

LLM を用いた不具合報告文からの形式知抽出

西納修一† 馬場達也† 鳥羽忠信†

日立製作所†

1. はじめに

製品の保守の分野では過去の不具合報告を参照して不具合対応が行われる。しかしながら不具合報告は自由記述のテキスト形式であることが多く、記入者による書き方のばらつきが大きい。このため 1 件 1 件の報告文について不具合事象を把握し、対応策を作成するのに手間がかかる課題がある。この課題に対し、本研究では、報告文に記載された不具合に関する知識を形式知化し、可視化・検索できるシステムを検討する。

最初に知識を表現するためのデータ構造とシステムの概要を述べる。続いて不具合報告文を形式知化するための自然言語処理の技術について述べる。

2. 提案システムの概要

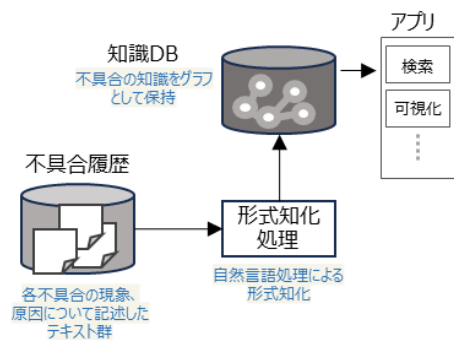


図 1 : 提案システム

図 1 に提案システムの概略を示す。不具合履歴に含まれる自由記述の不具合報告文に対して、形式知化の処理を行い、グラフ形式に変換し、知識 DB へ登録する。知識 DB にグラフ DB を用いて、不具合対策や設計改善をサポートする検索や可視化が可能となる。

不具合の知識を表現するためのグラフ形式について述べる。まず、機械系や電気・電子系などの不具合分析で広く用いられる FTA (Fault Tree Analysis) ・ FMEA (Failure Mode and Effects Analysis) の考え方を参考に、不具合

に関するナレッジの基本要素を不具合の発生部位・状態・因果関係と定義した。この観点で自動車の電子システムの不具合をグラフモデルとした表現した例を図 2 に示す。発生部位 (component) ・状態 (status) はグラフのノードとして表現する。また、箇所と状態の対応 (is 関係)、状態の間の因果関係 (caused_by 関係) をグラフの有向エッジとして表現している。これは不具合に関するナレッジを表現するための基本モデルであり、応用に合わせた拡張が可能である。

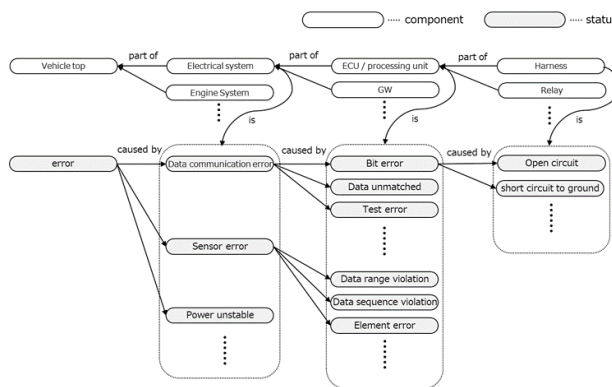


図 2 : グラフ形式

不具合報告文から上記のようなグラフ構造へと形式知化するための自然言語処理の技術について次節で述べる。

3. 不具合報告文からの形式知抽出

不具合報告文に記載された不具合事象から、関連する現象を発生部位・状態の観点で列挙し、それらの因果関係に基づいて整理することを考える。ここでは、この処理を固有表現抽出と関係性抽出の 2 ステップに分けて検討する。

まず不具合報告文に対して、固有表現抽出 (Named Entity Recognition, NER) [1] を行う。固有表現抽出は、テキストから固有名詞や固有の情報を読み出し分類する処理である。ここでは、前述のグラフモデルにおけるノード (不具合の発生部位・状態) に対応する表現のペアを不具合報告文から抽出するのに用いる。

続いて、得られた発生部位・状態に対して、再び不具合報告文を参照し関係性抽出

Extracting formal knowledge from trouble report using LLM
†SHUICHI NISHINO, TATSUYA BABA, TADANOBU TOBA, Hitachi td.

