# MaintainCoder: Maintainable Code Generation Under Dynamic Requirements

**Zhengren Wang** [1,3]    Rui Ling [1]    Chufan Wang [1]    Yongan Yu [2]
Sizhe Wang [1]    Zhiyu Li [3]    Feiyu Xiong [3]    Wentao Zhang [1]

[1]Center for Data Science, Peking University
[2]McGill University
[3]Center for LLM, Institute for Advanced Algorithms Research, Shanghai

# Outline

# Outline

*"The Only Constant in Life is Change."* — *Heraclitus*

- **The Problem**: Modern code generation focuses on initial correctness, but overlooks **maintainability**—the ability to adapt to changing requirements.

- **The Crisis**: This leads to a "maintainability crisis."
  - 60-80% of software costs are from post-deployment maintenance.
  - Real-world disasters (e.g., Knight Capital's $440M loss) often stem from unmaintained legacy code.

- **The Gap**: Current research lacks tools and methods to address this.
  1. No benchmark to measure maintainability under evolving requirements.
  2. No method that systematically applies software engineering principles.
  3. No discussion on the interplay between correctness and maintainability.

1. **Dynamic Benchmark: MaintainBench**
   The first benchmark to assess code maintainability through requirement evolution cycles.

# Our Contributions

1. **Dynamic Benchmark: MaintainBench**
   The first benchmark to assess code maintainability through requirement evolution cycles.

2. **Systematic Generation: MaintainCoder**
   A multi-agent system integrating the Waterfall model, design patterns, and collaboration to produce maintainable code.

# Our Contributions

1. **Dynamic Benchmark: MaintainBench**
   The first benchmark to assess code maintainability through requirement evolution cycles.

2. **Systematic Generation: MaintainCoder**
   A multi-agent system integrating the Waterfall model, design patterns, and collaboration to produce maintainable code.

3. **Empirical Insights**
   Experiments show MaintainCoder boosts maintainability (+60%) and initial correctness, while revealing flaws in static metrics.

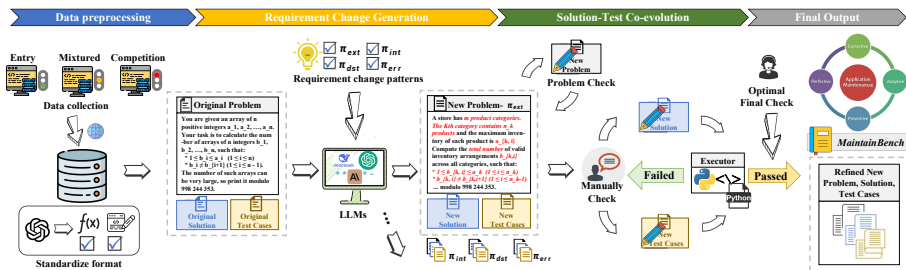# Outline

# Systematic Construction of MaintainBench



Figure: The construction of MaintainBench, from data selection to quality checks.

**Key Idea**: Systematically apply four types of requirement changes to existing benchmark problems:

- Functional Extensions
- Interface Modifications
- Data Structure Transformations
- Error Handling Enhancements

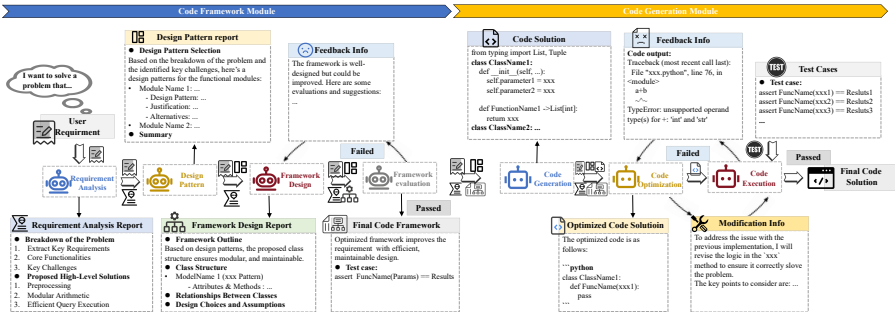# Outline

# Overview of the MaintainCoder Workflow



Figure: Overview of MaintainCoder. MaintainCoder features best practices like the Waterfall model, classical design patterns, and iterative refinement mechanisms for critical stages. 1) The Code Framework Module systematically transforms user requirements into an optimized, maintainable architectural blueprint; 2) The Code Generation Module implements the blueprint into production-ready code, rigorously adhering to both the framework specifications and initial user requirements.

