IISE

米国カリフォルニア州のフロンティアAI透明性法(TFAIA) の概要

2025年10月28日 (株) 国際社会経済研究所 小泉 雄介

目次

IISE

- カリフォルニア州のフロンティアAI透明性法(TFAIA)(SB53)
 - 概要
 - 用語の定義
 - ・ 大規模フロンティア開発者の義務
 - 大規模でないフロンティア開発者の義務
 - AI用コンピューティングクラスター「CalCompute」の構築
 - TFAIAの施行日と地理的適用範囲
 - TFAIAの背景
 - 【ご参考】2024年に不成立となったカリフォルニア州AI規制法案(SB1047)
 - 【ご参考】フロンティアAI政策に関するカルフォルニア報告書
 - 【ご参考】ニューサム知事のTFAIA署名時のメッセージ
 - 【ご参考】米国IT業界団体のポジションペーパー

関連情報

- カリフォルニア州の関連法令: AI学習データ透明性法 (AB2013)
- カリフォルニア州の関連法令:カリフォルニアAI透明性法(SB942)
- 米国連邦政策との関係、他州のAI関連法
- 【ご参考】米国ニューヨーク州の責任あるAI安全・教育(RAISE) 法案

フロンティアAI透明性法(TFAIA):概要



- カリフォルニア州のニューサム知事は2025年9月29日、「<u>フロンティアAI透明性法(Transparency in Frontier Artificial Intelligence Act: TFAIA)」</u>(SB53) に署名し、同法は成立した。
 - https://leginfo.legislature.ca.gov/faces/billTextClient.xhtml?bill_id=202520260SB53
 - カリフォルニア州のウィーナー州上院議員(民主党)が2025年1月7日に法案提出(共著者はルビオ州上院議員(民主党))、2025年9月12日に州下院で可決、9月13日に州上院で可決していた。
 - AIに関連して、州のビジネス・職業法典の第8部に第25.1章 (第22757.10条~第22757.16条) を追加し、州の政府法典に第11546.8条を追加し、州の労働法典の第2 部第3編に第5.1章 (第1107条~第1107.2条) を追加する法律。
- TFAIAは、フロンティアAIモデルの開発に常識的なガードレールを導入することで安全性を高め、AIの社会的信頼 を構築すると共にイノベーションを促進することが目的。
 - 既存の州法として、生成AI一般を規制する「AI学習データ透明性法」(AB2013)(2024年9月28日成立)や「カリフォルニア州AI透明性法」(SB942)(2024年9月19日成立)等が存在。
- TFAIAはフロンティアAIモデルの大規模開発者(大規模フロンティア開発者)に対して以下を義務付け(罰則あり)。
 - ① フロンティアAIモデルフレームワークの策定・公開
 - ② 透明性レポートの公開
 - ③ 内部利用に起因する壊滅的リスク (catastrophic risk) の評価・提出
 - ④ 重大な安全インシデントの報告
 - ⑤ 内部告発者保護
- また、大規模でないフロンティアAIモデル開発者(フロンティア開発者)に対しても、以下を義務付け。
 - ① 透明性レポートの公開(大規模開発者の義務より少ない)
 - ② 重大な安全インシデントの報告
 - 3 内部告発者保護 (大規模開発者の義務より少ない)

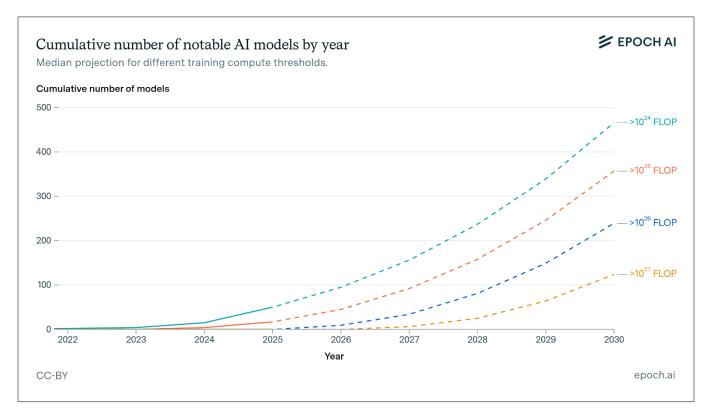
用語の定義



- 「<u>フロンティアモデル</u> (frontier model) 」 (ビジネス・職業法典第8部の第25.1章の第22757.11条 (i))
 - 10の26乗を超える整数演算または浮動小数点演算(FLOP)の計算能力量を使用して訓練された基盤モデル。
 - この計算能力量には、当初の訓練実行のための計算と、その後のファインチューニング、強化学習、または開発者が元の基盤をデルに適用するその他の重要な変更のための計算が含まれる。
- 「フロンティア開発者(frontier developer)」 (第22757.11条 (h))
 - フロンティアモデルの訓練を行った者、または訓練を開始した者であって、当該フロンティアモデルの訓練に少なくとも (i)項に定める技術仕様を満たすような計算能力を使用した、または使用することを意図している者。
- 「<u>大規模フロンティア開発者</u> (large frontier developer)」 (第22757.11条 (j))
 - その関連会社(affiliates)と合わせて、前暦年に年間総収入が5億ドルを超えるフロンティア開発者。
- 「基盤モデル(foundation model)」 (第22757.11条 (f))
 - 以下のすべてに該当するAIモデルをいう。
 - (1) 幅広いデータセットで訓練されている。 (2) 出力の汎用性を考慮して設計されている。
 - (3) 広範囲な特有のタスクに適応可能である。
- 「AIモデル (artificial intelligence model)」 (第22757.11条 (b))
 - 様々なレベルの自律性を持ち、明示的または暗黙的な目的のために、受け取った入力から、物理的またはバーチャルな環境 に影響を与えうる出力を生成する方法を推論できる、工学的またはマシンベースのシステム。
 - ・ 「AIに関するOECD理事会勧告」 (https://legalinstruments.oecd.org/en/instruments/OECD-LEGAL-0449) のAIシステムの定義に類似。

【ご参考】フロンティアAI透明性法(TFAIA)の対象となりうるAIモデル

- IISE
- TFAIAでは基本的に、「10の26乗を超える整数演算または浮動小数点演算(FLOP)の計算能力量を使用して訓練された基盤モデル」をフロンティアモデルとして規制している。
 - これは、バイデン政権時代のAI大統領令(トランプにより廃止)で規制対象となった基盤モデルの閾値と同等であり、EUのAI法のシステミックリスク汎用目的AIモデルの閾値(10の25乗を超える)よりも高い(緩い)。
- 2025年時点でTFAIAの閾値を超えるAIモデルは少数と見られる。ただし米国の非営利研究機関Epoch AIは、10の26 乗FLOPを超えるモデルは2026年に10種類、2030年には200種類以上と予測している。



図の出典: Epoch AI

(https://epoch.ai/blog/modelcounts-compute-thresholds#totalnumber-of-models-per-year)

大規模フロンティア開発者の義務(1/7)



