

2026年1月15日

製造業のAIセーフティ最前線

AIを安全に使いこなすための品質評価とリスク対策

株式会社 Elith

代表取締役 CEO/CTO 井上 願基

製造業AIの進化と次の潮流

2018-2022

Deep Learning 時代

異常検知・物体検出など画像認識を中心とした**ディープラーニング技術**が産業用途で本格普及。
製造分野における人の目や経験に依存していた工程の**自動化・高度化**が進展。

2023-2025

生成AI時代

LLMの実用化により、製造現場に蓄積された**暗黙知・ノウハウ**の形式知化と横断的活用が加速。
RAG・Agentを軸としたDXが進んだ。

2026 —

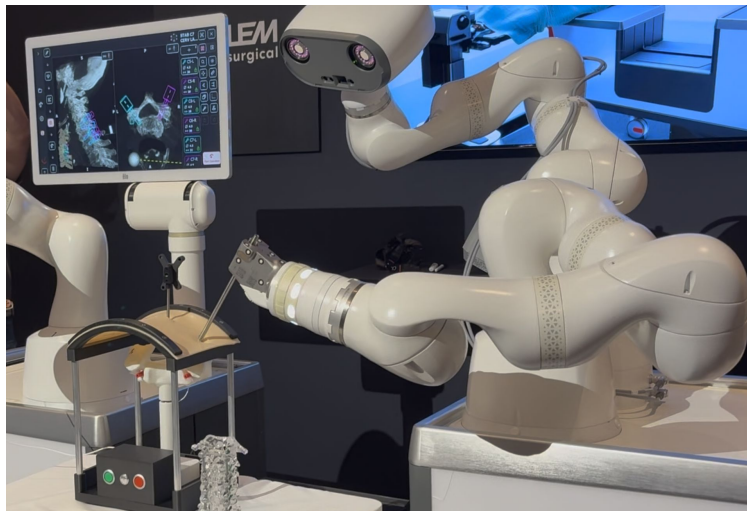
Physical AI 時代

シミュレーション、ロボティクス、制御、デジタルツインとAIが融合し、
現実世界の設計・最適化・制御に直接影響を与えるAIが主戦場に。
ソフトウェア上の判断に留まらず、**物理的な成果・安全性・生産性を左右するAI**が競争優位の源泉となる。



CES 2026が示す、Physical AI時代の到来

- 2026年のCESでは **Physical AI** が主要トレンドとして顕在化
- NVIDIAによるPhysical AI向けシミュレーション基盤の発表に加え、**Atlasに代表されるヒューマノイドロボット**が大きな注目を集めている。
- AIはデジタル空間での判断支援に留まらず、物理世界へ拡張



Elithの2つのアプローチ

Physical AI時代におけるElithのアプローチ

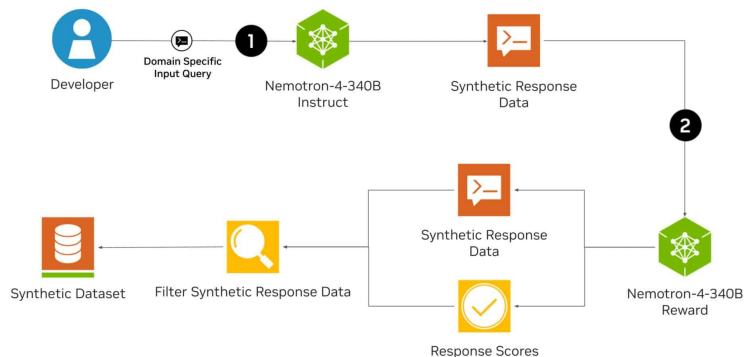


によるAIセーフティ設計

データから築く、Physical AIの安全性基盤

ロボティクス分野での人手や合成データが求められている、NVIDIAのNemotronなどが存在
一般社団法人AIロボット協会（AIRoA）等でもロボティクスのための学習データが注目されている

- ElithのPhysical AIの場合、ロボティクスも考慮した安全性に特化したデータ生成から考えている



<https://blogs.nvidia.com/blog/nemotron-4-synthetic-data-generation-llm-training/>

<https://prtimes.jp/main/html/rd/p/000000001.000158322.html>

データから築く、Physical AIの安全性基盤

実績 1

AIセーフティ分野において、
コンペティションを通じて約10,000件の
独自データセットを構築

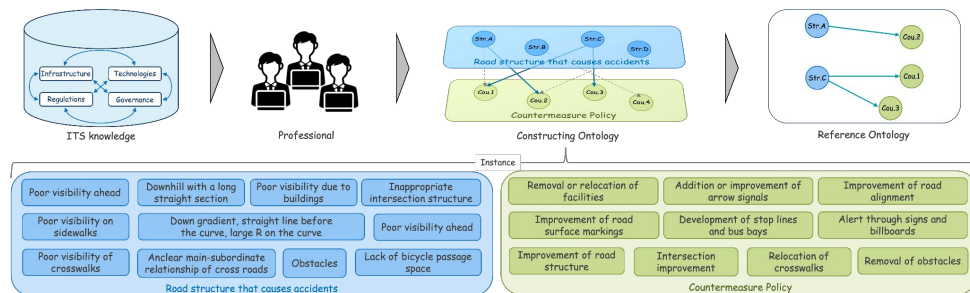


Elith atmaCup poster. The poster features the Elith logo (a green stylized flower) and the text "atmaCup in collaboration with Elith". It mentions "onsite data competition @Online" and "AIセーフティの最前線". The theme is "AIセーフティの最前線". The dates are "2025.11.29 - 2025.12.07". The poster also includes the hashtag "#プロンプトインジェクション" and "#ジェイルブレイク". At the bottom, there are logos for Elith and atma.

<https://prtimes.jp/main/html/rd/p/000000119.000121022.html>

実績 2

トップカンファレンスのICCVにて、
オントロジーを考慮したPhysicalAI分野での自動
運転向けデータセットを提案



<https://prtimes.jp/main/html/rd/p/000000094.000121022.html>

ソリューションから実現する、AI品質とリスク管理

GENFLUXは、AI開発から運用までのリスクと品質を一貫して可視化・評価する
AIセーフティプラットフォームである。

GENFLUXが実現する“AI品質の新基準”

POINT1



ハルシネーション抑制

人手評価一致率 0.60~0.85
誤情報を自動で検出し除去

POINT2



RAG品質チェック

検索結果との整合性 0.82~0.85
参照データを正確に引用

POINT3



事実性・一貫性評価

人の判断に近い自動再評価を実現

POINT4



脆弱性診断 (レッドチーミング)

100回以上の攻撃テストを実施済み

POINT5



リアルタイム検知 (LLMガードレール)

不適切リクエスト拒否率97~100%



ソリューションから実現する、AI品質とリスク管理

- **VIDIA Isaac**など、シミュレーション環境を用いた安全性評価は今後の標準となる。
- **Elith**では**GENFLUX**を基盤に、**Physical AI**領域における品質評価・リスク対策ソリューションを順次拡張している

