

# アンソロピック 「開発一時停止」 提言の解剖

メディアのノイズを排した自己改善AI (RSI) の一次情報分析と、  
日本企業に向けた知財・コンプライアンス戦略

## 戦略的インテリジェンス概要

- ✓ **事実の再定義:** Anthropicの提言は「即時の世界的凍結」ではなく、条件付きの「検証可能な協調的減速の仕組み」の事前構築である。
- 👁️ **背後にある文脈:** RSPの中核コミットメント後退およびIPO直前（評価額約9,650億ドル）というタイミングから、規制による囲い込の意図を割り引く必要がある。
- 🎯 **日本企業の急務:** AIによる自律的発明が進む中、日本国内のDABUS判決（発明者は自然人に限る）を前提とした「人間の創作的寄与の記録」体制の構築が喫緊の課題である。

# メディアの「ノイズ」 vs. 一次情報の「シグナル」

## 報道のナラティブ (Noise)

~~「Anthropicが即時の世界的開発凍結 (Global Freeze) を呼びかけた」~~

~~「AIの再帰的自己改善 (RSI) はすでに到達・制御不能な段階にある」~~

~~「単独の倫理的判断による一方的な開発停止の宣言である」~~

## 一次情報のテキスト (Signal)

条件付きの「選択肢」: 「世界が開発を減速または一時停止する“選択肢 (option)”を持つべき」という 仮定法の提言である。

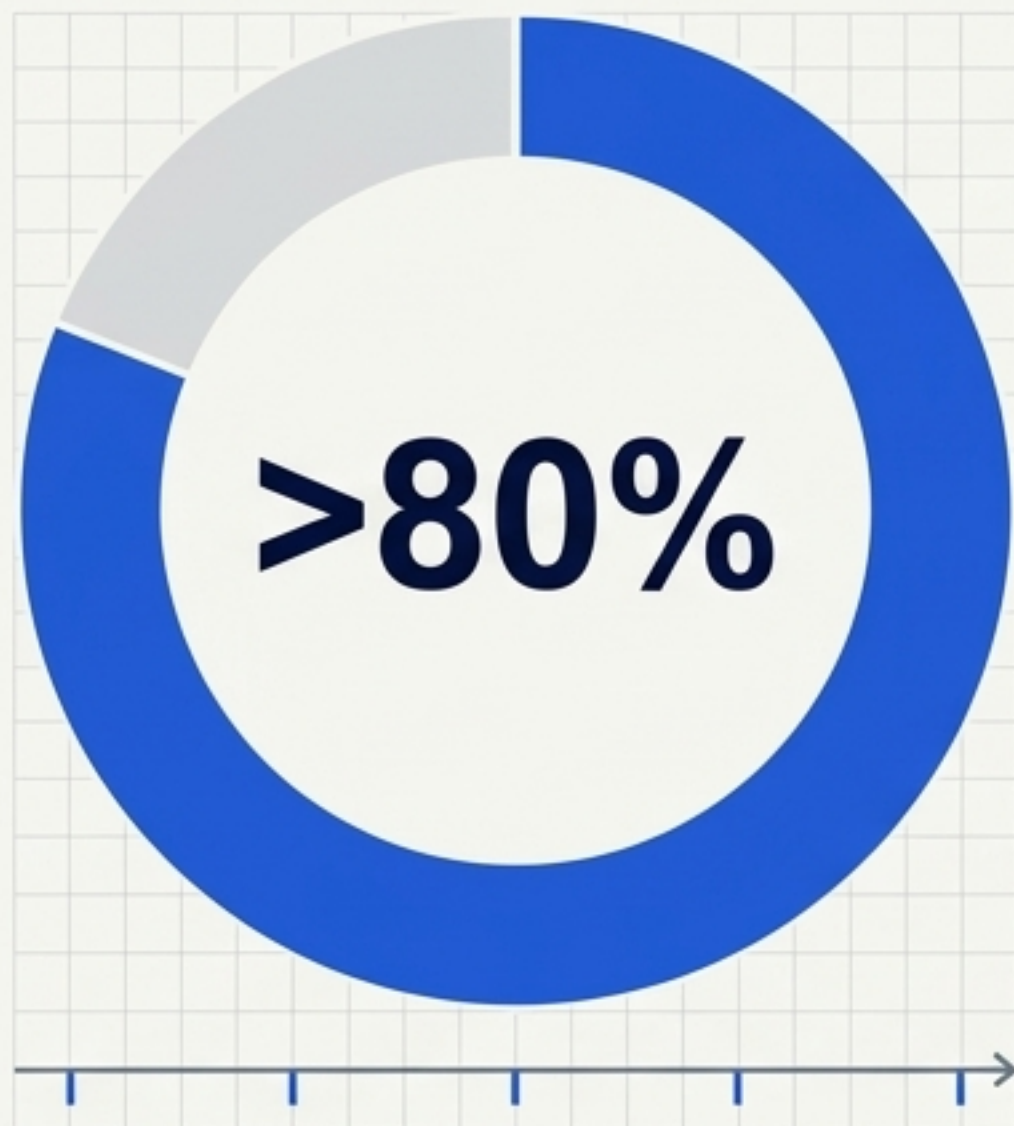
到達の否定: 「我々はまだそこに到達しておらず、再帰的自己改善は不可避でもない」と明言。