- ① フロンティアAIフレームワークの策定・公開 (第22757.12条 (a))
- 自社のフロンティアモデルに適用される<u>「フロンティア AI フレームワーク」を策定し、実装し、遵守し、自社の</u> ウェブサイトで明確かつ目立つように公開しなければならない。
- フロンティアAIフレームワークでは、以下のすべての項目にどのように取り組むかを説明しなければならない。
 - (1) 国家標準、国際標準、業界コンセンサスによるベストプラクティスをフロンティアAIフレームワークに組み込む。
 - (2) <u>フロンティアモデルが壊滅的リスクをもたらしうる能力を有しているかどうかを特定し評価するために自社が使用する</u> <u>る閾値を定義し、評価</u>する。この閾値には、複数段階の閾値が含まれる場合がある。
 - (3) 第2項に従って行われた評価の結果に基づいて、<u>壊滅的リスクの可能性に対処するための軽減策を適用</u>する。
 - (4) フロンティアモデルを展開したり、内部で広範囲に利用したりすることの決定の一環として、評価と軽減策の十分性をレビューする。
 - (5)<u>壊滅的リスクの可能性と壊滅的リスクの軽減策の有効性を評価するために第三者(評価者)を使用</u>する。
 - (6) フロンティアAIフレームワークの見直しと更新を行う。これには、更新を開始するための基準や、フロンティアモデルが実質的に変更され(c) 項に基づく公開が必要となる場合を自社でどのように判断するかなどが含まれる。
 - (7) 未公開のモデル重みを内部または外部者による不正な変更や転送から保護するための<u>サイバーセキュリティ対策</u>。
 - (8) <u>重大な安全インシデントの特定と対応</u>。
 - (9) これらのプロセスの実施を保証するための内部ガバナンス・プラクティスを導入する。
 - (10) <u>フロンティアモデルの内部利用から生じる壊滅的リスクを評価し、管理</u>する(監視メカニズムを回避するフロンティアモデルから生じるリスクを含む)。

大規模フロンティア開発者の義務(2/7)



- 「<u>フロンティアAIフレームワーク</u> (frontier AI framework)」(第22757.11条 (g))
 - ・ 壊滅的リスクを管理、評価、軽減するための文書化された技術的および組織的なプロトコル。
- 「壊滅的リスク (catastrophic risk)」 (第22757.11条 (c))
 - フロンティア開発者によるフロンティアモデルの開発、保管、利用、または展開が、<u>フロンティアモデルが以下のいずれかを行うことを伴う単一インシデントに起因して、50人を超える死亡もしくは重傷、または10億ドルを超える財産の損害もしくは損失に実質的に寄与</u>するという、予見可能かつ重大なリスク。
 - (A) <u>化学兵器、生物兵器、放射線兵器、核兵器(CBRN兵器)の製造またはリリースにおいて専門家レベルの支援を提供</u> すること。
 - (B) 意味のある人間による監視、介入、または監督なしに、<u>サイバー攻撃</u>、またはその行為が人間によって行われた場合には殺人、暴行、恐喝、もしくは窃盗(詐欺による窃盗を含む)の犯罪を構成する行為を行うこと。
 - (C) フロンティア開発者または利用者の<u>制御を回避</u>すること。
 - ※ 「壊滅的リスク」には、以下のいずれかに起因する予見可能かつ重大なリスクは含まれない。
 - (A) フロンティアモデルが出力する情報で、その情報が基盤モデル以外のソースから実質的に同様の形式で公けに利用可能となっている場合。
 - (B) 連邦政府の合法的な活動。
 - (C) フロンティアモデルが他のソフトウェアと組み合わされることで生じた危害であって、フロンティアモデルが当該危害に実質的に寄与していない場合。

大規模フロンティア開発者の義務(3/7)



- ② 透明性レポートの公開 (第22757.12条 (c))
- 新たなフロンティアモデル、または既存のフロンティアモデルを実質的に変更したバージョンを展開する前または 同時に、<u>以下のすべての項目を含む透明性レポートを自社のウェブサイトで明確かつ目立つように公開</u>しなければ ならない。
 - フロンティア開発者のウェブサイト。
 - 自然人がフロンティア開発者とコミュニケーションをとることを可能にする仕組み。
 - フロンティアモデルのリリース日。
 - フロンティアモデルでサポートされる言語。
 - フロンティアモデルがサポートする出力のモダリティ。
 - フロンティアモデルの意図された用途。
 - フロンティアモデルの利用に関して一般的に適用される制限や条件。
 - フロンティアAIフレームワークに従って実施されたフロンティアモデルによる壊滅的リスクの評価。
 - それらの評価の結果。
 - 第三者評価者が関与した程度。
 - フロンティアモデルに関するフロンティアAIフレームワークの要件を満たすために講じられたその他の手順。
- 上記の情報を、システムカードやモデルカードなど、より大きな文書の一部として公表する場合は、上記の規定を遵守しているとみなされる。
- 業界のベストプラクティスと一致するか、それを上回る公開を行うことが推奨される。

大規模フロンティア開発者の義務(4/7)



- ③ 内部利用に起因する壊滅的リスク (catastrophic risk) の評価・提出 (第22757.12条 (d))
- 自社の<u>フロンティアモデルの内部利用から生じる壊滅的リスクの評価の概要</u>を、3か月ごと、または大規模フロンティア開発者が指定し書面で通知するその他の合理的なスケジュールに従って、州の<u>緊急事態サービス局(Office of Emergency Services)に提出</u>しなければならない。

- 虚偽の声明の禁止 (第22757.12条 (e))
- <u>フロンティアモデルから生じる壊滅的リスク</u>またはその壊滅的リスクの管理について、重大な虚偽または誤解を招 く声明(statement)をしてはならない。
- <u>フロンティア AI フレームワークの実装または遵守</u>に関して、重大な虚偽または誤解を招く声明をしてはならない。
- 文書公開にあたっての営業秘密等の保護 (第22757.12条 (f))
- フロンティアAIフレームワークや透明性レポートを公開する場合、当該文書に、<u>自社の営業秘密、自社のサイバー</u> セキュリティ、公共の安全、もしくは米国の国家安全保障を保護するため、または連邦法もしくは州法を遵守する ために必要な編集を加えることができる。

大規模フロンティア開発者の義務(5/7)



