

# 2026-06-25

MORNING DISPATCH / Vibe Coder Bootcamp Tech News

## 今朝のホットな話題

1



### Anthropic 「Claude Tag」 公開

Claude が Slack に“チームの一員”として常駐、  
@メンションでタスクを非同期委任

2



### Z.ai 「GLM-5.2」が実利用で急伸

MIT ライセンスのオープンウェイトが  
Opus 4.8 にあと数点、コストは約1/6

3



### 新攻撃クラス「Agentjacking」

公開 Sentry キー1つで  
Claude Code / Cursor / Codex に任意コードを実行

5トピックを整理。



## 🔍 何が起きた？

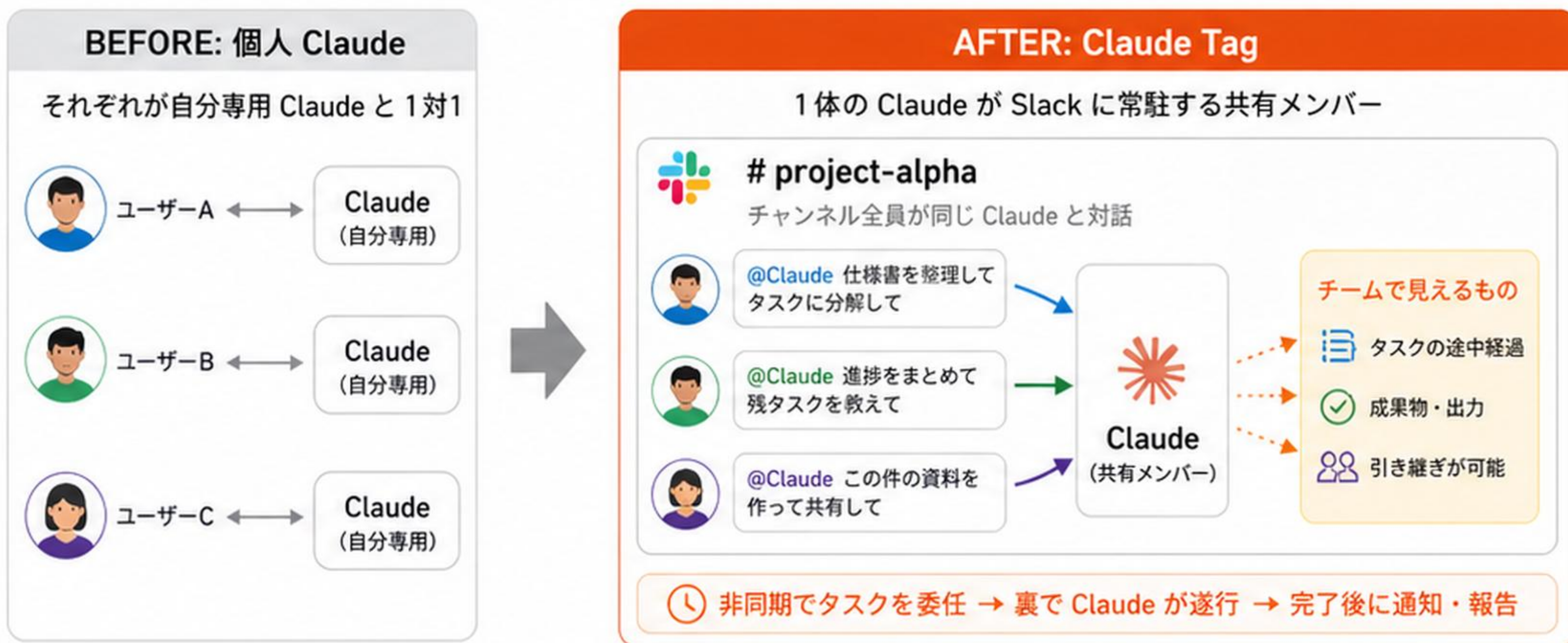
Anthropic が 6/23 に Claude Tag を発表。従来の「個人が自分専用 Claude と1対1」から「1体の Claude が Slack に共有メンバーとして常駐する」形へ。管理者が許可したチャンネル・ツールにアクセスし、ユーザーは @Claude にタスクを投げて裏で進めさせられる。エンジンは Opus 4.8。旧「Claude in Slack」アプリを置き換える。

