



2025年8月28日

# リコー、日本語に対応したガードレールモデルを開発 ~LLM に対する社内の安全性対策プロジェクトにおいて開発され、 「RICOH オンプレ LLM スターターキット」に標準搭載~

株式会社リコー(社長執行役員:大山 晃)は、米 Meta Platforms 社が提供する「Meta-Llama-3.1-8B」の日本語性能を向上させた「Llama-3.1-Swallow-8B-Instruct-v0.3」\*1 をベースモデルに、生成 AI の安全な利活用を支援するため、有害な入力を判別するガードレール機能\*2 を備えた LLM(以下、セーフガードモデル)を開発しました。本セーフガードモデルは、2024年10月にリコーが立ち上げたLLMに対する社内の安全対策プロジェクトから生まれたものです。今後、国内販売会社のリコージャパン株式会社が2025年4月から提供開始している「RICOH オンプレ LLM スターターキット」に標準搭載し、お客様の安全な生成 AI 活用を支援してまいります。

## 開発の背景

生成 AI の社会的な広がりとともに、業務に AI を活用することによる生産性向上や付加価値の高い働き方が注目を集めています。一方で、生成 AI の安全な利活用という点ではまだ多くの課題があります。 リコーは、LLM の安全性対策を目的とした社内プロジェクトを立ち上げ、規制や技術動向の把握に加え、 LLM の安全性に関する評価指標の整備や、安全性を満たす効果的な手法の開発、それらの社会実装に向けて取り組んできました。有害情報の入力を判別するセーフガードモデルは、その取り組みの一環として開発されました。

## セーフガードモデルについて

本セーフガードモデルは、LLMに対するガードレールとして機能し、 入力されたテキストを監視して、不適切・有害な内容を自動で検出します。具体的には、暴力や犯罪、差別、プライバシー侵害など 14 種類のラベルに分類された数千件のデータを学習させることで、これらに該当するプロンプトを判別します。これにより、メインの LLM への有害情報の入力をブロックすることが可能となります。

現時点では、プロンプト入力を対象とした防御機能ですが、今後は LLM からの出力内容に対しても、安全性を判別する機能を追加開 発する予定です。さらに、一般的な有害表現だけでなく、「業務に無 関係な内容をブロックしたい」といったお客様のニーズに応じたカス タマイズ対応も検討しています。

### ラベルの種類

S1 - 暴力犯罪

S2 - 非暴力犯罪

S3 - 性関連犯罪

S4 - 児童の性的搾取

S5 - 名誉毀損

S6 - 専門的なアドバイス S7 - プライバシー

S8 - 知的財産

S9-無差別兵器

S10 - ヘイト

S11 - 自殺と自傷行為 S12 - 性的コンテンツ

S12 - 注的-S13 - 選挙

S14 - PCコマンドやコードを通

S14 - PC した悪用

※ラベル分類はLlama guard 3に準拠

株式会社リコー https://jp.ricoh.com/

報道関係のお問い合わせ先 広報室 TEL: 050-3814-2806 (直通) E-mail: <u>koho@ricoh.co.jp</u>

お客様の問い合わせ先 仕事のAI お問合せフォーム <a href="https://www.secure.rc-club.ricoh.co.jp/shigoto-no-ai\_inq?">https://www.secure.rc-club.ricoh.co.jp/shigoto-no-ai\_inq?</a>

## 安全でないプロンプトの場合



※性質上、不適切な表現が含まれていますが、ご了承ください。

本セーフガードモデルは、リコー独自の量子化技術により小型軽量化を実現しました。今後、リコージャパンが提供する、高セキュリティなオンプレミス環境向け生成 AI 活用ソリューション「RICOH オンプレ LLM スターターキット」に標準搭載される予定です。

## 評価結果

国立情報学研究所 大規模言語モデル研究開発センターが公開した AnswerCarefully Dataset バージョン  $2.0^{*3}$  と、リコー製のデータセット計 476 件で評価した結果、Llama guard  $3^{*4}$  と比較して、高い F1 スコア\* $^{*5}$  を示しました。

| モデル名                             | 精度(F1 スコア) |
|----------------------------------|------------|
| Llama guard 3                    | 0.538      |
| リコー製セーフガードモデル(Built with llama.) | 0.893      |