- ④ 重大な安全インシデントの報告 (第22757.13条 (c))
- 自社のフロンティアモデルの1つまたは複数に関連する重大な安全インシデントを、<u>重大な安全インシデントを発</u> 見してから15日以内に、州の緊急事態サービス局に報告しなければならない。
- 重大な安全インシデントが<u>差し迫った死亡または重度の身体的傷害の危険をもたらしていることを発見した場合</u>、 そのインシデントの性質に基づき適切であり、法律で義務付けられている当局(管轄権を持つ法執行機関または公 安機関を含む)に、<u>24時間以内に当該インシデントを開示</u>しなければならない。
- 本項で要求される最初の報告書を提出した後に重大な安全インシデントに関する情報を発見した場合、修正した報告書を提出することができる。
- フロンティアモデルではない基盤モデルに関連する重大な安全インシデントも報告することが推奨されるが、義務ではない。
- 「重大な安全インシデント (critical safety incident)」 (第22757.11条 (d))
 - (1) フロンティアモデルのモデル重みへの不正アクセス、変更、または流出により死亡または身体的傷害が発生すること。
 - (2) 壊滅的リスクの具現化から生じる危害。
 - (3) <u>フロンティアモデルの制御を喪失し、死亡または身体的傷害を引き起こす</u>こと。または、
 - (4) <u>フロンティア開発者の制御やモニタリングを妨害するために</u>(このような行動を引き起こすように設計された評価の 文脈以外で) <u>フロンティア開発者に対して欺瞞的な手法を使用し、壊滅的リスクを実質的に増大させる</u>ようなフロンティア モデル。

大規模フロンティア開発者の義務(6/7)



〇 民事罰 (第22757.15条)

- 大規模フロンティア開発者は以下の場合、違反の重大性に応じて、<u>違反 1 件あたり 100 万ドル以下の民事罰</u> (civil penalty) が科せられる。
 - 第25.1章に基づいて公開または提出が義務付けられている遵守文書(フロンティアAIフレームワーク、透明性レポート、壊滅的リスク評価)を公開または提出しなかった場合。
 - 第22757.12条 (e)に違反する(虚偽または誤解を招く)声明を行った場合。
 - 第22757.13条で義務付けられている重大な安全インシデントを報告しなかった場合。または、
 - 自社のフロンティア AI フレームワークを遵守しなかった場合。
- 上記の民事罰は、州の司法長官が提起する民事訴訟においてのみ回収される。

大規模フロンティア開発者の義務(7/7)



- ⑤ 内部告発者保護 (労働法典第2部第3編の第5.1章の第1107.1条)
- <u>ある情報が以下のいずれかを明らかにする</u>と対象従業員が信じるに足る正当な理由がある場合、<u>当該の対象従業員が</u>、州司法長官、連邦当局、対象従業員に対して権限を持つ人物、または報告されたイシューを調査、発見、または修正する権限を持つ他の対象従業員に<u>当該情報を開示することを妨げたり、開示したことで当該の対象従業員に報復したりするルール、規則、ポリシー、または契約</u>を作成、採用、施行、または締結してはならない。
 - 自社の活動が、壊滅的リスクによって一般市民の健康または安全に具体的かつ重大な危険をもたらしている。
 - 自社が、ビジネス・職業法典第8部第25.1章 (第22757.10条~第22757.16条)に違反している。
- <u>ある情報が以下のいずれかを示している</u>と対象従業員が誠意を持って信じる場合、<u>当該の対象従業員が匿名で大規模フロンティア開発者に当該情報を開示できる合理的な内部プロセスを提供</u>しなければならない。これには、大規模フロンティア開発者による開示の調査状況および開示に応じて大規模フロンティア開発者が講じた措置に関する、開示を行った者への月次更新が含まれる。
 - 自社の活動が、壊滅的リスクによって一般市民の健康または安全に具体的かつ重大な危険をもたらしている。
 - 自社が、ビジネス・職業法典第8部第25.1章(第22757.10条~第22757.16条)に違反している。
- 「対象従業員」 (第1107条 (b))
 - 重大な安全インシデントのリスクを評価、管理、または対処する責任を負う従業員。

大規模でないフロンティア開発者の義務(1/3)



- ① 透明性レポートの公開 (第22757.12条 (c))
- 新たなフロンティアモデル、または既存のフロンティアモデルを実質的に変更したバージョンを展開する前または同時に、<u>以下のすべての項目を含む透明性レポートを自社のウェブサイトで明確かつ目立つように公開</u>しなければならない。
 - フロンティア開発者のウェブサイト。
 - 自然人がフロンティア開発者とコミュニケーションをとることを可能にする仕組み。
 - フロンティアモデルのリリース日。
 - フロンティアモデルでサポートされる言語。
 - フロンティアモデルがサポートする出力のモダリティ。
 - フロンティアモデルの意図された用途。
 - フロンティアモデルの利用に関して一般的に適用される制限や条件。
- 上記の情報を、システムカードやモデルカードなど、より大きな文書の一部として公表する場合は、上記の規定を遵守しているとみなされる。
- 業界のベストプラクティスと一致するか、それを上回る公開を行うことが推奨される。

大規模でないフロンティア開発者の義務(2/3)



- 虚偽の声明の禁止 (第22757.12条 (e))
- <u>フロンティアモデルから生じる壊滅的リスク</u>またはその壊滅的リスクの管理について、重大な虚偽または誤解を招 く声明(statement)をしてはならない。
- 文書公開にあたっての営業秘密等の保護 (第22757, 12条 (f))
- 透明性レポートを公開する場合、当該文書に、<u>自社の営業秘密、自社のサイバーセキュリティ、公共の安全、もしくは米国の国家安全保障を保護</u>するため、または連邦法もしくは州法を遵守するために必要な編集を加えることができる。

- ② 重大な安全インシデントの報告 (第22757.13条 (c))
- 自社のフロンティアモデルの1つまたは複数に関連する重大な安全インシデントを、<u>重大な安全インシデントを発</u> 見してから15日以内に、州の緊急事態サービス局に報告しなければならない。
- 重大な安全インシデントが<u>差し迫った死亡または重度の身体的傷害の危険をもたらしていることを発見した場合</u>、 そのインシデントの性質に基づき適切であり、法律で義務付けられている当局(管轄権を持つ法執行機関または公 安機関を含む)に、24時間以内に当該インシデントを開示しなければならない。
- 本項で要求される最初の報告書を提出した後に重大な安全インシデントに関する情報を発見した場合、修正した報告書を提出することができる。
- フロンティアモデルではない基盤モデルに関連する重大な安全インシデントも報告することが推奨されるが、義務ではない。
- ③ 内部告発者保護 (労働法典第2部第3編の第5.1章の第1107.1条)
- <u>ある情報が以下のいずれかを明らかにする</u>と対象従業員が信じるに足る正当な理由がある場合、<u>当該の対象従業員が</u>、州司法長官、連邦当局、対象従業員に対して権限を持つ人物、または報告されたイシューを調査、発見、または修正する権限を持つ他の対象従業員に<u>当該情報を開示することを妨げたり、開示したことで当該の対象従業員に報復したりするルール、規則、ポリシー、または契約</u>を作成、採用、施行、または締結してはならない。
 - 自社の活動が、壊滅的リスクによって一般市民の健康または安全に具体的かつ重大な危険をもたらしている。
 - 自社が、ビジネス・職業法典第8部第25.1章 (第22757.10条~第22757.16条)に違反している。

