

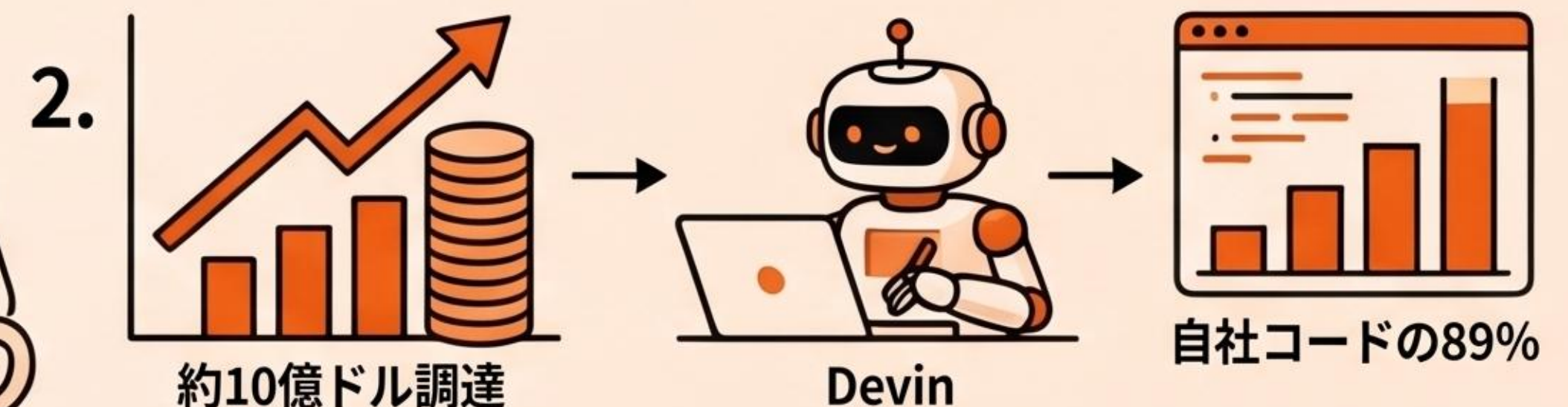
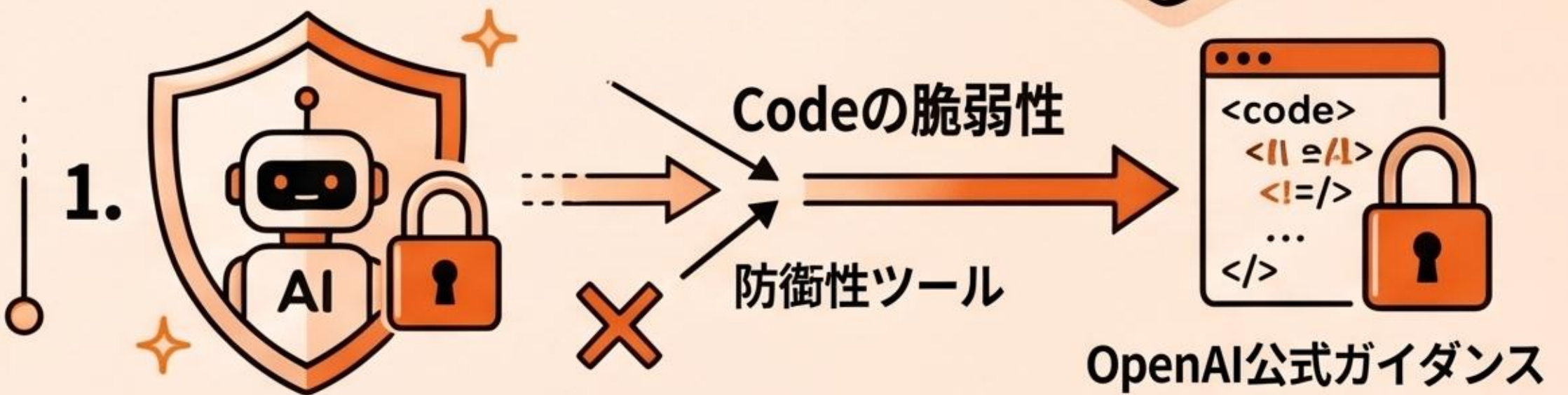


## 今朝のホットな話題

- 1. AIコーディングエージェントの「セキュリティ」が新たな主戦場に — OpenAIが公式ガイドス、脆弱性・防御ツールも続々
- 2. Cognition (Devin 開発元) が約10億ドルを調達、評価額260億ドルに — 自社コードの89%を Devin が記述
- 3. StepFun が「Step 3.7 Flash」をオープンウェイトで公開 — 198B MoE・Apache 2.0・最大400トークン/秒