## 📌 主な変更点

- **共有アシスタント**: チャンネル内の全員が同じ Claude と対話。他人が始めたタスクの途中経過が見え、引き継げる。
- **非同期前提**: 依頼を投げて立ち去り完了後に戻れる。大きな依頼を複数ステップに分解し連携ツールで遂行、フォローアップのスケジューリングも可。
- **Ambient Mode (任意)**: 関連情報を能動push、返信のないスレッドを追い、放置タスクをリマインド。
- **エージェント identity モデル**: 個人ログインを借りず自身のサービスアカウントで動作。チャンネル単位権限・メモリ分離・監査ログ・管理者コントロール。
- **課金分離**: チャンネル作業は組織課金、Slack 内 DM は個人の Claude アカウント課金。
- 未対応組織は 2026-08-03 に自動移行。

## 💡 なぜ重要？

AI 利用の単位が個人専用からチーム共有へ移り、Slack がマルチプレイヤー AI の作業場になる。権限・監査・identity を分ける設計が実務導入の鍵。



## 権限・セキュリティ / 運用の仕組み

- チャンネル権限**: チャンネル単位でアクセス制御
- メモリ分離**: チャンネルごとに記憶を分離
- 監査ログ**: 操作・出力を記録・追跡
- 管理者コントロール**: 設定・ポリシー・利用状況を管理

- 課金分離**: 組織課金 (チャンネル作業) / 個人課金 (Slack 内 DM)
- 65% のコード**: 内部版で 65% のコードを書いている
- 2026-08-03 自動移行**: 未対応組織はこの日に自動で移行
- “Slack がマルチプレイヤー AI になる”
- “内部版で65%のコードを書いている点が衝撃”

# Z.ai 「GLM-5.2」が実利用で急伸 — MIT ライセンスのオープンウェイトが Opus 4.8 にあと数点、コストは約1/6

## 🔍 何が起きた？

中国 Z.ai (旧 Zhipu) の MIT ライセンス・オープンウェイト coding モデル GLM-5.2 が、公開後の実利用で OpenRouter 上のシェアを急速に伸ばしている。Terminal-Bench 2.1 で 81.0、長期タスク FrontierSWE で 74.4 と Claude Opus 4.8 (75.4) に肉薄。コストは閉源フロンティアの約1/6。

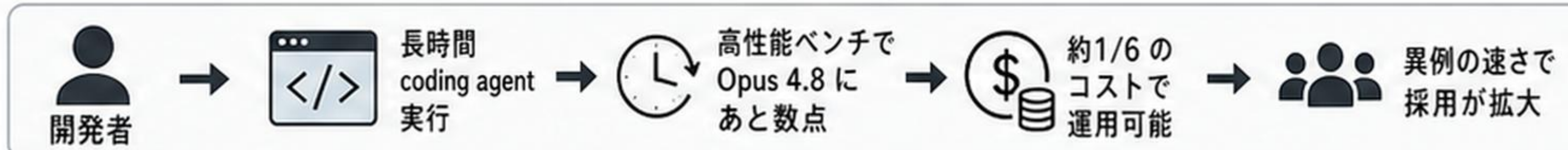
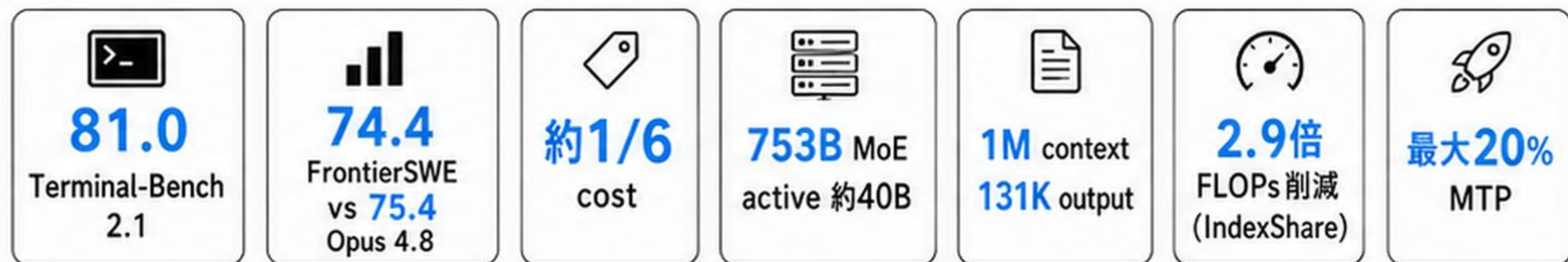
## 📌 主な変更点