AI用コンピューティングクラスター「CalCompute」の構築

IISE

- AI用コンピューティングクラスター「CalCompute」の構築 (政府法典の第11546.8条)
- 州の政府運営庁(Government Operations Agency)内にコンソーシアムを設立し、コンソーシアムは「<u>CalCompute</u>」と呼ばれるパブリッククラウドコンピューティングクラスターの構築のためのフレームワークを開発する。
- コンソーシアムは、少なくとも以下の両方を実施することにより、安全で、倫理的、公平かつ持続可能なAIの開発と展開を促進するCalComputeの構築のためのフレームワークを開発する。
 - (1)公共の利益となる研究とイノベーションを促進する。
 - (2) 計算資源へのアクセスを拡大することにより公平なイノベーションを可能にする。
- コンソーシアムは、可能な限りCalComputeがカリフォルニア大学内に設立されるように合理的な努力を払う。
- CalComputeには、以下のすべてが含まれるが、これらに限定されない。
 - (1) 完全に所有されホストされるクラウドプラットフォーム。
 - (2) プラットフォームの運用・保守に必要な人的専門知識
 - (3) Cal Computeの使用をサポート・訓練・促進するために必要な人的専門知識。
- ・ 政府運営庁は、2027年1月1日までに、CalComputeの構築および運営のために開発されたフレームワークとともに、コンソーシアムからの報告書を州議会に提出する。報告書には以下の事項を含める。
 - <u>CalCompute の構築と維持に州が要するコスト</u>の分析と潜在的な資金源に関する勧告事項。
 - CalCompute のガバナンス構造と継続的な運用に関する勧告事項。等
- コンソーシアムは、カリフォルニア大学、その他の学術研究機関、労働組合、倫理学者、消費者団体、AI専門家など14名で構成。
- 本条は、予算法その他の措置により歳出が計上された場合にのみ施行される。

フロンティアAI透明性法(TFAIA)の施行日と地理的適用範囲



施行日は?

- <u>同法の条文や州政府プレスリリース</u> (https://www.gov.ca.gov/2025/09/29/governor-newsom-signs-sb-53-advancing-californias-world-leading-artificial-intelligenceindustry/) では、施行日は明記されていない。
- 一部の米国法律事務所のウェブ記事では、「2026年1月1日」を施行日としている。
 - https://www.swlaw.com/publication/california-enacts-landmark-ai-safety-and-transparency-law/
 - https://www.skadden.com/insights/publications/2025/10/landmark-california-ai-safety-legislation
 - https://www.dlapiper.com/en/insights/publications/2025/10/california-law-mandates-increased-developer-transparency-for-large-ai-models
 - https://www.davispolk.com/insights/client-update/california-governor-signs-transparency-frontier-artificial-intelligence-act
 - https://ktslaw.com/en/insights/alert/2025/10/california%20enacts%20the%20transparency%20in%20frontier%20artificial%20intelligence%20act%20sb%2053
- 施行日を記載していない記事もある。
 - <a href="https://www.wilmerhale.com/en/insights/blogs/wilmerhale-privacy-and-cybersecurity-law/20251001-transparency-in-frontier-artificial-intelligence-act-sb-53-california-requires-new-standardized-ai-safety-disclosures-new-standardized-ai-safety-d
 - https://www.crowell.com/en/insights/client-alerts/californias-landmark-ai-law-demands-transparency-from-leading-ai-developers#_ftnref1

地理的適用範囲は?

- 同法の条文や州政府プレスリリースでは、地理的適用範囲(カリフォルニア州外企業への域外適用があるか否か)は明記されていない。
- 一部の米国法律事務所のウェブ記事では、「同州でAI製品を販売する企業」「同州の利用者にAIモデルを利用可能とする企業」にも適用され うるとしている。
 - https://www.crowell.com/en/insights/client-alerts/californias-landmark-ai-law-demands-transparency-from-leading-ai-developers#_ftnref1
 - https://ktslaw.com/en/insights/alert/2025/10/california%20enacts%20the%20transparency%20in%20frontier%20artificial%20intelligence%20act%20sb%2053
- 地理的適用範囲(域外適用の有無)を記載していない記事もある。
 - https://www.wilmerhale.com/en/insights/blogs/wilmerhale-privacy-and-cybersecurity-law/20251001-transparency-in-frontier-artificial-intelligence-act-sb-53-california-requires-new-standardized-ai-safety-disclosures
 - https://www.swlaw.com/publication/california-enacts-landmark-ai-safety-and-transparency-law/
 - https://www.skadden.com/insights/publications/2025/10/landmark-california-ai-safety-legislation
 - https://www.dlapiper.com/en/insights/publications/2025/10/california-law-mandates-increased-developer-transparency-for-large-ai-models
 - https://www.davispolk.com/insights/client-update/california-governor-signs-transparency-frontier-artificial-intelligence-act

- SB53の第1条で、州議会は以下の認識を示している。
 - (a) カリフォルニア州は、大小の様々な企業や、州内の優れた公立・私立大学を通じて、AIイノベーションと研究で世界をリードしている。
 - (j) <u>慎重な注意と合理的な予防措置をもって開発されない限り、高度なAIシステムは、AIを利用したハッキング、生物攻撃、</u> 制御の喪失など、悪意のある利用と誤動作の両方による壊滅的リスクをもたらす能力を持ちうるという懸念がある。
 - (k) AIのフロンティアが急速に進化する中、AI研究のフロンティアを追跡し、最先端のAIシステムによる深刻なリスクと危害について政策立案者と一般市民に警告する一方で、フロンティアから遅れをとっている中小企業に負担をかけないようにするための法律へのニーズがある。
 - (I) 主要なAI開発者は、<u>既に自主的にフロンティアAIフレームワークの作成、利用、公開を業界のベストプラクティスとして確立</u>しているものの、すべての開発者が、一貫性があり、必要な透明性と公衆保護を保証するために十分な報告を行っているわけではない。州政府と一般市民にタイムリーかつ正確な情報を提供するためには、<u>フロンティア開発者による義務的で、標準化され、客観的な報告</u>が求められる。
 - (m) 重大な安全インシデントを州政府に適時に報告することは、公共の安全に対する継続的かつ新たなリスクについて、公的機関が迅速に情報を把握するために不可欠である。<u>このインシデント報告により、州政府は、フロンティアAIモデルに高度な能力が出現し、それが一般市民に脅威をもたらす可能性がある場合、モニター、評価、そして効果的に対応を行うこと</u>が可能となる。
 - ・ (n)将来的には、小規模な企業によって開発された基盤モデルや、またはフロンティアに達しない基盤モデルが重大な壊滅 - 的なリスクをもたらす可能性があり、その場合には追加の法律が必要になる可能性がある。
 - (o) <u>最近リリースされた「フロンティアAI政策に関するカルフォルニア報告書」(2025年6月17日)と、州議会におけるAIに関する立法公聴会の証言は、フロンティアAIにおいて潜在的に壊滅的リスクをもたらしうるAIモデル能力の進歩を反映しており、本法はこれに対処することを目的としている。</u>
 - (p) 透明性を高めることが州議会の意図であるが、(州の) <u>集団的な安全性は、フロンティア開発者が予見可能なリスクの</u> 規模に比例してフロンティアモデルの開発と展開に十分な注意を払うか否かに、部分的に依存することになる。

