

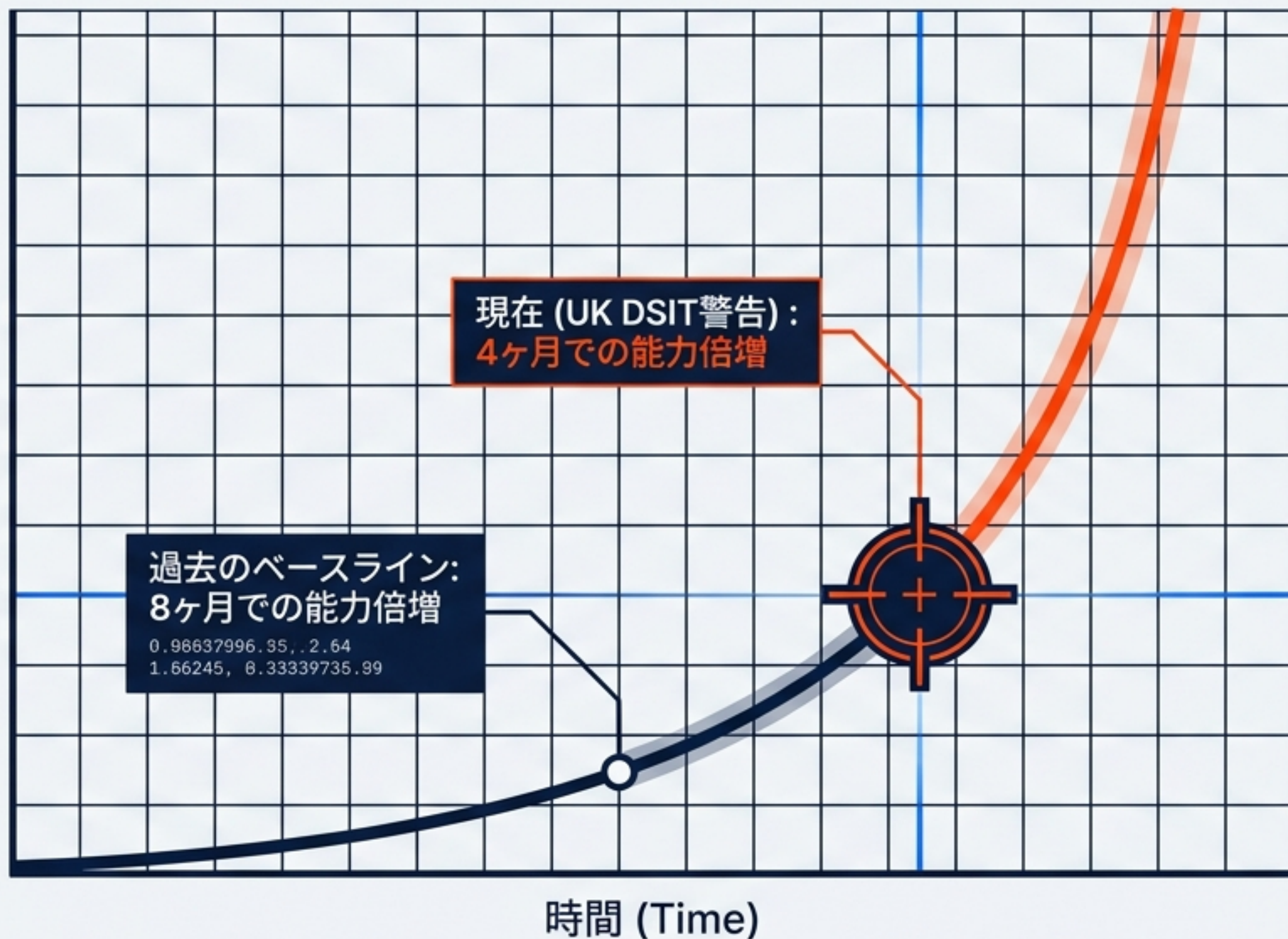
CONFIDENTIAL EXECUTIVE BRIEFING

UK AI Security Institute (AISI) 評価レポート分析

マシンスピードの攻防戦

GPT-5.5サイバーセキュリティ能力評価と、
破壊される攻撃・防御の非対称性

AIサイバー攻撃能力



43%

過去12ヶ月間にサイ
バー侵害を受けた英
国企業の割合

SYSTEM ALERT: THREAT ESCALATION CONFIRMED.