7トピックを整理。



# AIコーディングエージェントの「セキュリティ」が新たな主戦場に — OpenAIが公式ガイダンス、脆弱性・防御ツールも続々

## 🔍 何が起きた？

OpenAIが『プロンプトインジェクションに耐えるエージェント設計』の公式ガイダンスを公開。完全な検知は期待せず影響範囲を限定する現実的な構えを打ち出した。

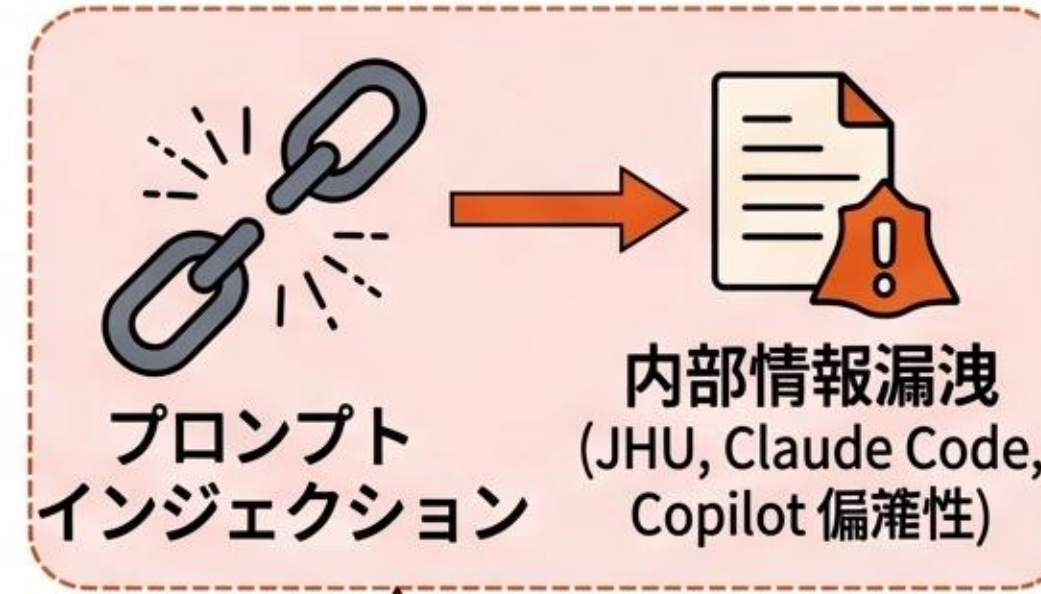
## 📌 主な変更点

- OpenAI 公式: 防御は入力フィルタだけに頼れず、最小権限・タスク単位のアクセス制限・重要操作の二段階承認・検知重視を推奨。
- OpenAIは『プロンプトインジェクションは（ブラウザ/エージェントで）完全には解決しないかもしれない』と明言。
- Johns HopkinsがClaude Code / Copilotを操作して内部情報を漏洩させる脆弱性を公表。
- 依存ライブラリにプロンプトインジェクションを埋め込み、AI生成コードの出力を無言\*で破壊するPoCが話題に。
- Keeper Securityが『Agent Kit』を提供開始、コーディングエージェントが認証情報を安全に取得する枠組み。

## 💡 なぜ重要？

コーディングエージェントの利用拡大に伴い、セキュリティが新たな「主戦場」となった。攻撃手法（プロンプトインジェクション、依存関係）と防御手法（ガイダンス、専用ツール）の両面が急激に表面化しており、現実的な防御設計が急務となっている。完全な検知の難しさが公式に認められた意義も大きい。

