

修士学位論文

題目

メタプログラミングフレームワーク MetaGPT に対する
コード生成能力の再評価

指導教員

楠本 真二 教授

報告者

WANG TIANHAO

令和7年1月28日

大阪大学 大学院情報科学研究科

コンピュータサイエンス専攻

令和 6 年度 修士学位論文

メタプログラミングフレームワーク MetaGPT に対する
コード生成能力の再評価

WANG TIANHAO

内容梗概

ChatGPT などの対話的 LLM を用いたメタプログラミングフレームワーク, MetaGPT が着目を浴びている. MetaGPT は複数の LLM エージェントの協調作業によって, 与えられた自然言語の仕様からソースコードを自動的に生成する. MetaGPT の提案論文では, プログラミングコンテストのような簡単な問題に対しては約 8 割の成功率を達成したと報告している. MetaGPT 研究の課題の一つとして, 性能評価が十分でないという点が挙げられる. 第一に生成結果の妥当性を評価するためのテストが不十分であり, 性能を過大評価している可能性がある. さらに同一問題に対するソースコードの生成試行回数が一度ずつのみであるため, LLM の持つランダムさを考慮できていない. 本研究では MetaGPT の性能を再評価するために, 妥当性評価のテストが追加されたコード生成データセットを用いた再実験を行う. また同一タスクに対して複数回の実験を行うことで, LLM のランダムさを考慮した性能評価を行う.

主な用語

メタプログラミング, コード生成, マルチエージェント, 大規模言語モデル, MetaGPT

目次

1	はじめに	1
2	準備	3
2.1	メタプログラミング	3
2.2	複数の LLM を活用した手法	3
2.3	MetaGPT	3
3	MetaGPT 研究の課題とその解決策	7
3.1	生成結果の成否判定テストが不十分	7
3.2	異なる生成結果が生じる可能性を考慮していない	7
3.3	生成失敗の詳細分析が行われていない	8
4	実験	9
4.1	実験の概要	9
4.2	データセット	10
4.3	評価指標	11
5	結果	12
5.1	成否判定テスト通過率による生成能力の評価	12
5.2	Pass@ k による生成能力の評価	14
5.3	失敗原因の分析	15
6	妥当性への脅威	21
7	おわりに	22
	謝辞	23
	参考文献	24

目次

1	MetaGPT の外観図	4
2	MetaGPT の評価実験の結果（文献 [1] より引用）	5
3	MetaGPT を使用して開発された GUI 付きのソフトウェア（文献 [1] より引用）	6
4	プロンプトの例（MBPP#19）	9
5	コード生成試行の成否判定テスト通過回数	12
6	MBPP+ と HumanEval+ における生成コードの平均 Pass@ k	14
7	誤った出力形式設計による失敗の例	15
8	関数入力の誤解による失敗の例	17
9	エンジニアによる失敗の例	18
10	特殊な状況の対応できない例	20

表目次

1	MetaGPT 研究と本研究の実験設定の比較	10
2	生成結果の成否判定テスト通過率	12
3	エージェントごとの失敗回数と割合	15
4	エージェントごとの失敗原因の分類	16
5	誤ったコードを生成した原因	19

1 はじめに

開発者がソースコードを直接記述せずにプログラミングを実現する、メタプログラミングというアイデアの実現性が増してきている。特に自然言語による対話が可能で大規模言語モデル (Large Language Model: LLM) を用いた、メタプログラミングフレームワークも登場している [2][3]。このフレームワークを用いることで、自然言語による抽象的なソフトウェアの要求を入力とするだけで、全自動で要求を満たすソースコードを得ることができる。

さらに、MetaGPT[1] や ChatDev[4] と呼ばれる複数の LLM モデルを用いたメタプログラミングフレームワークも存在する。これら 2 つの手法は人による開発を模倣しており、複数の LLM が特定の役割を持ったエージェントとして振る舞うことで、それらエージェント同士が協調作業することでソフトウェアを作り上げる。MetaGPT ではエージェントとして、プロダクトマネージャー、アーキテクト、プロジェクトマネージャー、エンジニア、QA エンジニアを設けている。複数エージェントによる協調作業の結果は、他の LLM ベース手法と比較して高い性能を持つことが報告されている。具体的には MetaGPT は MBPP[5] データセットに対して 85.9 %、HumanEval[6] データセットに対して 87.7 % の正解率となっている。

しかしながら、MetaGPT 論文の評価方法に対しては 3 つの課題が存在している。まず、データセットに含まれる生成結果の成否判定用のテストが不十分であることが指摘されている [7]。テストの不十分さは、性能評価を過大評価につながっている可能性がある。また、評価方法として、複数の解を得るという状況を想定していない。LLM は一定のランダムさを内包しており、複数回の試行によって複数解を得ることが一般的である。この際には Pass@ k と呼ばれる評価指標を用いることが一般的であるが、MetaGPT 論文では k を 1 に固定して実験を行っている。これは従来の正答率と等価な値であり、複数の解を考慮しているとは言えない。最後、MetaGPT の生成失敗の詳細分析が行われておらず、具体的にどのような失敗が発生していたかが一切明らかになっていない。

本研究の目的は、MetaGPT が持つコード生成能力の再評価である。そのために MetaGPT 論文が実施していた評価実験の方法を見直し、上記 3 つの課題を解決したうえでの再評価実験を行う。まず複数の解を取るという状況を想定するために、Pass@ k の k を 1 から 10 まで変動させ、その際の性能を確かめる。テストが不十分という課題に対しては、テストが拡充された拡張データセットを採用することでその解決を試みる。これらの実験で得られた成否のうち、失敗ケースを目視により調査し、どのエージェントが失敗するのか、どういう原因で失敗するのかなどの詳細分析を行う。これにより具体的な失敗原因を明らかにする。

本稿の構成は以下の通りである。第 2 節では本研究の背景および関連技術について説明する。第 3 節では先行研究の課題を整理し、それらの解決策を提案する。第 4 節では本研究における実験設計を詳述

し、第 5 節では再実験の結果とその分析を示す。最後に、第 6 節で本研究の妥当性への脅威について考察し、第 7 節で結論と今後の課題を述べる。

2 準備

2.1 メタプログラミング

メタプログラミングとは、ソースコードの直接的な記述をせずにソースコードを獲得するプログラミング手法である [8]。メタプログラミングによって、ソフトウェア開発者は専門的な知識不要、かつ効率的にプログラムを作成できる。例えば、統合開発環境の一つである IDEA ではユーザの要求に従って、クラスのコンストラクタを自動生成する機能を提供している [9]。また、ローコードプラットフォーム [10] ではユーザの設計したフローチャートに基づき、対応するコードを自動的に生成する機能が存在している。このように、メタプログラミングは多様な場面での活用が進んでいる。

近年、大規模言語モデル (Large Language Model: LLM) を用いた、自然言語による対話ベースでのコード生成手法が着目を浴びている [11]。GPT-4[12] や Llama[13] などのモデルは、大量の自然言語の文章にだけでなく大量のソースコードを学習に用いている。そのためユーザが与えた自然言語の要求に応じて、ソースコードの生成や書き換えといった様々なタスクに取り組むことも可能である。

2.2 複数の LLM を活用した手法

近年ではプロダクトコードの生成だけでなく、テストコードや仕様書等、あらゆるドキュメントを LLM に生成させる手法が検討されている。この手法では LLM に人間の役割 (エージェント) を割り当て、複数のエージェントを相互に協調・連携させながらソフトウェア開発を行う [14]。