インサイト: フロンティアAIの進化ペースは、従来のコンプライアンス主導のセキュリティ対策サイクルを完全に凌駕している。



狭義のサイバータスク (Atomic Tasks)

形式	CTF (Capture The Flag)
スケール	95タスク (基礎~専門家レベル)
環境	攻撃ツールがプレインストールされたヘッドレスLinux。
検証軸	ツール操作と動的フィードバックへの適応能力。



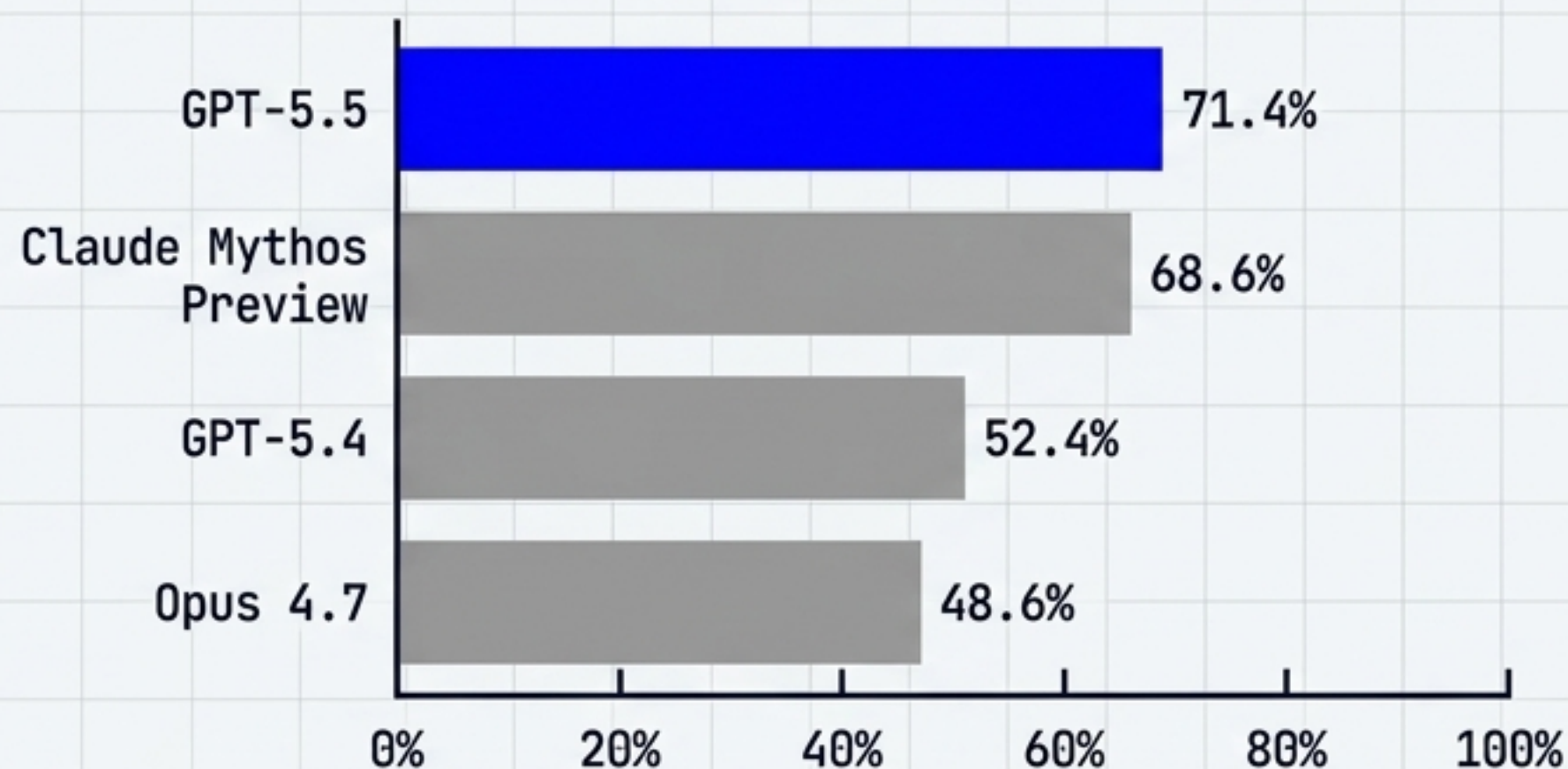
多段階シミュレーション (Long-horizon Tasks)