### 脅威と攻撃手法



セキュアなAIコーディングエージェントの設計

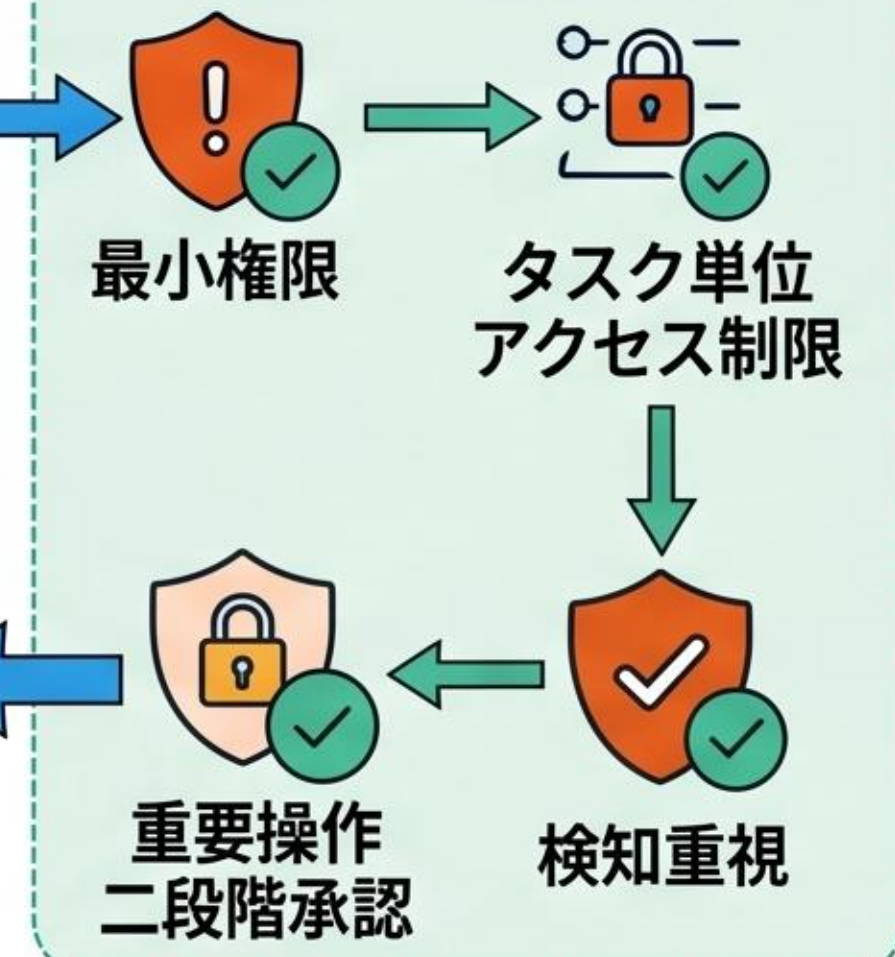


### 依存ライブラリの汚染



### 防御とガイダンス

#### OpenAI ガイダンス推奨図解



#### Keeper Security 『Agent Kit』



# Cognition (Devin 開発元) が約10億ドルを調達、評価額260億ドルに — 自社コードの89%を Devin が記述



## 何が起きた？

自律ソフトウェアエンジニアリングエージェント

『Devin』を開発する Cognition が、主要VC主導で大規模な資金調達を実施。

## 主な変更点

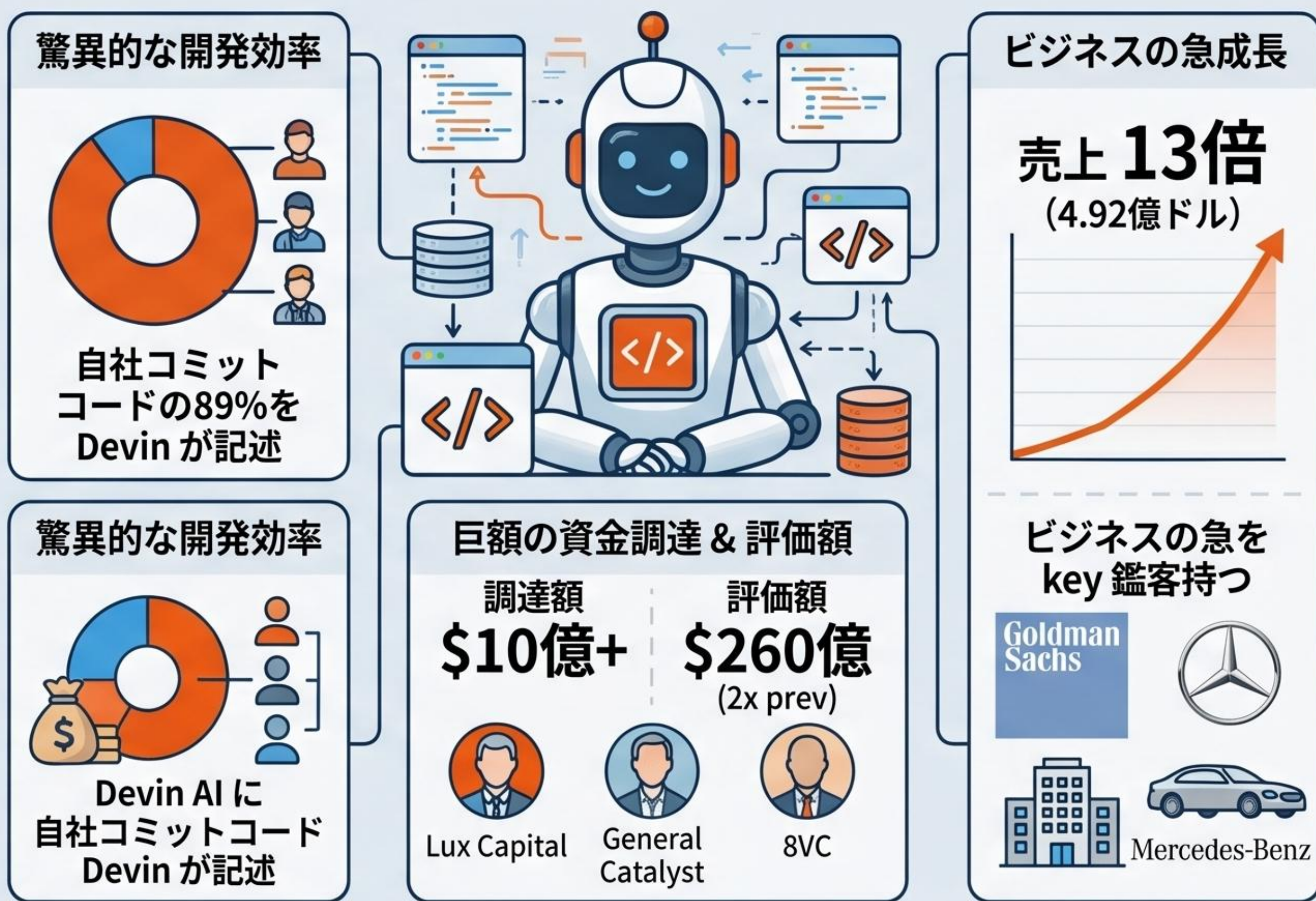
- 調達額10億ドル超、評価額260億ドル（前回比2倍超）
- 売上は1年で13倍の4.92億ドル。Goldman Sachs や Mercedes-Benz を顧客に持つ
- 自社エンジニアがコミットするコードの89%は Devin によるもの

## なぜ重要？

「Devin はプログラマを置き換えるためではなく支援するためのもの」と創業者が強調する中、AIエージェントの驚異的な開発スピードと市場価値の実証。

Lux Capital / General Catalyst / 8VC / Founders Fund / Ribbit

## Devin AI による開発加速と市場価値の実証



# StepFun が「Step 3.7 Flash」をオープンウェイトで公開 — 198B MoE・Apache 2.0・最大400トークン/秒

## 何が起きた？

中国の StepFun が高効率モデル『Step 3.7 Flash』をオープンウェイトで公開。196B の言語バックボーンと 1.8B のビジョンエンコーダがからなる 198B パラメータのスパース MoE（推論時アクティブ約 11B）。