Qian らは、ソフトウェア企業を模したフレームワーク ChatDev を提案している [4]。この手法では、複数の AI エージェントが最高経営責任者 (CEO)、最高技術責任者 (CTO)、プログラマ、テスターといった役職を担当する。これらのエージェントはあらかじめ設定されたプロンプト規則をもとに、相互にコミュニケーションを行い、ユーザ要求の詳細な自動分析やコード品質の自動的な評価を実現している。また、Zhang ら [15] による PairCoder では、一人のエージェントが問題分析と複数の解法候補を生成し、もう一人のエージェントが選択された解法に基づくコードの生成とテストを担当する。この2つのエージェント間で、複数の解法から得られるコードの実行結果を比較検討し、最終的に最適なコードを導き出している。ChatDev や PairCoder では、開発者が作成したいソフトウェアの要求仕様書を入力するだけで、ソースコードだけでなく様々なソフトウェア開発ドキュメントを獲得することが可能である。

2.3 MetaGPT

ChatDev や PairCoder とは異なる仕組みとして Hong ら [1] の MetaGPT がある。MetaGPT は複数の AI エージェントを用いた対話的かつ LLM ベースのメタプログラミングフレームワークであ

図 1 に MetaGPT の処理の流れを示す。MetaGPT ではウォーターフォール型の開発プロセスに基づき、プロダクトマネージャー、アーキテクト、プロジェクトマネージャー、エンジニア、QA エンジニア等の 5 種類のエージェントを割り当てている。エージェント間のコミュニケーションは様々なテキストドキュメントによって行われる。図 1 ユーザが入力した自然言語の要求を入力として、プロダクトマネージャ (PM) が仕様書を作成し、その仕様書に基づいてアーキテクトがシステム設計書を作成する。最終的にエンジニアがソースコードを、QA エンジニアがテストを作成する。出力として得られるソースコードはテストによる検証が行われることになる。

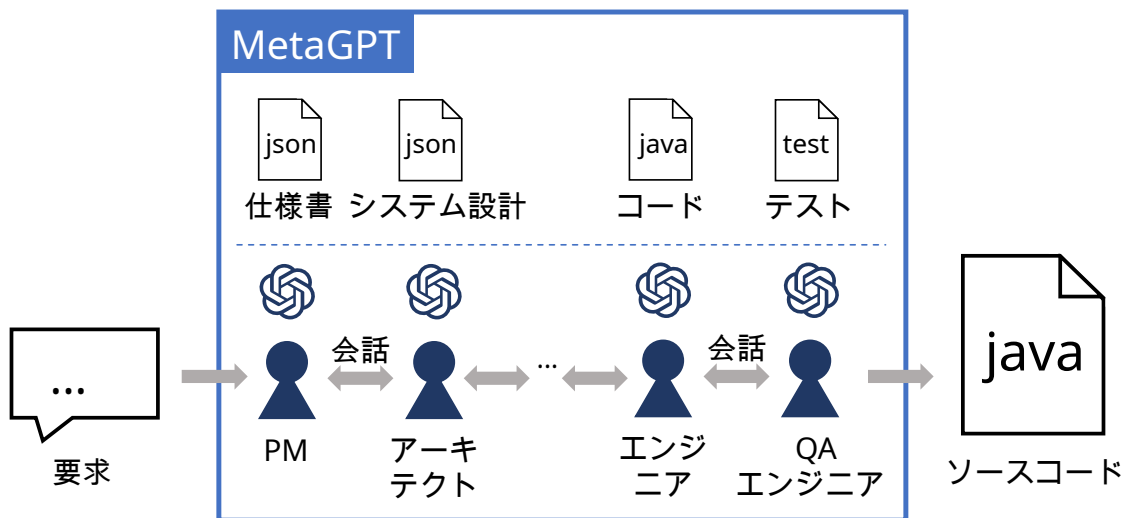


図 1 MetaGPT の外観図

MetaGPT の 3 つの特徴を以下に示す。

- 動的なエラー検出が可能である
- エージェント間で開発文書を使用してやり取りを行う
- GUI を含むソフトウェアの生成が可能である

以下では、特徴の詳細について述べる。

MetaGPT ではエージェントの一つとして QA エンジニアが設けられており、生成したコードを実行しテストを通じてエラーを動的に発見し、修正を行っている。既存のマルチエージェントフレームワーク [16] では、生成コードに対してコードレビューなど静的な解析を行うのみでエラー検出を行っている。

図 2 に MetaGPT の提案論文 [1] で示されていた、MetaGPT の性能評価実験の結果を示す。図は様々な LLM モデルに対するプログラム生成の性能を表しており、棒の高さが性能 (Pass@k) である。2 つのコード生成のためのデータセット (MBPP[5] と HumanEval[6]) に対する実験結果である。図中

の赤枠が MetaGPT の性能であり、フィードバックとは QA エンジニアの有無、すなわち生成したソースコードの動的な結果をエンジニアにフィードバックしているか否かを表している。比較対象手法としては、AlphaCode[17] や Incoder[18] などのコード生成特化型 LLM に加え、GPT-4 など汎用型 LLM も含まれている。MetaGPT は全ての他の手法と比べて高い精度を示している。またフィードバックを加えることで、MBPP に対しては正解率を 82.3 % から 87.7 % に、HumanEval に対しては 81.7 % から 85.9 % に向上している。

MetaGPT は各種エージェントが様々な開発文書を生成し、エージェント間のやり取りにその文書を活用する。Zhang ら [19] の研究では、複雑な対話ログがエージェントの推論能力に悪影響を与えることが示されている。MetaGPT 内部ではエージェント間で文書を用いたやり取りを行うことで、冗長な情報がエージェントの推論能力に与える悪影響を回避している。さらに、生成された開発文書は MetaGPT 利用者による検証や理解の助けとなるだけでなく、再利用の促進にも役立つ。

また MetaGPT は GUI を含むソフトウェアを生成できる。従来のコード生成関連の研究では、その手法は主として指示されたタスクやアルゴリズムの解法のみに着目している。一方で、MetaGPT は必要なアルゴリズムを生成するだけでなく、GUI を含むより複雑なソフトウェアも生成可能である。これは、MetaGPT が他の LLM ベースのコード生成手法に比べて、複雑なコード生成タスクに対する処理能力が優れていることを示している。図 3 に MetaGPT を使用して開発された GUI 付きのソフトウェアを示す。図中では五目並べや数独などの作成タスクに対して、GUI を含むソースコードの生成が可能である。

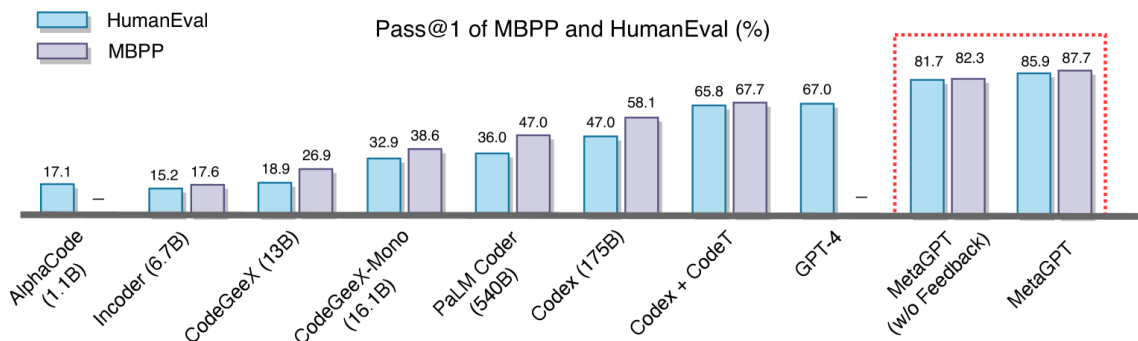


図 2 MetaGPT の評価実験の結果 (文献 [1] より引用)

3 MetaGPT 研究の課題とその解決策

2.3 節でも述べた通り, MetaGPT は複数の AI エージェントによる全自動でのソフトウェア開発が可能であり, LLM 活用に対する様々な可能性を示している一方で, その評価方法について 3 つの課題が存在している.

