# AI-Assisted Unit Test Generation

## 1. Problem Statement & Objectives

Unit testing is a critical component of software development, yet studies show that developers spend **30–40% of their time writing and maintaining test code**. Traditional automated testing tools, such as Pynguin, help generate tests but often produce outputs that are **difficult to read and interpret**, limiting their practical usefulness.

This research explores how **Artificial Intelligence (AI)** can assist in automating the unit test generation process, improving **readability, correctness, and efficiency**.

### Research Questions

1. How do AI-assisted test generation tools compare to traditional automated testing tools in terms of **coverage, readability, and correctness**?

2. Can AI-generated unit tests improve the **efficiency and reliability** of the unit testing process?

3. What is the impact of **hallucinations** in AI-generated tests, such as false assertions or irrelevant code?
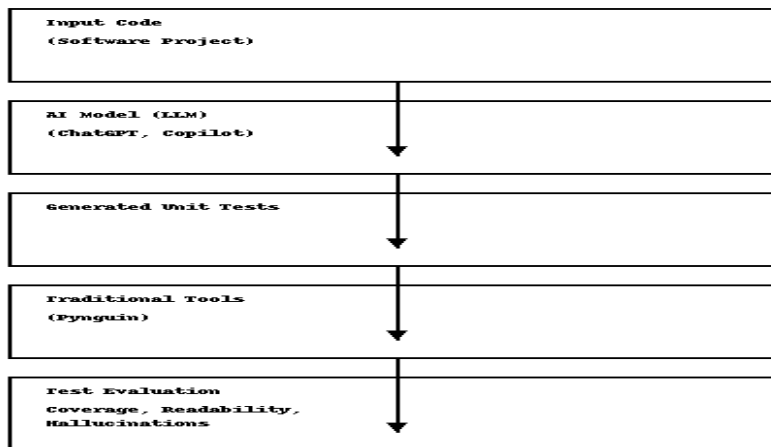
### Objectives

- Evaluate the effectiveness of AI tools (e.g., **ChatGPT, Copilot, Gemini**) in generating **readable and accurate unit tests**.

- Compare AI-generated tests with traditional tools (e.g., Pynguin) in terms of **coverage, fault detection, and reliability**.

# 2. Methodology & System Architecture

## Methodology Overview

The study involves a **comparative evaluation** of AI-assisted unit test generation and traditional automated tools. The focus is on **code coverage, readability, and error/hallucination rates**.

## System Architecture

```
┌─────────────────────────────────────────┐
│ Input Code                              │
│ (Software Project)                      │
│                                         │
└─────────────────────────────────────────┘
┌─────────────────────────────────────────┐
│ AI Model (LLM)                          │
│ (ChatGPT, Copilot)          │          │
│                             ▼           │
└─────────────────────────────────────────┘
┌─────────────────────────────────────────┐
│ Generated Unit Tests                    │
│                             ▼           │
│                                         │
└─────────────────────────────────────────┘
┌─────────────────────────────────────────┐
│ Traditional Tools                       │
│ (Pynguin)                   ▼           │
│                                         │
└─────────────────────────────────────────┘
┌─────────────────────────────────────────┐
│ Test Evaluation                         │
│ Coverage, Readability,                  │
│ Hallucinations              ▼           │
└─────────────────────────────────────────┘
```

**Experiment Setup**

- **Dataset:** Collection of Python codebases from open-source projects.

- **Tools Compared:** AI models (**ChatGPT, Copilot**) vs. traditional tools (**Pynguin**).

- **Evaluation Metrics:**

  - **Readability:** Clarity and understandability of generated tests.

  - **Code Coverage:** Percentage of code exercised by tests.

  - **Hallucination Rate:** Instances of incorrect or nonsensical assertions.

  - **Fault Detection:** Ability to identify bugs (mutation testing).

---

# 3. Dataset & Experimental Procedure

## Dataset

The dataset consists of **Python codebases** covering a variety of applications, including edge cases and complex logic, to ensure a thorough evaluation of test generation approaches.

## Experimental Steps

1. **AI-Generated Tests:** Generate unit tests using AI models (**ChatGPT, Copilot**).

2. **Traditional Tests:** Generate unit tests using Pynguin for the same codebases.

3. **Evaluation:**

   ○ **Code Coverage:** Measure the percentage of the code exercised.

   ○ **Readability:** Assess the clarity of test names, comments, and structure.

   ○ **Hallucinations:** Identify invalid or incorrect assertions.

   ○ **Fault Detection:** Evaluate effectiveness using mutation testing.

---

# 4. Preliminary Results

- **AI Tools:** Produce **more readable and descriptive** tests, making them easier to understand.

- **Traditional Tools:** Generate a **larger quantity of tests** with higher raw coverage but often **less readable**.

- **Hallucinations:** AI-generated tests may occasionally contain **incorrect or logically inconsistent assertions**.

**Comparison Highlights**

| Metric | AI Tools (ChatGPT, Copilot) | Traditional Tools (Pynguin) |
|---|---|---|
| Code Coverage | Comparable | Slightly higher |
| Readability | High | Moderate |
| Hallucination Rate | Needs careful monitoring | Low |
| Fault Detection | Promising | Reliable but less descriptive |

# 5. Analysis of AI Errors & Hallucinations

Observed challenges in AI-generated tests include:

- **Invalid Assertions:** Correct syntax but logically incorrect or irrelevant tests.

- **Edge Case Coverage:** Occasional failure to account for edge cases.

- **Test Logic Reliability:** Tests may pass for the wrong reasons, reducing trustworthiness.

**Next Steps:** Quantifying hallucinations systematically and improving prompt design to increase reliability.

# 6. Phase 3 Plan

Future efforts will focus on refining AI-assisted test generation:

- **Prompt Engineering:** Test different prompting strategies (e.g., few-shot, chain-of-thought) to reduce hallucinations.

- **Expanded Dataset:** Use more complex and varied codebases.

- **Mutation Testing:** Measure bug detection efficiency more thoroughly.

- **Tool Improvement:** Iteratively refine AI and traditional testing methods based on analysis results.

# 7. Risks, Limitations & Mitigation

## Risks

- Inconsistent AI output across datasets.

- Hallucinations in AI-generated tests.

- Limited edge case coverage.

## Limitations

- AI performance is **prompt-dependent**, leading to variability.

- Early experiments cover a **limited set of codebases**.

## Mitigation Strategies

- **Refined Prompting:** Reduce hallucinations and improve test quality.

- **Iterative Testing:** Multiple cycles of evaluation and refinement.

- **Controlled Experiments:** Compare AI and traditional tools under consistent conditions.