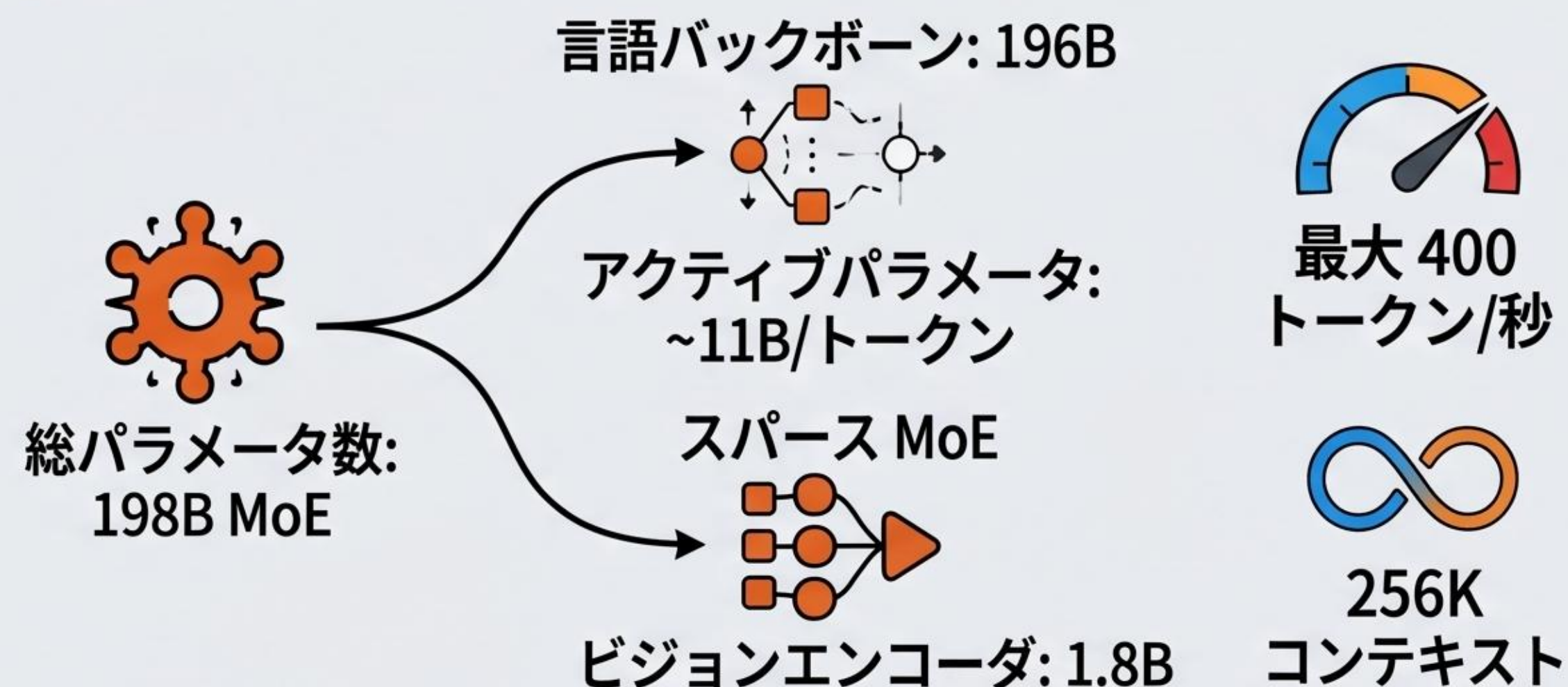
## 主な変更点

- 198BパラメータMoE、アクティブ約11B/トークン、最大400トークン/秒
- 256Kコンテキスト、画像・動画理解に対応するビジョン言語モデル
- ライセンスは Apache 2.0（商用利用可のオープンウェイト）
- エージェント型コーディング・ツール利用・長文推論・検索を主用途に設計
- 提供先: StepFun Open Platform / OpenRouter / NVIDIA NIM、vLLM・SGLang・llama.cpp で動作

## なぜ重要？

SWE-bench Pro などのベンチマークでオープンモデル最高水準を主張しており、商用利用可能な高性能オープンウェイトモデルの新たな選択肢となる。

## モデル構成と性能



## 主な特徴と展開

 Apache 2.0 (商用OK)	 マルチモーダル (画像・動画)	 エージェント機能 (コーディング・ ツール利用)	 ベンチマークリーダー (SWE-bench Pro 最高水準主張)
 StepFun Open Platform	 OpenRouter	 NVIDIA NIM	 vLLM・SGLang・llama.cpp

## 🔍 何が起きた？

Opus 4.8 と同時に発表された Claude Code の『Dynamic Workflows』（計画→数百の並列サブエージェント→自己検証→報告を動的に組む研究レビュー機能）が、公開直後から開発者の間で『久々に最も価値あるアップデート』と高評価を集めている。あわせて『タスク途中で Claude の指示を差し替えてもプロンプトキャッシュを壊さない』改善も実装の妙として注目された。