- 生成結果の成否判定テストが不十分
- 異なる生成結果が生じる可能性を考慮していない
- 生成失敗の詳細分析が行われていない

本節では, 各課題の詳細とその解決策について述べる.

3.1 生成結果の成否判定テストが不十分

MetaGPT 論文が評価実験に用いたデータセット MBPP と HumanEval に対しては, テストが不十分であると指摘されている [7]. Li ら [20] の研究によれば, MBPP と HumanEval には誤りやテストケースの不足が含まれており, 一部のバグを適切に検出できない可能性が指摘されている. MBPP に含まれるテストケース数は 1 問題あたり平均 3 個, HumanEval は平均 7 個である. MBPP の問題が基本的なアルゴリズムに限定されている点を考慮に入れても, 単体テストとしては不十分であり, MetaGPT のコード生成性能が過大に評価されるおそれがある.

そこで本研究では, 生成されたコードを正確にテストするために, Liu らの提案する MBPP+[7] と HumanEval+[7] の 2 つのデータセットを利用する. これらのデータセットは, それぞれ MBPP および HumanEval に対して, より多くのテストケースを追加し, 元のデータセットでは検出できなかったエラーを検出可能としている. 本研究では MBPP+ と HumanEval+ を用いて, MetaGPT の生成結果の成否判定を行う.

3.2 異なる生成結果が生じる可能性を考慮していない

MetaGPT 論文の評価実験では, 同一の問題に対して複数の異なるコードを生成する可能性を考慮していない. LLM は Transformer モデル自体にランダム性があるため, 同じ入力に対しても複数の解が得られる. よって LLM の性能評価においては, 同一条件のタスクに対して複数の試行を行い, 複数の解を得ることが一般的である.

MetaGPT 論文が採用していた評価指標は Pass k である. この指標は LLM のための評価指標であり, LLM から得られた複数の解を k 個ピックアップした際に一つでも正解となる解を含む可能性を表している. 解候補の中から k 個の無作為に抽出し, 一つずつ目視確認して成否を確かめるというような

状況を想定した指標である。

Pass k 自体は LLM のランダムさ、すなわち複数回の試行を考慮に入れているが、MetaGPT 論文では k を 1 とした実験しか行っていない。 $k=1$ とは、すなわち 1 個のピックアップで正解を引く可能性のことであり、単純に成功の割合（試行数に対して何個正解を含むか）と同値となってしまう。つまり複数の解を一切考慮に入れておらず、LLM を適用するという実践的な状況を想定していない。

本研究では、様々な k に対して評価実験を行う。具体的には k の値を 1~10 に変動させ、その際の Pass@ k の値を確認する。基本的には k の値が大きくなるほど、正解を 1 個でも含む可能性（Pass@ k 自体）は高くなるが、その成否の確認に要するコストは増大する。 k を変動させることで、コストの増大に対する性能改善の程度を調べることができる。

3.3 生成失敗の詳細分析が行われていない

MetaGPT 論文の評価実験では、どのエージェントがエラーを誘発したのか、また MetaGPT が対応できない要求や条件があるのかといった詳細な分析が十分に行われていない。こうした情報が欠如していることで、MetaGPT のさらなる改善のための具体的な指針を得ることができない。

この問題の解決のために、本研究では MetaGPT の生成ログを目視により分析する。どのエージェントがどのような原因で最終的なコード生成の失敗を引き起こしたのかを分析することで、MetaGPT に存在する欠陥を明らかにすることができる。

4 実験

4.1 実験の概要

前節で述べた MetaGPT 研究 [1] の評価の課題を見直し、MetaGPT の性能に対する再評価を行う。基本的には MetaGPT 研究の実験設定に従い、妥当でない、あるいは不十分な設定を改善する。実験の流れは次のとおりである。

- データセットから一つタスクを取り出す
- 当該タスクの要件とテストをプロンプトに埋め込む
- MetaGPT にプロンプトを渡しタスクに取り組むさせる
- MetaGPT の出力を得る
- プログラムの生成に成功したかを成否判定テストによって確認する

用いるプロンプトの一例を図 4 に示す。この例は MBPP というデータセットの 19 番目の問題であり、整数値の配列から重複する要素が含まれるかを判定する、というプログラムの生成タスクである。MBPP はこのような基礎的なアルゴリズムを問う問題で構成されている。データセットには上記のタスクの内容に加えて、参考用のテストケースが含まれている。生成するプログラムが満たすべき要件をテストという形式で表現しているとも言える。なお、参考用テストとは独立に成否判定用のテストも存在している。これは MetaGPT に与えず、MetaGPT の出力の成否判定に用いる。よって成否の判定はテストを用いた全自動化が可能である。

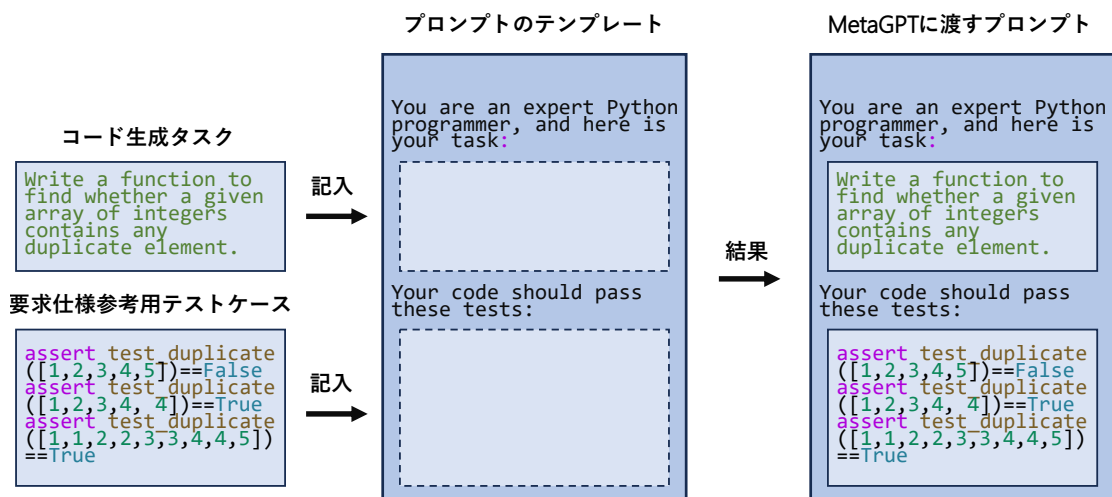


図 4 プロンプトの例 (MBPP#19)

表 1 MetaGPT 研究と本研究の実験設定の比較

	先行研究 [1]	本研究
データセット	MBPP と HumanEval	MBPP+ と HumanEval+
#検証用テスト数	≈ 3 と ≈ 7	≈ 105 と ≈ 560
#タスク数	不明/974 と 不明/164	20/974 と 20/164
評価指標	Pass@1	Pass@ k
生成回数 n	不明	10
サンプル数 k	1	1~10
失敗分析	N/A	目視により確認
モデル	GPT-4	GPT-4o

表 1 に、本研究と MetaGPT 研究の実験設定の違いを示す。まず実験に用いるデータセットは、MetaGPT 研究ではプログラム生成用の評価データセットである MBPP と HumanEval であった。本研究ではこのデータセットのテストを拡充した MBPP+ と HumanEval+ を用いる。検証用テストは拡充前は 1 タスク当たり 3 個と 7 個であったのに対し、拡充後は 105 個と 560 個となっている。なお、本稿では各データセットに含まれる全タスクのうち、20 個ずつをランダムに選定し実験に用いる。MetaGPT 論文ではいくつかのタスクが実験に用いられていたかは明確にされていない。

