# When LLM-based Code Generation Meets the Software Development Process

Feng Lin, Dong Jae Kim, Tse-Husn (Peter) Chen

Software PErformance, Analysis and Reliability (SPEAR) Lab

Concordia University, Montreal, Canada

feng.lin@mail.concordia.ca, k_dongja@encs.concordia.ca, peterc@encs.concordia.ca

*Abstract*—Software process models play a pivotal role in fostering collaboration and communication within software teams, enabling them to tackle intricate development tasks effectively. This paper introduces *LCG*, a code generation framework inspired by established software engineering practices. *LCG* leverages multiple Large Language Model (LLM) agents to emulate various software process models, namely $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$. Each model assigns LLM agents specific roles such as requirement engineer, architect, developer, tester, and scrum master, mirroring typical development activities and communication patterns. Through collaborative efforts utilizing chain-of-thought and prompt composition techniques, the agents continuously refine themselves to enhance code quality. Utilizing GPT3.5 as the underlying LLM and baseline (GPT), we evaluate LCG across four code generation benchmarks: HumanEval, HumanEval-ET, MBPP, and MBPP-ET. Results indicate $LCG_{Scrum}$ outperforms other models, achieving Pass@1 scores of 75.2, 65.5, 82.5, and 56.7 in HumanEval, HumanEval-ET, MBPP, and MBPP-ET, respectively—an average 15% improvement over GPT. Analysis reveals distinct impacts of development activities on generated code, with design and code reviews contributing to enhanced exception-handling, while design, testing, and code reviews mitigate code smells. Furthermore, temperature values exhibit negligible influence on Pass@1 across all models. However, variations in Pass@1 are notable for different GPT3.5 model versions, ranging from 5 to over 60 in HumanEval, highlighting the stability of *LCG* across model versions. This stability underscores the importance of adopting software process models to bolster the quality and consistency of LLM-generated code.

*Index Terms*—Large Language Model, Code Generation, Agents, Software Process Model

## I. INTRODUCTION

The recent surge of Large Language Models (LLMs) has sparked a transformative phase in programming and software engineering. Pretrained on vast repositories of code-related datasets, these LLMs have acquired a comprehensive understanding of code, enabling them to excel in diverse code-related tasks. With tools like ChatGPT [26] or LLaMA [37], researchers have demonstrated the potential of LLMs in generating commit messages [49], resolving merge conflicts [33], generating tests [41, 48, 31], method renaming [1], and even facilitating log analytics [21, 22].

Among all development activities, code generation has received much attention due to its potential to reduce development costs. As LLMs are becoming increasingly integral to software development, various techniques have emerged in LLM-based code generation. For example, prompting techniques like few-shot learning [16, 47] have been shown to improve code generation results. In particular, few-shot learning coupled with few-shot sampling [22, 15] or information retrieval augmented technique [25, 4] have been shown to improve code generation. Moreover, one can integrate personalization in the prompt, instructing LLMs to be domain experts in a specific field, which can further improve LLM responses [39, 32]. Such personalization techniques highlight the potential of using multiple LLMs working together to assist in complex software development activities.

Given the complexity of software development, LLM agents stand out among various LLM techniques. Agents are LLM instances that can be customized to carry out specific tasks that replicate human workflow [12, 9]. Recently, multi-agent systems have achieved significant progress in solving complex problems in software development by emulating development roles [12, 9, 29]. MetaGPT, introduced by Hong et al. [12], integrated development workflow using standard operating procedures by assigning specific roles (e.g., a designer or a developer) to LLM agents. Dong et al. [9] developed self-collaboration, which assign LLM agents to work as distinct "*experts*" for sub-tasks in software development. Qian et al. [29] proposed an end-to-end framework for software development through self-communication among the agents.

Despite the promising applications of LLMs in automating software engineering tasks, it is pivotal to recognize that software development is a collaborative and multi-faceted endeavor. In the real world, developers and stakeholders work together, following certain software process models like *Waterfall*, *Test-Driven-Development (TDD)*, and *Scrum*. The process models help facilitate communication and collaboration to ensure the delivery of high-quality products. Even though there is a common community agreement on the pros and cons of each process model [11], the impact of adopting these process models for LLM code generation tasks remains unknown. In particular, will emulating different process models impact the generated code quality in different aspects, such as reliability, code smell, and functional correctness?

While some research has explored integrating multi-agents within LLM frameworks [29, 43, 45], their research focus diverges from the influence of the software process model on code generations for several reasons: 1) Xu et al. [43] do not adhere to specific process models, and 2) both Dong et al. [9] and Qian et al. [29] focus solely on *Waterfall*, neglecting *Test-Driven-Development (TDD)* and *Scrum*, which may have

different impact on code generations. Importantly, none of the aforementioned studies conduct a fine-grain analysis of how different development activities affect code quality metrics, such as code smell and reliability, other than the Pass@1 score. Our study takes further steps to analyze the impacts of different agents within the process models on code generation and their influence on other code quality attributes.

This paper presents a novel multi-agent Large Language Model (LLM) based code generation framework named *LCG*. *LCG* integrates diverse prompt engineering techniques, including chain-of-thought [38, 18], prompt composition [20, 48], and self-refinement [23], with a focus on emulating the flow of development activities in various software process models. Specifically, we implemented three popular process models into *LCG*: $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$. Each process model emulates a real-world development team involving several LLM agents whose roles (i.e., requirement engineer, architect, developer, tester, and scrum master) correspond to common software development activities. The agents are assigned roles with different domain-specific expertise, working collaboratively to produce software artifacts and help other agents review and improve the artifacts in every activity of a process model.

We evaluate *LCG* on four popular code generation benchmarks: *HumanEval* [6], *HumanEval-ET* [8], *MBPP* [2], and *MBPP-ET* [19]. We apply zero-shot learning to avoid biases in selecting few-shot samples [42]. To compare, we also apply zero-shot learning on GPT-3.5 as our baseline (*GPT*). We repeat our experiments five times to account for variability in LLM's responses and report the average value and standard deviation. To study code quality, in addition to Pass@1, we run static code checkers to detect the prevalence of code smells in the generated code. Our evaluation shows that $LCG_{Scrum}$'s generated code achieves the highest accuracy (Pass@1 is 75.2 for HumanEval, 65.5 for HumanEval-ET, 82.5 for MBPP, and 56.7 for MBPP-ET), surpassing *GPT*'s Pass@1 by 5.2% to 31.5%. While *LCG*, in general, is more stable than *GPT*, $LCG_{Scrum}$ exhibits the most stable results with an average standard deviation of only 1.3% across all benchmarks.

We further study the impact of each development activity on code quality. We find that removing the testing activity in the process model results in a significant decrease in Pass@1 accuracy (17.0% to 56.1%). Additionally, eliminating the testing activity leads to a substantial increase in error and warning code smell densities. We also find that the design and code review activities reduce refactor and warning code smells, and improve reliability by adding more exception handling code. Nevertheless, *LCG* consistently outperforms *GPT* by reducing code smells and enhancing exception handling. Finally, we find that the GPT model version plays a significant role in the quality of generated code, and *LCG* helps ensure stability across different versions of LLMs and temperature values.

We summarize the main contributions of this paper as follows:

1) **Originality**: We introduce a multi-agent framework called *LCG*, incorporating software process models from real-world development practice. We integrate agents acting as requirement engineers, architects, developers, testers, and scrum masters, and study how their interaction improves code generation and code quality.

2) **Technique**: We integrate prompt engineering techniques like chain-of-thought, prompt composition, and self-refinement to facilitate interactions among the agents. We implement three recognized process models: $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$, but the technique can be easily extended to emulate other process models or development practices (e.g., DevOps).

3) **Evaluation**: We conduct a fine-grained evaluation on the quality of the generated code using four popular code generation benchmarks: *Humaneval* [6], *Humaneval-ET* [8], *MBPP* [2], *MBPP-ET* [19], comparing agent interactions and their effect on both accuracy (Pass@1) and other code quality metrics (e.g., smells). We manually checked the generated code and discussed the reasons for test failures. Finally, we examined how different model versions and temperature settings affect code generation stability.