## 📌 主な変更点

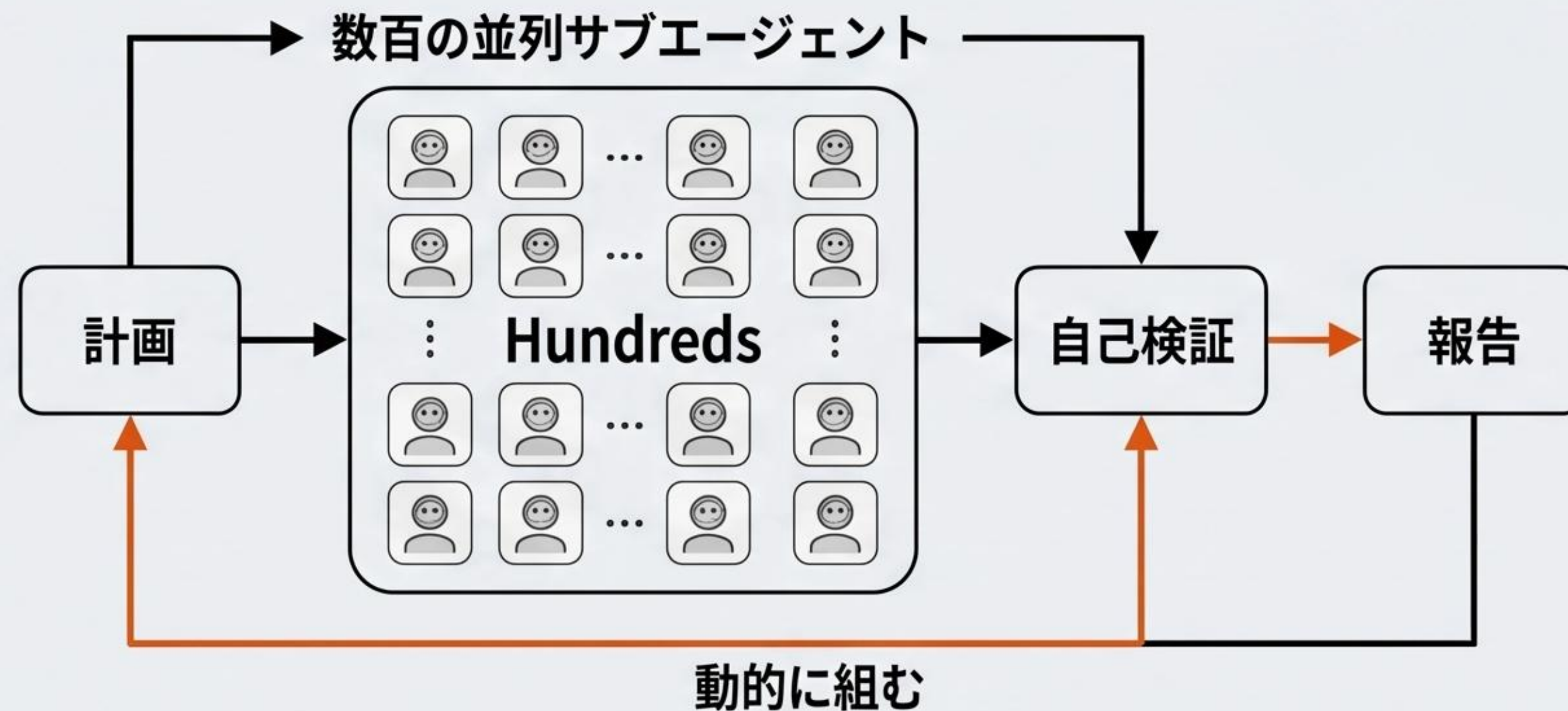
- Dynamic Workflows は計画→並列サブエージェント→自己検証→報告を動的に構成する研究レビュー機能
- 『ここ最近の Claude Code で最も価値あるアップデート』との実戦評価が複数
- タスク途中で system 指示を差し込んでも prompt cache が壊れない設計が好評
- effort control（思考量の制御）と組み合わせた長時間タスクの安定運用が狙い

## 💡 なぜ重要？

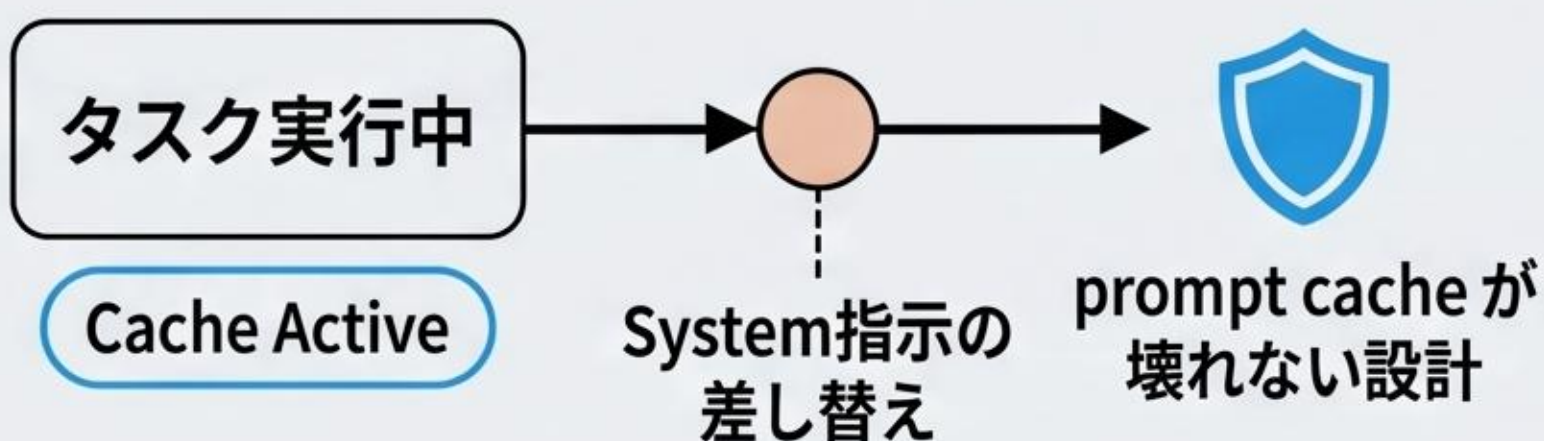
Why it matters の開発者のでデベロセをとに、複数をわれるタスク、複数突なタスクはまえは合的な価値値えび、ダイミックコントロール、また、技術術的な優化ココンク住姿定し存保籍などプロンプトキャッシュなどの prompt caching prompt preservation, ctrcs の保存可能に保存された。