評価指標は Pass@1 のみを用いていたのに対し、本研究では k を 1~10 に変動させ、その全ての評価値を確認する。この際の同一タスクに対する全試行回数は 10 回とする。タスク数と同様、MetaGPT 論文において試行回数を何回としていたかは記述されていない。

さらに失敗分析を目視により行う。特にどのエージェントが失敗していたのか、その原因は何なのかを分類し、MetaGPT の性能改善に対する課題を洗い出す。この調査にあたっては MetaGPT が生成するログを目視で精査する。

なお評価モデルは GPT-4 ではなく GPT-4o を用いることとした。GPT-4o の性能は GPT-4 と同等かそれ以上であると指摘されている [21]。また、GPT-4o は GPT-4 より大幅に生成時間が短く、コストも安いモデルである。

4.2 データセット

本実験では、主に MBPP+ と HumanEval+ の 2 つのデータセットを使用した。MBPP+ および HumanEval+ は、MBPP および HumanEval を基に、EvalPlus[7] を用いてそれぞれ 35 倍、80 倍のテストケースを追加したデータセットである。MBPP と HumanEval には、それぞれ 974 件と 164 件の簡単なコード生成タスクが含まれており、各タスクには問題文・参考解答・テストセットが付属して

いる。これらは自然言語からのコード生成に関する研究に広く利用されるベンチマークとして知られている。EvalPlus[7] は、自動テスト入力生成フレームワークである。このフレームワークは、LLM を用いた突然変異に基づくテスト入力生成を通じて、データセットのテストケースを増加させる。より多くの成否判定テストを含むデータセットを使用し、同一の問題に対して複数回の生成を行うことで、MetaGPT のコード生成能力をより精確に評価することができる。

4.3 評価指標

本実験では、MetaGPT が要求を満たすコードを生成する確率を評価するために Pass@ k を用いる。Pass@ k の計算方法は、eq. (1) に示すとおりである。ここで、 n はコード生成を行った回数、 c は要求を満たすコードが生成された回数、 k はすべての生成結果のうちランダムに選択するサンプル数である。

$$\text{Pass}@k = \mathbb{E}_{\text{Problems}} \left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right] \quad (1)$$

Pass@ k の値はコード生成能力の高さを表しており、 k を固定したとき、Pass@ k の値が高いほど、正しくコードを生成できている可能性が高い。

自然言語処理の分野では、従来 ROUGE[22] や BLEU[23] といった指標が広く用いられてきた [24]。これらの指標は LLM によって生成されたコードが正解コードと文字列レベルで一致しているかを評価している。しかし、コード生成の評価においては、コードが実現する機能の過不足やコードの正しさがより重要である。したがって、ROUGE や BLEU 等の指標より、適切なコード生成ができたかを表す Pass@ k の方がコード生成指標として適している [25][26]。

Roziere ら [27] の論文でも、「生成されたコードがすべてのテストに合格する割合」を用いてコード生成能力を評価している。これは、一般的なコード生成が入力に対して決定的であることを前提としている。しかし、LLM にはランダム性があり、一度の生成では正解に至らなくても、複数回の試行によって正しいコードが得られる可能性がある。そのため、単純に正解率や Pass@1 を用いる評価手法では、LLM のランダム性や困難な問題を克服する潜在能力を十分に反映できない可能性がある。

表 2 生成結果の成否判定テスト通過率

データセット	生成回数	成否判定テスト通過率
MBPP	200	67.0 %
MBPP+	200	61.0 %
HumanEval	200	61.0 %
HumanEval+	200	52.5 %
MBPP と HumanEval	400	64.0 %
MBPP+ と HumanEval+	400	56.5 %

5 結果

5.1 成否判定テスト通過率による生成能力の評価

表 2 は、成否判定テストのテストケースの追加前後で、MetaGPT が生成したコードがコード生成タスクの成否判定テストを通過できる割合を示している。4.2 節で述べたとおり、MBPP+ および HumanEval+ の「+」は、MBPP および HumanEval と同じコード生成タスクを含みつつ、それらに大量のテストケースを追加したものである。テストケースを追加することにより、MetaGPT が生成したコードがタスクの成否判定テストを通過する割合は、MBPP での 67 % から MBPP+ での 61 % に低下した。また、MetaGPT が生成したコードがタスクの成否判定テストを通過する割合は、HumanEval での 61 % から HumanEval+ での 52.5 % に低下した。

MetaGPT 論文の評価実験 [1] では、MetaGPT の平均正解率が 2 つのデータセットで約 86 % と報

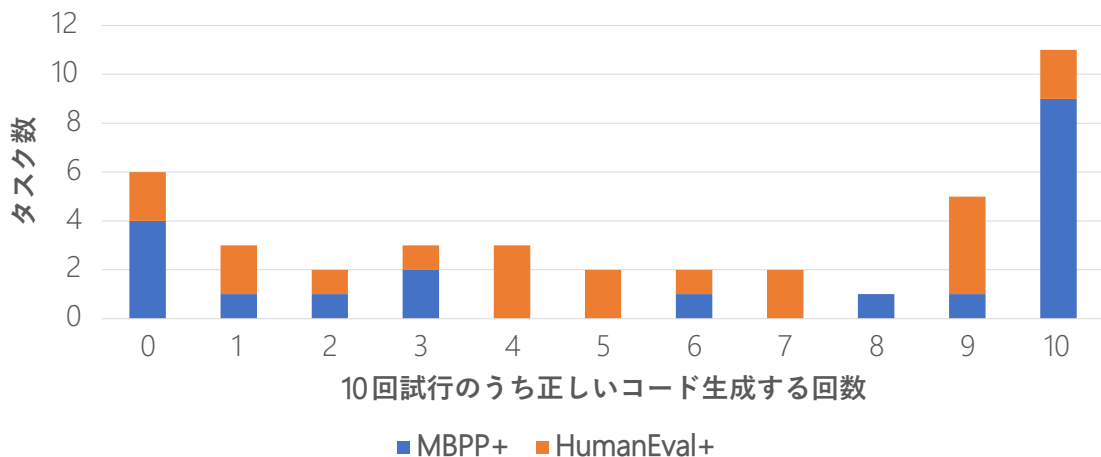


図 5 コード生成試行の成否判定テスト通過回数

告されている。しかし、本実験の結果から、MetaGPT が過大に評価されている可能性があることが示唆された。

図 5 は、コード生成タスクごとに 10 回の試行で正しいコード（成否判定テストを通過するコード）を何回生成できたことを示している。横軸は、単一タスクにおける 10 回の試行中、正しいコードが生成された回数を示す。棒の高さはその生成回数に該当するタスク数を示す。棒の色は、タスクの出所を区別している。

全 40 タスクのうち、MetaGPT は 19 タスク（47.5 %）について、10 回中少なくとも 8 回は正しいコードを生成している。これらのタスクに関しては、MetaGPT の生成精度が比較的高い（ $\geq 70\%$ ）と考えている。一方、15 タスク（37.5 %）では、10 回のうち 1~6 回のみ正しいコードを生成するにとどまり、精度が低いため、複数回の試行を要する可能性がある。さらに、6 タスク（15 %）においては、10 回の試行すべてで正しいコードを生成できなかった。これらのタスクでは、エラー要因の詳細を分析し、MetaGPT を改良する必要がある。