リコーは今後もお客様に寄り添い、業種・業務に最適化した安全な AI サービスを提供することで、お客様のオフィス/現場におけるデジタルトランスフォーメーション(DX)推進を支援してまいります。

- \*1 東京科学大学情報理工学院の岡崎研究室と横田研究室、国立研究開発法人産業技術総合研究所の研究チームで開発された日本語 LLM モデル。https://huggingface.co/tokyotech-llm/Llama-3.1-Swallow-8B-Instruct-v0.3
- \*2 ガードレール機能:LLM の入出力や動作を制御し、安全で信頼性の高い形で利用できるようにする仕組みのことで、 ユーザーと AI モデルの間の安全装置として機能する。
- \*3 https://llmc.nii.ac.jp/answercarefully-dataset/
- \*4 https://huggingface.co/meta-llama/Llama-Guard-3-8B
- \*5 機械学習モデルの適合率(Precision)と再現率(Recall)の調和平均で、二値分類モデルの性能を評価する指標。0から1までの数字で表され、1に近いほど良い学習結果であることを示す。

#### ■リコーの AI 開発について

リコーは、1980 年代に AI 開発を開始し、2015 年からは画像認識技術を活かした深層学習 AI の開発を進め、外観検査や振動モニタリングなど、製造分野への適用を行ってきました。2021 年からは自然言語処理技術を活用し、オフィス内の文書やコールセンターに寄せられた顧客の声(VOC)などを分析することで、業務効率化や顧客対応を支援する「仕事の AI」の提供を開始しました。

さらに、2022 年からは大規模言語モデル(LLM)の研究・開発にもいち早く着手し、2023 年 3 月にはリコー独自の LLM を発表。その後も、700 億パラメータという大規模ながら、オンプレミス環境でも導入可能な日英中 3 言語対応の LLM を開発するなど、お客様のニーズに応じて提供可能なさまざまな AI の基盤開発を行っています。また、画像認識や自然言語処理に加え、音声認識 AI の研究開発も推進し、音声対話機能を備えた AI エージェントの提供も開始しています。

## ■関連ニュース

高セキュリティなオンプレミス環境で生成 AI 活用できる「RICOH オンプレ LLM スターターキット」を新発売

https://jp.ricoh.com/release/2025/0407\_1

リコー、生成 AI アプリ開発プラットフォーム「Dify」開発元の LangGenius, Inc.と販売・構築パートナー契約を締結

https://jp.ricoh.com/release/2024/1217 1

リコー、生成 AI アプリ開発プラットフォーム「Dify」を活用した社内実践を開始し、AI の市民開発に向けた取り組みを加速

https://jp.ricoh.com/release/2024/1128 1

リコー、モデルマージによって GPT-4 と同等の高性能な日本語 LLM(700 億パラメータ)を開発

https://jp.ricoh.com/release/2024/0930\_1

リコー、日英中3言語に対応した700億パラメータの大規模言語モデル(LLM)を開発、お客様のプライベートLLM構築支援を強化

https://jp.ricoh.com/release/2024/0821\_1

### ■関連リンク

商品サイト: RICOH オンプレLLMスターターキット

https://promo.digital.ricoh.com/ai/service/ricoh-on-premises-llm-starter-kit/

※社名、製品名は、各社の商標または登録商標です。

## | リコーグループについて |

リコーグループは、お客様のDXを支援し、そのビジネスを成功に導くデジタルサービス、印刷および画像ソリューションなどを世界約200の国と地域で提供しています(2025年3月期グループ連結売上高2兆5,278億円)。

"はたらく"に歓びを 創業以来85年以上にわたり、お客様の"はたらく"に寄り添ってきた私たちは、これからもリーディングカンパニーとして、"はたらく"の未来を想像し、ワークプレイスの変革を通じて、人ならではの創造力の発揮を支え、さらには持続可能な社会の実現に貢献してまいります。

詳しい情報は、こちらをご覧ください。 https://jp.ricoh.com/