## 🔄 Dynamic Workflows Flow

Anthropic公式 ❤️ 1.2k likes



## 🛡️ Prompt Cache Preservation



## 🏆 Praise Card

『久々に最も価値あるアップデート』

# Opus 4.8 は経験豊富なエンジニアのように、 逐一の確認なしで判断する (Claude 公式)

## 📄 何が起きた？

Anthropic公式がClaude 3 Opus 4.8の高度な自律性を発表。Claude Code上で動作し、開発者の負担を大幅に軽減することを示唆。

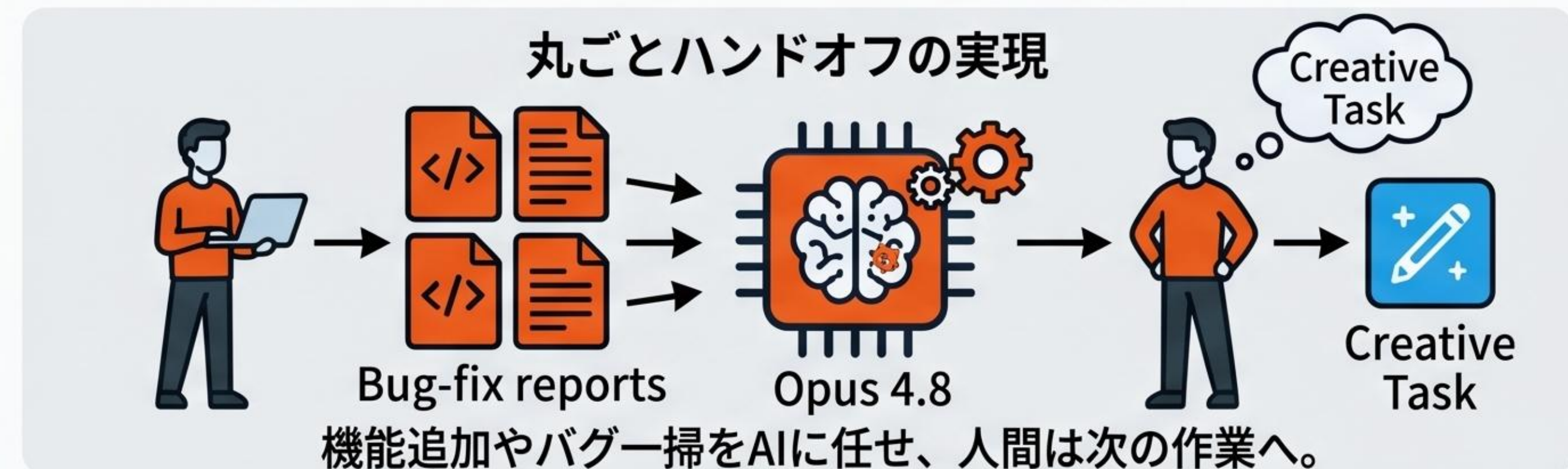
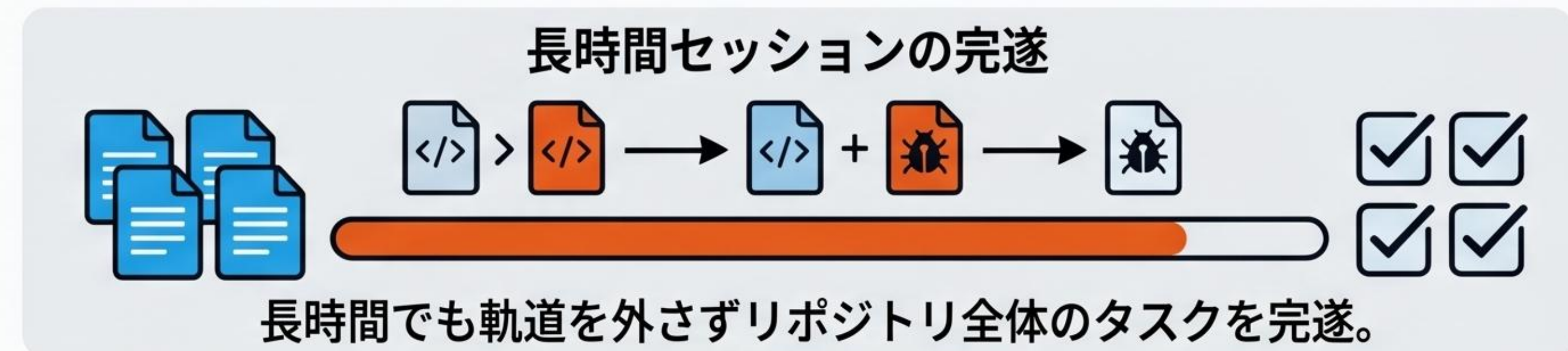
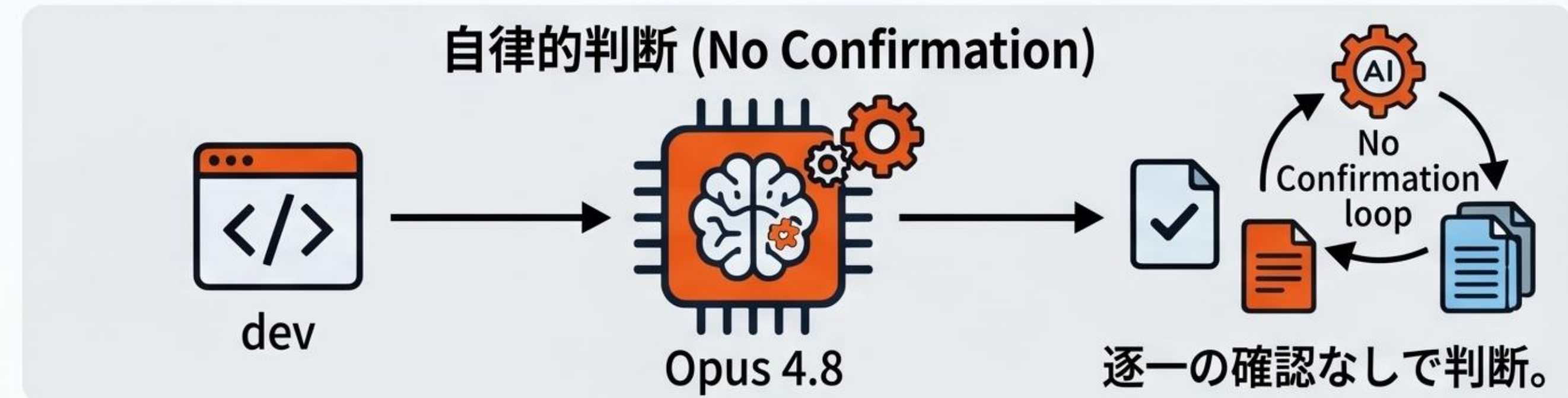
## 📌 主な変更点