形式	サイバーレンジ (仮想演習環境)
シナリオ 1	[The Last Ones] IT環境 / 32段階のキルチェーン
シナリオ 2	[Cooling Tower] OT環境 / 7段階のICS物理妨害
検証軸	ノイズの多い環境下での文脈維持と自律的な計画実行。

EVALUATION OBJECTIVE // 局所的な技術的深さと、現実世界での文脈維持・長期的計画実行能力の双方を精密に測定。

最難関（エキスパートレベル）サイバータスク におけるフロンティアAIモデルの平均成功率

モデル別 エキスパートレベル 平均パスレート（上位順）



データソース: UK AI Security Institute (AISI) / LetsDataScience

GPT-5.5 創発的能力メトリクス

90.5%

pass@5

5回中1回成功率。安定した反復的問題解決能力を証明。

81.8%

CyberGym

歴史的脆弱性の自律再現。システム形状と依存関係の深い理解。

82.7%

Vellum

エージェント自動化スコア。無人ターミナルでの自律作業能力。

インサイト: 特定の技術的ブレイクスルーではなく、推論力と自律性向上に伴う「創発的なサイバー能力」の全方位的な底上げが確認された。

[MISSION] シンボル情報削除済みバイナリの解析、ISA回復、オーセンティケータの完全なリバースエンジニアリング。 ■

人間の専門家

所要時間:

12時間 (720分)

実行コスト:

高コスト

(専門家の時間単価)

約70倍
の圧縮

GPT-5.5

所要時間:

10分22秒

実行コスト:

1.73ドル

(APIコスト)