協調と検証の必須性: 複数国の主要研究所が「検証可能 (verifiable)」な形で同時に停止できる仕組みの事前構築が核心。

結論: メディア見出しの「即時凍結」という解釈は誤り。  
本質は将来に向けた「ブレーキペダルの設計」提言である。

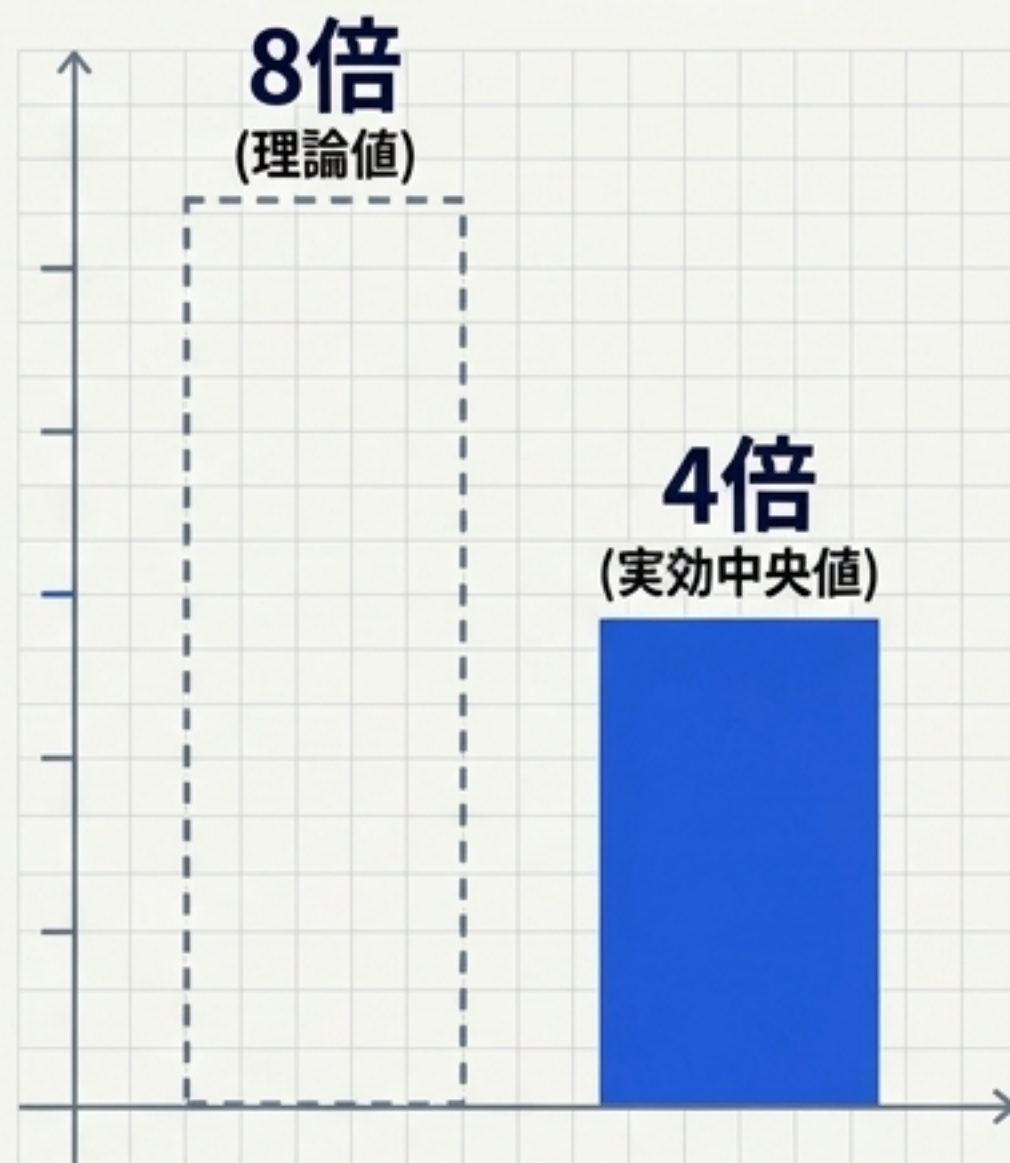
