

AIミユトス:金融システムを脅かす「今そこにある危機」と共同防御の道筋

高度なサイバー能力を持つAIの台頭と、官民一体で取り組むべき防御策「日本版Project Glasswing」

2026年4月7日



Anthropicが「Mythos」と「Glasswing」を発表。高度なサイバー能力を持つAIモデル「Claude Mythos Preview」と、それに対抗するための共同防御枠組み「Project Glasswing」が公開。

2026年4月22日



片山金融相による異例の会見。4月24日に自局、中経、3メガバンク、取引所が集まる初の官民会議を開催することを事前告知。

2026年4月24日



官民連携会議の開催と作業部会の設置。「インシデントへの備え」の重要性を共有し、実務レベルの検討を行う事務レベル作業部会が立ち上がり。

AIミユトスが引き起こす5つの深刻な懸念

ゼロデイ脆弱性探索の高速化

主要OSやブラウザの脆弱性をAIが複数回で発見・悪用可能になり、共通の基盤を使う金融機関へ同時に被害が波及する恐れがあります。

第三者依存・集中のリスク

特定のAIサービスやクラウドプロバイダーに依存することで、そこが「米一轉差点」となり、システム全体が停止するリスクが高まります。

AIそのものへの攻撃

プロンプトインジェクションやデータ漏えいにより、AIに接続された顧客情報や内部ルールが悪意ある者に奪われる危険があります。

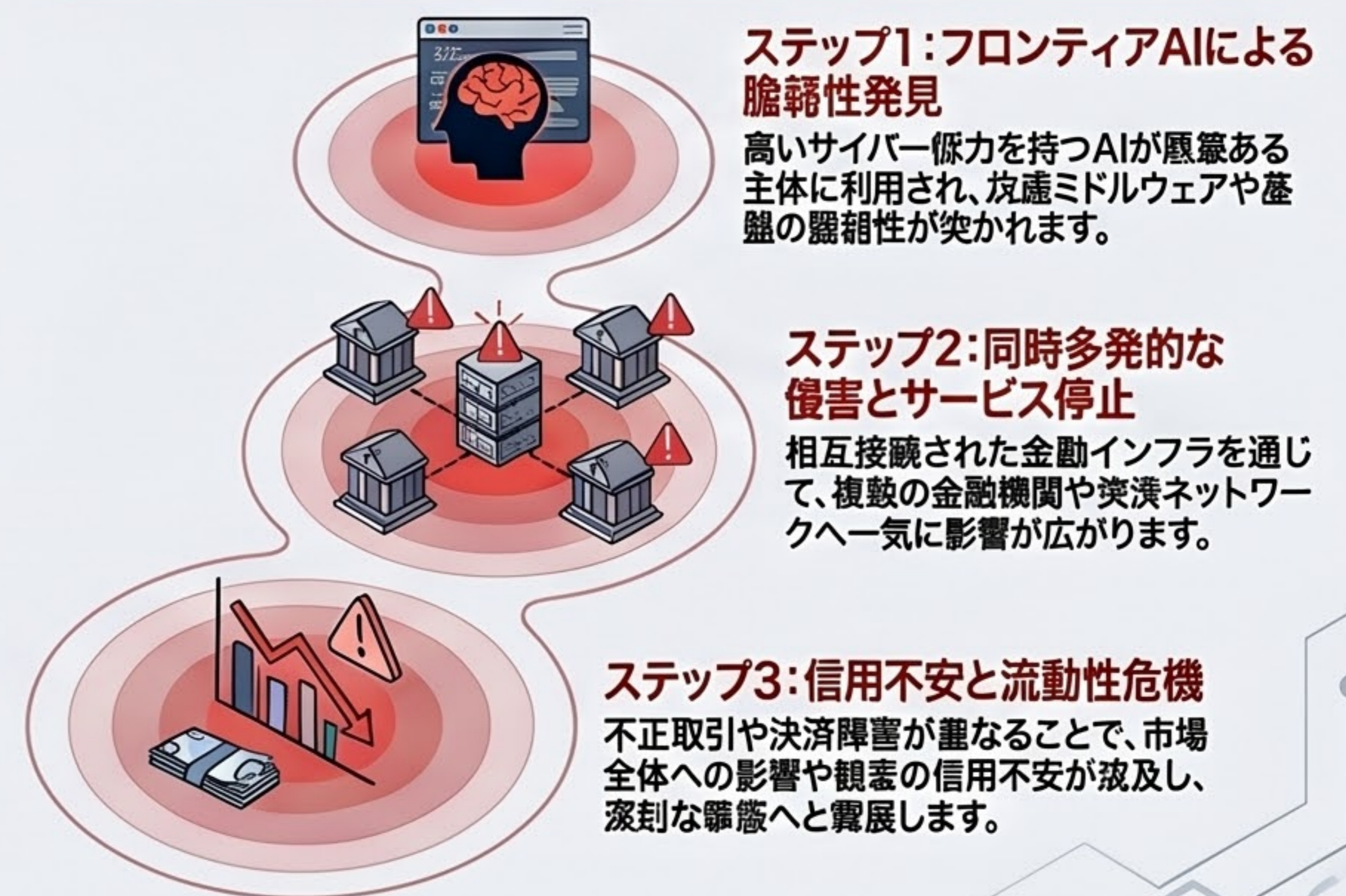
顧客接点・不正検知の逆用

攻撃者がAIを悪用して本人確認を突破したり、偽情報を混入させたりすることで、顧客の信頼が低下するシナリオが想定されます。

人材・監督の非対称性

AIの進化激速に駆逐のモニタリングや専門人材の確保が追いつかず、初動の遅れや誤判断を担う「後継者層」が割の機関で指摘されています。

リスク波及のメカニズム



海外の動向と日本への示唆

地域	中心的な特組み	特徴	日本への示唆
米国	AIEOG / 金融向けAI RMF	実装可能な用語集や管理フレームワークを先行確立	事故分類や報告テンプレートの共通化が必要
英国	AI Consortium / BoE監視	監督制と実地テスト(Live Testing)を重視	継続的なサーベイと安全な運用テストの仕組み
EU	DORA(デジタル運用レジリエンス法)	第三者リスク管理やレジリエンス演習を法制化	AI単独だけでなく、金融機関の実務ルールの積み増し

実務提言:今すぐ取り組むべき優先アクション

「日本版Project Glasswing」の制度化
厳格な管理下で閉鎖目的のAIを活用し、秘密保持を前提とした脆弱性診断や共同点検を行う体制を年間に設計する必要があります。

AI時代版「重要資産台帳」の整備
外部露出資産(ASM)や委託先、AI API接続先を一体で調査し、代替調達手帳まで含む台帳を直ちに作成すべきです。

AI事故シナリオに基づく共同演習
生成AIの停止やコード生成の汚染を想定した、セクター横断的なサイバー演習を3~6か月以内に実施することが推奨されます。

ガイドラインのAI固有リスクへの拡張
プロンプト管理やデータ導管、AI事故報告の町文化など、既存の監督指針をAI的視点で再設計することが急務です。