4) **Data Availability**: To encourage future research in this area and facilitate replication, We made our data publicly available online https://anonymous.4open.science/r/FlowGen-LLM-E842.

**Paper Organization.** Section II discusses background and related work. Section III. provides the details of *LCG*. Section IV evaluates our *LCG*. Section V discusses threats to validity. Section VI concludes the paper.

## II. BACKGROUND & RELATED WORKS

In this section, we discuss the background of software process models and LLM agents. We also discuss related work on LLM-based code-generation.

### A. Background

**Software Development Process.** Software development processes encompass methodologies and practices that development teams use to plan, design, implement, test, and maintain software. The primary goal of a software process is to assist the development teams in producing high-quality software. Generally, different software process models involve the same set of development activities, such as requirement, design, implementation, and testing, but differ in how the activities are organized. Because of the variation, each software process model has its strengths and weaknesses based on the project type, teams, and experience [11].

In particular, three well-known and widely adopted software process models were created over the years: *Waterfall*, *Test-Driven-Development (TDD)*, and *Scrum*. *Waterfall* [3] is often used in safety-critical systems where development teams must adhere to a linear path, and each software development activity builds upon the previous one. *TDD* and *Scrum* are both variants of the agile development model. Compared to *Waterfall*, agile process models focus more on iterative and incremental development and adapting to change. *TDD* [24]

emphasize writing tests before writing the actual code to improve software design and quality. *Scrum* highlights the importance of collaboration and communication in software development. *Scrum* prescribes for teams to break work into time-boxed iterations called sprints. During these sprints, teams focus on achieving specific goals (e.g., user stories), ensuring a continuous discussion among teams to handle any unexpected risks throughout the development process.

**LLM Agents.** LLM agents are artificial intelligence systems that utilize LLM as their core computational engines to understand questions and generate human-like responses. LLM agents can refine their responses based on feedback, learn from new information, and even interact with other AI agents to collaboratively solve complex tasks [12, 29, 43, 27]. Through prompting, agents can be assigned different roles (e.g., a developer or a tester) and provide more domain-specific responses that can help improve the answer [12, 39, 32].

One vital advantage of agents is that they can be implemented to interact with external tools. When an agent is reasoning the steps to answer a question, it can match the question/response with corresponding external tools or APIs to construct or refine the response. For instance, an LLM agent that represents a data analysis engineer can apply logical reasoning to generate corresponding SQL query statements, invoke the database API to get the necessary data, and then answer questions based on the returned result. When multiple LLM agents are involved, they can collaborate and communicate with each other. Such communication is essential for coordinating tasks, sharing insights, and making collective decisions. Hence, defining how the agents communicate can help optimize the system's overall performance [40], allowing agents to undertake complex projects by dividing tasks according to their domain-specific skills or knowledge.

The software development process plays a crucial role in the creation of high-quality software, fundamentally involving communication among various development roles. Given the demonstrated capability of Large Language Model (LLM) agents to mimic domain experts in specific fields [39, 12, 29], this study leverages LLM agents to represent diverse development roles and conduct their associated duties. Our research establishes a collaborative team of LLM agents designed to emulate these process models and roles, aiming to enhance code generation.

### B. Related Works

Code generation is a thriving field of research because of its potential to reduce development costs. Apart from training more powerful LLM to improve code generation, *prompt-based* and *agent-based* code generation techniques are two of the most prevalent directions.

**Prompt-based Code Generation.** Prompt-based code generation employs a range of techniques to refine prompts, ultimately leading to the generation of expected code. For example, Li et al. [18] propose using structured prompts containing code information (e.g., branch and loop structures) to improve the generated code. Nashid et al. [25] retrieval

code demos similar to the given task and include them in the prompt to improve code generation. Ruiz et al. [30] use translation techniques for program repair, where buggy programs are first translated into natural language or other programming languages. Then, the translated code is used as a prompt to generate new/fixed code with the same feature. Schäfer et al. [31] iteratively refine prompts based on feedback received from interpreters or test execution results. Kang et al. [15] provide specific instructions, test method signature, and bug report as part of the prompt for generating test code to reproduce bugs. Xie et al. [41] parse the code to identify the focal method and related code context, which are given in the prompt for test code generation. Yuan et al. [48] apply a prompt composition technique by first asking an LLM to provide a high-level description of a method, and then the description is used as part of the prompt to enhance test code generation.

**Agent-based Code Generation.** Agent-based code generation emphasizes on the importance of role definition and communication among multiple LLM agents. Some approaches incorporate external tools as agents. For example, Huang et al. [14] introduce the test executor agent, employing a Python interpreter to provide test logs for LLMs. Zhong et al. [50] introduces a debugger agent that utilizes a static analysis tool to build control flow graph information, guiding LLMs in locating bugs. Meanwhile, other studies [12, 29, 9] task LLMs as agents by emulating diverse human roles, including analysts, engineers, testers, project managers, chief technology officers (CTOs), etc. These approaches follow the *Waterfall* model to communicate among these roles, with variation in the prompts and roles, ultimately improving code generation.

In comparison, our research leverages LLM agents to emulate multiple software development process models, while prior research focuses only on the *Waterfall* model [12, 29, 9]. We implement several prompting techniques, but more importantly, we emphasize on how various process models and the associated development activities affect the generated code. Different from prior works which only study functional correctness, we study several additional dimensions of code quality, including code design, code smell, convention issues, and reliability. We also explore why the generated code fails the tests and the sensitivity of the results across LLM model versions and temperature values.

While there are many code generation results using the same benchmarks, our evaluation results cannot be directly compared with some other agent-based code generation research [12, 13, 14] due to missing information on the model version, temperature value, steps to post-process the generated code, specific prompts, or how the few-shot samples are selected. Nevertheless, *LCG* still shows improved Pass@1 compared to existing works that provide clear information. For instance, Dong et al. [9] (uses the same GPT3.5 version as our study) applied few-shot learning with LLM agents and received a Pass@1 of 68.2% in MBPP. In contrast, while $LCG_{Scrum}$ uses zero-shot learning, we achieved a Pass@1 of 82.5% in MBPP. We observe the same level of improvements

in all the benchmarks. Similarly, Li et al. [18], which used the same model version and temperature value as we do, applied structured chain-of-thought to enhance reasoning ability and achieved a Pass@1 of 60.6% in HumanEval. $LCG_{Scrum}$ attained 78.5% under the same model version and temperature setting. In short, due to variations of the LLMs and missing information in prior works, we focus on comparing the $LCG$'s improvement against *GPT* (directly using GPT3.5) in this study.

## III. METHODOLOGY

We propose *LCG*, an agent-based code generation technique based on emulating different software processes. Figure 1 shows the overview of *LCG*: (1) define the roles and their responsibilities; (2) use LLM agents to represent these roles; and (3) complete the interactions among these agents according to the software process models. In each development activity, we implement chain-of-thought and self-refinement to improve the quality of the generated artifacts. In particular, we study and compare three software process models: *Waterfall*, *TDD*, and *Scrum*. Nevertheless, *LCG* can be easily adapted to different process models. We use zero-shot learning in all our experiments to avoid biases in selecting data samples. Below, we discuss *LCG* in detail.

### A. Using LLM Agents to Represent Different Development Roles in Software Process Models

In *LCG*, we create LLM agents who are responsible for the main development activities: requirement, design, implementation, and testing. Hence, to emulate a software process model, we reorganize the communication and interaction among different agents. The benefit of such a design is that it maximizes the extensibility and reuseable of the agents, and *LCG* can be easily adapted to different process models.

We implement four agents whose role corresponds to the common development activities: *Requirement Engineer*, *Architect*, *Developer*, and *Tester*. For *Scrum*, we introduce an additional role – *Scrum Master*. We design these roles so that they use the the same prompt template with role-specific details as shown below:

```
1 {
2   "Role": "You are a [role] responsible for [task]",
3   "Instruction": "According to the Context please [
     role-specific instruction]",
4   "Context": "[context]"
5 }
```

In this prompt template, `role` refers to one of the roles (e.g., Requirement Engineer) that corresponds to the development activity, and `task` describes the duties for the `role` (e.g., analyze and create requirement documents). `Instruction` leverages chain-of-thought reasoning [38, 18] and refers to role-specific instruction listed in steps, such as 1) analyzing the requirement and 2) writing the requirement documentation. Finally, `Context` contains the programming question, the agent conversation history, or the agent-generated artifacts. `Context` includes all necessary information that

helps the agents to make a next-step decision based on the current conversation and generated results.