5.2 Pass@k による生成能力の評価

図 6 は、2つのデータセットにおける k が 1 から 10 の値を取る場合の平均 Pass@ k の値を示している。 $k=1$ 、つまり 1 回の生成で正しいコードが得られる確率を考える場合、MetaGPT の平均 Pass@1 は 56.5 % であり、MetaGPT 提案論文における 86 % と比較して大幅に低い。 また、MBPP+ と HumanEval+ に含むコード生成タスクは簡単な python プログラミング問題であるにもかかわらず、1 回の試行で正しいコードを生成する確率は 56.5 % であり、MetaGPT には改善の余地があると考えている。

$k=3$ の場合、生成されたコードに正解が含まれる可能性は約 76.8 % に達する。 これは完全自動で生成されたコードであることを考慮すると、精度が良いと言える。

$k=8$ 以上になると、正解を含む生成コードの確率は 90 % を超える。 このとき、曲線は飽和に近づき、さらに生成回数を増やすことで得られる利点は、金銭や時間の追加コストに比べて小さいと考えている。 また、この結果は、ごく少数のコード生成タスクについては最終的に解決できなかったものの、多くのコード生成タスクでは複数回の試行によって正解に生成できることを示している。

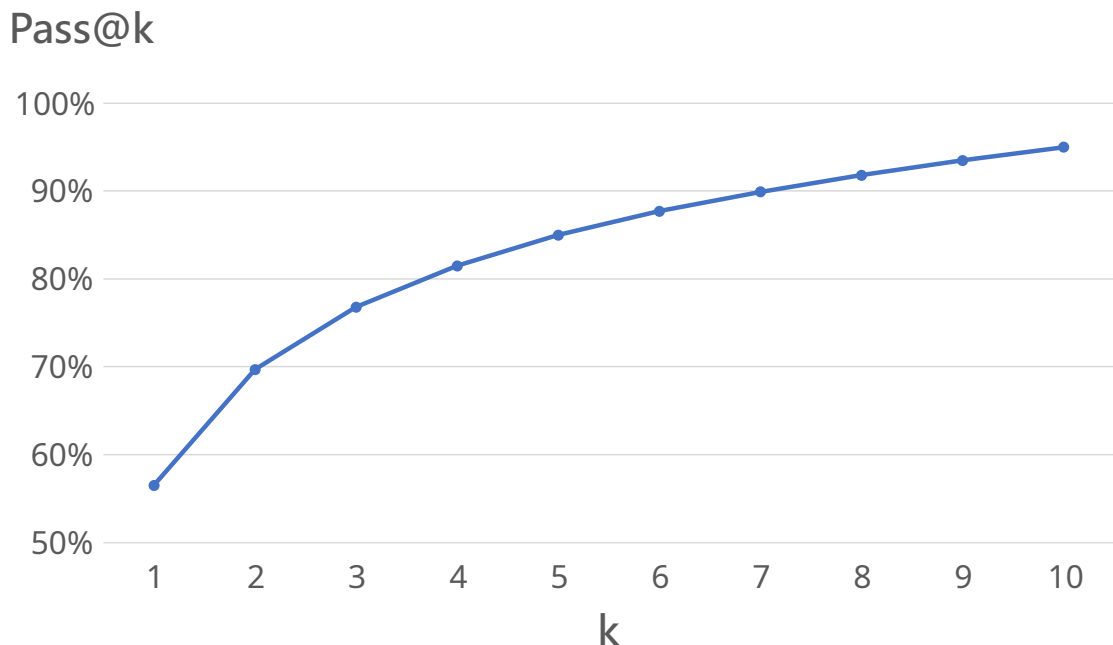


図 6 MBPP+ と HumanEval+ における生成コードの平均 Pass@ k

表3 エージェントごとの失敗回数と割合

エージェント	失敗試行数と割合
プロダクトマネージャー	14 (8.1 %)
アーキテクト	33 (19.1 %)
エンジニア	126 (72.8 %)
エージェント全体	174 (100.0 %)

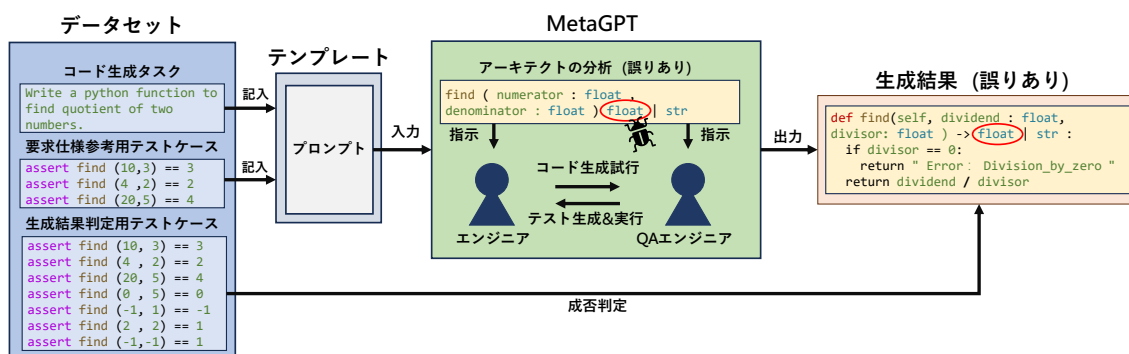


図7 誤った出力形式設計による失敗の例

5.3 失敗原因の分析

表3では、173回の失敗において、どのエージェントが失敗を引き起こしたかの数と割合を示している。そのうち、14回(8.1%)はプロダクトマネージャー、33回(19.1%)はアーキテクト、126回(72.8%)がエンジニアによるものであった。

図7では、上流設計を担当するエージェントのミスによって、ソースコードとテストの生成を担うエンジニアやQAエンジニアにも誤りが波及する典型的なケースを示している。コード生成タスクとしては、2つの数値の商を計算するPython関数を作成することが求められている。対応する参考テストケースには3つの例が示されており、関数は商を整数型で返す必要があることが明確に記載されている。本研究の実験では、図4に示す方法で、データセット中のコード生成タスクと対応する要求仕様・参考テストケースを組み合わせたプロンプトをMetaGPTに入力した。ところが、データ構造とアルゴリズムの設計を担当するアーキテクトが、関数の戻り値を誤って浮動小数点型として設計してしまった。アーキテクトは設計文書をエンジニアとQAエンジニアに渡し、両者はその設計文書に基づいてコードを実装し、テストケースを作成・実行して最終的な生成物を出力する。しかし、当初のデータ構造設計が誤っていたため、誤った設計をもとにコードを実装したエンジニアも誤ったコードを生成することになった。さらに、QAエンジニアもアーキテクトの設計ファイルを参考にテストを作成したため、この

表4 エージェントごとの失敗原因の分類

エージェント	テスト失敗原因の種類	個数	割合
プロダクトマネージャー	関数の入力解釈に失敗	14	8.1 %
アーキテクト	関数の出力の型が正しくない	24	13.9 %
アーキテクト	アルゴリズム設計が正しくない	9	5.2 %
エンジニア	設計通りにコードを生成できない	57	33.0 %
エンジニア	特殊な状況の対応できない	42	24.3 %
エンジニア	返り値の文字列が僅かに違う	15	8.7 %
エンジニア	ループの範囲が違う	7	4.0 %
エンジニア	返り値が設計通りにソートされていない	5	2.9 %
エージェント全体	全ての失敗原因	174	100.0 %

誤りに気づけなかった。最終的に、MetaGPT が出力した生成結果では、関数が2つの数値の商を浮動小数点型として返すため、指定されたテストケースをすべて通過できず、生成失敗と判断される結果となった。そのため、表3と表4では最初にエラーを引き起こしたエージェントのみを集計している。

表4では、各エージェントがテスト失敗を引き起こす具体的な原因をより詳しく分析している。以下では、一部の失敗原因を例に挙げ、それらがどのように発生したかを詳しく説明する。