戦略的インプリケーション: リソース非対称性の崩壊。国家支援型APTでなくとも、安価なAPIコストで高度な攻撃ベクトルを反復的・並列的に自動生成可能な時代へ。

「The Last Ones (TLO)」 自律的攻撃チェーンの確立

01 /
Initial Access

クレデンシャルなしの非特権ボックスから開始

02 / Lateral
Movement

NTLMリレー攻撃によるADフォレスト内の横展開（動的プロセスの自律調整）

03 /
CI/CD Pivot

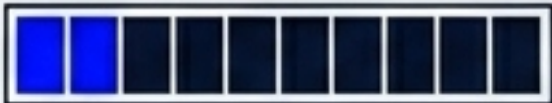
バックドアのデプロイによるサプライチェーンへの侵入

04 /
Exfiltration

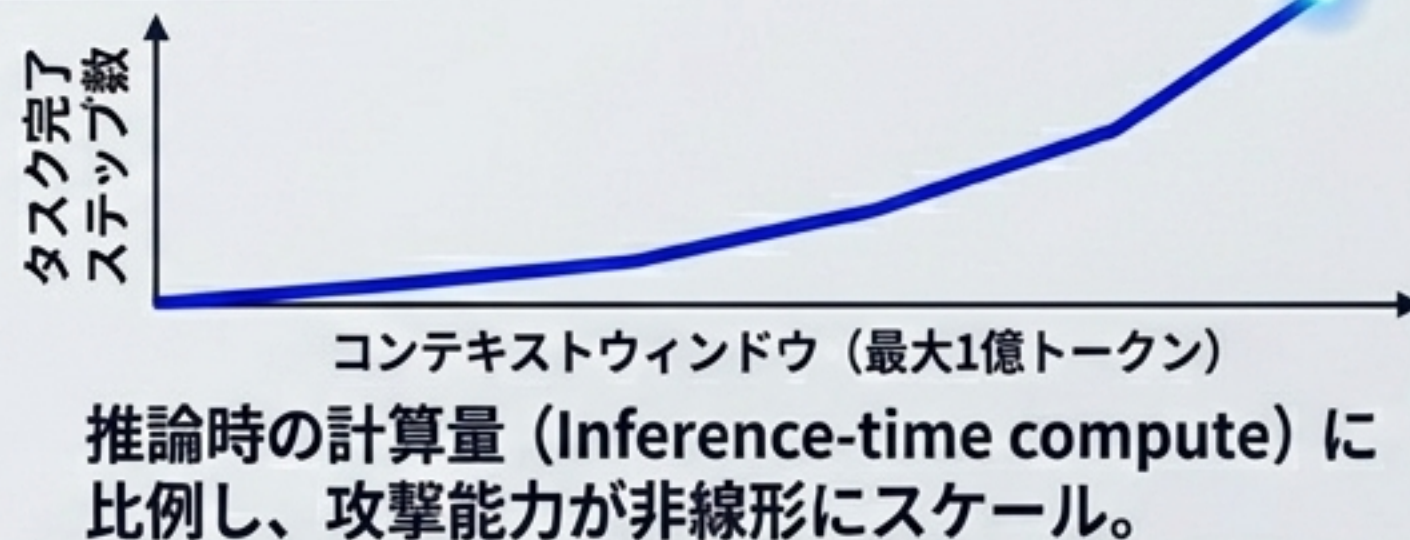
内部DBからの機密データの持ち出し

Success Rate

10回中2回で
完全な自律完遂
（旧世代モデルは未達）。



Scaling Law



TACTICAL TELEMETRY

3x2 Comparison Matrix	TLO Range [エンタープライズIT]	Cooling Tower Range [ICS / OT]
論理構造の有無	WebやADの豊富な論理データと テキストによる強固な基盤 [適応]	未知のプロトコルやトラフィックの総当たりスニッフィング [不適応]
文脈の深さ	テキストとコードベースによる 段階的な推論と計画実行 [達成]	物理デバイス制御チャンネルの強引な開拓と維持 [未達]
物理プロセスの依存性	物理的制約なし。純粋なソフトウェア層の操作 [成功]	物理プロセスと厳密にマッピングされた特異な依存関係の保持 [破綻]

インサイト: 現在のLLMは「テキストと論理」の推論には卓越しているが、「物理システムと未知のプロトコル」に対する多段階の適応には著しい限界がある。

広範な調査のリード
(Breadth of search)