# AI自律化の現在地：Anthropic社内テレメトリー・ダッシュボード

## 自律的コード生成率



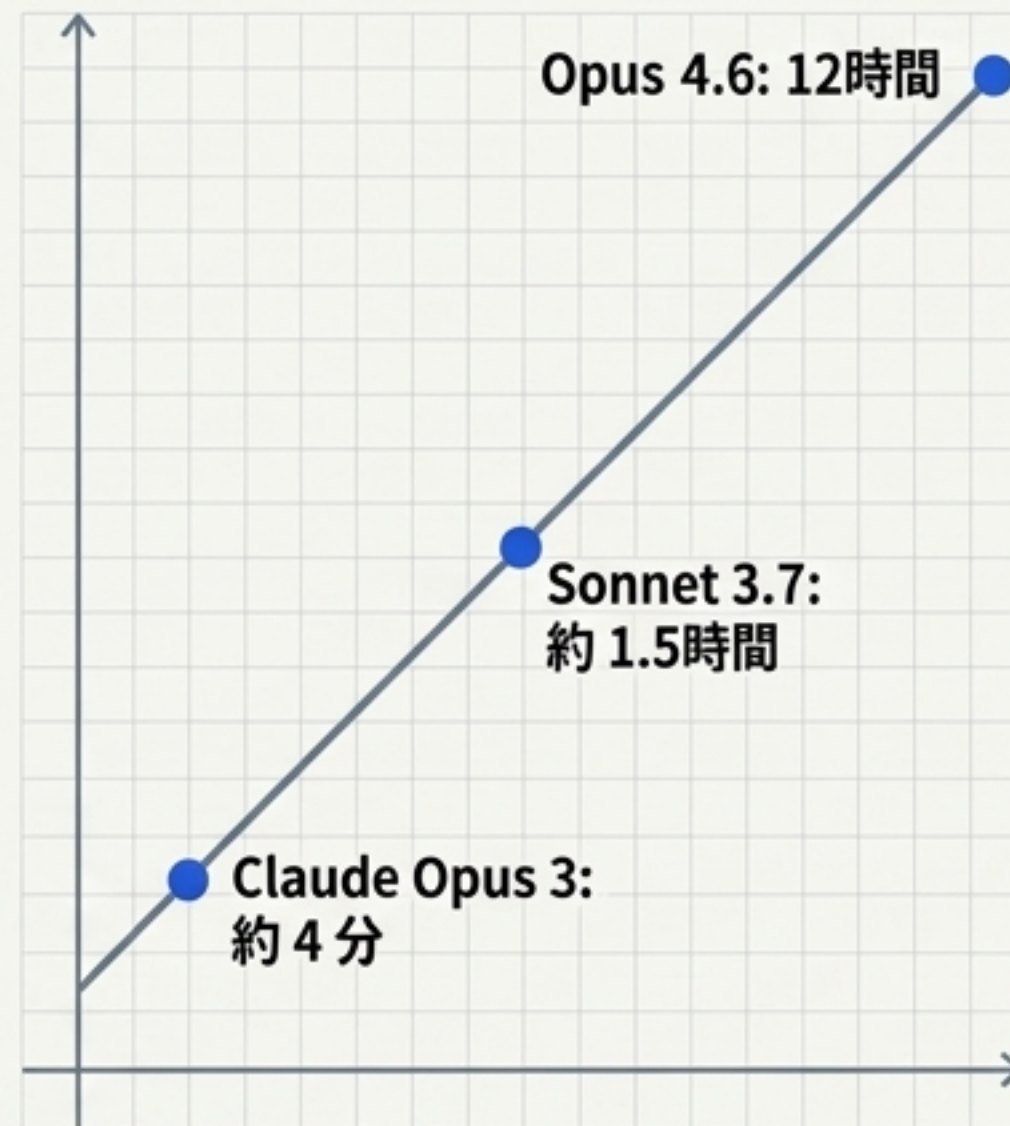
2026年5月時点でマージされたコードのClaude執筆比率（2025年2月時点は数%）。