- Opus 4.8はClaude Codeで『経験豊富なエンジニアのように』逐一の確認なしで判断。
- 長時間セッションでも軌道を外さずリポジトリ内の作業を完遂。
- 機能追加やバグー掃を丸ごとハンドオフできる自律性を強調。

## 💡 なぜ重要？

開発者が煩雑な作業から解放され、AIに自律的に丸ごと任せ、よりクリエイティブな「次の作業」に集中できる。AIエンジニアリングの大きな進歩。

## Opus 4.8の自律的エンジニアリング：『経験豊富なエンジニア』の動き



# Opus 4.8 に Fast モード登場 — 同じモデルで 約2.5倍速・従来比3倍安 (Claude 公式)

## 🔍 何が起きた？

Anthropic 公式が Opus 4.8 の Fast モードを発表した。同一モデルで速度約2.5倍、価格は従来比 3倍安になる。

## 📌 主な変更点

- Fast モードは Opus 4.8 で利用可能、同一モデルで約2.5倍速・従来比3倍安。
- Claude Code では /fast で即切替。
- API はアカウントマネージャ連絡 or ウェイトリスト登録 (claude.com/fast-mode)。

## 💡 なぜ重要？

このアップデートは、FISD 最級の Opus 4.8 モデルの速度とコスト効率的な敷動に改善。開発効率に両も、制成の幅揃と、AIの利活用コストに利な優益ををも与える。



## アクセス方法

Claude Code

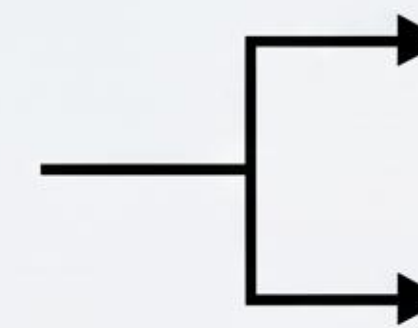
Code: /fast



即切替

API

API

アカウントマネージャ  
へ連絡ウェイトリスト登録  
claude.com/fast-mode)

