

大規模言語モデル (LLM) マルチエージェントの開発 エンジニアのための真のアシスタントを目指して

佐藤亨太[†], 佐橋功一^{*}, 高橋憲司[‡], 高橋雅裕[§]

[†]富士通株式会社, ^{*}株式会社 NTT データグループ, [‡]株式会社 NTT データアイ, [§]株式会社デンソー

開発における問題点

ソフトウェア開発の現場へのLLM（大規模言語モデル）活用が進んでいる。しかし、コード生成タスクをLLMに依頼した際、生成されたコードに正当性がない場合がある

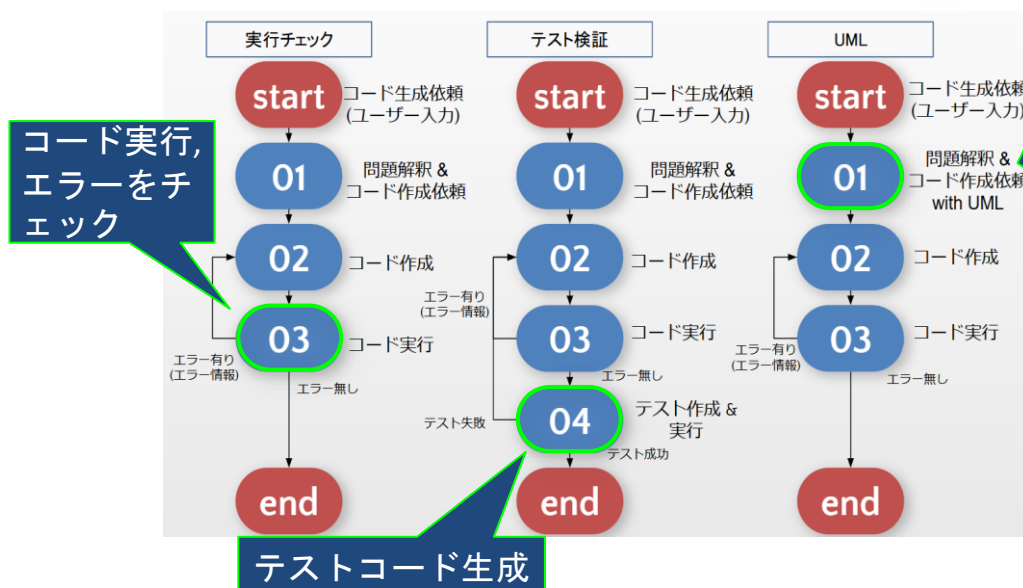
- 実行するとエラーが発生
- 要件を満たしていない

手法の適用による解決

- エラーが発生⇒生成されたコードを実行チェックする
- 要件を満たされていない⇒テスト検証を行い、要件を満たすかチェックする
- 扱う問題が難しい⇒モデリング言語を使って問題を構造化する

提案アーキテクチャ

LangGraph フレームワークを用いて3種類の LLM マルチエージェントアプリケーションを作成(下図)。以下では特徴的なエージェントを ○ で表示。



アーキテクチャ評価

AtCoderの問題をMarkdown形式で読み込み, LLMでコードを生成した. 生成したコードをAtCoderに提出することで得られた採点結果から, 以下の評価指標を算出した。

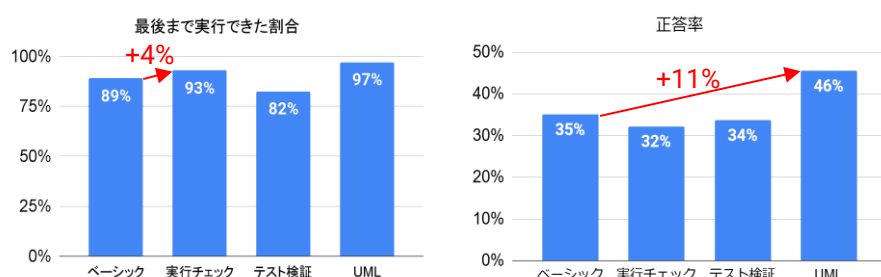
$$\text{最後まで実行できた割合} = 1 - \frac{\text{WA} + \text{TLE} + \text{RE}}{\text{AC} + \text{WA} + \text{TLE} + \text{RE}}$$

$$\text{正答率} = \frac{\text{AC}}{\text{AC} + \text{WA} + \text{TLE} + \text{RE}}$$

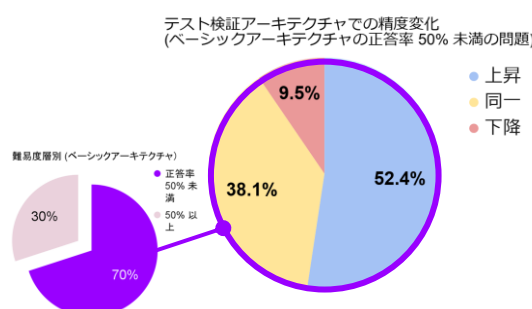
AC 正解 TLE 実行時間オーバー
WA 不正解 RE 実行時エラー

結果

- ベーシックアーキテクチャ(シンプルなLLMへの問い合わせ)を含む, 4種類のアーキテクチャを評価。
- 実行チェックアーキテクチャは最後まで実行できた割合がベーシックアーキテクチャと比較して4%向上。
- UMLアーキテクチャはベーシックアーキテクチャと比較して正答率が11%向上。



- テスト検証アーキテクチャでは全体の性能は低下したが, ベーシックアーキテクチャで正答率の低い問題に対して精度向上。



考察

UMLアプローチ

No.1!

正答率 ↑

最後まで実行できた割合と正答率が大幅(約10%↑)に向上。問題を構造化し, コーディングしやすく記号化することで精度が向上したのではないかと。

最後まで実行できた割合 ↑

実行チェックアプローチ

テスト検証アプローチ

最後まで実行できた割合 / 正答率が下がる結果が得られた。プロンプトの微妙な差異によって初期生成コードの品質に差が生まれその結果が性能に影響していることが分かった。

今後の課題 / 展望

- ① 正答率が減少した原因の分析と改善
- ② さらに実行エラーの減少と正答率の向上
- ③ 中・大規模なアプリケーションの作成
⇒汎用的なソースコード生成ソリューションの構築
⇒エージェントアプリを使った自己再帰的な改善