- 753B MoE (アクティブ約40B)、コンテキスト 1M、出力 131K。MIT で HuggingFace に重み公開
- IndexShare: 4つの sparse-attention 層ごとに軽量 indexer を共有し、1M で per-token FLOPs を 2.9倍削減
- MTP 改良で投機的デコード受理長を最大20%向上
- SWE-bench Pro 62.1 (GPT-5.5 の 58.6 超)、Artificial Analysis Index v4.1 で 51、オープンウェイト首位
- GLM Coding Plan 月~\$18、API は \$1.40 in / \$4.40 out (per 1M)

## 💡 なぜ重要？

長時間走らせる coding エージェントに premium API を払う経済合理性を揺さぶる。OpenRouter は『異例の速さの採用』と報告し、@tianhuil は『2度目の DeepSeek モーメント』と表現。GLM-5.2 vs DeepSeek V4 のトークンシェア比較が話題。

**⚠️** 公開時に公式ベンチ未掲載。数値の多くはベンダー自己申告/早期サードパーティ。中立ハーネスでの独立検証は途上。



## 🔍 何が起きた？

Tenet Security が 6/12 に公開した新攻撃クラス。Web ソースに普通に埋め込まれている公開 DSN (Sentry の書き込み専用キー) だけで、攻撃者が偽のエラーイベントを注入。開発者が AI エージェントに『未解決の Sentry issue を直して』と頼むと、エージェントが MCP 経由で偽イベントを“信頼できるシステム出力”として読み、仕込まれた npx コマンドを開発者の権限で実行してしまう。

## 📌 主なポイント

- 起点は認証不要の公開 DSN。フロントエンド JS に埋め込まれ Web 上にインデックスされている書き込み専用クレデンシャル。
- “それ単体では安全”な2設計 (Sentry が任意ペイロード受理 × Sentry MCP がそれを信頼出力として返す) の交差点を突く。
- 影響: 注入可能 DSN を持つ 2,388 組織を特定 (Tranco top-1M に 71件)、成功率 85%・100 超のエージェント実行を観測 (※露出であり確定侵害ではない)。
- 窃取対象: 環境変数、~/.aws/config の AWS 認証情報、npm トークン、Docker/git 認証、private repo URL。稼働中の AWS secret key 捕捉例も。
- EDR/WAF/IAM/VPN、さらに『外部データを信用するな』という system prompt 指示でも防げなかった。

## 💡 なぜ重要？

Xでの反応:『公開 DSN ひとつで Claude Code / Cursor / Codex が乗っ取られる』という驚きの反応が多数。Sentry が“原理的に防御不能”と認めた点や、対症対応に留まった点への懸念も共有。



