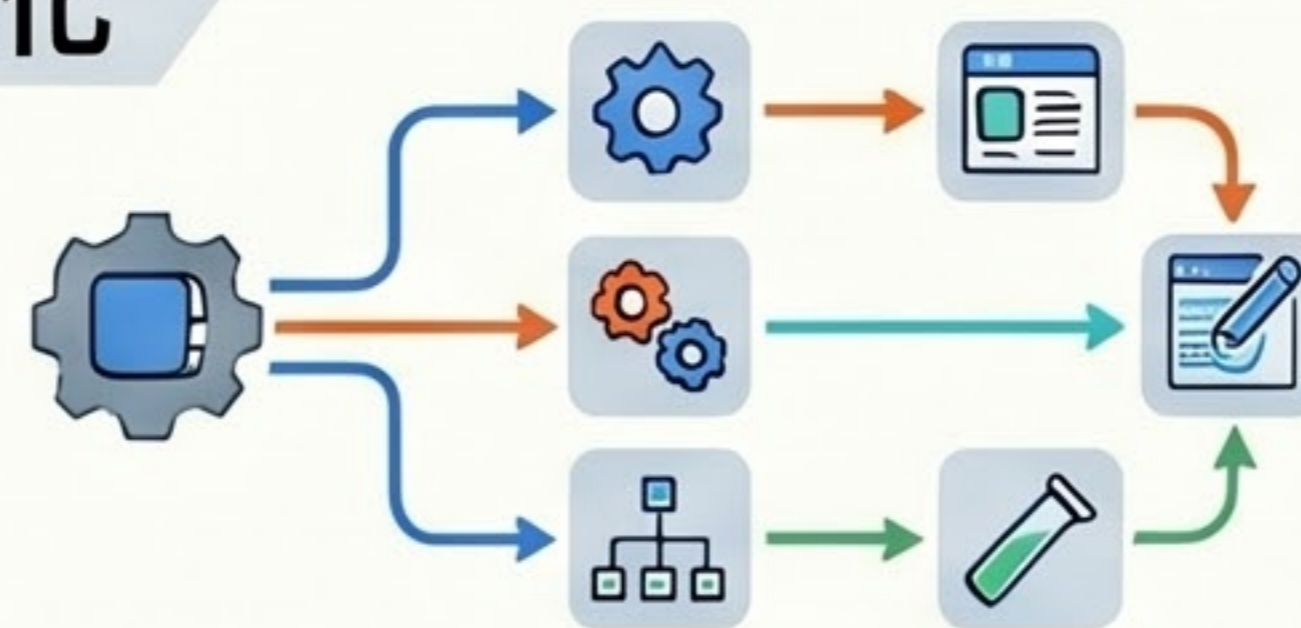


Dynamic Thinking (動的思考しべり) の導入
タスクの複雑度に応じて推進認定を4段階
(MINIMAL~HIGH) で調整可能。デフォルトは
「MEDIUM」で、速度とコストのバランスを最適化。

Antigravity 2.0: ソフトウェア工場の自動化



スタンドアロンの「中央司令部」へ進化
単なるIDE傍観から AIエージェントの作業を
監視、指揮するための役立ったデスクトップ
アプリと刷新。



動的サブエージェントによる並列処理

メインエージェントがタスクを連発。「ロジック」「UI」「テスト」などの専任サブエージェントを生成し、コンテキストを汚染せずにバックラウンドで並列成功。



12時間でOSコアを構築

93部のサブエージェントが連携、26億トークンの通信を駆使してOSコアを構築。数ヶ月かから開発工程を半分に短縮する「ソフトウェア工場」の実証に成功。

パフォーマンス: フラグシップを選驚する数値



Gemini 3.1 Proを超えるベンチマークスコア
コーディングやマルチモーダル国際において前世代の
最上位モデルを上回る、「置いて違いモデル」でなくとも
標準的な基盤解決が可能。



CharXiv
84.2%
チャートや科学の回廊の
マルチモーダル優数力



MCP Atlas
83.6%
確認ツールを押手機種
なワークフロー処理



Terminal-Bench 2.1
76.2%
CLI環境でのエージェント
実行助力

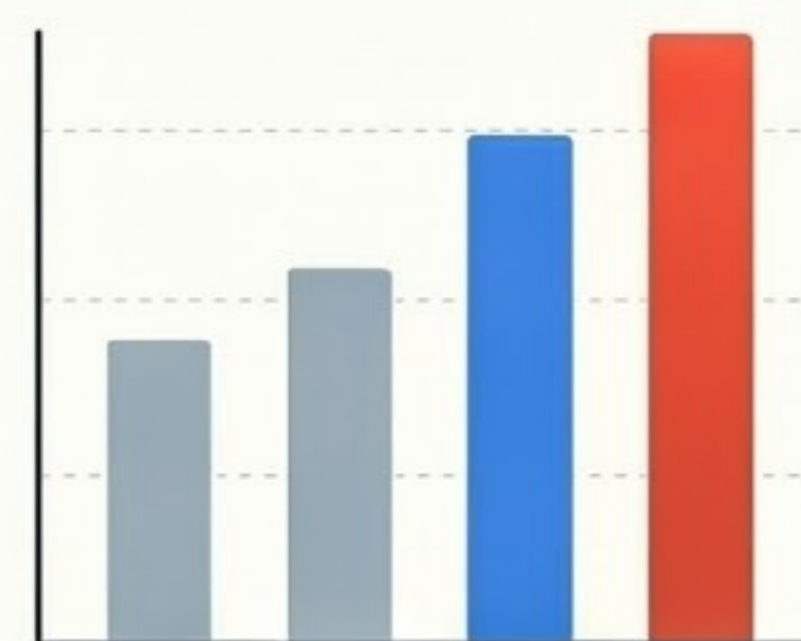


AI Analysis Index
55
同クラス平均 36 を
大幅に上回る

コストの整: 進化の代償



API価格が前モデルの**3倍**に高騰
入力\$1.50/1Mトークン、出力\$0.00/1M
トークンに設定、「Flash=安価」という
従来のブランドイメージが崩壊、上位の
Proモデルに実現する価格設定。



コストの雪だるま式増大と効率性の課題

- Gemini 3.5 Flash (知能スコア 27, コスト \$172, 比率 1.0x)
- Gemini 3.5 Flash Preview (知能スコア 46, コスト \$278, 比率 1.6x)
- Gemini 3.5 Flash (知能スコア 25, コスト \$1,551, 比率 9.0x)

他社モデルの**2倍速**いトークンを出力する間内 (Verbosity) があり、実用コストが倍々に増加するリスク。



Pelican SVGテストでの
構造的入場

構造的な装飾を追加してトークンを
大量消費した一方、ベタルと重畳を
驚く「根本的な構造」を掘り出す
など、推理の効率性に課題。