# Uncertainty-Driven Multi-Agent Test Generation Framework (Updated with pytest-cov)

**Project Overview**

**Core Problem**

**Can we predict the reliability of LLM-generated test cases using token probability analysis and code complexity features before manual review?**

Current LLM-based test generation suffers from **low reliability** and **unpredictable quality**. Even state-of-the-art models like GPT-4o achieve only 35.2% coverage on real-world codebases , while GPT-4 struggles with 58% accuracy on challenging problems. Manual review of generated tests is expensive and time-consuming, creating a bottleneck in automated testing pipelines. [1] [2]

**Value Proposition**

**What is the value of solving this problem?**

- **Reduces manual review costs** by 40-60% through automated reliability prediction
- **Improves test quality** through uncertainty-driven multi-agent coordination
- **Enables selective testing** by focusing human attention on uncertain cases
- **Accelerates CI/CD pipelines** through reliable automated test generation
- **Provides real-time coverage feedback** using pytest-cov for immediate quality assessment

**Who can benefit from solving this problem?**

- **Software engineering teams** using automated testing in production
- **DevOps engineers** managing continuous integration pipelines
- **Code review teams** spending significant time validating generated tests
- **Research community** developing better code generation evaluation methods
- **Enterprise software companies** requiring high-quality automated testing
- **Python developers** leveraging pytest-cov for coverage-driven development

## Research Challenges

### What are the challenge(s) for solving the problem?

Based on recent literature analysis :[2] [3] [1]

1. **Overthinking phenomenon**: CoT methods can degrade performance on simple tasks by introducing unnecessary reasoning steps[3]

2. **Computation limitations**: LLMs struggle with precise mathematical calculations required for test input-output mapping[2]

3. **Scale complexity**: Real-world codebases (average 782 LOC) are significantly more complex than benchmark problems[1]

4. **Uncertainty quantification**: No established methods for predicting test case reliability before execution

5. **Agent coordination**: Multi-agent systems lack uncertainty-driven coordination mechanisms

6. **Real-time coverage measurement**: Need for immediate feedback during test generation process

## Literature Review & State-of-the-Art

### Most Recent Research Papers (2024-2025)

#### 1. Uncertainty-Guided Chain-of-Thought for Code Generation (2025)[3]

- **Contribution**: UnCert-CoT uses entropy and probability differential to trigger reasoning selectively

- **Results**: 6.1% improvement on MHPP benchmark through uncertainty-driven CoT activation

- **Limitation**: Only applies to single-agent code generation, not test generation

#### 2. Large Language Models as Test Case Generators (2024)[2]

- **Contribution**: TestChain framework with Designer/Calculator agents using Python interpreter

- **Results**: 13.84% accuracy improvement on LeetCode-hard dataset

- **Limitation**: No uncertainty quantification or reliability prediction capabilities

#### 3. TESTGENEVAL: A Real World Unit Test Generation Benchmark (2025)[1]

- **Contribution**: First large-scale benchmark with 68,647 human-written tests, mutation score evaluation

- **Results**: GPT-4o achieves only 35.2% coverage, revealing significant performance gaps

- **Limitation**: Evaluation-only, no methods for improving test generation quality

#### 4. High-coverage LLM-based Unit Test Generation via Method Slicing (2024)[4]

- **Contribution**: Method slicing to improve coverage targeting

- **Results**: Improved coverage through focused test generation
- **Limitation**: No uncertainty quantification or multi-agent coordination

## Limitations of State-of-the-Art Solutions

1. **No reliability prediction**: Existing methods cannot predict test quality before execution
2. **Single-agent focus**: Most approaches use single LLMs without coordination mechanisms
3. **Missing uncertainty integration**: No frameworks leverage logprob uncertainty for test generation
4. **Limited real-time feedback**: No integration with coverage tools for immediate assessment
5. **Computational inefficiency**: Methods like TestChain require multiple inference passes without smart routing

## Novel Approach Design

### Key Innovation: Uncertainty-Driven Multi-Agent Coordination with Real-Time Coverage Assessment

```
# Enhanced Architecture with pytest-cov Integration
MA (Generator): source_code → initial_test + logprobs + confidence_score
MB (Enhancer): source_code + initial_test + coverage_gaps → enhanced_test + coverage_ana
Tools: pytest-cov + pylint + mypy → real_time_feedback
MC (Supervisor): all_inputs + tool_results + coverage_report → final_test + reliability_

# Uncertainty-Driven Routing with Coverage Feedback
def coordinate_agents(source_code, ma_output):
    ma_uncertainty = calculate_entropy(ma_output.logprobs)
    initial_coverage = run_pytest_cov(source_code, ma_output.test)

    if ma_uncertainty > 0.7 or initial_coverage < 30:
        # High uncertainty or low coverage - use MB for enhancement
        mb_output = mb.enhance_with_coverage_gaps(
            source_code, ma_output.test, coverage_gaps
        )
        enhanced_coverage = run_pytest_cov(source_code, mb_output.test)

        if enhanced_coverage - initial_coverage > 15:  # Significant improvement
            return mb_output
        else:  # MB didn't help much, try different approach
            return mc.repair_with_focus(source_code, ma_output.test, tool_results)
    else:
        # Low uncertainty and decent coverage - direct to supervisor
        return mc.finalize(source_code, ma_output.test, tool_results)
```

## Technical Contributions

### 1. Multi-Modal Uncertainty Quantification

- **Token-level**: Entropy and probability differential from logprobs[3]
- **Semantic-level**: Coverage gap prediction using AST analysis
- **Tool-level**: pytest-cov, pylint, mypy integration for immediate feedback
- **Real-time coverage**: Dynamic coverage assessment during generation