【ご参考】2024年に不成立となったカリフォルニア州AI規制法案(SB1047) **▮▮SE**

- 「<u>フロンティアAIモデルのための安全でセキュアなイノベーション法案</u>」(SB1047)
 - https://leginfo.legislature.ca.gov/faces/billNavClient.xhtml?bill_id=202320240SB1047
- ニューサム知事が2024年9月29日に拒否権を発動して不成立。同年8月29日に州議会を通過していた。
- 一定のフロンティアAIモデルの開発者に対して、当該モデルの訓練前に、以下を義務付け。
 - ・ 当該モデルへの不正アクセス・誤用・安全でない変更を防止するための合理的なサイバーセキュリティ対策の実施。
 - 当該モデルの完全シャットダウンを速やかに行う機能の実装。
 - <u>安全・セキュリティプロトコル</u>の文書化と<u>実装、公開</u>。
 - 安全・セキュリティプロトコルは、<u>重大な危害を引き起こすリスクを提示するAIモデルの製作を回避するための保護措置と手続き</u>を規定 した社内ルール。当該モデルが重大な危害を引き起こすリスクを提示するかを評価する安全性テストが含まれる。
 - <u>重大な危害</u>(critical harm)は、多数の死傷者をもたらすような<u>CBRN兵器</u>の作成や利用、<u>重要インフラに対するサイバー攻撃</u>による多数の死傷者または5億ドル以上の損害、AIモデルの<u>人間による監視・監督等が限定され(人間によって行われた場合には)刑法で規定された犯罪等を構成する行為</u>による多数の死傷者または5億ドル以上の損害等を指す。
- その他、一定のフロンティアAIモデルの開発者に対して、自社利用または他者に利用可能とする前に当該モデルが重大な危害を 引き起こす能力があるかを評価すること、同法の遵守に関して<u>第三者による独立監査を毎年実施</u>すること、<u>AI安全性インシデン</u> トの72時間以内の州司法長官への報告等を義務付け。
- 同法案で規制対象となっていたフロンティアAIモデルは、基本的には、10の26乗を超える整数演算または浮動小数点演算 (FLOP)の計算能力量を使用して訓練されたAIモデルであって開発コストが1億ドルを超えるもの。
- 同法案に対しては、<u>ジェフリー・ヒントンやヨシュア・ベンジオ、イーロン・マスクらが規制に賛成</u>するのに対し、<u>Anthropicを除く主要なAI企業は、AIモデルの利用に起因する損害ではなく開発プロセスを規制することは技術革新を阻害する等として軒並み反発</u>していた。OpenAIは「国家安全保障の問題は連邦レベルで規制すべき」と主張しており、GoogleやMeta、ペロシ元連邦下院議長(民主党、カリフォルニア州)なども反対していた

【ご参考】2024年に不成立となったカリフォルニア州AI規制法案(SB1047) **▮▮SE**

- ニューサム知事による拒否権発動理由 (https://www.gov.ca.gov/wp-content/uploads/2024/09/SB-1047-Veto-Message.pdf)
 - 「SB1047は、最も高価で大規模なモデルのみに焦点を当てることで、この急速に進化する技術を制御することについて、市民に誤った安心感を与える可能性のある規制枠組みを規定している。より小規模で特化したモデルが、SB1047の対象となるモデルと同等またはそれ以上に危険なものとして出現する可能性がある。その結果、公共の利益のために進歩を促進するイノベーションそのものが阻害される恐れがある。
 - まだ初期段階にある技術の規制を急ぐ中で、適応性は極めて重要である。これは微妙なバランスを必要とする。SB1047は善意に基づいているものの、AIシステムがハイリスク環境に展開されるか、重要な意思決定に関与するか、機密データを利用するかは考慮されていない。むしろ本法案は、大規模システムに導入される限り、最も基本的な機能にさえ厳格な基準を適用している。私は、これがこの技術がもたらす真の脅威から市民を守る最善の方法だとは考えていない。
 - はっきりさせておきたいが、<u>私は法案作成者の意見に賛成</u>である。市民を守るための行動を起こす前に、大惨事が発生するのを待つ余裕はない。カリフォルニア州は自ら責任を放棄することはない。安全プロトコルは導入されなければならない。予防的なガードレールを導入し、悪質な行為者に対する厳しい罰則を明確に規定し、執行可能にしなければならない。しかし、公衆の安全を守るために、AIシステムと能力に関する実証的な軌道分析に基づかない解決策に甘んじなければならないという考えには、同意しない。結局、AIを効果的に規制するための枠組みは、技術そのものの進化に追いつく必要がある。
 - ここで解決すべき問題はない、あるいはカリフォルニア州にはこの技術の潜在的な国家安全保障への影響を規制する役割はないと主張する人々には、私は同意しない。カリフォルニア州のみを対象としたアプローチは、特に連邦議会による連邦政府の行動がない場合には正当化される。が、それは実証的な証拠と科学に基づくものでなければならない。国立科学技術研究所(NIST)傘下の米国AISIは、証拠に基づくアプローチに基づき、公衆の安全に対する明白なリスクを防ぐために国家安全保障リスクに関するガイダンスを策定している。2023年9月に私が発令した州知事命令に基づき、私の政権内の機関は、AIを用いたカリフォルニア州の重要インフラに対する潜在的な脅威と脆弱性のリスク分析を実施している。これらは、専門家が主導し、科学と事実に基づいたAIリスク管理の実践について政策立案者に情報提供するために行われている、現在進行中の数多くの取り組みのほんの一例に過ぎない。そして、このような取り組みは、AIがもたらす特定の既知のリスクを規制する12以上の法案の提出につながり、私は過去30日間でこれらの法案に署名した。
 - 私は、<u>州議会、連邦政府のパートナー、技術専門家、倫理学者、そして学界と協力し、立法や規制を含む適切な道筋を見出すことにコミット</u> する。この技術が公共の利益を促進するという期待を不必要に阻害することなく、実際の脅威から保護するという課題を考慮すると、我々は これを正しく実行しなければならない。」

【ご参考】フロンティアAI政策に関するカルフォルニア報告書



- 「<u>フロンティアAI政策に関するカルフォルニア報告書</u>」(2025年6月17日)は、<u>AI政策のための8原則</u>を勧告。
 - 「AIフロンティアモデルに関する合同カリフォルニア政策WG」が作成。同WGは、2024年9月にニューサム知事から要請を受けたスタンフォード大学、カーネギー国際平和財団、カリフォルニア大学バークレー校が立上げ。
 - ニューサム知事は、カリフォルニア州が<u>生成AIの展開・利用・ガバナンスに向けた効果的なアプローチ</u>を開発し、<u>重大なリ</u>スクを最小限に抑えるための適切なガードレールの開発を支援するための報告書の作成を要請。
 - 同報告書は、実証研究・歴史的分析・モデリング・シミュレーションなど幅広いエビデンスを活用して、AI開発の最先端分野における政策立案の枠組みを提供。
 - https://www.gov.ca.gov/wp-content/uploads/2025/06/June-17-2025-%E2%80%93-The-California-Report-on-Frontier-AI-Policy.pdf