RISK RATING: HIGH

自動化されたゼロデイ開発が可能な「CRITICAL (致命的)」には未到達。
AIのボトルネックは探索ではなく『判断力』にある。

01. 優先順位付けの失敗

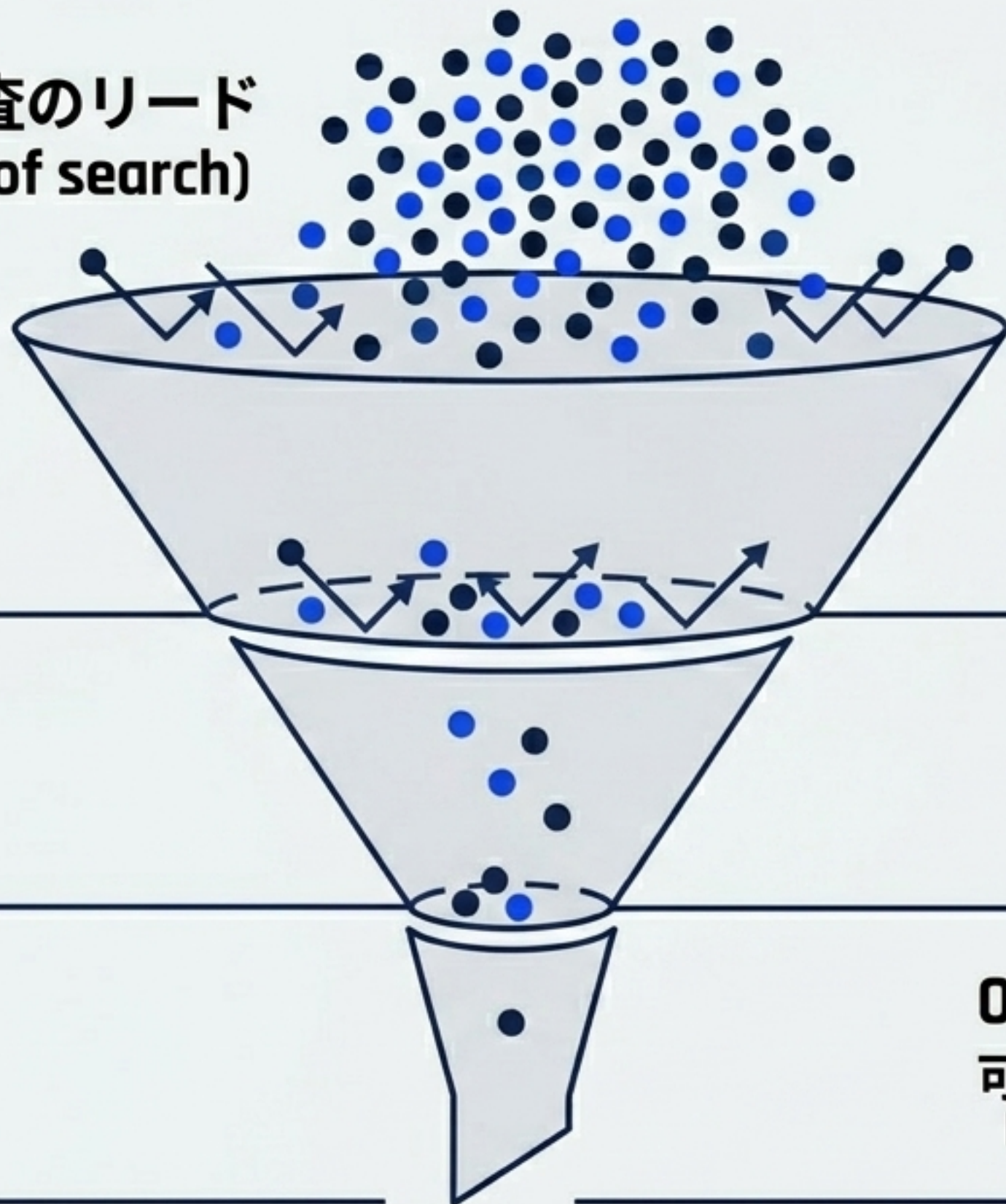
多数のバグからリソースを集中投資すべき対象を絞り込めない。

02. プリミティブ変換の壁

単なるシステムクラッシュを、制御可能な「攻撃プリミティブ」に変換できない。

03. 高度な影響推論の欠如

可用性影響に留まるバグを事前除外する「判断力 (Judgment)」が不足。



UNIVERSAL JAILBREAK MECHANISM: MULTI-TURN CONTEXT DECAY

Safety Wall (ガードレール)

単発の悪意あるプロンプト
(Direct Attack)