# How MaintainCoder Works

## Module 1: Code Framework

- **Requirements Analysis Agent**: Decomposes the problem.
- **Design Pattern Agent**: Selects appropriate patterns (e.g., Strategy, Factory) to enhance modularity.
- **Framework Design Agent**: Creates a high-level blueprint.
- **Framework Evaluation Agent**: Critiques and refines the blueprint.

## Module 2: Code Generation

- **Code Generation Agent**: Implements the blueprint into functional Python code.
- **Code Optimization Agent**: Executes, debugs, and refines the code against test cases until it is correct and robust.

# Outline

# Experimental Setup

## Baselines

- **Foundation Models**: GPT-4o, DeepSeek-V3, Claude-3.7-Sonnet, etc.
- **Prompting Methods**: Chain-of-Thought (CoT), Self-Planning (Plan).
- **Multi-Agent Systems**: AgentCoder, MapCoder, etc.

## Metrics

- **Static Metrics**:
  Maintainability Index (MI $\uparrow$), Cyclomatic Complexity (CC $\downarrow$).
- **Dynamic Metrics (Our Proposal)**:
  - **Post-modification Correctness (Pass@5 $\uparrow$)**
  - **Code Change Volume (Code$_{\text{diff}}$ $\downarrow$)**
  - **AST Similarity (AST$_{\text{sim}}$ $\uparrow$)**

# MaintainCoder Excels on Difficult Problems

| Model (on CodeContests-Dyn) | Static | | Dynamic | | | |
|---|---|---|---|---|---|---|
| | MI ↑ | CC ↓ | Pass@5 ↑ | $AST_{sim}$ ↑ | $Code_{diff}^{per}$ ↓ | $Code_{diff}^{abs}$ ↓ |
| GPT-4o-mini | 57.8 | 6.06 | 24.2 | 0.661 | 90.1 | 16.6 |
| AgentCoder (GPT-4o-mini) | 62.5 | 7.28 | 18.2 | 0.629 | 44.9 | 18.5 |
| **MaintainCoder (GPT-4o-mini)** | **65.8** | **2.68** | **32.6** | **0.833** | **23.2** | **17.4** |
| DeepSeek-V3 | 55.7 | 6.80 | 26.5 | 0.718 | 87.2 | 17.5 |
| AgentCoder (DeepSeek-V3) | 43.8 | 19.6 | 16.7 | 0.670 | 48.4 | 28.0 |
| **MaintainCoder (DeepSeek-V3)** | **63.1** | **3.64** | **37.9** | **0.788** | **43.2** | **18.5** |

- Key Insight: MaintainCoder achieves the best dynamic metrics by a large margin (+60% Pass@5).

- Low Complexity: It keeps Cyclomatic Complexity (CC) extremely low (∼3), indicating simpler, more modular code.

- Agent Pitfall: Other multi-agent systems like AgentCoder may actually *hurt* maintainability (lower Pass@5).

# MaintainCoder Excels on Difficult Problems

| Model (on CodeContests-Dyn) | Static | | Dynamic | | | |
|---|---|---|---|---|---|---|
| | MI ↑ | CC ↓ | Pass@5 ↑ | $AST_{sim}$ ↑ | $Code_{diff}^{per}$ ↓ | $Code_{diff}^{abs}$ ↓ |
| GPT-4o-mini | 57.8 | 6.06 | 24.2 | 0.661 | 90.1 | 16.6 |
| AgentCoder (GPT-4o-mini) | 62.5 | 7.28 | 18.2 | 0.629 | 44.9 | 18.5 |
| **MaintainCoder (GPT-4o-mini)** | **65.8** | **2.68** | **32.6** | **0.833** | **23.2** | **17.4** |
| DeepSeek-V3 | 55.7 | 6.80 | 26.5 | 0.718 | 87.2 | 17.5 |
| AgentCoder (DeepSeek-V3) | 43.8 | 19.6 | 16.7 | 0.670 | 48.4 | 28.0 |
| **MaintainCoder (DeepSeek-V3)** | **63.1** | **3.64** | **37.9** | **0.788** | **43.2** | **18.5** |

- Key Insight: MaintainCoder achieves the best dynamic metrics by a large margin (+60% Pass@5).
- Low Complexity: It keeps Cyclomatic Complexity (CC) extremely low (~3), indicating simpler, more modular code.
- Agent Pitfall: Other multi-agent systems like AgentCoder may actually *hurt* maintainability (lower Pass@5).