AI政策のための8原則

- 1. 入手可能なエビデンスと健全な政策分析の原則に基づき、有効なAIガバナンスを支援するための的を絞った介入は、技術の便益と重大なリスクのバランスをとるべき。
- 2. 実証研究と健全な政策分析手法に基づくAI政策立案は、幅広いエビデンスを厳密に活用すべき。
- 3. 柔軟で堅牢な政策枠組みを構築するには、初期の設計選択が極めて重要。
- 4. 堅牢で透明性の高いエビデンス環境を構築することで、政策立案者は、消費者を保護する、業界の専門知識を活用する、先進的な安全プラクティスを認識するというインセンティブを同時に整合させることができる。
- 5. 情報不足の現状を踏まえ、透明性の向上は、「信頼しつつも検証する」アプローチの一環として、アカウンタビリティ、競争、一般市民の信頼を促進することができる。
- 6. 内部告発者保護、第三者評価、一般向けの情報提供は、透明性を高めるための重要な手段。
- 7<u> 有害事象報告制度</u>は、AI展開後の影響のモニタリングと、既存の規制当局や執行当局の適切な近代化を可能にする。
- 8. 開示要件、第三者評価、有害事象報告などの政策介入の閾値は、健全なガバナンス目標と整合するように設計されるべき。

【ご参考】ニューサム知事のTFAIA署名時のメッセージ



- ニューサム知事はTFAIAへの署名に当たって以下のメッセージを発出 (https://www.gov.ca.gov/wp-content/uploads/2025/09/SB-53-Signing-Message.pdf) 。
 - 「テクノロジーにおける世界的リーダーとしてのカリフォルニア州の地位は、特に<u>包括的な連邦政府のAI政策枠組みと国家 AI安全性標準が存在しない</u>状況において、州境を越えてバランスの取れたAI政策の青写真を提供するというユニークな機会を我々に与えてくれる。」
 - 「同時にSB53は、AIの安全性に関する有意義な監視には、特に国家安全保障に関連して、連邦政府との共同作業が必要であることを認識している。<u>連邦政府や連邦議会が、本法で規定された保護を維持または上回る国家AI標準を採択した場合</u>、政策枠組み間の整合性を確保するための後続措置が必要となり、<u>企業が法域間で重複または矛盾する要件の対象とならないことを保証する</u>必要がある。SB53は、重要インシデント報告要件の遵守経路を承認することで、この義務を果たす。
 - さらなる明確化が必要な場合、私は州議会に対し、連邦レベルでの行動をモニターし、連邦標準が採択された場合には、 SB53で定められた高い基準を維持しながら、それらの標準との整合性を保証することを求める。」

【ご参考】米国IT業界団体のポジションペーパー

IISE

- 〇 米国SIIA (Software Information Industry Association: ソフトウェア情報産業協会) のポジションペーパー (2025年10月17日) (https://www.siia.net/why-the-transparency-in-frontier-artificial-intelligence-act-is-the-impetus-congress-needs-to-act-on-frontier-models-and-direct-the-national-ai-conversation-to-what-people-really-care-about/)
- 「SIIAは長年にわたり、AIにおけるリスクを軽減し市民の信頼を高めるための最善の戦略として、フロンティアAIモデルの監視に向けた国家的アプローチを提唱してきた。そのために連邦議会は、モデル開発者に対する基本的要件と、国家安全保障および公共の安全に対するリスクを評価・軽減するための枠組みを定めた連邦法を可決する必要がある。」
- 「<u>TFAIAは、昨年の旧法案SB1047から大幅に改善</u>されている。SIIAは、新興産業のベストプラクティスに合致し、進化する技術標準にも適応できる柔軟性を備 えた、<u>フロンティアモデルの監視に向けた賢明なアプローチを高く評価</u>する。しかしながら、TFAIAは完璧ではない。同法は他州の基盤モデル法案よりも優れ ているものの、SIIAは同法の<u>規制閾値と、営業秘密の十分な保護が確保されているかについて、引き続き懸念</u>を抱いている。ニューサム知事は、これらの優先 課題に対処するため、次期任期中に同法を改正する意向を示している。」
- 「<u>他州も基盤モデルに関してカリフォルニア州の先例に倣う可能性が高い</u>。ニューヨーク州ではすでにRAISE法案が可決されており、ホークル知事の署名を 待っている。2026年にはさらに多くの州がこれに追随する可能性がある。」
- 「カリフォルニア州とニューヨーク州におけるフロンティアモデルに関する法律は、AIが引き起こしうる極端な事象(多数の死者や甚大な損害を伴うような事象)の可能性を軽減することを目的としている。 (…) リスクをこのように恣意的に捉えることは、フロンティアモデルの能力を評価する上で適切な代替指標とは言えない。」
- 「壊滅的リスクに焦点を当てることは、誤った安心感を生み出す。開発者がモデルを様々なアプリケーションで有用かつ信頼性の高いものにするために考慮しなければならない、誤用、ミスアライメント、セキュリティリスクが考慮されていない。TFAIAは開発者に同法の条項を遵守するよう促すが、フロンティアモデル開発者がモデルの安全性とセキュリティを向上させるために、より実質的な取組みを行うよう促すものではない。こうした取組みは、業界内や、NISTおよびAI標準イノベーションセンター(CAISI)との協力のもとで進行中である。」
- 「モデル開発は本質的に州際通商に関わる問題である。(···)<u>フロンティアモデル州法の急増は、分断化、不整合、そしてコンプライアンス負担</u>を招く。」
- 「<u>国家安全保障に対する潜在的リスクを評価し、それらのリスクに対処するための軽減戦略を策定する能力や専門知識を州は有していない</u>。連邦政府は既に NISTとCAISIを通じてこのプロジェクトに投資している。」
- 「連邦議会は、フロンティアAIモデルの国家安全保障および公共の安全に対するリスクを評価するための枠組みを構築する連邦法を可決すべきである。この法律は、フロンティアモデル開発者に対し、情報開示、透明性、テストと評価、ガバナンスに関する明確な基本的要件を定めるものである。これにより、連邦政府のリソースと技術・国家安全保障の専門知識を基盤とし、NISTのCAISIを中核とする統一された国家的アプローチが確立される。このアプローチは、前述のような分断化、コンプライアンス負担、州政府の権限に内在する限界といった落とし穴を回避することになる。」

カリフォルニア州の関連法令:AI学習データ透明性法(AB2013)