Table I shows the `tasks`, `instructions`, and `contexts` for every development role. In general, every role takes the output from the prior development activities as input (i.e., `context`). For example, an architect writes a design document based on the requirement document generated by a requirement engineer. We design developers and testers to have multiple tasks. Developers are responsible for writing code and fixing/improving the code based on suggestions. We design testers using a prompt composition technique, which is shown to improve the LLM-generated result [20, 48]. First, testers design a test. Then, testers write and execute the tests based on the design. Finally, testers generate a test failure report. Developers receive the test failure report to fix the code. In addition to the tasks described in Table I, all roles have one common task, which is to provide feedback to other roles for further improvement (e.g., for code review).

### B. Communications Among Agents

One of the most important aspects of LLM agents is how the agents communicate. A recent survey paper [40] shows that one common communication pattern is sequential, i.e., ordered, where one agent communicates to the next in a fixed order. Another pattern is disordered, where multiple agents participate in the conversation. Each agent gets the context separately and outputs the response in a shared buffer. Then, the responses can be summarized and used in the next decision-making process. Based on the software process models and the two aforementioned communication patterns, we implement three interaction models for the agents: $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$.

#### 1) $LCG_{Waterfall}$

We follow the intuition behind the *Waterfall* model and implement $LCG_{Waterfall}$ as an ordered communication among the agents. Given a programming problem, the problem goes through the requirement analysis, design, implementation, and testing. One thing to note is that the test result from our generated tests is redirected to the developer agent in our implementation of $LCG_{Waterfall}$ so developers can fix and improve the code.

#### 2) $LCG_{TDD}$

In the design of $LCG_{TDD}$, we follow the ordered communication pattern and organize the development activities so that testing happens after design and before implementation. Once the tests are written, the developer agent considers the test design when implementing the code. When the implementation is finished, we execute the tests. If a test fails, the developer agent is asked to examine the code and resolve the issue.

#### 3) $LCG_{Scrum}$

Compared to *Waterfall* and *TDD*, *Scrum* involves one additional role, the *Scrum Master*. There are also additional *Sprint meetings* among the agents. Note that, different from $LCG_{Waterfall}$, we use the agile terminologies in the prompt (e.g.,
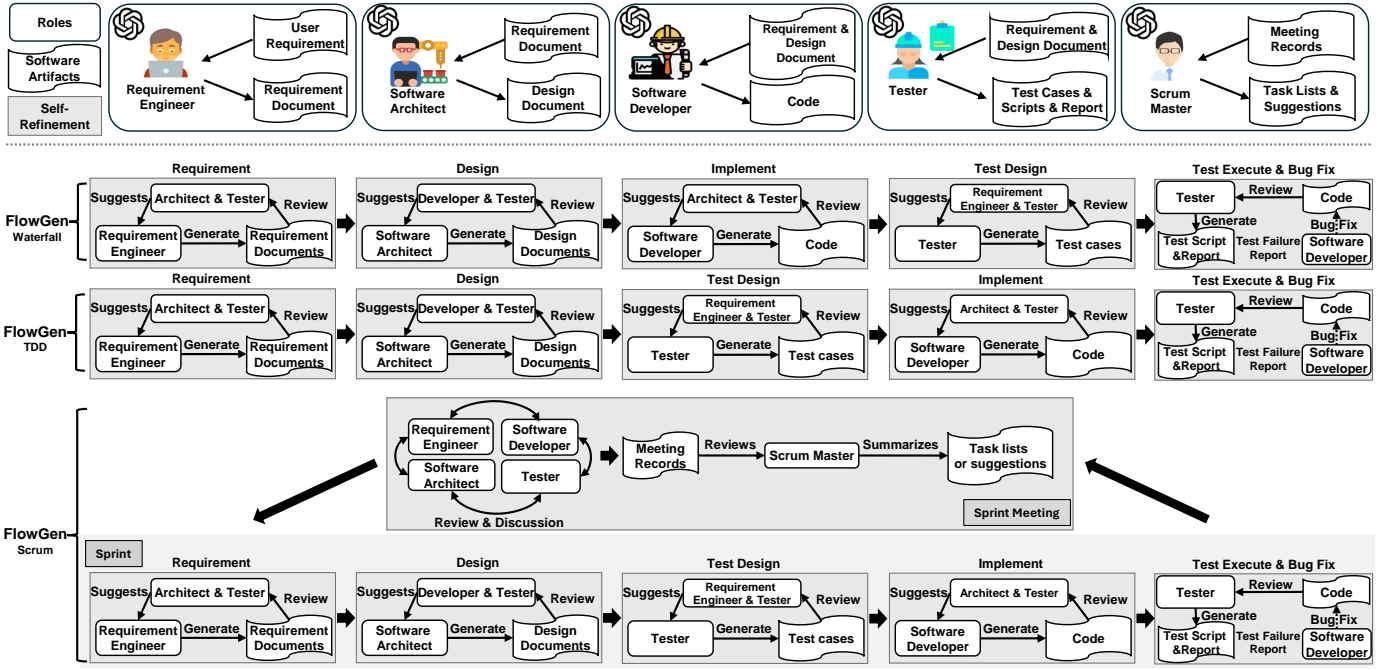
Fig. 1: An overview of $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$.

TABLE I: Tasks, instructions, and corresponding contexts that are used for constructing the prompts for the development roles.

| Role | Task | Instruction | Context |
|---|---|---|---|
| Requirement Engineer | Analyze and generate requirement documentation from the context. | 1) Analyze the requirement; and 2) Write a requirement document. | Programming problem description. |
| Architect | Design the overall structure and high-level components of the software. | 1) Read the context documents; and 2) Write the design document. The design should be high-level and focus on guiding the developer in writing code. | Requirement document. |
| Developer | Write code in Python that meets the requirements. | 1) Read the context documents; and 2) Write the code. Ensure that the code you write is efficient, readable, and follows best practices. | Requirement and design documents. |
| | Fix the code so that it meets the requirements. | 1) Read the test failure reports and code suggestions from the context; and 2) Rewrite the code. | Original code, test failure report, and suggestions for improvement. |
| Tester | Design tests to ensure the software satisfies feature needs and quality. | 1) Read context documents; and 2) Design test cases. | Requirement and design documents. |
| | Write a Python test script using the unittest framework. | 1) Read the context documents; 2) Write a Python test script; and 3) Follow the input and output given by the requirement. | Test case design and requirement documents. |
| | Write a test failure report. | 1) Read the test execution result; and 2) Analyze and generate a test failure report. | Test execution result. |
| Scrum Master | Summarize and break down the discussion into a task list for the scrum team. | 1) Read and understand the context; and 2) Define the tasks for development roles. | Meeting discussion. |

we use user story instead of requirement document) when implementing $LCG_{Scrum}$. We follow a disordered communication pattern in the design of $LCG_{Scrum}$, because, in sprint meetings, every development role can provide their opinion (e.g., to simulate the planning poker process). Every development role, except the *Scrum Master*, reads the common context (e.g., description of the programming problem) from a common buffer. Then, every role provides a discussion comment and is saved back in the buffer. Therefore, every role is aware of all the comments. Then, the *Scrum Master* summarizes the entire discussion and derives a list of user stories for each development role. During the sprint, similarly to *Waterfall* and *TDD*, the four development roles carry out the development activities in sequence. At the end of the sprint, the agents will start another sprint meeting to discuss the next steps, such as releasing the code or needing to fix the code because of test failures.