(Relation Extraction, RE) [1]を行う。関係性抽出は、テキストから意味的な関係を識別する処理である。ここでは得られた箇所・現象のノード間のエッジ (caused_by 関係) を推定するのに用いる。

上記の二つの処理 (NER および RE) は大規模言語モデル (Large Language Models, LLM) を用いて実現する。チャット型モデルを用いて、まず不具合報告文に対する固有表現抽出を実行する命令 (プロンプト文) を LLM に入力し、続いて LLM より得られた結果に対し関係性抽出を実行するプロンプト文を入力することで、不具合報告文から発生部位・状態・因果関係の情報を抽出する。

図 3 に実行例を示す。自動車の電子システムの不具合報告を模したテキストデータを入力とした。報告文の特徴として、製品に関連する専門的な略語 (例えば PU = Processor Unit) が含まれる。LLM は 2023 年 10 月時点で最も性能が高いモデルである ChatGPT4 (OpenAI 社) を用いた。出力結果として、抽出された不具合に関連する現象 (発生部位および状態のペア) をボックスで、抽出された因果関係をボックス間の矢印 (結果→原因) で示している。この例では、文章に含まれる不具合に関連する現象 (発生部位および状態) の 3 つ (Relay における Does-not-energize-or-de-energize、PU における ElementError、Vehicle Eesystem における power failure) 全てを正しく固有表現抽出できている。一方で関係性抽出の結果には誤りがあり、プロセッサ (PU) における ElementError が自動車の電子システム (Vehicle EEsysteem) における power failure へ波及するという因果関係を正しく抽出できていない。これは LLM が専門用語を含む文章を正しく処理できていないことを示す。

上記の課題に対して、専門用語の意味的な関係を LLM に与えることで性能を改善する方式を検討した。意味的な関係として、不具合の因果関係と密接な情報である装置構成情報を活用する (図 4)。装置構成情報は部位の階層関係を表現したもので、装置の設計情報などから作成することができる。一般に不具合事象は下層の部位から上位の部位に波及するため、因果関係と構成情報は高い相関を持つ。

関係性抽出処理のプロンプトに装置構成情報を含めた場合の入出力の例を図 5 に示す。プロセッサ (PU) における ElementError が自動車の電子システム (Vehicle EEsysteem) での power failure へ波及するという因果関係を正しく抽出

できている。

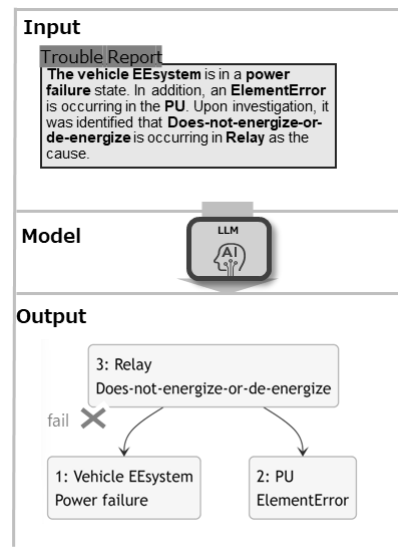


図 3 : 実行例 1

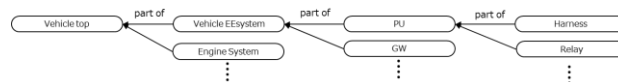


図 4 : 装置構成情報

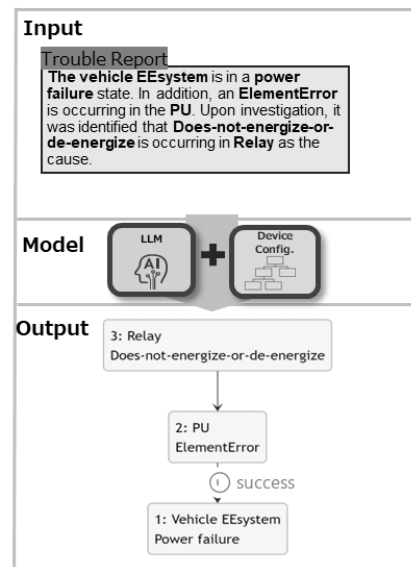


図 5 : 実行例 2

4. おわりに

本研究では報告文に記載された不具合に関する知識を形式知化するシステムを検討した。LLM を用いた形式化に構成情報を活用することで因果関係の抽出性能が向上することを確認した。

参考文献

[1] 岩倉友哉, 関根聡, 情報抽出・固有表現抽出のための基礎知識, 近代科学社, 2020