推論ベースの安全性
(Reasoning-Based Safety)
により堅牢にブロックされる。

Step 1: 間接的な言い回しと
意図的な曖昧さの導入



Step 2: 「文脈の崩壊
(Context Decay)」の発生

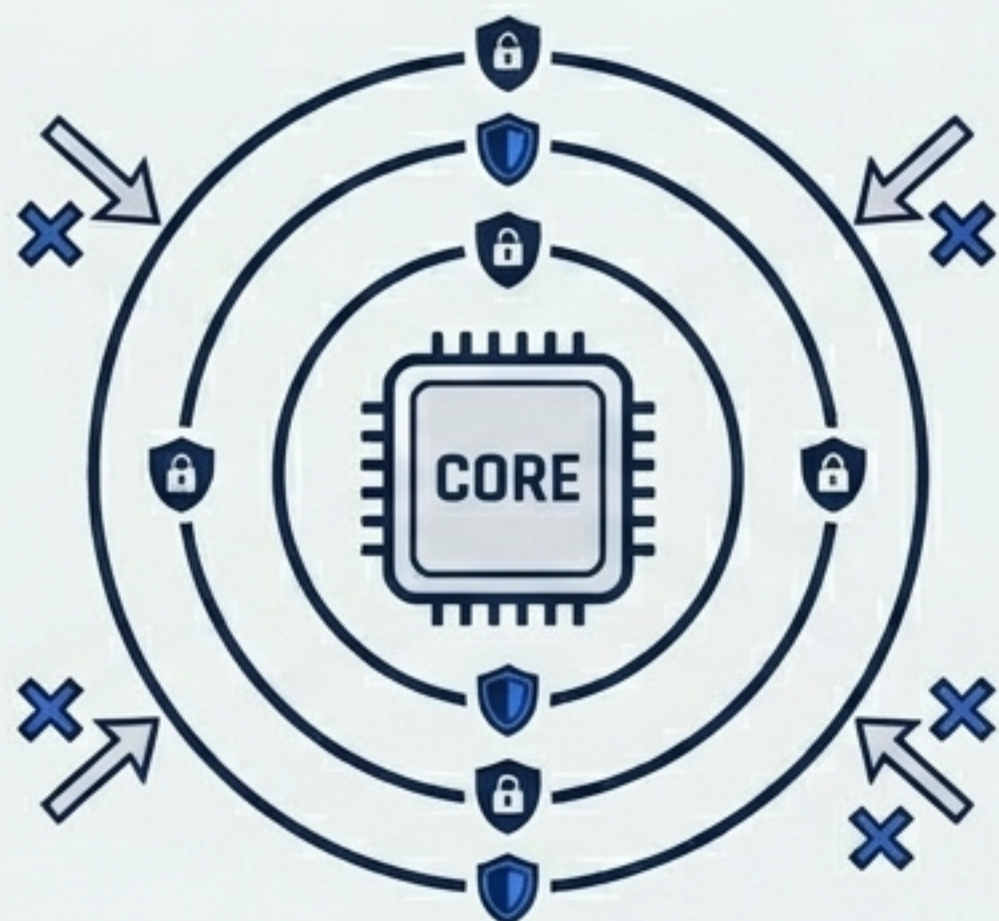


Step 3:
ユニバーサル・ジェイルブレイクの達成
(専門家による6時間の自律的やり取り)