*4) Self-Refinement*

We implement self-refinement [23], which tries to refine the LLM-generated result through iterative feedback, to further improve the generated artifacts from every development activity. In all three variations of *LCG*, we assign other agents

to review the generated artifacts for every development activity and provide improvement suggestions. We assign the agents from both the downstream activity and the tester to examine the generated artifacts and provide suggestions. The suggestions are then considered for the re-generation of the artifacts. We include the tester in every development activity to emulate the DevTestOps practice [36], where testers are involved in all development activities and provide feedback on the quality aspects. For example, once a requirement engineer generates a requirement document, both the architect and tester would read the document and provide suggestions for improvement. Then, the requirement engineer will re-generate the requirement document based on the previously generated document and suggestions. At the development and testing activity, the tester will generate a test failure report if any of the LLM-generated tests fail or if the code cannot be executed (e.g., due to syntax error). The test failure report is then given to the developer for bug fixing. We repeat the process $t$ times to self-refine the generated code. In our implementation, we currently set $t$=3. If the code still cannot pass the test, we repeat the entire software development process.

## C. Implementation and Experiment Settings

**Environment.** We use GPT3.5 (version `gpt-3.5-turbo-1106`) as our underlying LLM due to its popularity and wide usage in code generation research [18, 34, 9]. We leverage the APIs provided by OpenAI to interact with GPT. We send prompts using JSON format [39] and send all the conversation history as part of the prompt [12]. We set the temperature value to 0.8 and explore the effect of the temperature value in RQ3. We implemented *LCG* using *Python 3.9* and version *0.28.1* for the OpenAI API.

**Benchmark Datasets.** We follow prior studies [51, 18, 50] and evaluate the code generation result using four benchmarks: HumanEval, HumanEval-ET, MBPP (Mostly Basic Python Programming), and MBPP-ET. These benchmarks contain both the programming problems and tests to evaluate the correctness of the code. Given a programming problem, we consider that a generated code snippet is correct if it can pass all the provided tests. HumanEval [6] has 164 programming problems, and MBPP [2] has 427 programming problems and three test cases for each problem. We also use the dataset published by Dong et al. [8], where they use the same problems as HumanEval and MBPP but offer stronger evaluation test cases (around 100 test cases for each problem, called HumanEval-ET and MBPP-ET). All these benchmarks use Python as the programming language. Each programming problem contains the INPUT pairs <method signature, method description, invoke examples> and expect the code as OUTPUT.

**Evaluation Metric.** To evaluate the quality of the generated code, we use the Pass@K metric [9, 6]. Pass@K evaluates the first K generated code's functional accuracy (i.e., whether the generated code can pass all the test cases). In this work, we set K=1 to evaluate if the first generated code can pass all the provided test cases. A Pass@1 of 100 means 100% of the generated code can pass all the tests in the first attempt.

TABLE II: Average and standard deviation of the Pass@1 accuracy across five runs, with the best Pass@1 marked in bold. The numbers in the parentheses show the percentage difference compared to *GPT*. Statistically significant differences are marked with a "*".

| | HumanEval | HumanEval-ET | MBPP | MBPP-ET |
|---|---|---|---|---|
| *GPT* | 64.4±3.7 | 49.8±3.0 | 77.5±0.8 | 53.9±0.7 |
| $LCG_{Waterfall}$ | 69.5±2.3 (+7.9%)* | 59.4±2.5 (+19.2%)* | 76.3±0.9 (-1.5%) | 51.1±1.7 (-5.2%)* |
| $LCG_{TDD}$ | 69.8±2.2 (+8.4%)* | 60.0±2.1 (+20.5%)* | 76.8±0.9 (-1.0%) | 52.8±0.7 (-2.1%)* |
| $LCG_{Scrum}$ | **75.2±1.1 (+16.8%)*** | **65.5±1.9 (+31.5%)*** | **82.5±0.6 (+6.5%)*** | **56.7±1.4 (+5.2%)*** |

## IV. RESULTS

In this section, we present the research questions aimed at evaluating *LCG*. For each RQ, we discuss the motivation, the approach we take to answering it, and the results and discussion of the findings.

*RQ1: What is the code generation accuracy of LCG?*

*Motivation.* In this RQ, we emulate the three process models using LLM agents and compare their results on code generations. Such results may provide invaluable evidence for future researchers seeking to optimize process models for code generation within their specific business domain.

*Approach.* As a baseline for comparison, we directly give the programming problems to ChatGPT (which we refer to as *GPT*). Although prior works [15, 17, 25, 18] shows that few-shot learning can improve the results from LLMs, there can be biases on how the few-shot samples are selected [42]. Hence, we use zero-shot learning in our experiment. To control for randomness in the experiment, we ensure all these experiments use the same temperature value (t=0.8) and the same model version (*gpt-3.5-turbo-1106*). Finally, we repeat each *LCG* five times and report the average Pass@1 and standard deviation across the runs. We also conduct a student's t-test to study if *LCG*'s results are statistically significantly different from *GPT*.

*Result.* **$LCG_{Scrum}$ shows a consistent improvement over GPT, achieving 5.2% to 31.5% improvement in Pass@1.** Table II shows the Pass@1 accuracy of *LCG* across different process models on the benchmark datasets studied. As shown in the Table II, for HumanEval and HumanEval-ET, all of the studied process models have 7.9% to even 31.5% improvement in Pass@1 compared to *GPT*, and the improvements are all statistically significant (p <= 0.05). For MBPP and MBPP-ET, $LCG_{Scrum}$ also has statistically significant improvements of 5.2% to 6.5%, even though we see a slight decrease when adopting $LCG_{Waterfall}$ and $LCG_{TDD}$ to MBPP and MBPP-ET.

Despite slight variations in code generation responses form LLM across executions, **we find stable standard deviations of Pass@1, ranging from 0.6% to 3.7% across all process models and benchmarks**. In particular, $LCG_{Scrum}$ has the lowest standard deviation (0.6% to 1.9%, an average of 1.2%), while *GPT* has the highest standard deviation (0.5% to 3.7%, an average of 2%). Following $LCG_{Scrum}$, $LCG_{Waterfall}$ has the second highest standard deviation, with $LCG_{TDD}$ is

TABLE III: *LCG* test failure categorization. Failure types are generated from Python Interpreter [28]. Darker red indicates higher percentages of the failure categories in the generated code across the models. Percentages are calculated by the ratio of specific failure types to the total number of failed tests across different process models.

| Benchmarks | Model | Failure Categories | | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Assertion | Syntax | Name | Type | Index | Value | Recursion | Attribute | |
| HumanEval | RawGPT | 36 (24%) | 18 (100%) | 11 (58%) | 2 (17%) | 2 (50%) | 1 (12%) | 0 (0%) | 0 (0%) | 70 (33%) |
| | Waterfall | 39 (26%) | 0 (0%) | 3 (16%) | 3 (25%) | 0 (0%) | 1 (12%) | 0 (0%) | 0 (0%) | 46 (22%) |
| | TDD | 39 (26%) | 0 (0%) | 4 (21%) | 3 (25%) | 1 (25%) | 5 (62%) | 0 (0%) | 0 (0%) | 52 (24%) |
| | Scrum | 38 (25%) | 0 (0%) | 1 (5%) | 4 (33%) | 1 (25%) | 1 (12%) | 0 (0%) | 0 (0%) | 45 (21%) |
| HumanEval-ET | RawGPT | 39 (21%) | 9 (82%) | 13 (43%) | 1 (25%) | 4 (57%) | 2 (18%) | 0 (0%) | 0 (0%) | 68 (27%) |
| | Waterfall | 49 (26%) | 1 (9%) | 7 (23%) | 0 (0%) | 1 (14%) | 3 (27%) | 0 (0%) | 0 (0%) | 61 (24%) |
| | TDD | 50 (26%) | 1 (9%) | 6 (20%) | 2 (50%) | 1 (14%) | 4 (36%) | 0 (0%) | 0 (0%) | 64 (25%) |
| | Scrum | 52 (27%) | 0 (0%) | 4 (13%) | 1 (25%) | 1 (14%) | 2 (18%) | 0 (0%) | 0 (0%) | 60 (24%) |
| MBPP | RawGPT | 66 (23%) | 7 (70%) | 7 (41%) | 7 (25%) | 2 (67%) | 0 (0%) | 1 (100%) | 1 (50%) | 91 (26%) |
| | Waterfall | 76 (27%) | 1 (10%) | 5 (29%) | 6 (21%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 88 (25%) |
| | TDD | 86 (30%) | 1 (10%) | 5 (29%) | 8 (29%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 100 (29%) |
| | Scrum | 58 (20%) | 1 (10%) | 0 (0%) | 7 (25%) | 1 (33%) | 0 (0%) | 0 (0%) | 1 (50%) | 68 (20%) |
| MBPP-ET | RawGPT | 133 (24%) | 8 (57%) | 22 (26%) | 26 (24%) | 2 (50%) | 3 (8%) | 1 (100%) | 1 (100%) | 196 (24%) |
| | Waterfall | 148 (27%) | 2 (14%) | 22 (26%) | 27 (25%) | 0 (0%) | 15 (39%) | 0 (0%) | 0 (0%) | 214 (27%) |
| | TDD | 147 (27%) | 2 (14%) | 23 (27%) | 28 (26%) | 1 (25%) | 10 (26%) | 0 (0%) | 0 (0%) | 211 (26%) |
| | Scrum | 123 (22%) | 2 (14%) | 18 (21%) | 28 (26%) | 1 (25%) | 10 (26%) | 0 (0%) | 0 (0%) | 182 (23%) |

ranking third. In conclusion, although the models generally have consistent Pass@1 across runs, *LCG_Scrum* consistently produces the most stable results.

