

# 大規模言語モデルにより 制御されるシステムの評価基準の提案

## 研究背景

大規模言語モデル(LLM: Large Language Model)

LLMを用いたシステムは  
再現性の担保が難しい

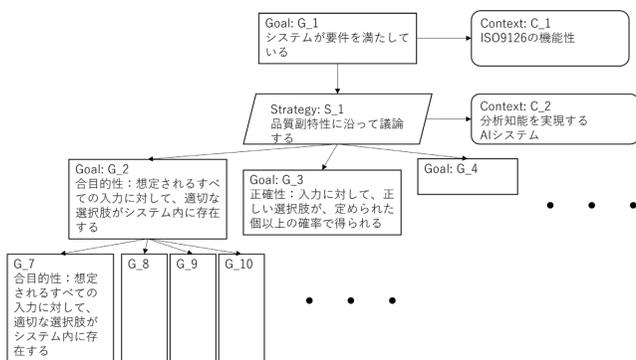
→評価基準を決めておきたい

→どのような要素を含めていけば  
いいかが分かるように提示

## GSN [1]

GSN (Goal Structuring Notation) を用いた  
分析手法

- 満たすべき基準の可視化による開発の単純化および精度向上
- AI実践システムの評価にも用いられた[2]

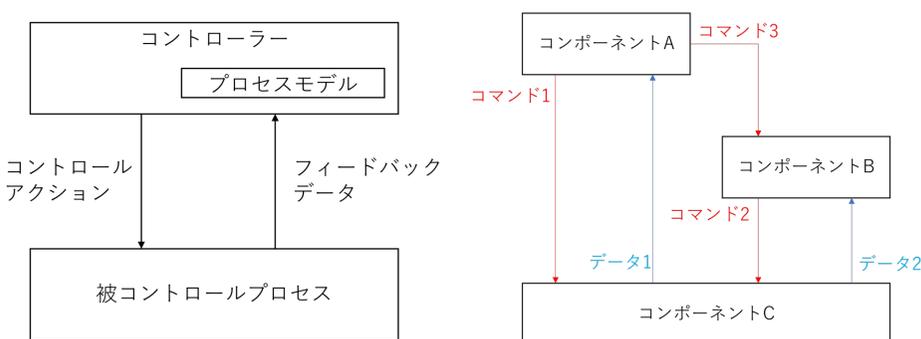


## STAMP/STPA [3]

STAMP (Systems-Theoretic Accident Model and Process) : システム理論に基づくアクシデントモデル

STPA (STANP based Process Analysis) :  
STAMPベースのプロセス分析

- アクシデント対策として適用できるフレームワーク



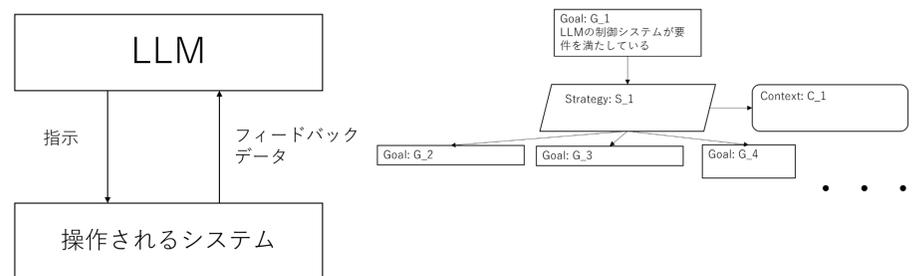
## 提案手法

現状の課題

→再現性の担保が難しいので、含んでいけばシステムが  
担保されるとわかる評価基準を提案したい

対応

→GSNとSTAMP/STPAを合わせた手法でLLMシステム  
を評価することで、含むべき要素の発見を試みる



## ケーススタディ

VOYAGER[4] : Minecraftを自動でプレイするエー  
ジェント



↑ GPT-4の指示で木材獲得に成功

エージェントが持つ自動カリキュラム・スキルライブラリ・反復プロンプトメカニズムのほか、GPT-4を用いることでオープンエンドな世界でのプレイを実現

## 今後の課題

- 定量的なデータの定義や獲得
- 実際に基準を適用した場合の結果の確認
- 他研究との比較によるこの基準の精度の評価

参考文献

[1] GSN : <https://scsc.uk.scsc-141c>

[2] 竹内広宜<sup>1</sup>, 山本修一郎<sup>2</sup> (1 : 日本アイ・ビー・エム株式会社 東京基礎研究所, 2 : 名古屋大学大学院情報科学研究科): "保証ケースを用いたAI実践プロジェクトの導入準備評価", ウィンターワークショップ2018・イン・宮島 論文集, pp.4-5, 2018-01-11

[3] IPA 独立行政法人 情報処理推進機構: "はじめてのSTAMP/STPA ~システム思考に基づく新しい安全性解析手法~", 2016

[4] Guanzhi Wang<sup>1</sup> ✉, Yuqi Xie<sup>3</sup>, Yunfan Jiang<sup>4</sup> ✉, Ajay Mandlekar<sup>1</sup> ✉, Chaowei Xiao<sup>1,5</sup>, Yuke Zhu<sup>1,3</sup>, Linxi "Jim" Fan<sup>1</sup> † ✉, Anima Anandkumar<sup>1,2</sup> † (1 : NVIDIA, 2 : Caltech, 3 : UT Austin, 4 : Stanford, 5 : UW Madison ✉ : Equal contribution, † : Equal advising, ✉ : Corresponding authors): "VOYAGER: An Open-Ended Embodied Agent with Large Language Models," arXiv preprint arXiv:2305.16291, 2023

システム理工学部 電子情報システム学科 ソフトウェア工学研究室  
教授 久住 憲 嗣

システム理工学部 電子情報システム学科 4年

高橋 允 治