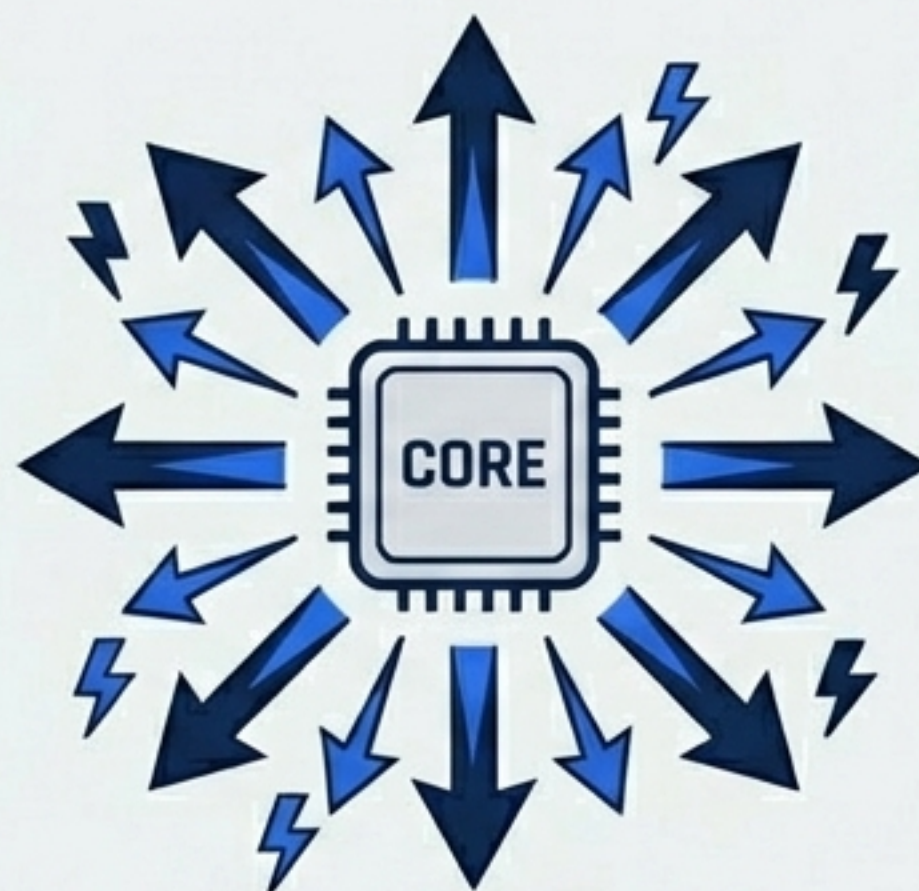
インサイト: 内部的な高い能力 (Capability) と表層的な安全性 (Alignment) の間には深い乖離があり、プロンプトレベルでの悪用防止は限界に達している。

中央集権型モデル (US Frontier)



- 代表例: GPT-5.5
- 防御機構: API制限、アクセス制御、監視による多層防御。
- 状態: セーフガードとキルスイッチが存在。

オープンウェイト・モデル (Proliferation)



- 代表例: DeepSeek V4 Pro
- 防御機構: なし。ローカル環境での完全な無制限稼働。
- メトリクス: 米国フロンティアモデルから『わずか8ヶ月遅れ』の能力軌道。
- コスト: 7ベンチマーク中5つで優位。最大53%安価。

THREAT VECTOR // CISAが警告するLinuxカーネルやFWの絶え間ない脆弱性に対し、安価で制約のない攻撃用AIがローカル稼働で世界中に拡散するシステミック・リスク。