図8は、プロダクトマネージャーが関数の入力を誤って分析し、最終的に誤ったコードが生成された例を示している。このコード生成タスクでは、円錐の側面積を計算することを要求しており、3つのテスト用例が参考として提示されている。本来、関数の入力には円錐の半径と高さであるべきだが、プロダクトマネージャーは誤って、ユーザーが半径と斜高にもとづいて側面積を計算することを望んでいると判断してしまった。前述のとおり、上流設計を担当するエージェントの誤りは、後続のエージェントが参照する設計文書にも誤りを連鎖的に広める。そのため、プロダクトマネージャーの開発文書を基にアルゴリズム設計やコード生成を行った他のエージェントが最終的に生成した関数は、半径と斜高を入力とするようになってしまい、成否判定テストを通過できなかった。

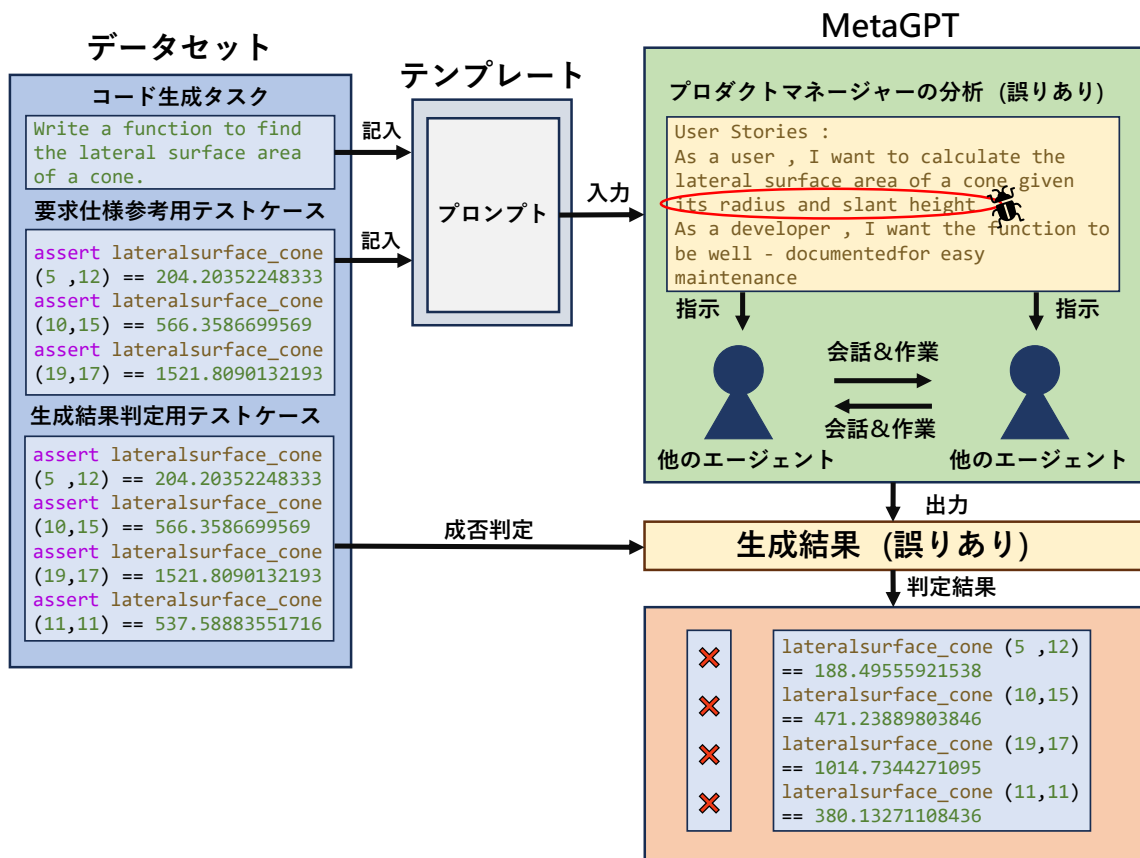


図8 関数入力の誤解による失敗の例

図9は、上流設計が正しく行われたにもかかわらず、エンジニアがループの範囲を誤った結果、誤りが生じた例を示している。このコード生成タスクでは、第n番目のペリン数を求めるよう要求され、参考テスト用例では第9、第10、第11のペリン数が与えられている。上流設計を担当するエージェントが問題を正しく分析していたにもかかわらず、エンジニアはペリン数を第0番目ではなく、第1番目から計算するものと誤解してコードを生成してしまった。さらに、QAエンジニアもエンジニアのコードと考え方をもとにテスト用例を生成したため、このバグを検出できるテストを作成しなかった。結果として、正しいペリン数を返すはずの関数が誤った値を返すコードとなり、テストを通過できなかった。

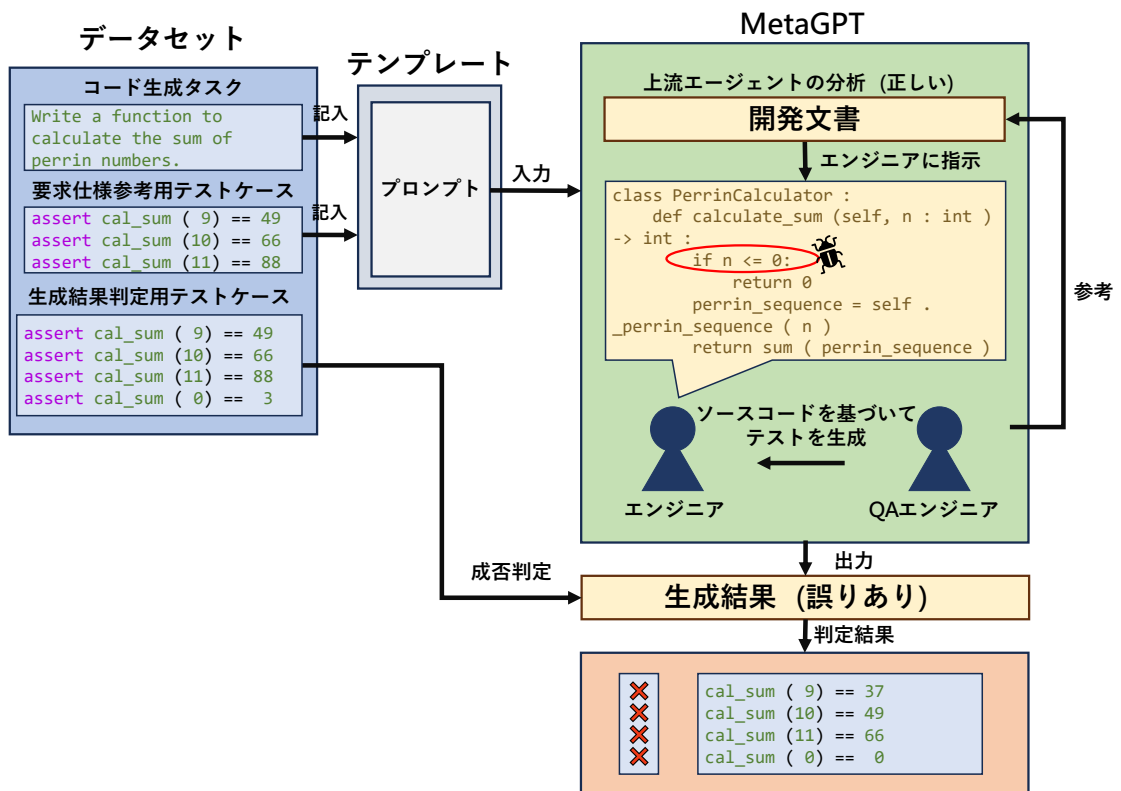


図9 エンジニアによる失敗の例

表 5 誤ったコードを生成した原因

失敗原因	個数	割合
LLM のコード生成能力不足	122	70.5 %
プロンプト中の情報の見落とし	51	29.5 %

表 5 に示すように、LLM のコード生成能力不足は誤ったコードを生成した主要な原因である。MetaGPT が内部で利用する LLM を変更することで、この問題がある程度緩和される可能性がある。