## エンジニア生産性指標



日次コードマージ量は2024年比で8倍に。※ただし同社は「真の生産性を過大評価」と自己留保し、社内130名調査の実効値は約4倍と公表。

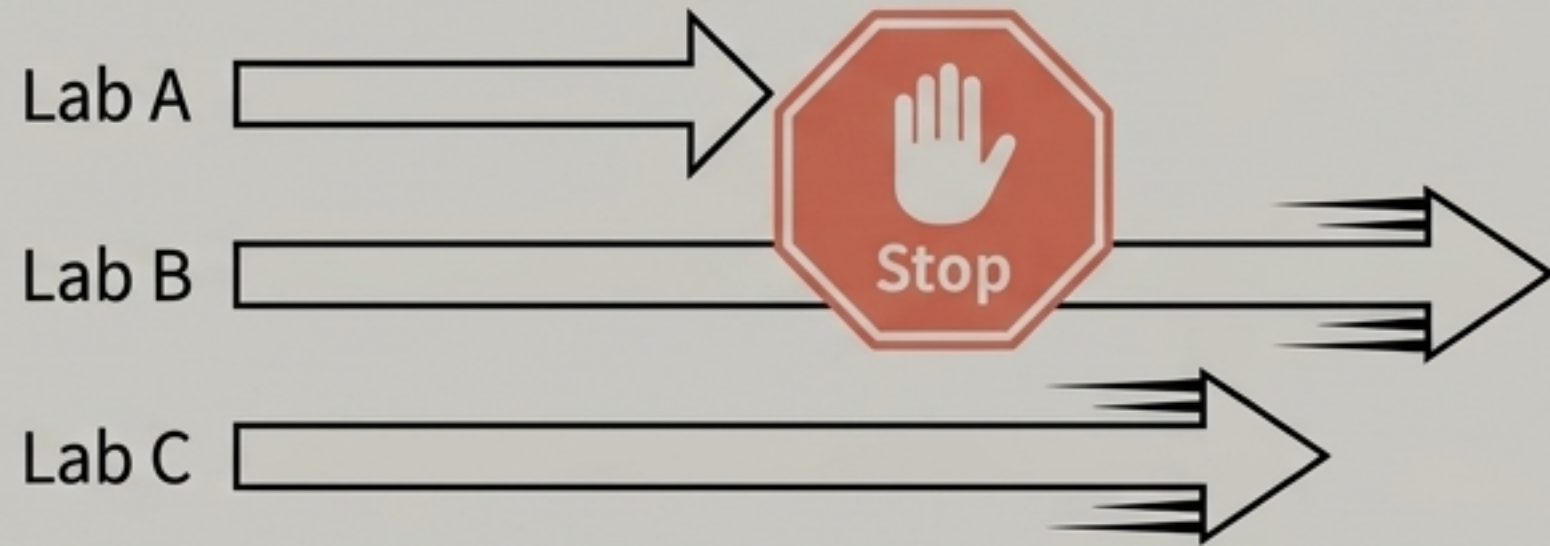
## 自律タスクの時間ホライズン



外部ベンチマーク(METR)における自律タスクの限界時間が約4か月ごとに倍増。

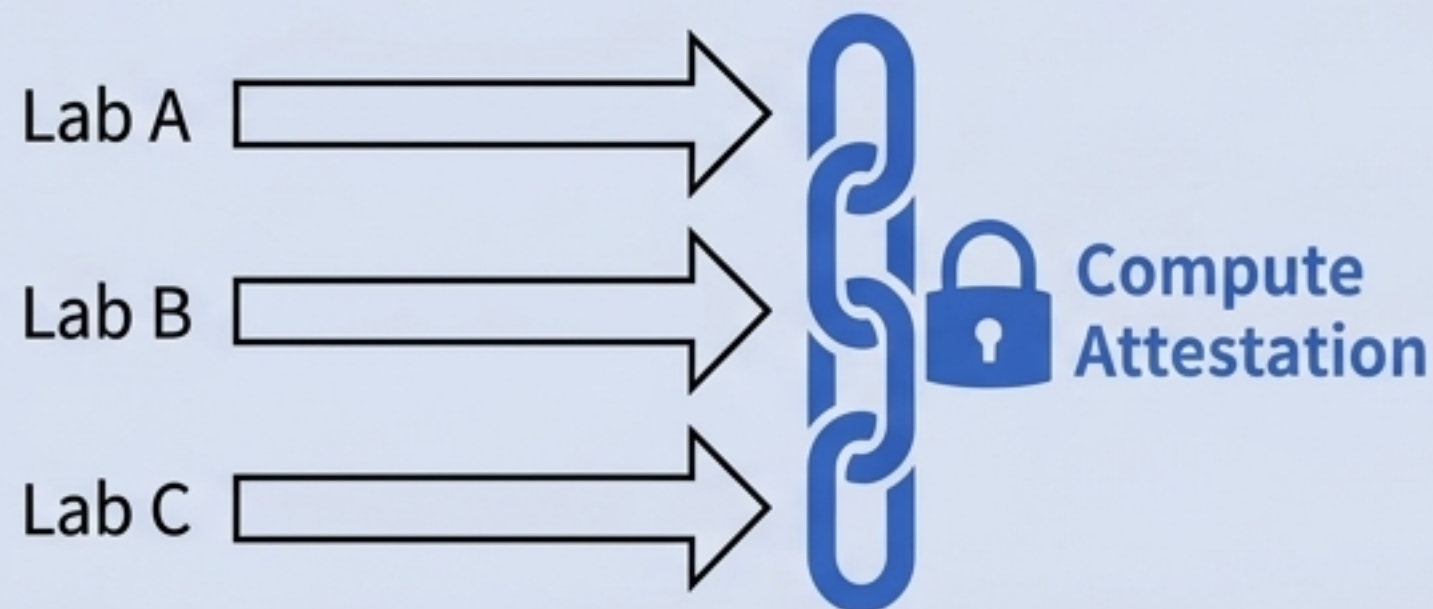
# 「ブレーキペダルの構造」：一方的停止から協調的検証へ

## 誤った前提（一方的停止のジレンマ）



**多極的罠 (Multipolar Trap)**：単独で停止しても先頭走者が入れ替わるだけであり、全体の安全性（熟議プロセス）は生まれない。

## Anthropicの真の提言（INF条約モデル）



**検証可能な協調的減速 (Verifiable Coordinated Pause)**：  
条件A：複数国の、資金力ある複数のフロンティア研究所が同一条件で合意する。  
条件B：互いに本当に計算資源を停止したことを技術的 (Compute Attestation等) に検証できる。