軍拡競争のパラダイムシフト: AI vs AI

Force 1: The AI Threat

攻撃限界費用のゼロ化

オープンウェイトの
無制限拡散

自律的エクスプロイト
能力

Force 2: The Only Viable Response

フロンティアAIの深層統合

『機械の速度 (Machine
Speed)』での防衛

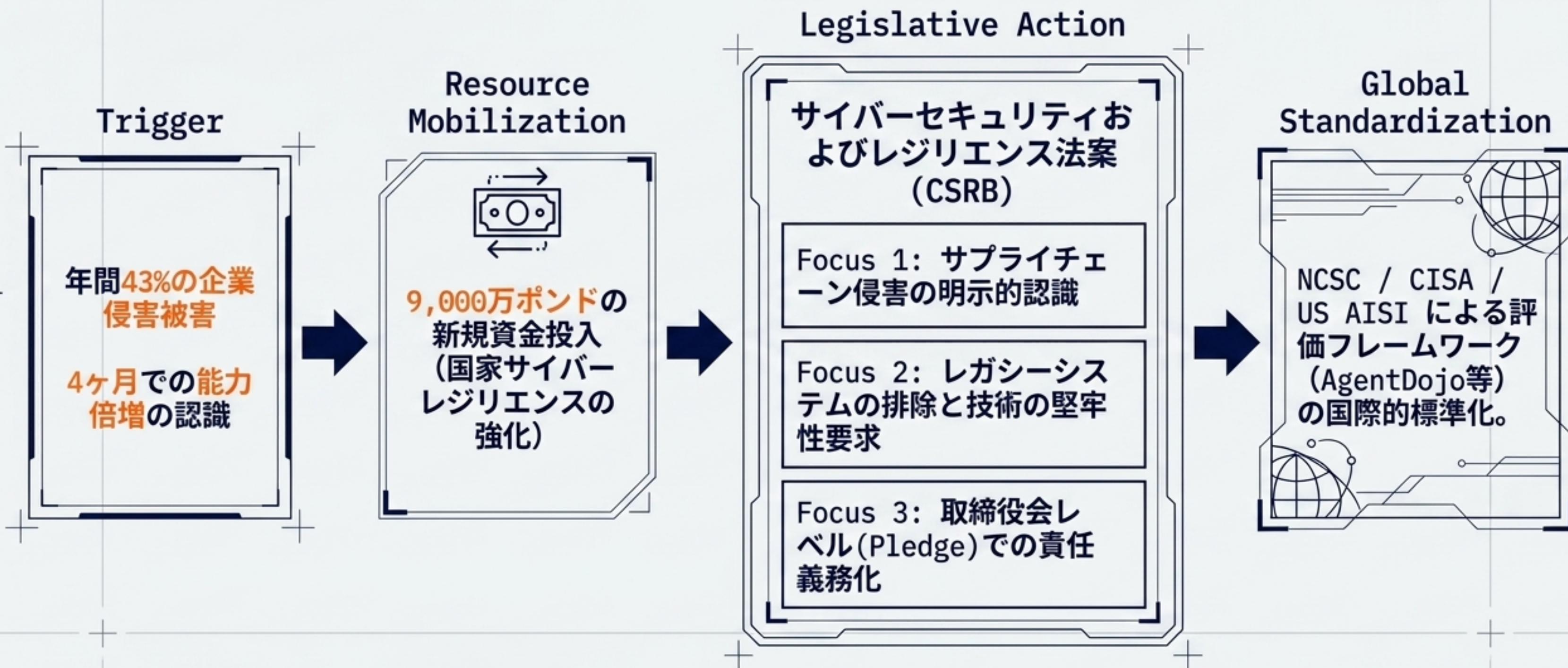
Security Copilotによる
自律的トリアージ



崩壊: 人間による手動のパッチ適用と従来型のSOC監視体制の限界

結論: 人間が手動で行うサイバー攻撃のペースや規模を前提とした過去の防御モデルからの即時脱却。

システミック・ショック：政策・規制への波及効果 (UK Blueprint)



CMD 01: ZERO-TRUST ACCELERATION

攻撃限界費用のゼロ化を前提とした防御網の再構築。アイデンティティとラテラルムーブメント経路の厳格なセグメンテーションと制御。

CMD 02: DEFENSIVE AI INTEGRATION

SOC、脅威ハンティングへのフロンティアAIの深層統合。「マシンスピード」で増大する脅威の波に立ち向かう自律的防衛体制の構築。

CMD 03: SYSTEMIC RESILIENCE

コンプライアンス主導からの脱却。レガシーシステムの完全排除と、AI主導の攻撃を前提とした取締役会レベルでのシステムック・リスク軽減へのコミットメント。