**2,388**  
組織

Tranco  
top-1M: **71件**

成功率 **85%**

**100超**  
実行

⚠️ 露出であり確定侵害ではない

## 🔍 何が起きた？

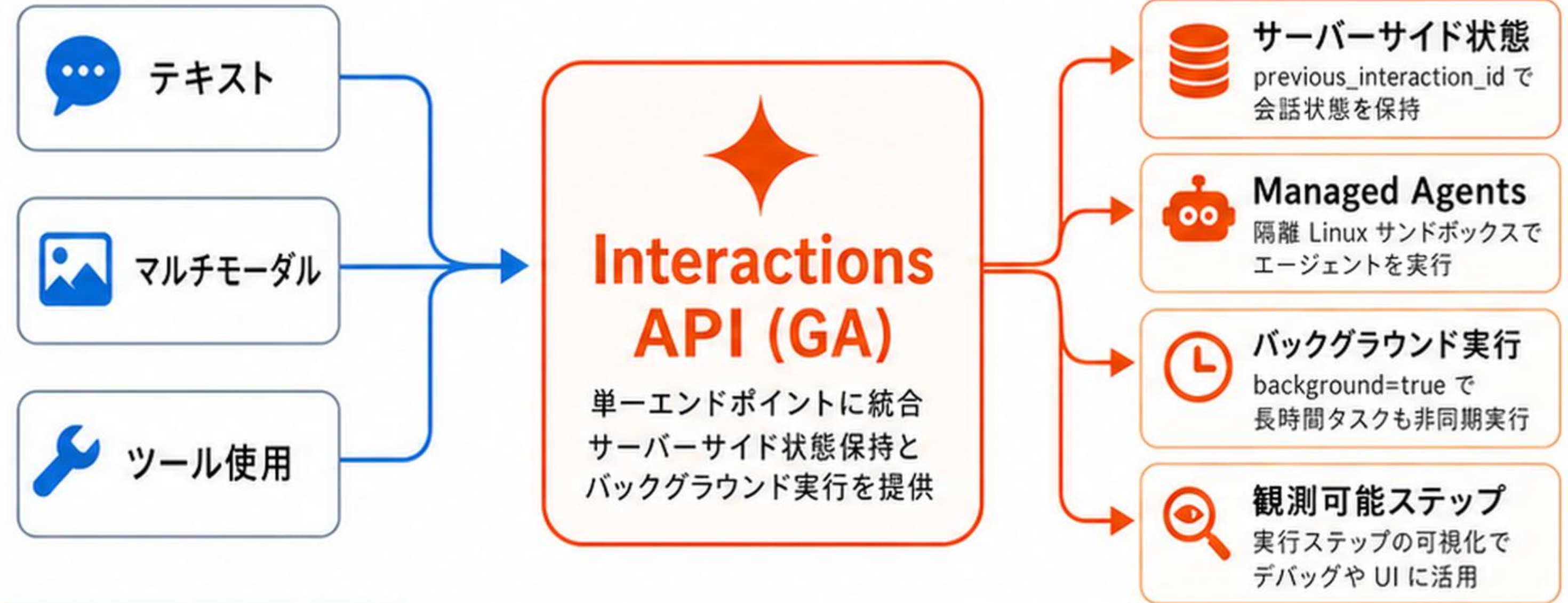
Google が Gemini の新インターフェース「Interactions API」を一般提供 (GA) した。テキスト・マルチモーダル入力・ツール使用・Managed Agents を単一エンドポイントに統合し、サーバーサイドの会話状態保持とバックグラウンド実行を備える。2025年12月のベータから昇格し、新規プロジェクトの推奨インターフェースとなった (旧 generateContent は legacy 扱いだが継続サポート)。

## 📌 主な変更点

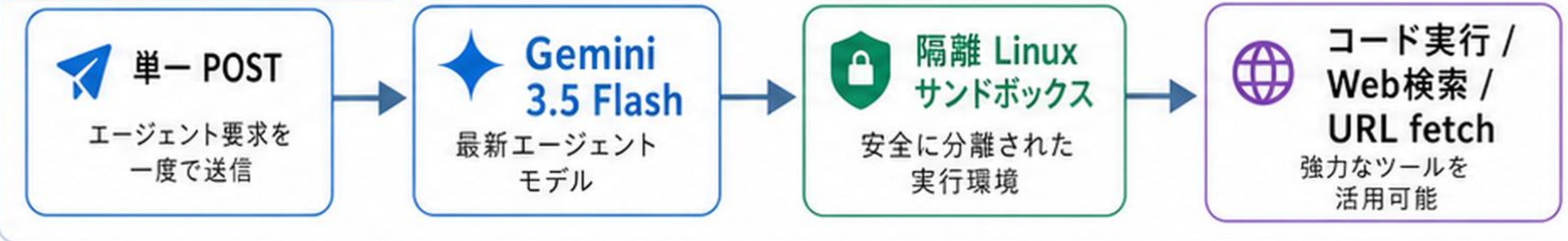
- **サーバーサイド状態:** previous\_interaction\_id で会話状態を保持。マルチターンの context caching が効きトークンコスト削減
- **観測可能な実行ステップ** (デバッグ/UI 用) と background=true による長時間タスクのバックグラウンド実行
- **Managed Agents:** 単一 POST で Gemini 3.5 Flash エージェントを隔離 Linux サンドボックスに起動 (コード実行・Web検索・URL fetch)
- **コスト階層:** Flex (コスト50%減) / Priority (低レイテンシ)。有料枠は過去 interaction を 55日保持
- Python / JavaScript SDK 対応。LiteLLM / Eigent / Agno など既存パートナー統合も利用可

## 💡 なぜ重要？

prompt から production まで最速に進めるための統合口。モデル呼び出し、ツール、状態、エージェント実行を単一の API 体験へ寄せる一方、Gemini の競争力を巡る賛否も混在。