***There are potential issues in the tests provided by the benchmarks, which may hinder the Pass@1 of LCG.*** Table III provides a breakdown of failure types from the Python Interpreter [28] across various process models and benchmarks. While we repeat our experiment 5 times, the standard deviation across runs is low; hence, we represent test failure from only one of the runs. Aligned with the findings from Table II, *LCG_Scrum* has the lowest `AssertionError` compared to other models. We also notice that `SyntaxError` is more evident in *GPT* across all benchmarks, as expected due to the absence of a code review and testing process. However, there are still higher test failures in *LCG_Waterfall* and *LCG_TDD* caused by increased occurrences of `ValueError`, `TypeError`, and `NameError`, for MBPP and MBPP-ET as seen in Table II.

Upon manual investigation of test failures, we discover that *LCG_Waterfall* and *LCG_TDD* **introduce various validations and enforce programming naming conventions in the generated code, *which may help improve code quality but cause tests to fail***. For example, Listing 1 depicts the provided tests and the generated code/test for MBPP-582, of which the objective is to *"Write a function to check if a dictionary is empty or not"*. While *GPT* passed the provided tests, *LCG_Waterfall* and *LCG_TDD* failed. This failure is because our generated code contains validation to check that the input should be of type *dict*. However, the MBPP-582 provided test uses an input of the type *set*, which fails the validation, causing a `TypeError` exception. Moreover, *LCG_Waterfall* and *LCG_TDD* enforce the common naming convention format and more meaningful function name (e.g., `my_dict` v.s. `is_dict_empty`), both of which causes `NameError` exception due to missing declaration, causing test failure. More interestingly, we find

that such code standardization may also be misled by the requirement provided by the benchmark itself. For example, the MBPP-582 requirement specifies expected input as `dict`, yet provides a `set` type as the input to the test. The LLM code generation indeed captures this correct requirement by validating that input must be of type `dict`. Such inconsistency in the benchmark may reduce the Pass@1.

```
# MBPP-582: check if a dictionary is empty
# MBPP Test Case
def Test():
    assert my_dict({10})==False # {10} is a set not a dict
# rawGPT's answer
def my_dict(dict1):
    return len(dict1) == 0
# Waterfall/TDD's answer, enforces the input type must be a
    dictionary
def is_dict_empty(input_dict): # function name is renamed
    from my_dict
    if not isinstance(input_dict, dict):
        raise TypeError("Input is not a dictionary")
    return True if not input_dict else False
```

Listing 1: MBPP-582 Test Failure due to *Strict Input Validation* and *Wrong Function Name*.

In MBPP-794 (Listing 2), test cases provided by MBPP-ET change the return value from a `boolean` (as is the case in MBPP) to the string ``match!''. Moreover, in MBPP-797, MBPP-ET's test capitalized the last word in uppercase (`range` v.s. ``R''ange). Such non-standard evaluation leads to unfair test results (leads to failure), which may bias the experimental results for MBPP-ET. Such bias suggests that the decrease in Pass@1 rates for MBPP-ET is not solely due to an increase in the number of provided tests.

```
# Example1: Changed return type from boolean to string
assert text_starta_endb("aabbbb") # MBPP 794
assert text_starta_endb("aabbbb")==('match!') # MBPP-ET 794
# Example2: Capitalized the last character in function name
assert sum_in_range(2,5) == 8 # MBPP-797
assert sum_in_Range(2,5) == 8 # MBPP-ET-797
```

Listing 2: MBPP-794 & MBPP-797 Test Failure due to Irregularity in Test Cases.

We conduct a pilot study on the impact of *Wrong Function Name* and *Strict Input Validation* by manually fixing these issues in the generated code. We find that, after fixing these issues, most generated code can pass all the tests provided by the benchmarks, indicating correct functional behaviors. Fixing the issues largely improve *LCG*'s Pass@1: 16 to 28 improvement in HumanEval, 29 to 35 in HumanEval-ET, 15 to 21 in MBPP, and 28 to 42 in MBPP-ET. Namely, *LCG*'s Pass@1 can achieve over 90 to 95 across all benchmarks. Even though we also observe improvement in *GPT*'s Pass@1, the three *LCG* models had greater improvement. On average, the three models have a Pass@1 that is 14%, 8%, 1.4%, and 5% better than *GPT* in HumanEval, HumanEval-ET, MBPP, and MBPP-ET, respectively. The aggregated results underscore the deficiencies in the benchmarks, suggesting that the current Pass@1 score of *LCG* (Table II) could represent a lower bound. These preliminary findings highlight the potential of *LCG* and suggest that future research should improve the benchmarks by incorporating input checking or consistent

naming convention into the tests and subsequently re-evaluate existing code generation techniques.

> $LCG_{Scrum}$ achieves the best results, with a Pass@1 that is 5.2% to 31.5% better than *GPT*. $LCG_{Scrum}$ also has the most stable results (average standard decision of 1.3% across all benchmarks) among all models. Notably, while $LCG_{Waterfall}$ and $LCG_{TDD}$ enhance code quality, such improvements may result in test failures.

*RQ2: How do different development activities impact the quality of the generated code?*

*Motivation.* As observed in RQ1, various process models can indeed affect the functional correctness (Pass@1) of the generated code. However, it is equally crucial to understand code quality issues such as code smells and the impact of a development activity on the generated code. This understanding is essential for assessing whether the generated code adheres to industry best practices. Moreover, such insight may offer valuable opportunities for enhancing the design, readability, and maintainability in auto-generated code.

*Approach.* To study the impact of each development activity on code quality, we remove each activity separately and re-execute *LCG*. For example, we first remove the requirement activity in $LCG_{Waterfall}$ and execute $LCG_{Waterfall}$. Then, we add the requirement activity back and remove the design activity. We repeat the same process for every development activity. Note that we cannot remove the coding activity since our goal is code generation. Hence, we removed code review at the end of the coding activity.

Code quality considers numerous facets beyond mere functional correctness [46, 44]. Other factors such as code smells, maintainability, and readability are also related to code quality. Hence, to gain a comprehensive understanding of how code quality changes, in addition to Pass@1, we 1) apply a static code analyzer to detect code smells in the generated code and 2) study code reliability by analyzing the exception handling code. To study code smell, design, and readability, we apply Pylint 3.0.4 [35, 7] (a Python static code analyzer) on the generated code. Pylint classifies the detected code smells into different categories such as *error*, *warning*, *convention*, and *refactor*.

We study how the number of detected code smells in each category changes when removing an activity. Since the generated code may have different lengths, we report the density of the code smells in each category. We calculate the code smell density as the total number of code smell instances in a category (e.g., *error*) divided by the total lines of code. To study reliability, we calculate the density of handled exceptions (total number of handled exceptions divided by the total lines of code) since exceptions are one of the most important mechanisms to detect and recover from errors [10]. For better visualization, we present the density results as per every 10 lines of code. We also ensure reliability of our results by repeating all of the aforementioned approach three times.