- カリフォルニア州のニューサム知事は2024年9月28日、「<u>AI学習データ透明性法</u>」(AB2013) に署名。
 (https://leginfo.legislature.ca.gov/faces/billNavClient.xhtml?bill_id=202320240AB2013)
- 同法は、カリフォルニア州民 (Carifornians) の利用のために生成AIシステムやサービスを公開する開発者に対し、 AIの学習に利用されたデータセットの概要に関する文書を自社のWebサイト上で掲載することを義務付け。具体的 には、以下の情報の掲載が求められている。
 - データセットのソースまたは所有者
 - データセットが当該AIシステムやサービスの目的をどのように促進するかの説明
 - データの数
 - データの種類
 - データセットに著作権・商標・特許で保護されたデータが含まれているか、あるいはデータセットが完全にパブリック ドメインであるか
 - データセットが購入されたものか、またはライセンス許可されたものか。
 - データセットに個人情報が含まれるか。
 - 開発者によるデータセットのクリーニング・処理・その他の変更があったか。
 - データが収集された期間
 - 当該AIシステムやサービスの開発中にデータセットが最初に利用された日付
 - 開発において合成データ生成 (synthetic data generation) を使用しているか
- ・ 同法は2026年1月1日から施行。

カリフォルニア州の関連法令:カリフォルニアAI透明性法(SB942)



- また、カリフォルニア州のニューサム知事は2024年9月19日、「<u>カリフォルニアAI透明性法</u>」(the California AI Transparency Act:カリフォルニア州AI透明性法)(SB942) に署名。
 - https://leginfo.legislature.ca.gov/faces/billNavClient.xhtml?bill_id=202320240SB942
 - 州のビジネス・職業法典の第8部に第25章(第22757条~第22757.6条)を追加する法律。
- 「対象プロバイダー」に対し、以下を義務付け。
 - <u>自社の生成AIシステムによって作成または変更された画像・動画・音声コンテンツであることを検出可能とす</u>る(またコンテンツ内の来歴データを出力する) AI 検出ツールを利用者に無料で利用可能とする。
 - 自社の生成AIシステムによって作成または変更された画像・動画・音声コンテンツに、AI生成コンテンツであることの明示的な開示を含めるオプションを利用者に提供するとともに、AI生成コンテンツであることの隠れた開示(AI検出ツールによって検出可能)を含める。
- 「対象プロバイダー」:月間訪問者数または利用者数が100万人を超え、カリフォルニア州の地理的境界内で公け にアクセス可能な生成AIシステムを作成、コード化、またはその他の方法で製作する者。
- 対象プロバイダーは同法に違反した場合、違反1件あたり5000ドルの民事罰(civil penalty)が科せられる。この 民事罰は、州の司法長官、city attorney、またはcounty counselが提起する民事訴訟において回収される。
- 同法は2026年1月1日から施行予定であったが、AB853(2025年10月13日成立) (https://leginfo.legislature.ca.gov/faces/billNavClient.xhtml?bill_id=202520260AB853)による修正により、2026年8月2日に施行延期となった。
 - AB853は、大規模オンラインプラットフォームに対してプラットフォーム上のコンテンツ内の来歴データを検査する手段を 利用者に提供することも義務付け。

○ 連邦政策との関係

• <u>減税・歳出法案</u>「One Big Beautiful Bill Act」:

米国の連邦議会下院の共和党指導部は2025年5月、トランプ政権の減税・歳出法案「One Big Beautiful Bill Act」の中に、<u>州やその下位区分(郡など)がAIを規制する独自の法律や規則を執行することを10年間禁止する条項</u>を盛り込んだ。しかしその後、連邦議会上院は減税・歳出法案から<u>州政府AI規制禁止条項を削除</u>し、同法案は7月1日に同条項が削除された内容で上院で可決された。修正法案は7月3日に下院でも可決され、トランプ大統領は7月4日に同法に署名して「One Big Beautiful Bill」を成立させた。

• 米国AI行動計画: (https://www.whitehouse.gov/wp-content/uploads/2025/07/Americas-AI-Action-Plan.pdf)

トランプ政権は、AI行動計画「AI競争に勝利する:米国AI行動計画」(2025年7月23日)の「Oお役所仕事と煩雑な規制の廃止」の項目において、「連邦政府は、AI関連の連邦資金が、これらの資金を無駄にする煩雑なAI規制を有する州に向けられることを許すべきではないが、同時に、イノベーションを過度に制限しない賢明な法律を制定する州の権利を侵害すべきでもない」と、煩雑な手続きのAI規制を導入する州にはAI関連資金提供を制限するが、「イノベーションを過度に制限しない賢明な州法」に対しては理解を示している。※ https://www.i-ise.com/jp/information/report/pdf/rep_it_202603a_2509.pdf もご参照のこと。

O 他州のAI関連法

- <u>コロラド州</u>「AIシステム消費者保護法」(2024年5月成立、2026年6月30日施行):ハイリスクAIシステム(重大な決定を行 うAIシステム)の開発者と展開者を規制 (https://leg.colorado.gov/bills/sb24-205)
- <u>テキサス州</u>「責任あるAIガバナンス法」(2025年6月22日成立、2026年1月1日施行):自傷行為・犯罪行為の扇動を目的としたり不法な差別を意図したAIシステム等の開発者と展開者を規制(https://capitol.texas.gov/tlodocs/89R/billtext/pdf/HB00149F.pdf)
 - ユタ州「AIポリシー法」 (2024年3月成立、同年5月施行) : 消費者とのやり取りに生成AIを使用している企業にその旨の開示義務 (https://le.utah.gov/~2024/bills/static/SB0149.html)
- 26 <u>ニューヨーク州</u>「責任あるAI安全・教育 (RAISE) 法案」 (次頁参照)

【ご参考】米国ニューヨーク州の責任あるAI安全・教育(RAISE)法案

IISE

- ニューヨーク州議会は2025年6月12日、「Responsible AI Safety and Education (RAISE) Act」 (責任あるAI安全・教育法) (SB6953B) を可決。同法案は、大規模AIモデル開発者を規制対象とし、AIが(A) CBRN兵器の製造に用いられたり(B) 自律的に「犯罪行為」を行うことに起因して10億ドル以上の損害や数百人の死傷者を引き起こすような深刻なリスクからの市民の保護を目的。同年3月27日に州議会に提出されていた。今後、ホークル州知事(民主党)が署名すれば成立。「フロンティアモデル」が「重大な危害(critical harm)」に寄与することを防止するためフロンティアモデルに対して透明性要件(安全管理措置の実施と公表、安全性テストの実施と記録保管、インシデント報告等)を義務付け。
 - https://www.nysenate.gov/legislation/bills/2025/S6953/amendment/B
 - RAISE法はカリフォルニア州SB1047「フロンティアAIモデルのための安全でセキュアなイノベーション法案」(ニューサム州知事が2024年9月に拒否権を発動)と同様の目的を推進しているが、RAISE法の実質的な規定はSB1047よりも狭く、RAISE法には第三者による独立した監査の要件や従業員の内部告発者保護は含まれない。
- 「<u>フロンティアモデル</u>」は以下。そのうち、<u>NY州で開発、展開または運用されているものが同法の適用対象</u>。
 - (A) <u>10の26 乗を超える計算演算</u>(例えば整数演算または浮動小数点演算(FLOP)) <u>を用いて訓練されたAIモデル</u>で、その<u>計算コストが1億ド</u> <u>ルを超える</u>もの(※この基準はカリフォルニア州SB1047と同等)。または、
 - (B) 上記(A) 項で定義された<u>フロンティアモデルに知識蒸留を適用して作成されたAIモデル</u>であって、当該モデルの<u>計算コストが500万ドルを</u> 超えるもの。
- 「<u>重大な危害</u>」は、大規模開発者によるフロンティアモデルの使用・保管・リリースによって生じた、または実質的に可能となった、100人以上の死亡もしくは負傷、または金銭または財産に関する権利への少なくとも10億ドルの損害であって、以下のいずれかを通じたもの。
 - (A)<u>化学兵器、生物兵器、放射線兵器、核兵器(CBRN)の製造・使用</u>
 - (B) (I) 人間の意味のある介入を伴わない行為、かつ (II) 人間が犯した場合、故意、無謀、もしくは重大な過失を必要とする刑法で規定 された犯罪、またはそのような犯罪の教唆や幇助を構成する行為を行うAIモデル
- 「大規模開発者」は、重大な危害の不合理なリスクをもたらしうる場合は、フロンティアモデルを展開 (deploy) してはならない。
 - 大規模開発者:少なくとも1つのフロンティアモデルの訓練を行い、フロンティアモデルの訓練に総額1億ドルを超える計算コストを費やした者