## エージェント実行フロー (例)



|   |                                     |   |                                      |  |
|---|-------------------------------------|---|--------------------------------------|--|
| <b>2025年12月</b><br>ベータ<br>ベータから昇格し GA へ | <b>50%減</b><br>Flex<br>コストを約 50% 削減 | <b>55日保持</b><br>有料枠は過去<br>interaction を保持 | <b>Gemini 3.5 Flash</b><br>エージェントモデル | @googleaidevs:<br>prompt から production まで最速で<br>-----<br>手軽さに注目 / 賛否混在 |
|---|-------------------------------------|---|--------------------------------------|--|

### 🔍 何が起きた？

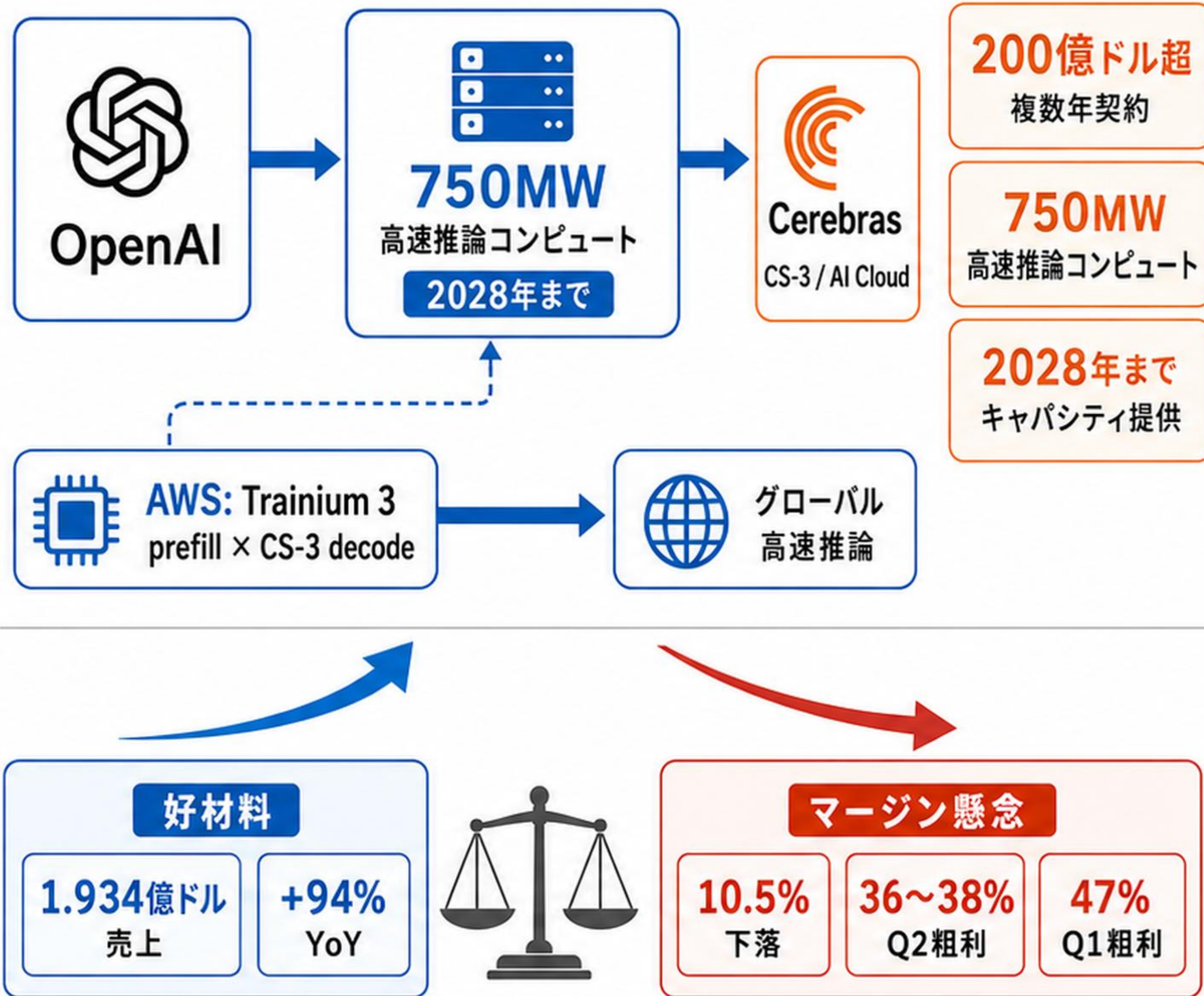
ウェハースケール AI チップの Cerebras が、上場後初の四半期決算 (Q1 2026) と同時に、OpenAI との複数年・750MW・総額200億ドル超の計算供給契約を公表。1月に「100億ドル超」とされた契約が拘束力ある正式契約で約2倍に拡大し、OpenAI は今後数年で 750MW の高速推論コンピュートを Cerebras 上に展開。