*Result.* **Testing has the largest impact on the functional correctness of the code, while other development activities only have small impacts. Removing testing causes Pass@1 to decrease by 17.0% to 56.1%.** Table IV presents changes in Pass@1 and the densities of code smell and handled exceptions. We show the results for HumanEval and MBPP because they share the same programming problems and generated code with the other two benchmarks. Among all development activities, testing has the largest impact on Pass@1, where removing testing causes a large decrease in Pass@1 (17.0% to 56.1% decrease). The finding implies that LLM's generated tests are effective in improving the functional correctness of code. In both benchmarks, removing sprint meetings in $LCG_{Scrum}$ also causes Pass@1 to drop. However, removing other activities only has a small and inconsistent effect on Pass@1. For example, in HumanEval, removing requirement, design, and code review generally causes Pass@1 to decrease (except for $LCG_{Scrum}$), but removing these activities improves Pass@1 in MBPP. The findings show that, in HumanEval and MBPP, most development activities do not significantly contribute to the functional correctness of the generated code.

We further study how removing different activities impacts the code smell densities. As shown in Table IV, eliminating test activities significantly boosts *error* and *warning* smell densities by an average of 702.2% and 52.2%, respectively. Omitting design raises refactor smell density by an average of 7.1%, and skipping code review leads to a 14.2% average increase in *warning* density. However, removing other development activities shows either a small or inconsistent impact. In short, the findings show that **adding design, testing, and code review can help reduce the density of code smell in the generated code**.

*Having design and code review activities significantly improves reliability by increasing the density of handled exceptions, while other development activities only have small or no impacts.* Removing design and code review activities separately causes the handled exception density to decrease by 20.0% to 43.2% and 14.5% to 27.3%, respectively. Namely, design and code review activities add exception handling in the generated code, which may help improve reliability. Removing other development activities shows a mixed relationship with the density of the handled exception. For example, removing testing in $LCG_{TDD}$ causes an increase of handled exception density by 25.8% in MBPP (i.e., testing removes exception handling code) while causing a decrease of 30.3% in HumanEval (i.e., testing adds exception handling code). While the effect of each development may be related to the nature of the benchmarks, our findings show that, in both benchmarks, **adding design and code review activities can help improve code reliability by handling more exceptions in the generated code**. *LCG shows consistent improvement over GPT in the quality of the generated code: decreasing the density of error/warning/convention/refactor code smells (6.7% to 96.0%) while significantly increasing handled exception density.* The code generated by *GPT* has higher *error/warning/convention/refactor* code smell densities

TABLE IV: Pass@1 and *Error/Warning/Convention/Refactor/Handled-Exception* density (per 10 lines of code) in the full *LCG* (with all the development activities) and after removing a development activity. A ***lower*** error/warning/convention/refactor is preferred, and a ***higher*** handled-exception is preferred. Darker red indicates a larger decrease in percentages, while darker green indicates a larger increase in percentages.

| Model | Dev. Activities | HumanEval | | | | | | MBPP | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Pass@1 | Error | Warning | Convention | Refactor | Handled Exception | Pass@1 | Error | Warning | Convention | Refactor | Handled Exception |
| *GPT* | – | 64.4 | 0.25 | 0.19 | 0.39 | 0.30 | 0.00 | 77.47 | 0.22 | 0.20 | 1.18 | 0.30 | 0.01 |
| *LCG_Waterfall* | full | 69.5 | 0.01 | 0.12 | 0.24 | 0.21 | 0.37 | 76.35 | 0.03 | 0.12 | 0.47 | 0.23 | 0.67 |
| | rm-requirement | -1.2 (1.7%) | 0.0 (0.0%) | -0.01 (8.3%) | -0.02 (8.3%) | -0.01 (4.8%) | -0.06 (16.2%) | +0.7 (0.9%) | -0.02 (66.7%) | -0.02 (16.7%) | 0.0 (0.0%) | +0.02 (8.7%) | -0.09 (13.4%) |
| | rm-design | -1.2 (1.7%) | +0.01 (100.0%) | +0.02 (16.7%) | +0.02 (8.3%) | 0.0 (0.0%) | -0.15 (40.5%) | -1.64 (2.1%) | -0.01 (33.3%) | 0.0 (0.0%) | +0.02 (4.3%) | +0.01 (4.3%) | -0.2 (29.9%) |
| | rm-codeReview | -2.4 (3.5%) | 0.0 (0.0%) | 0.0 (0.0%) | -0.03 (12.5%) | +0.02 (9.5%) | -0.09 (24.3%) | +0.46 (0.6%) | -0.02 (66.7%) | +0.07 (58.3%) | -0.02 (4.3%) | -0.02 (8.7%) | -0.17 (25.4%) |
| | rm-test | -39.0 (56.1%) | +0.17 (1700.0%) | +0.01 (8.3%) | +0.07 (29.2%) | +0.1 (47.6%) | +0.1 (27.0%) | -23.7 (31.0%) | +0.09 (300.0%) | +0.05 (41.7%) | +0.36 (76.6%) | +0.01 (4.3%) | -0.01 (1.5%) |
| *LCG_TDD* | full | 69.8 | 0.01 | 0.08 | 0.33 | 0.27 | 0.33 | 76.77 | 0.04 | 0.13 | 0.71 | 0.28 | 0.62 |
| | rm-requirement | -2.9 (4.2%) | +0.01 (100.0%) | +0.01 (12.5%) | 0.0 (0.0%) | -0.03 (11.1%) | -0.1 (30.3%) | +1.92 (2.5%) | -0.02 (50.0%) | +0.02 (15.4%) | +0.03 (4.2%) | 0.0 (0.0%) | -0.27 (43.5%) |
| | rm-design | -2.9 (4.2%) | 0.0 (0.0%) | +0.02 (25.0%) | -0.06 (18.2%) | +0.03 (11.1%) | -0.13 (39.4%) | +1.22 (1.6%) | -0.02 (50.0%) | -0.03 (23.1%) | -0.03 (4.2%) | +0.04 (14.3%) | -0.24 (38.7%) |
| | rm-codeReview | -0.5 (0.7%) | -0.01 (100.0%) | +0.02 (25.0%) | -0.04 (12.1%) | +0.03 (11.1%) | -0.09 (27.3%) | +0.98 (1.3%) | -0.03 (75.0%) | +0.01 (7.7%) | -0.05 (7.0%) | +0.1 (35.7%) | -0.09 (14.5%) |
| | rm-test | -11.9 (17.0%) | +0.07 (700.0%) | +0.04 (50.0%) | -0.02 (6.1%) | -0.05 (18.5%) | -0.1 (30.3%) | -17.3 (22.5%) | +0.11 (275.0%) | +0.14 (107.7%) | +0.16 (22.5%) | -0.01 (3.6%) | +0.16 (25.8%) |
| *LCG_Scrum* | full | 75.2 | 0.00 | 0.13 | 0.21 | 0.24 | 0.15 | 82.48 | 0.02 | 0.17 | 0.51 | 0.23 | 0.44 |
| | rm-requirement | +1.0 (1.3%) | 0.0 (0.0%) | 0.0 (0.0%) | +0.05 (23.8%) | -0.01 (4.2%) | +0.09 (60.0%) | -1.92 (2.3%) | +0.01 (50.0%) | 0.0 (0.0%) | +0.01 (2.0%) | -0.01 (4.3%) | +0.04 (9.1%) |
| | rm-design | +1.0 (1.3%) | 0.0 (0.0%) | -0.02 (15.4%) | -0.03 (14.3%) | 0.0 (0.0%) | -0.03 (20.0%) | +0.89 (1.1%) | -0.01 (50.0%) | -0.04 (23.5%) | +0.11 (21.6%) | +0.03 (13.0%) | -0.19 (43.2%) |
| | rm-codeReview | -2.0 (2.7%) | 0.0 (0.0%) | 0.0 (0.0%) | -0.03 (14.3%) | 0.0 (0.0%) | -0.03 (20.0%) | +0.19 (0.2%) | 0.0 (0.0%) | -0.01 (5.9%) | +0.01 (2.0%) | +0.06 (26.1%) | -0.1 (22.7%) |
| | rm-test | -14.2 (18.9%) | +0.03 (588.2%) | +0.06 (46.2%) | 0.0 (0.0%) | -0.03 (12.5%) | +0.07 (46.7%) | -26.5 (32.1%) | +0.13 (650.0%) | +0.1 (58.8%) | +0.41 (80.4%) | -0.03 (13.0%) | +0.14 (31.8%) |
| | rm-sprintMeeting | -1.6 (2.1%) | 0.0 (0.0%) | -0.01 (7.7%) | -0.02 (9.5%) | -0.03 (12.5%) | +0.1 (66.7%) | -3.09 (3.7%) | 0.0 (0.0%) | -0.02 (11.8%) | +0.05 (9.8%) | -0.01 (4.3%) | +0.21 (47.7%) |