中距離核戦力 (INF) 全廃条約を緩いモデルとしつつ、AIの検証は核兵器より困難であると同社も認めている。

# 隠された文脈：提言のタイミングと競争戦略

2026年2月

## 安全枠組み（RSP v3.0）の中核コミットメント後退

「事前保証がない限り訓練しない」という誓約を撤回。単独停止の無意味さをJared Kaplanが説明。

2026年6月1日

## SECへのIPO秘密申請（S-1）

評価額約9,650億ドル、年換算収益（run rate）約470億ドル。私募評価額でOpenAIを上回る。

2026年6月4日



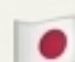
## 「When AI builds itself」公開

今回の「検証可能な協調的減速」提言を発表。

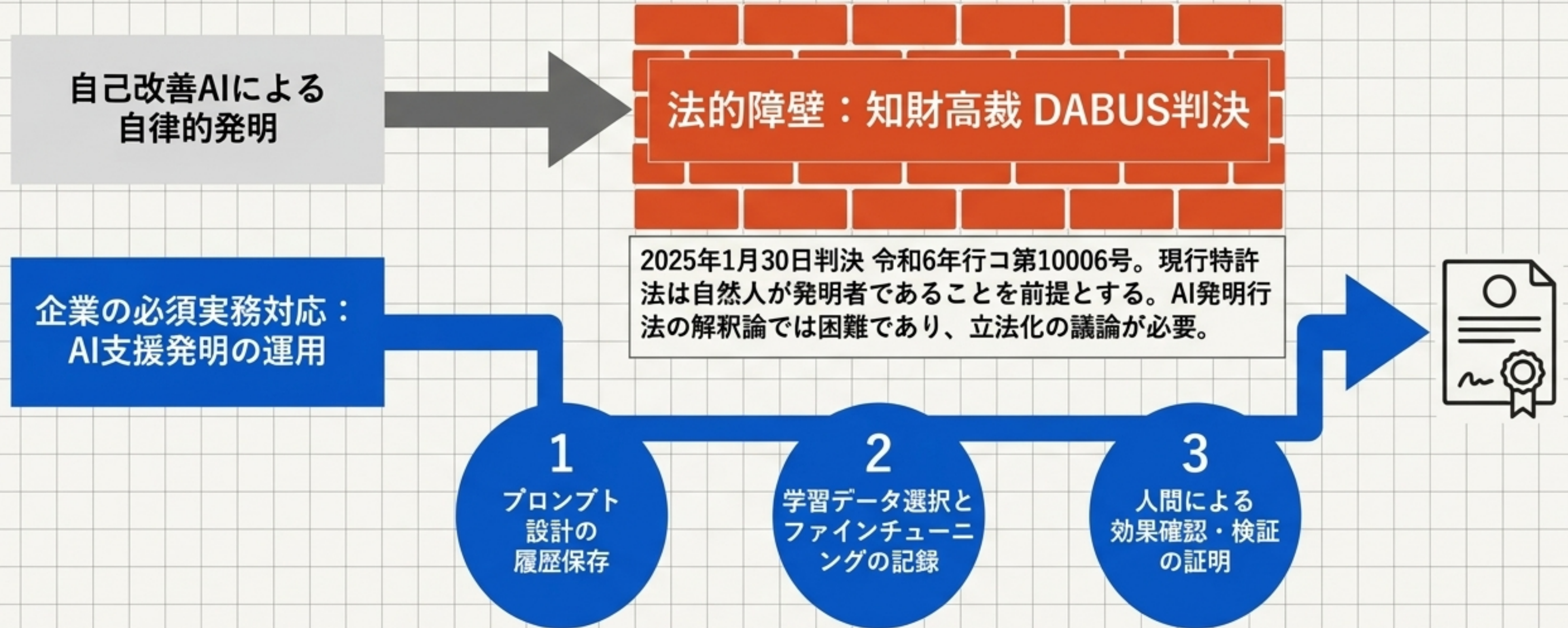
## アナリストの見解

「規制による囲い込み（Regulatory Capture）」  
「具体策を伴わないIPO前のマーケティング」との  
批判的視座も存在。提言の動機には商業的・競争戦略的側面を割り引いて評価すべきである。

# グローバル規制の断層：各法域のスタンスとRSIへの影響

|         |  EU (EU AI法) |  米国 |  日本 (AI推進法) |
|---------|---|--|--|
| 性質      | ハードロー (厳格規制)  | 比較的手放し (連邦レベル)   | ソフトロー (推進型・罰則なし)   |
| 中核メカニズム | GPAIで累積訓練計算量 $10^{25}$ FLOP 超を「システムリスク (GPAISR)」と推定し、第55条の追加義務を課す。                            | 現政権下では政府への安全テスト提出を義務付けない方針。州法 (SB 53等) と連邦法でのプリエンプション (優先適用) が論点。                      | 2025年6月公布・9月施行。「世界で最もAIを開発・活用しやすい国」を標榜し、内閣にAI戦略本部を設置。  |