# MaintainCoder Excels on Difficult Problems

| Model (on CodeContests-Dyn) | Static | | Dynamic | | | |
|---|---|---|---|---|---|---|
| | MI ↑ | CC ↓ | Pass@5 ↑ | $AST_{sim}$ ↑ | $Code_{diff}^{per}$ ↓ | $Code_{diff}^{abs}$ ↓ |
| GPT-4o-mini | 57.8 | 6.06 | 24.2 | 0.661 | 90.1 | 16.6 |
| AgentCoder (GPT-4o-mini) | 62.5 | 7.28 | 18.2 | 0.629 | 44.9 | 18.5 |
| **MaintainCoder (GPT-4o-mini)** | **65.8** | **2.68** | **32.6** | **0.833** | **23.2** | **17.4** |
| DeepSeek-V3 | 55.7 | 6.80 | 26.5 | 0.718 | 87.2 | 17.5 |
| AgentCoder (DeepSeek-V3) | 43.8 | 19.6 | 16.7 | 0.670 | 48.4 | 28.0 |
| **MaintainCoder (DeepSeek-V3)** | **63.1** | **3.64** | **37.9** | **0.788** | **43.2** | **18.5** |

- Key Insight: MaintainCoder achieves the best dynamic metrics by a large margin (+60% Pass@5).

- Low Complexity: It keeps Cyclomatic Complexity (CC) extremely low (∼3), indicating simpler, more modular code.

- Agent Pitfall: Other multi-agent systems like AgentCoder may actually *hurt* maintainability (lower Pass@5).

# Good Design Improves Both

## Pass@5 on Original Datasets

| Model | APPS | CodeContests | xCodeEval |
|---|---|---|---|
| GPT-4o-mini | 44% | 18% | 46% |
| MaintainCoder | 48% | 23% | 57% |
| DeepSeek-V3 | 66% | 48% | 75% |
| MaintainCoder | 69% | 51% | 77% |

- **Hypothesis Confirmed**: A focus on maintainability also improves the correctness of the *initial* code. The structured MaintainCoder leads to more robust and correct solutions from the start.
- The benefit is larger for complex problems (e.g., +11% on xCodeEval).
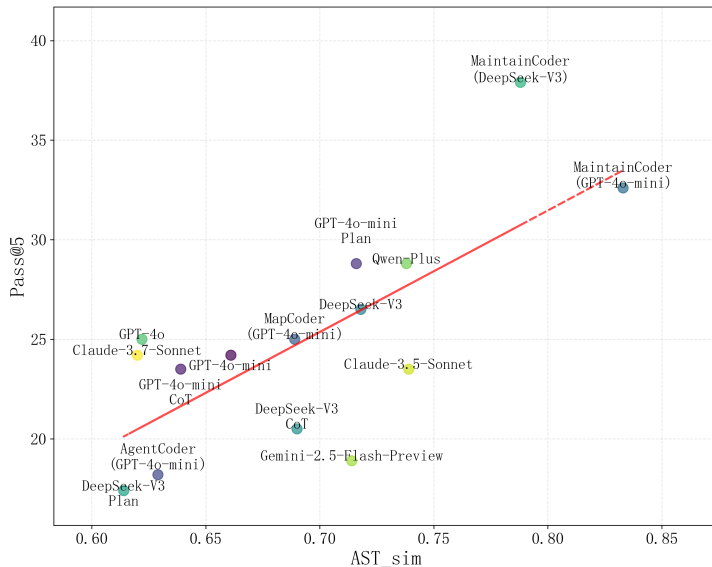
# Good Design Improves Both

## Pass@5 on Original Datasets

| Model | APPS | CodeContests | xCodeEval |
|-------|------|--------------|-----------|
| GPT-4o-mini | 44% | 18% | 46% |
| MaintainCoder | 48% | 23% | 57% |
| DeepSeek-V3 | 66% | 48% | 75% |
| MaintainCoder | 69% | 51% | 77% |

- **Hypothesis Confirmed**: A focus on maintainability also improves the correctness of the *initial* code. The structured MaintainCoder leads to more robust and correct solutions from the start.
- The benefit is larger for complex problems (e.g., +11% on xCodeEval).

## Static vs. Dynamic Metrics

Our experiments show static metrics (MI, CC) can be inconsistent and misleading. In contrast, our proposed dynamic metrics (Pass@k, $Code_{diff}$, $AST_{sim}$) are highly consistent and better reflect maintenance effort.

# Outline

# Conclusion and Future Work

## Summary

- We introduced **MaintainBench**, the first benchmark for evaluating code maintainability under dynamic requirements.
- We proposed **MaintainCoder**, a multi-agent system that leverages software engineering principles to generate highly maintainable code.
- Experiments show MaintainCoder significantly outperforms existing methods in both maintainability and initial correctness.

## Future Work

- Expand MaintainBench with more diverse modification types.
- Extend to more complex, repository-level tasks (e.g., SWE-Bench).
- Explore the trade-offs between initial computational cost and long-term maintenance savings.

# Thank You!

Resources available at:
https://github.com/IAAR-Shanghai/MaintainCoder