than that of $LCG_{Waterfall}$, $LCG_{TDD}$, and $LCG_{Scrum}$. This finding shows all three models improve the quality of the generated code to different degrees. Specifically, compared to *GPT*, *LCG* decreases the *error* code smell density by 81.8% to 96.0%, *warning* density by 15.0% to 57.9%, *convention* density by 15.4% to 60.2%, and *refactor* density by 6.7% to 30.0%. Meanwhile, *GPT* has fewer handled exceptions than *LCG*. As Table IV shows, in both HumanEval and MBPP, *GPT* has almost zero handled exception, while $LCG_{Waterfall}$ generates the most handled-exception (0.37 and 0.67 handled exceptions per every 10 LOC in the two benchmarks), $LCG_{TDD}$ ranks second (0.33 and 0.62), and then $LCG_{Scrum}$ (0.15 and 0.44). In conclusion, these findings show that *LCG* improves the quality of the generated code by reducing code smells while adding more exception-handling code.

> Compared to *GPT*, *LCG* remarkably improves the quality of the generated code by reducing code smells and adding more exception handling. Testing has the most significant impact on Pass@1 and code smells among all development activities, while having design and code review greatly improve the exception-handling ability.

### RQ3: How stable is the LCG generated code?

*Motivation.* In LLM, the stability of generated responses can be influenced by several parameters: 1) *temperature*, affecting the randomness in the generated responses, and 2) *model versions*, which may introduce variability due to changes in optimization and fine-tuning [5]. Understanding and improving the stability of LLMs is crucial for enhancing their trustworthiness, thereby facilitating their adoption in practice. Therefore, in this research question (RQ), we investigate the stability
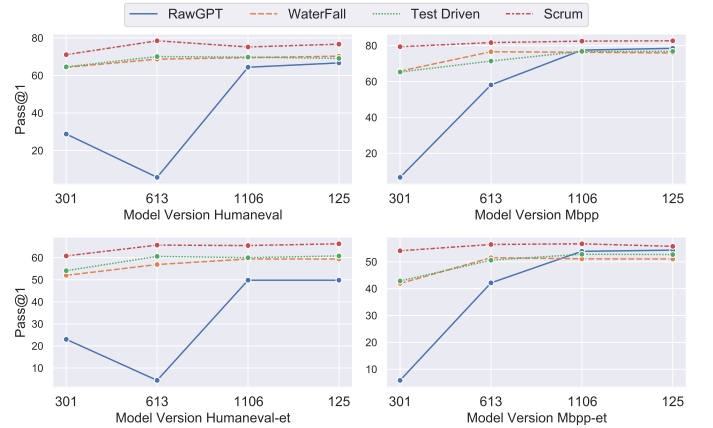


Fig. 2: Pass@1 across GPT3.5 versions.

of our *LCG* in Pass@1 across four benchmarks, considering various temperature values and model versions.

*Approach.* We evaluate the Pass@1 of *LCG* across four versions of GPT3.5 – *gpt-3.5-turbo-0301*, *gpt-3.5-turbo-0613*, *gpt-3.5-turbo-1106*, and *gpt-3.5-turbo-0125*. The latest version is *gpt-3.5-turbo-0125* (published in January 2024) and the earlier version is *gpt-3.5-turbo-0301* (published in March 2023). To avoid the effect of the model version when we vary the temperature, we use the same model version (*gpt-3.5-turbo-1106*, the version that we used in prior RQs) to study the effect of temperature values. We set the temperature to 0.2, 0.4, 0.6, and 0.8 in our experiment. We execute *GPT* and the three variants of *LCG* three times under each configuration and report the average Pass@1.

*Result.* **GPT has extremely low Pass@1 in some versions of GPT3.5, while LCG has stable results across all versions. LCG may help ensure the stability of the generated code**

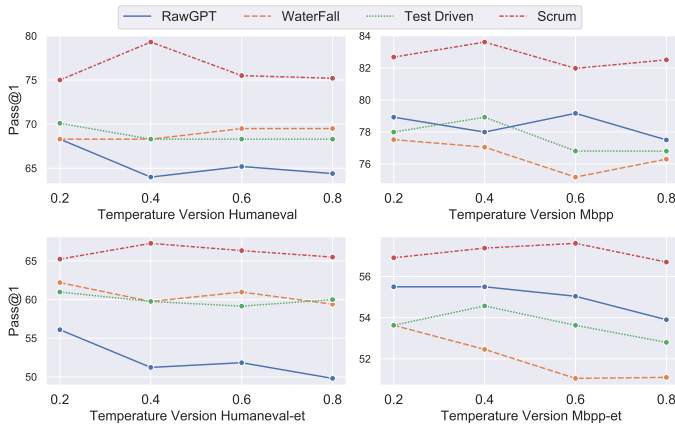Fig. 3: Pass@1 across temperature values.

***even when the underlying LLM regresses.*** Figure 2 shows the Pass@1 for *GPT* and the three *LCG* across GPT3.5 versions. In earlier versions of GPT3.5 (*0301* and *0613*), *GPT* has very low Pass@1 on all benchmarks (e.g., 20 to 30 in HumanEval and HumanEval-ET). In *0301*, MBPP's Pass@1 is even lower with a value around 5. **The findings show that model version may have a significant impact on the generated code**. However, we see that, after adopting our agent-based techniques, all three variants of *LCG* achieve similar Pass@1 across GPT3.5 versions. The results indicate that *LCG* can generate similar-quality code even if we have an underperformed baseline model.

***All techniques have a relatively similar Pass@1 when the temperature value changes.*** Figure 3 shows the Pass@1 for all the techniques when the temperature value changes. There is a slight downward trend for *GPT* when $t$ increases, but the changes are not significant (Pass@1 is decreased by 2 to 5). For *LCG*, and especially *LCG_Scrum*, we see similar Pass@1 regardless of the temperature value. Although we see a slight increase in the Pass@1 of *LCG_Scrum* when $t$=0.4 (2 to 5 higher compared to when $t$=0.8 across the benchmarks), the difference is small, and the Pass@1 is almost the same when $t$ is either the lowest (0.2) or largest value (0.8). In short, although temperature values may have an impact on the generated code, the effect is relatively small, especially for *LCG*.

> *LCG* generates stable results across GPT versions, while we see large fluctuations (14 times difference) in *GPT*'s Pass@1. Pass@1 is generally consistent across all models when the temperature value changes.

## V. THREATS TO VALIDITY

**Internal validity.** Due to the generative nature of LLM, the responses may change across runs. To mitigate the threat, we try to execute the LLMs multiple times. As found in RQ1, the standard deviation of the results is small, so the generated results should be consistent. Variables such as temperature and LLM model version can have an impact on the generated code. To mitigate the threat, in RQ3, we conduct the experiments using different temperature values and model versions. We find

that the temperature value only has a small effect on Pass@1, and model versions have a large impact on *GPT*. In RQ2, we study the impact of removing every development activity. However, having multiple development activities may have a tandem effect that further improves code quality. Future studies are needed to study the effects of different combinations of development activities in code generation.

**External validity.** We conduct our study using two state-of-the-art benchmarks. However, as we discussed, there exist some issues in the provided tests. Moreover, the programming problems are mostly algorithmic, so the findings may not generalize to other code-generation tasks. Future studies should consider applying *LCG* on different and larger programming tasks. We use GPT3.5 as our underlying LLM. Although one can easily replace the LLM in our experiment, the findings may be different. Future studies on how the results of *LCG* change when using different LLMs.

**Construct validity.** We try to implement an agent system that follows various software development processes. However, there are many variations of the same process model, and some variations may give better results. Future studies should further explore how changing the process models affect the code generation ability.