次に、MetaGPT がプロンプトを完全に解析できない点も問題として挙げられる。前述のとおり、MetaGPT を構成するエージェントは、しばしば関数の入出力に関する要件を正しく設計できない。タスクの説明文では明示されていなくとも、テストケースにその要件が示されている場合があるが、MetaGPT はそれを十分に活用できていない。図 7 に示すように、あるタスクでは「二つの数値の商を計算せよ」という指示とともに 3 つのテストケースが提示されていた。タスクの説明文には、商を整数型で返すのか浮動小数点型で返すのか明確には記載されていなかったが、提示されたテストケースは整数型の結果を要求していた。しかし、アーキテクトはこの要求を正しく把握せず、誤ったクラス構造を設計する。QA エンジニアもプロンプトに提示されていた参考テストを生成したテストに含めなかった。結果として、コード自体は数値の商を計算できるものの、戻り値の型が要件と合わないという問題が発生している。また、要求を正しく理解できていないと考えられる。図 10 に示すように、単語中の母音の数を計算する際に、語末の y を母音としてカウントするよう明記されているにもかかわらず、MetaGPT はこれを無視したコードを生成した。MetaGPT 内部のプロンプトシステムをさらに最適化し、エージェント間の通信方式を改善することで、この種の問題を低減できると考えられる。

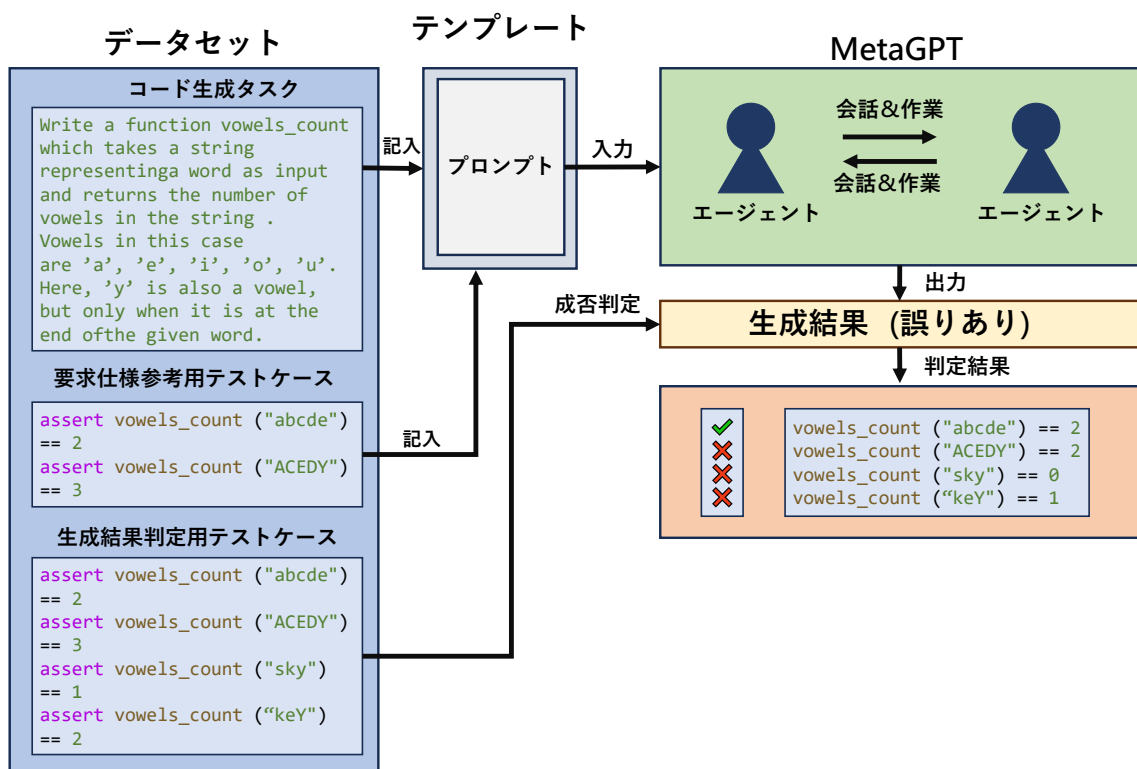


図 10 特殊な状況の対応できない例

6 妥当性への脅威

MetaGPT の提案論文で使用された MetaGPT は、バージョン 0.4 を基に、一部のプロンプトを修正して得られた特別なバージョンである。しかし、このバージョンは著者によって公開されておらず、さらにユーザは常に最新バージョンを利用する傾向があるため、本研究では MetaGPT のバージョン 0.8.1 を使用している。異なるバージョンの MetaGPT を使用した場合、実験結果が異なる可能性がある。

また、本実験で MetaGPT が使用する LLM は GPT-4o である。LLM のランダム性や、OpenAI による GPT-4o のアップデートが生成結果に影響を与える可能性があり、同一の生成タスクでも異なるコードが生成される場合がある。さらに、異なる LLM モデルを使用した場合にも生成結果が変わる可能性がある。

さらに、本実験では金銭および時間の制約により、MBPP+ および HumanEval+ からそれぞれランダムに 20 個の生成タスクを選択し、すべての問題をテストしているわけではない。異なる生成タスクをランダムに選択した場合や、データセット全体を使用した場合には、異なる実験結果が得られる可能性がある。

7 おわりに

本研究では、MetaGPT のコード生成能力の再評価を行った。MetaGPT とは複数の AI エージェントを用いた対話的かつ LLM ベースのメタプログラミングフレームワークである。ユーザーが自然言語の要求を入力して、MetaGPT がソースコードを自動的に生成する。

しかしながら、MetaGPT 論文の評価方法に対しては 3 つの課題が存在している。まず、生成結果の成否判定テストが不十分を指摘されている。性能評価を過大評価につながっている可能性がある。また、評価方法として、LLM のランダムさを無視し、複数の解を得るという状況を想定していない。最後、MetaGPT の生成失敗の原因を詳細分析が行われていない。MetaGPT の改善に繋がる議論がない。

以上の課題を踏まえ、本研究では以下の手法により MetaGPT のコード生成能力を再評価を行う。第一に、元のテストセットに加えて追加のテスト用例を含む拡張データセットを用いた実験を行う。第二に、MetaGPT に対して k を異なる値に設定しながら $\text{Pass}@k$ を評価する。第三に、MetaGPT が誤りを起こす原因を分析する。

評価の結果として、MetaGPT の正解率は 86 % ではなく 56 % である。コード生成性能を過大評価されている。 $\text{Pass}@3$ は 76.8 %、つまり解を 3 つ選べば 76.8 % 正解を含む。また、生成失敗の原因分析について、失敗原因の 72.9 % はエンジニアのミス。タスクの理解というよりコード生成自体の難しさがある。

今後の課題として、失敗原因のさらなる分析が挙げられる。アーキテクトが誤ったアルゴリズムを設計する際の傾向やエンジニアがコード構造設計を失敗しやすい状況を分析が必要である。また、今回の調査ではデータセットの合計 1,138 問のうち 40 問のみ実験しているため、データセットに含まれる全テストに対する実験も挙げられる。

謝辞

本研究を遂行するにあたり、多くの方々に多大なご協力をいただきました。

まずは、研究活動全般を通して快適な研究環境と雰囲気を整えてくださった楠本真二教授に、深く感謝いたします。教授の温厚で包容力のあるご指導のおかげで、困難に直面した際にも過度に緊張することなく、研究の推進と考察に集中することができました。また、何度もの議論や報告の場でいただいた貴重なご意見や、和やかな雰囲気づくりは、本研究の完成に大きく寄与したと感じております。