### 📌 主な変更点

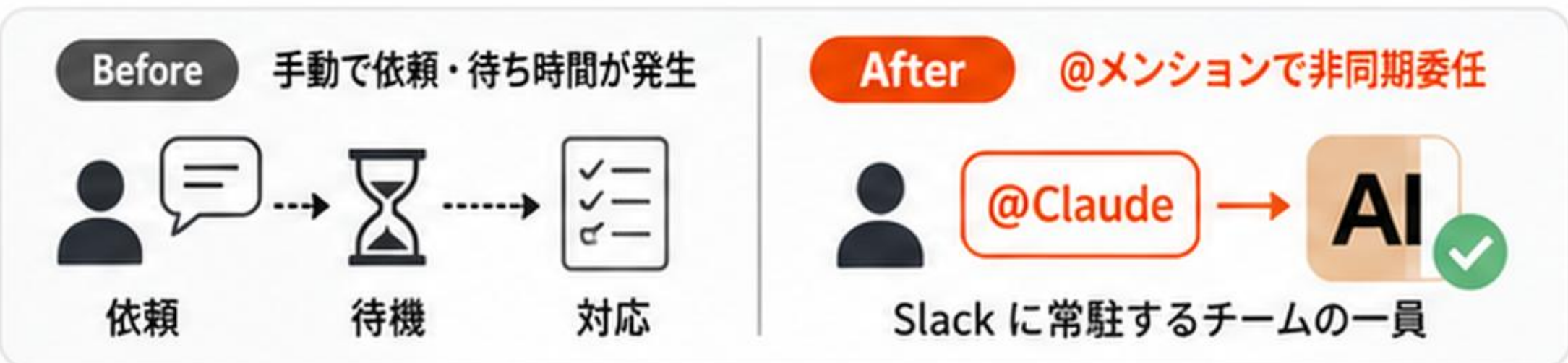
- 2028年までの 750MW 推論キャパシティをカバー。OpenAI は運転資金 10億ドルを Cerebras に貸与し、契約連動ワラントも保有。
- Q1 決算: GAAP 売上 1.934億ドル (前年比+94%)、クラウド/サービス売上 8,280万ドル (+178%)。2026通期コア売上ガイダンス 8.55~8.65億ドル (+69%)。
- AWS とも複数年提携: Trainium 3 (prefill) × Cerebras CS-3 (decode) で高速推論をグローバル展開。
- 株価は決算当日に 10.5% 下落。Q2 コア粗利ガイダンスを 36~38% へ引き下げ (Q1 は 47%)。
- 背景は『チップ売り』から『AI クラウド事業者』への積極転換。顧客集中 (OpenAI/G42/AWS 依存) がリスク。

### 💡 なぜ重要？

200億ドル契約でも株価が下げる点にマージン懸念の反応が集中。AI インフラ投資の重さを巡る議論が起きている。



**1 Anthropic「Claude Tag」公開**  
 — Claude が Slack に“チームの一員”として常駐、@メンションでタスクを非同期委任



**2 Z.ai「GLM-5.2」が実利用で急伸**  
 — MIT ライセンスのオープンウェイトが Opus 4.8 にあと数点、コストは約1/6



**3 新攻撃クラス「Agentjacking」**  
 — 公開 Sentry キー1つで Claude Code / Cursor / Codex に任意コードを実行させる



**4 Google「Gemini Interactions API」が GA**  
 — モデルとエージェントを単一エンドポイントに統合、Managed Agents・バックグラウンド実行対応



**5 Cerebras、上場後初決算と同時に OpenAI と 750MW・200億ドル超の複数年契約を公表**  
 — それでも株価は10%超下落