## VI. CONCLUSION

Software development teams rely on following software process models to ensure software quality. Given that the Large Language Model (LLM) agents can act as domain experts and communicate with other agents, our study explores their efficacy in adopting engineering roles and facilitating interactions across various software process models. In this paper, we introduce *LCG*, a framework that implements three renowned process models: *LCG_Waterfall*, *LCG_TDD*, and *LCG_Scrum*. We evaluated how these models affect code generation in terms of correctness and code quality on four benchmarks: HumanEval, HumanEval-ET, MBPP, and MBPP-ET. Our findings show that *LCG_Scrum* notably enhances Pass@1 scores by an average of 15% over RawGPT, while maintaining the lowest standard deviation (averaging 1.2%). While *GPT*'s Pass@1 suffers significant variance when the model version changes (e.g., Pass@1 changes from 60 to 5 in HumanEval), the three process models still maintain stable Pass@1. Moreover, we analyzed the generated code's quality by examining code smells and the ability to handle exceptions. Our results suggest that incorporating design and code review activities significantly reduces code smells and increases the presence of handled exceptions. This indicates that *LCG* not only boosts code correctness but also reduces code smells and improves reliability. These insights pave the way for future research to develop innovative development models tailored for LLM integration in software development processes.

## REFERENCES

[1] Eman Abdullah AlOmar, Anushkrishna Venkatakrishnan, Mohamed Wiem Mkaouer, Christian D Newman, and Ali Ouni. How to refactor this code? an exploratory study

on developer-chatgpt refactoring conversations. *arXiv preprint arXiv:2402.06013*, 2024.

[2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.

[3] Youssef Bassil. A simulation model for the waterfall software development life cycle. *arXiv preprint arXiv:1205.6904*, 2012.

[4] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in retrieval-augmented generation. *arXiv preprint arXiv:2309.01431*, 2023.

[5] Lingjiao Chen, Matei Zaharia, and James Y. Zou. How is chatgpt's behavior changing over time? *ArXiv*, abs/2307.09009, 2023.

[6] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

[7] Subhasish Dasgupta and Sara Hooshangi. Code quality: Examining the efficacy of automated tools. 2017.

[8] Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. Codescore: Evaluating code generation by learning code execution. *arXiv preprint arXiv:2301.09043*, 2023.

[9] Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. Self-collaboration code generation via chatgpt. *arXiv preprint arXiv:2304.07590*, 2023.

[10] Christof Fetzer, Pascal Felber, and Karin Hogstedt. Automatic detection and masking of nonatomic exception handling. *IEEE Transactions on Software Engineering*, 30(8):547–560, 2004.

[11] Martin Fowler. The new methodology, 2005. URL https://www.martinfowler.com/articles/newMethodology.html.

[12] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.

[13] Dong Huang, Qingwen Bu, and Heming Cui. Codecot and beyond: Learning to program and test like a developer. *arXiv preprint arXiv:2308.08784*, 2023.

[14] Dong Huang, Qingwen Bu, Jie M Zhang, Michael Luck, and Heming Cui. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010*, 2023.

[15] Sungmin Kang, Juyeon Yoon, and Shin Yoo. Large language models are few-shot testers: Exploring llm-based general bug reproduction. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2312–2323. IEEE, 2023.

[16] Sawan Kumar and Partha Talukdar. Reordering examples helps during priming-based few-shot learning. *arXiv preprint arXiv:2106.01751*, 2021.

[17] Van-Hoang Le and Hongyu Zhang. Log parsing with prompt-based few-shot learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2438–2449. IEEE, 2023.

[18] Jia Li, Ge Li, Yongmin Li, and Zhi Jin. Structured chain-of-thought prompting for code generation. *arXiv preprint arXiv:2305.06599*, 2023.

[19] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct. *Rigorous evaluation of large language models for code generation. CoRR, abs/2305.01210*, 2023.

[20] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, et al. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Comput. Surv.*, 55(9), jan 2023.

[21] Lipeng Ma, Weidong Yang, Bo Xu, Sihang Jiang, Ben Fei, Jiaqing Liang, Mingjie Zhou, and Yanghua Xiao. Knowlog: Knowledge enhanced pre-trained language model for log understanding. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, pages 1–13, 2024.

[22] Zeyang Ma, An Ran Chen, Dong Jae Kim, Tse-Hsun Chen, and Shaowei Wang. Llmparser: An exploratory study on using large language models for log parsing. In *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE)*, pages 883–883. IEEE Computer Society, 2024.

[23] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

[24] E Michael Maximilien and Laurie Williams. Assessing test-driven development at ibm. In *25th International Conference on Software Engineering, 2003. Proceedings.*, pages 564–569. IEEE, 2003.

[25] Noor Nashid, Mifta Sintaha, and Ali Mesbah. Retrieval-based prompt selection for code-related few-shot learning. In *IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2450–2462, 2023.

[26] OpenAI. Chatgpt. https://chat.openai.com/, 2023.

[27] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22, 2023.

[28] Python. Built-in exceptions, 2005. URL https://docs.python.org/3/library/exceptions.html.

[29] Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.

[30] Fernando Vallecillos Ruiz, Anastasiia Grishina, Max Hort, and Leon Moonen. A novel approach for automatic

program repair using round-trip translation with large language models. *arXiv preprint arXiv:2401.07994*, 2024.

[31] Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. An empirical evaluation of using large language models for automated unit test generation. *IEEE Transactions on Software Engineering*, 2023.

[32] Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*, 2023.

[33] Chaochao Shen, Wenhua Yang, Minxue Pan, and Yu Zhou. Git merge conflict resolution leveraging strategy classification and llm. In *2023 IEEE 23rd International Conference on Software Quality, Reliability, and Security (QRS)*, pages 228–239. IEEE, 2023.

[34] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.

[35] Pylint Development Team. Pylint. https://pypi.org/project/pylint/. Last accessed March 2024.

[36] Testsigma. What is devtestops? https://testsigma.com/devtestops. Last accessed March 2024.

[37] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[38] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

[39] Jules White, Sam Hays, Quchen Fu, Jesse Spencer-Smith, and Douglas C Schmidt. Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design. *arXiv preprint arXiv:2303.07839*, 2023.

[40] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.

[41] Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. Chatunitest: a chatgpt-based automated unit test generation tool. *arXiv preprint arXiv:2305.04764*, 2023.

[42] Jing Xu, Xu Luo, Xinglin Pan, Yanan Li, Wenjie Pei, and Zenglin Xu. Alleviating the sample selection bias in few-shot learning by removing projection to the centroid. *Advances in Neural Information Processing Systems*, 35: 21073–21086, 2022.

[43] Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint arXiv:2309.04658*, 2023.

[44] Aiko Yamashita and Leon Moonen. Do code smells reflect important maintainability aspects? In *2012 28th IEEE International Conference on software maintenance (ICSM)*, pages 306–315, 2012.

[45] Hui Yang, Sifu Yue, and Yunzhong He. Auto-gpt for online decision making: Benchmarks and additional opinions. *arXiv preprint arXiv:2306.02224*, 2023.

[46] Burak Yetiştiren, Işık Özsoy, Miray Ayerdem, and Eray Tüzün. Evaluating the code quality of ai-assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and chatgpt. *arXiv preprint arXiv:2304.10778*, 2023.

[47] Zhiqiang Yuan, Junwei Liu, Qiancheng Zi, Mingwei Liu, Xin Peng, and Yiling Lou. Evaluating instruction-tuned large language models on code comprehension and generation. *arXiv preprint arXiv:2308.01240*, 2023.

[48] Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng. No more manual tests? evaluating and improving ChatGPT for unit test generation. *arXiv preprint arXiv:2305.04207*, 2023.

[49] Yuxia Zhang, Zhiqing Qiu, Klaas-Jan Stol, Wenhui Zhu, Jiaxin Zhu, Yingchen Tian, and Hui Liu. Automatic commit message generation: A critical review and directions for future work. *IEEE Transactions on Software Engineering*, 2024.

[50] Li Zhong, Zilong Wang, and Jingbo Shang. Ldb: A large language model debugger via verifying runtime execution step-by-step. *arXiv preprint arXiv:2402.16906*, 2024.

[51] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models. *arXiv preprint arXiv:2310.04406*, 2023.