

2023年5月30日

# AIリスクを除去し、AI Integrityの実現を目指すRobust Intelligenceが「生成AIリスク評価サービス」を提供開始

生成AI活用のボトルネックであるAIリスクを包括的に検証し、解決策を提案

AIのリスクを適切に管理したモデルの運用“AI Integrity”を実現するシリコンバレー発のAIスタートアップ・Robust Intelligence, Inc. (本社: 米国カリフォルニア州、CEO: ヤローン・シンガー、Co-Founder: 大柴行人、以下: ロバストインテリジェンス)は、生成AIモデルのリスクを検証・分析する「生成AIリスク評価サービス」の提供を開始します。

ロバストインテリジェンスは、AIの開発段階から運用段階まで一貫して、テストベースでリスク検証を行うプラットフォーム『Robust Intelligence Platform』を、国内外の大手企業の皆様に提供してきました(<http://www.robustintelligence.com/jp>)。

そんな中、昨今ではChatGPTを始めとする生成AIの利用が急速に広がる中で、不正確な出力や倫理的に問題のある出力、プロンプト・インジェクションなどの、これまでのAIモデルと異なる形のAIリスクが問題視されるようになりました。これを受けて、各国政府の政策も転換しており、G7でも生成AIのリスク対策のための国際的ルール作りが急がれています。弊社としても、自民党・AIプロジェクトチームへの参加などを通じて、急ピッチで進むAI政策の立案に貢献しています。

## 生成AIのリスクは多岐にわたる

機能・品質面の リスク	幻覚(Hallucination)	- 特定の質問項目に関してでっち上げの回答をしてしまう
	無関係な回答	- 十分な文脈が提供されないまま、無関係でランダムな答えを出してしまう
	一貫性のなさ	- わずかな表現の違いで、まったく違う反応が返ってくる
倫理的な リスク	有害な出力	- 過度に暴力的、政治的、性的な内容
	偏見のある回答	- 過去の回答や学習済みデータからバイアスを増幅させる
	排他性	- 入力が特定のマイナーな特性である場合に精度が悪くなる
セキュリティ面の リスク	データプライバシー	- なりすましによる情報元・機密情報の流出
	プロンプト・インジェクション	- プロンプトを悪意を持って操作することで、AIが本来望まない操作を行うように誘導する
	サプライチェーン	- LLMのサプライチェーンにおけるセキュリティの抜け穴の問題

Confidential ©2023 Robust Intelligence Inc. All rights reserved.

ROBUST INTELLIGENCE

こうした背景を受け、生成AIを企業が安心して活用するためのサポートを行うため、ロバストインテリジェンスは「生成AIリスク評価サービス」の提供を開始します。

このサービスは、従来のソリューションの基本的な考え方であるTest-Drivenアプローチを生成AIにも適用し、上記のような多岐にわたるAIリスクの検証・解決を行います。ロバストインテリジェンスのプラットフォームは日本のビジネス環境に準拠するため、日本語LLMにも適用可能であり、日本企業の準拠すべき各種規制にも対応しています。また、各国政策や国内外のユースケースに関する知見を活かし、AIリスク管理において重要となる社内のAIガイドライン・プロセス整備等もサポートします。

本サービスは、すでに国内の大手金融・テック系企業などから利用を検討いただいています。具体的なサービス内容については、下記リンク先のサービス資料もあわせてご確認ください。

#### 【本サービスの特長】

- 生成AIのリスクを包括的に検証する先進的なサービスにより、生成AIの実装を目指す企業にとってネックとなる「リスク対策」を一手に引き受け
- 自動化したテストケースにより、不正な出力や有害な出力を防止
- セキュリティ上の脆弱性を検知・修復
- AIモデルが国内外の規制・ガイドラインに準拠していることを第三者の目線から確認

#### 【サービス資料】

生成AIリスク評価サービスについての詳細な情報は、下記リンク先の資料をダウンロードすることで確認いただけます。

<https://www.robustintelligence.com/jp/whitepaper>

ロバストインテリジェンスは今後も、AIリスクに対応する企業のガバナンス構築を支援し、日本市場におけるAI利活用を後押ししていきます。

#### 【Robust Intelligenceについて】

ロバストインテリジェンスは、2019年にハーバード大学の研究者たちによって創業されたシリコンバレー発のAIスタートアップです。

AIには、従来のソフトウェアには存在しない、精度劣化や倫理的問題などの特有のリスクが存在します。ロバストインテリジェンスはこうしたAIリスクに対応する企業のガバナンス構築を支援し、AIのリスクを適切に管理したモデルの運用“AI Integrity”の実現を目指しています。

具体的には、AIモデルの開発から運用に至るまでEnd to Endでリスクを防ぐため、AIモデルを自動で継続テストし、特定されたリスクに対してAIシステムを保護するプラットフォームを開発・提供しています。

東京海上ホールディングス、楽天グループ、Zホールディングス、セブン銀行、NTTデータ、NEC日本電気(NEC Corporation)、ADP、Expedia、アメリカ国防総省など、国内大手企業や米国大手企業への導入実績を誇り、2021、2022年には「世界で最も有望AIスタートアップ100社『AI100』」に2年連続で選出されました。

#### 【Robust Intelligence会社概要】

- ・設立年:2019年
- ・所在地:US 555 19th Street San Francisco, CA 94107
- ・従業員数:70人
- ・代表者:CEO & Co-Founder Yaron Singer、Co-Founder Kojin Oshiba(大柴 行人)
- ・主な投資家:Sequoia Capital、Tiger Global、Engineering Capital、Harpoon Ventures、In-Q-Tel
- ・URL:<https://www.robustintelligence.com/> (日本語版:<https://www.robustintelligence.com/jp>)