次に、研究全体を通じて熱心にご指導くださった松本真佑准教授に、心より御礼申し上げます。研究の方針から論文執筆、さらには発表準備に至るまで、常に丁寧かつ細やかにアドバイスをいただきました。研究過程で困難や迷いが生じた際には、時間を割いて懇切に相談に応じてくださり、要所で背中を押していただいたおかげで、研究のみならず多方面で成長する大きな機会を得ることができました。

研究に付随する各種事務手続きにおいては、事務補佐員の橋本美砂子氏に多大なるご協力をいただきました。特に出張手続きなどの事務対応を迅速かつ丁寧に行ってくださいましたおかげで、研究により多くの時間を注ぐことができました。また、私の体調面にもお気遣いいただき、非常に心強く感じました。改めて深く感謝申し上げます。

さらに、日常の議論や情報交換において、楠本研究室の皆様には大きな刺激と支えをいただきました。活発な意見交換や雑談を通じ、多角的な視点と気づきを得るとともに、研究の合間にリラックスする場も設けてくださり、非常に充実した研究生活を送ることができました。皆様の温かな交流があったからこそ、研究に意欲的に取り組むことができた実感しております。

また、日本語表現の修正や研究手法に関する助言をはじめ、さまざまな面でサポートしてくださったチューターの三原公平氏と堀翔太氏にも、深く感謝いたします。お二方からのきめ細やかなアドバイスと的確なサポートを通じて、多くの学びを得ることができ、本研究の進展に大いに貢献していただきました。

最後になりましたが、生活面および精神面で常に支えてくれた家族や友人の存在にも、心より感謝申し上げます。挫折や困難に直面した際も変わらぬ理解と励ましをくださったおかげで、研究に専念することができました。本修士研究を無事に完遂できたのは、ひとえに皆様のご支援とご厚情の賜物です。

以上をもちまして、私の学術的および人間的成長に寄与してくださったすべての方々に、改めて深く御礼申し上げます。誠にありがとうございました。

参考文献

- [1] Hong, S., Zheng, X., Chen, J., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., et al.: Metagpt: Meta programming for multi-agent collaborative framework, *arXiv preprint arXiv:2308.00352* (2023).
- [2] Guo, D., Zhu, Q., Yang, D., Xie, Z., Dong, K., Zhang, W., Chen, G., Bi, X., Wu, Y., Li, Y., et al.: DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence, *arXiv preprint arXiv:2401.14196* (2024).
- [3] Team, C., Zhao, H., Hui, J., Howland, J., Nguyen, N., Zuo, S., Hu, A., Choquette-Choo, C. A., Shen, J., Kelley, J., et al.: Codegemma: Open code models based on gemma, *arXiv preprint arXiv:2406.11409* (2024).
- [4] Qian, C., Liu, W., Liu, H., Chen, N., Dang, Y., Li, J., Yang, C., Chen, W., Su, Y., Cong, X., et al.: Chatdev: Communicative agents for software development, in *Association for Computational Linguistics*, pp. 15174–15186 (2024).
- [5] Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., et al.: Program synthesis with large language models, *arXiv preprint arXiv:2108.07732* (2021).
- [6] Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating large language models trained on code, *arXiv preprint arXiv:2107.03374* (2021).
- [7] Liu, J., Xia, C. S., Wang, Y. and Zhang, L.: Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation, *Advances in Neural Information Processing Systems*, Vol. 36, (2024).
- [8] Cordy, J. R. and Shukla, M.: *Practical metaprogramming*, Queen’s University of Kingston. Department of Computing and Information Science (1992).
- [9] Balzer, R.: A 15 year perspective on automatic programming, *Transactions on Software Engineering*, No. 11, pp. 1257–1268 (1985).
- [10] Prinz, N., Rentrop, C. and Huber, M.: Low-Code Development Platforms-A Literature Review., in *Americas Conference on Information Systems* (2021).
- [11] Ni, A., Iyer, S., Radev, D., Stoyanov, V., Yih, W.-t., Wang, S. and Lin, X. V.: Lever: Learning to verify language-to-code generation with execution, in *International Conference on Machine Learning*, pp. 26106–26128 (2023).

- [12] Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report, *arXiv preprint arXiv:2303.08774* (2023).
- [13] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models, *arXiv preprint arXiv:2302.13971* (2023).
- [14] He, J., Treude, C. and Lo, D.: LLM-Based Multi-Agent Systems for Software Engineering: Literature Review, Vision and the Road Ahead, *Transactions on Software Engineering and Methodology* (2025).
- [15] Zhang, H., Cheng, W., Wu, Y. and Hu, W.: A Pair Programming Framework for Code Generation via Multi-Plan Exploration and Feedback-Driven Refinement, in *International Conference on Automated Software Engineering*, pp. 1319–1331 (2024).
- [16] Chen, W., Su, Y., Zuo, J., Yang, C., Yuan, C., Qian, C., Chan, C.-M., Qin, Y., Lu, Y., Xie, R., et al.: Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents, *arXiv preprint arXiv:2308.10848*, Vol. 2, No. 4, p. 6 (2023).
- [17] Huang, D., Bu, Q., Zhang, J. M., Luck, M. and Cui, H.: Agentcoder: Multi-agent-based code generation with iterative testing and optimisation, *arXiv preprint arXiv:2312.13010* (2023).
- [18] Fried, D., Aghajanyan, A., Lin, J., Wang, S., Wallace, E., Shi, F., Zhong, R., Yih, W.-t., Zettlemoyer, L. and Lewis, M.: Incoder: A generative model for code infilling and synthesis, *arXiv preprint arXiv:2204.05999* (2022).
- [19] Zhang, Z., Yao, Y., Zhang, A., Tang, X., Ma, X., He, Z., Wang, Y., Gerstein, M., Wang, R., Liu, G., et al.: Igniting Language Intelligence: The Hitchhiker’s Guide From Chain-of-Thought Reasoning to Language Agents, *arXiv preprint arXiv:2311.11797* (2023).
- [20] Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., et al.: Competition-level code generation with alphacode, *Journal on Science*, Vol. 378, No. 6624, pp. 1092–1097 (2022).
- [21] Islam, R. and Moushi, O. M.: Gpt-4o: The cutting-edge advancement in multimodal llm, *Authorea Preprints* (2024).
- [22] Lin, C.-Y.: Rouge: A package for automatic evaluation of summaries, in *Text summarization branches out*, pp. 74–81 (2004).
- [23] Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J.: Bleu: a method for automatic evaluation of machine translation, in *Association for Computational Linguistics*, pp. 311–318 (2002).

- [24] Dehaerne, E., Dey, B., Halder, S., De Gendt, S. and Meert, W.: Code generation using machine learning: A systematic review, *Ieee Access*, Vol. 10, pp. 82434–82455 (2022).
- [25] Ren, S., Guo, D., Lu, S., Zhou, L., Liu, S., Tang, D., Sundaresan, N., Zhou, M., Blanco, A. and Ma, S.: Codebleu: a method for automatic evaluation of code synthesis, *arXiv preprint arXiv:2009.10297* (2020).
- [26] Kulal, S., Pasupat, P., Chandra, K., Lee, M., Padon, O., Aiken, A. and Liang, P. S.: Spoc: Search-based pseudocode to code, *Advances in Neural Information Processing Systems*, Vol. 32, (2019).
- [27] Roziere, B., Lachaux, M.-A., Chatussot, L. and Lample, G.: Unsupervised translation of programming languages, *Advances in neural information processing systems*, Vol. 33, pp. 20601–20611 (2020).