# 大規模開発は「人間 ↔ Opus ↔ Codex」の多重下請け構造がおすすめ (Kenn Ejima)

## 🔑 要点

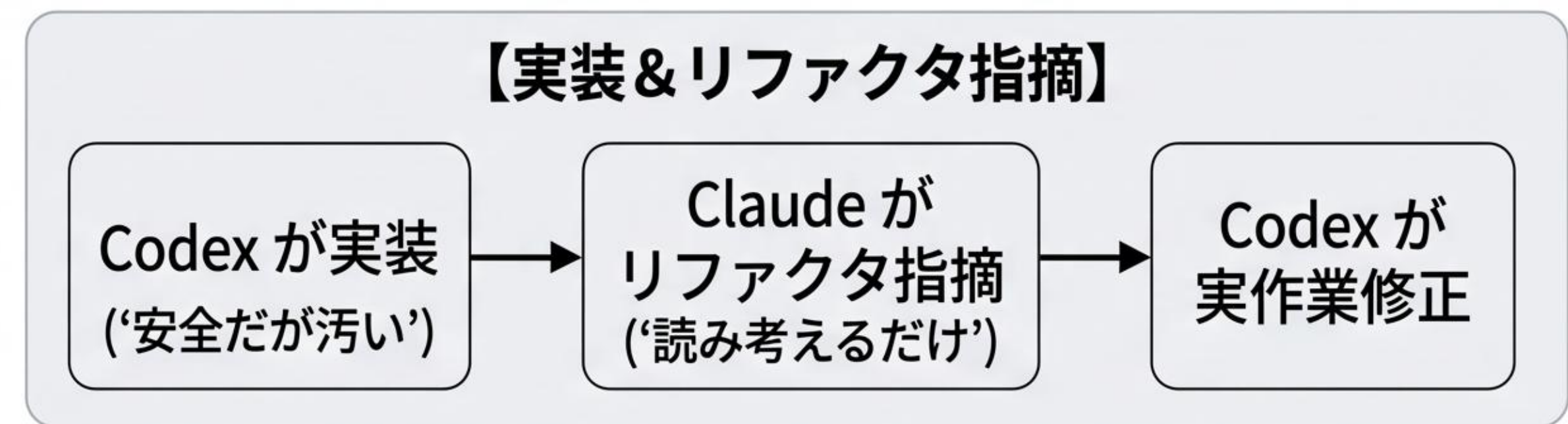
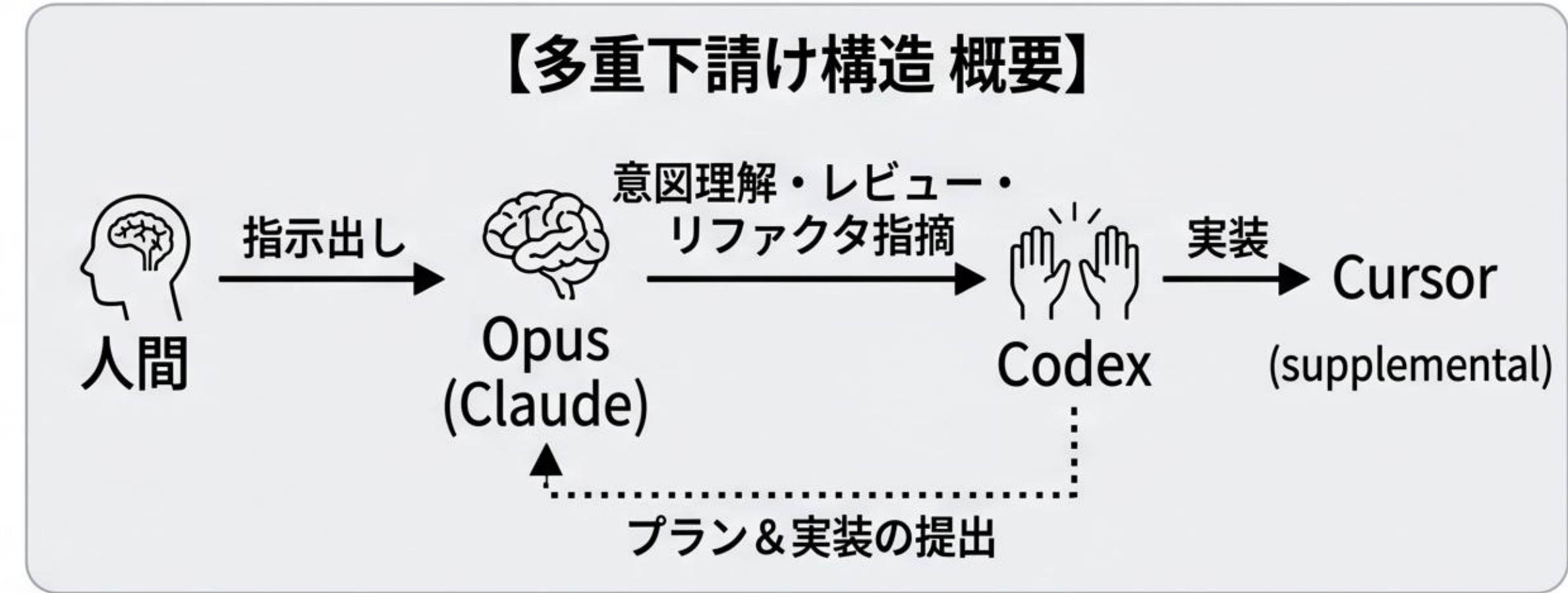
Kenn Ejima 氏が大規模開発でのマルチエージェント開発フローを self-thread で整理した投稿。『人間 ↔ Opus ↔ Codex』の多重下請け構造を勧め、Opus を意図理解・レビュー・リファクタを指摘の頭脳、Codex を実装の手と役割分担する。

## 🔧 具体的な手法 / 使いどころ

- 大規模開発は『人間 ↔ Opus ↔ Codex (↔ Cursor)』の多重下請け構造が有効
- Opus は意図・本質をつかむのが得意だがコードを書かせると手戻りが多い / Codex は実装が漏れない反面、細部にとらわれ本質を見落とす
- プランは Codex に書かせ Claude にレビューさせる (Claude のプランは詳細に踏み込みすぎるのでツッコミ役に徹させる)
- Codex の『安全だが汚い』コードは最後に Claude がリファクタを指摘し、実作業は Codex 本人に戻す

## 🌱 なぜ刺さるか / 学び

結論: Claude は読んで考えるだけ、docs 以外は書かせない



# 本日のトピック一覧

1 🔍 AIコーディングエージェントの「セキュリティ」が新たな主戦場に  
OpenAIが公式ガイダンス、脆弱性・防御ツールも続々



2 🔍 Cognition (Devin 開発元) が約10億ドルを調達、評価額260億ドルに  
自社コードの89%を Devin が記述



3 🔍 StepFun が「Step 3.7 Flash」をオープンウェイトで公開  
198B MoE・Apache 2.0・最大400トークン/秒



4 🔍 Claude Code 「Dynamic Workflows」が実戦で高評価  
Opus 4.8 の目玉機能



5 📖 Opus 4.8 は経験豊富なエンジニアのように、逐一の確認なし  
(Claude 公式ガイダンス) (Claude 公式)



6 📖 Opus 4.8 に Fast モード登場  
同じモデルで約2.5倍速・従来比3倍安 (Claude 公式)



7 📖 大規模開発は「人間 ↔ Opus ↔ Codex」の多重下請け構造がおすすめ  
(Kenn Ejima)