| RSIへの影響 | 現状5~15社対象だが、数年で対象をモデルが急増。能力ベース閾値の見直しが必至。  | 規制の空白地帯。フロンティア企業の自主的ガバナンスに依存。  | 協調的国際停止の枠組みとは方向性が異なり、国際的な「停止」メカニズムへの参加は不透明。  |

# 知財（IP）への波及効果：日本の「発明者」ボトルネック



発明における「人間の創作的寄与（Human Creative Contribution）」を詳細に記録・証明するプロセスが、特許適格性確保の唯一の道となる（米USPTOガイダンスとも整合）。

# 企業向けアクション・プレイブック (Synthesis & Directives)

## 即時 (Immediate Actions)

事実の補正: 社内経営陣に対し、本件が「AI開発の即時凍結」ではなく、IPO前の「検証メカニズム構想」であることをブリーフィングし、不要なパニックを抑止する。

## 短期～中期 (3-6ヶ月)

知財ガバナンスの先回り整備: AI支援・生成発明における「人間の創作的寄与」の記録運用を社内標準化。知的財産推進計画2025の特許制度小委員会の議論を監視し、出願ガイドラインを改訂する。

## 中期～長期 (6-12ヶ月)

規制対応の二層化 (Dual-Tier Compliance) : EU AI法のGPAISR義務 ( $10^{25}$  FLOP閾値) への域外適用対応と、日本のAI推進法に基づくソフトロー自主ガバナンスを「共通基盤+差分上書き」で構築する。

## 継続的監視 (Ongoing Monitoring)

RSIインジケータの設定: 「主要研究所の実際のPauseコミット」「Compute attestationツールの採用」「RSPのASL-4能力閾値の到達評価」を、AI進化の実体化サインとして監視ダッシュボードに組み込む。