

科学技術の潮流

JST研究開発戦略センター

290

AI（人工知能）技術の発展が目覚ましは「人間中心のAI社会原則」を発表し、高度な処理の自動化による生産性向上、新たな事業機会の可能性に期待が高まっている。一方で、AIの脆弱（ぜいじやく）性や悪用・誤用に起因するリスクが深刻化し、安全性確保のためのAIセキュリティ技術が不可欠になった。

AI（人工知能）技術の発展が目覚ましは「人間中心のAI社会原則」を発表し、高度な処理の自動化による生産性向上、新たな事業機会の可能性に期待が高まっている。一方で、AIの脆弱（ぜいじやく）性や悪用・誤用に起因するリスクが深刻化し、安全性確保のためのAIセキュリティ技術が不可欠になった。

AIセキュリティ強化を

原則から実践へ

CD)がまとめたAIの攻撃法など、AIの原則に42カ国が署名した。公平性、透明性、説明責任、プライバシーが強く要請されるようは、技術開発に裏付け

シムから19年に、産業技術総合研究所から20年に、それぞれガイドラインの初版が公開

安全性に焦点

AIについて、安全性を含む社会的要請が国際的に議論されるようになったのは201



科学技術振興機構(JST)研究開発戦略センター
フェロー(システム・情報科学技術ユニット) **福島 俊一**

東京大学理学部物理学科卒、IT企業にて自然言語処理・情報検索の研究開発に従事後、16年から現職。工学博士。11-13年東大大学院情報理工学研究所客員教授、情報処理学会フェロー。

保護など、AIに倫理性を求めた。22年に人間と区別困難な応答を返す生成AIが登場し、偽・誤情報生成AIを誤動作させ

報の生成・拡散による世論誘導や犯罪の巧妙化といった悪用問題、生成AIを誤動作させ

になり、23年に英国主権でAI安全サミット開催されたのを契機は、AIの安全性・信頼性に関する実践的方針に、安全性に関する評価手法や基準の検討・イドライン作りと技術整備への取り組みが、国際的にも早い時期に

され、開発現場での活用も広がっている。24年には経済産業省と総務省による「AI事業者ガイドライン」も公開された。これらは開

取り組みを促し、利用者をだましたり誘導したりする攻撃、悪意にも備えなければならぬ。そのような攻撃・悪用を見破って防御するAIセキュリティ技術の開発が急がれている。

「AI利用・運用時のリスク」



- (a) AI特有の脆弱性を突いて誤動作させたり情報搾取したりする
- (b) フェイク生成等AIを使って人をだましたり誘導したりする
- (c) 安全なAIだけでなく粗悪なAIや邪悪なAIが潜み利用者を欺く

【出典】研究開発戦略センター(CRDS)「人工知能研究の新潮流2025～基盤モデル・生成AIのインパクトと課題～」を基に筆者作成

AIはさらに発展し、自ら計画を立てて実行したり、他のAIと連携・交渉したりするようになりつつある。そこで生じる新たなリスクに先んじて対処するためにも、AIセキュリティ技術の一層の強化が急務である。(金曜日掲載)