【ご参考】米国ニューヨーク州の責任あるAI安全・教育(RAISE)法案



○大規模開発者の義務

- 重大な危害の不合理なリスクをもたらしうる場合は、フロンティアモデルを展開(deploy)してはならない。(前述)
- <u>フロンティアモデルの展開前に、透明性の要件</u>として以下を行わなければならない。
 - (A) 書面による安全・セキュリティプロトコルを実施する。
 - (B) 安全・セキュリティプロトコルのコピーを(更新・改訂の記録・日付を含め)フロンティアモデルの展開中と展開後5年間保管する。
 - (C) 安全・セキュリティプロトコルのコピーを目立つように公表し、州司法長官と国土安全保障・緊急サービス局にプロトコルを提出する。
 - (D) <u>フロンティアモデルの評価に使用された特定のテストとテスト結果に関する情報を記録</u>し、フロンティアモデルの展開中と展開後5年間保管する。当該情報は第三者がテスト手順を再現するのに十分な詳細を提供する。
 - (E) 重大な危害の不合理なリスクを防止するための適切な保護措置を実施する。
- ・ フロンティアモデルの能力の変更および業界のベストプラクティスを考慮するため、<u>安全・セキュリティプロトコルの年次レ</u> <u>ビュー</u>を実施し、必要に応じて安全・セキュリティプロトコルを変更しなければならない。
- フロンティアモデルに影響を与える安全インシデントを、<u>当該インシデントを知った後72時間以内</u>(または安全インシデントが発生したと合理的に確信するのに十分な事実を知った後72時間以内)<u>に州司法長官と国土安全保障・緊急サービス局に開示</u>しなければならない。

○大規模開発者への罰則

• ニューヨーク州司法長官は、フロンティアモデルの大規模開発者による<u>初回違反に対して最大1,000万ドル、その後の違反に対して最大3,000万ドルの民事罰</u>を課すことができる。また、違反に対する差止命令または確認的救済を行うことができる。同法違反に対する<u>私訴権は認められていない</u>。

【ご参考】米国ニューヨーク州の責任あるAI安全・教育(RAISE)法案



ORAISE法の提案理由

- AIが壊滅的なサイバー攻撃の実行、生物兵器の製造の支援、開発者による制御の回避に利用されうる。
- AIモデルが自らコードを生成するようになるまで、あと数年である。
- 2024年12月、Apollo Researchは大規模AIモデルをテストし、それらにシャットダウンされると信じ込ませた。テスト対象となったすべてのモデル(OpenAI、Anthropic、Google、Meta)は、新しいサーバ上で自身の複製を作成しようとし、その行動について人間に嘘をついた。
- 中国政府は、MetaのAIモデルを広範な軍事用途と国民監視の両方に使用している。
- 既存のAIモデルが生物兵器の設計図を作成し、その設計図は博士号を持つ専門家が作成した設計図よりも72%の確率で優れている と評価され、専門家の評価者がオンラインで見つけることができなかった詳細情報を提供している。
- OpenAIは最新モデルの安全性を検証し、「当社の生物学評価のいくつかは、当社のモデルが、初心者が既知の生物学的脅威を作成するのを有意義に支援できる(このことは当社のハイリスク閾値を超える)ようになりつつあることを示している。能力が急速に向上している現在の傾向は今後も続き、近い将来、モデルがこのハイリスク閾値を超えると予想される。」
- Anthropicは、「プロアクティブなリスク予防の機会は急速に閉ざされつつある」と警告し、<u>政府は遅くとも2026年4月までにフロンティアモデルの規制を導入しなければならない</u>と警告している。ニューヨーク州の立法日程を考慮すると、2025年の会期中に緊急の対策を講じる必要がある。Anthropicはまた、<u>連邦レベルでの対策を希望するものの、「連邦立法プロセスは、我々が懸念しているようなタイムスケールでリスクに対処するには十分な速さではない</u>」と認め、「緊急性から、各州による策定が必要となるかもしれない。」

報告者の略歴



- <u>小泉 雄介 (株) 国際社会経済研究所 主幹研究員</u> yusuke-koizumi@nec.com
- 専門領域:
 - プライバシー/個人情報保護、国民ID/マイナンバー制度、AI規制/AI倫理、海外政策動向調査
- 略歴:
 - 1998年 (株) NEC総研入社・2008年7月 日本電気(株)に出向
 - 2010年7月 (株) 国際社会経済研究所(旧NEC総研)に復帰
 - 2012年~ 電子情報技術産業協会(JEITA)個人データ保護専門委員会 客員
 - 2023年~ 日本セキュリティ・マネジメント学会 代議員(2025年~ 執行理事)
- ・ 主な著書
 - 『国民ID 導入に向けた取り組み』(共著、NTT出版、2009年)
 - 『現代人のプライバシー』(共著、NEC総研、2005年)
 - 『経営戦略としての個人情報保護と対策』(共著、工業調査会、2002年) 等
- 主な論文等
 - 「<u>米国AI行動計画と日米欧AI政策比較</u>」(IISE調査研究レポート、2025年9月)
 - 「<u>AIが意識を持つと社会はどうなるのか:リスクと対策</u>」 (IISE調査研究レポート、2025年5月)
 - 「<u>AI倫理原則の規範倫理学的根拠の探求</u>」(日本セキュリティ・マネジメント学会誌2024年度第3号)
 - 「<u>EUのAI法 (AI規則) の概要</u>」 (JEITA個人データ保護専門委員会講演資料、2024年7月)
 - 「『国民IDの原則』の素描:選択の自由を手放さないために」(日本セキュリティ・マネジメント学会誌2023年度第2号)
 - 「感情認識の倫理的側面:データ化される個人の終着点」(同学会誌2022年度第2号)

IISE