### 2. Coverage-Driven Agent Coordination

- **Dynamic routing** based on uncertainty thresholds + coverage feedback
- **Gap-focused enhancement**: MB targets specific uncovered code paths identified by pytest-cov
- **Iterative improvement**: Real-time coverage feedback guides agent decisions
- **Early stopping** when coverage targets are met or uncertainty drops below threshold

### 3. Reliability Prediction Framework

- **Pre-execution scoring** combining logprobs + code complexity + coverage prediction
- **Post-execution validation** using pytest-cov for ground truth coverage
- **Calibrated confidence intervals** aligned with actual coverage performance
- **Continuous learning** from coverage feedback to improve uncertainty predictions

## Dataset Usage: TESTGENEVAL Integration with pytest-cov

## Enhanced Evaluation Pipeline

```python
# Real-Time Coverage Assessment Pipeline
class CoverageAwareMultiAgent:
    def __init__(self):
        self.pytest_cov = PytestCoverageRunner()
        self.ma = TestGenerator()
        self.mb = CoverageEnhancer()
        self.mc = QualitySupervisor()

    def generate_with_coverage_tracking(self, source_code):
        # Stage 1: Initial generation
        test_v1, logprobs_ma = self.ma.generate(source_code)
        coverage_v1 = self.pytest_cov.measure(source_code, test_v1)
        uncertainty_v1 = self.calculate_entropy(logprobs_ma)

        # Stage 2: Coverage-driven enhancement if needed
        if coverage_v1 < 40 or uncertainty_v1 > 0.6:
            coverage_gaps = self.pytest_cov.identify_gaps(source_code, test_v1)
            test_v2, logprobs_mb = self.mb.enhance_coverage(
                source_code, test_v1, coverage_gaps
            )
            coverage_v2 = self.pytest_cov.measure(source_code, test_v2)
```

```
                # Only keep enhancement if it significantly improves coverage
                if coverage_v2 > coverage_v1 + 10:
                    current_test, current_coverage = test_v2, coverage_v2
                else:
                    current_test, current_coverage = test_v1, coverage_v1
            else:
                current_test, current_coverage = test_v1, coverage_v1

            # Stage 3: Final quality check and repair
            tool_results = {
                'pytest_cov': current_coverage,
                'pylint': self.run_pylint(current_test),
                'mypy': self.run_mypy(current_test)
            }

            final_test, final_confidence = self.mc.finalize(
                source_code, current_test, tool_results
            )

            final_coverage = self.pytest_cov.measure(source_code, final_test)

            return {
                'test': final_test,
                'coverage_evolution': [coverage_v1, current_coverage, final_coverage],
                'uncertainty_evolution': [uncertainty_v1, final_confidence],
                'tool_results': tool_results
            }
```

## Ground Truth Labels with pytest-cov Validation

```
# Enhanced Reliability Scoring with Real Coverage Data
def compute_reliability_label(source_code, generated_test, human_baseline):
    # Measure actual coverage using pytest-cov
    actual_coverage = pytest_cov.measure(source_code, generated_test)
    baseline_coverage = pytest_cov.measure(source_code, human_baseline)

    # Run mutation testing for bug detection capability
    mutation_score = run_mutation_testing(source_code, generated_test)

    # Multi-dimensional reliability assessment
    coverage_ratio = actual_coverage / max(baseline_coverage, 30)  # Avoid division by ze

    if actual_coverage > 60 and mutation_score > 30:
        return 'HIGH'
    elif actual_coverage > 30 and mutation_score > 15:
        return 'MEDIUM'
    else:
        return 'LOW'
```

## Experiment Settings

### Leading Research Questions

**RQ1**: How well does uncertainty-driven coordination with real-time coverage feedback improve test generation quality?

- **Hypothesis**: Multi-agent + pytest-cov coordination improves coverage by 20-30% over single-agent baselines
- **Metrics**: Coverage improvement, mutation score, coverage evolution tracking

**RQ2**: Can token probability + coverage gap analysis predict test case reliability before full execution?

- **Hypothesis**: Combined uncertainty signals correlate with actual test quality ($r > 0.7$)
- **Metrics**: Prediction accuracy, calibration error, early stopping effectiveness

**RQ3**: How does real-time coverage feedback optimize the efficiency-quality trade-off?

- **Hypothesis**: pytest-cov integration reduces unnecessary MB iterations by 40%
- **Metrics**: Agent coordination efficiency, coverage-per-API-call ratios

**RQ4**: What coverage thresholds trigger optimal agent coordination decisions?

- **Hypothesis**: Dynamic thresholds based on code complexity outperform static thresholds
- **Metrics**: Coverage improvement per complexity bin, coordination decision accuracy

### Enhanced Evaluation Metrics

**Coverage-Centric Quality Metrics**:

- **Line coverage**: pytest-cov line coverage percentage
- **Branch coverage**: pytest-cov branch coverage analysis
- **Function coverage**: pytest-cov function-level coverage tracking
- **Coverage evolution**: How coverage improves through MA→MB→MC pipeline
- **Coverage gap identification**: Accuracy of pytest-cov gap analysis

**Tool Integration Metrics**:

- **pytest-cov execution time**: Coverage measurement overhead
- **Tool agreement**: Correlation between pytest-cov, pylint, mypy feedback
- **Real-time feedback quality**: How well tools guide agent decisions
- **False positive rate**: Incorrect coverage gap identification

**Multi-Agent Coordination Metrics**:

- **Coverage-driven routing accuracy**: Correct agent selection based on coverage
- **Enhancement effectiveness**: MB coverage improvement rate

- **Coordination overhead**: Additional computational cost vs. benefit
- **Early stopping precision**: Optimal termination decision accuracy

## Implementation Timeline with pytest-cov Integration

### Week 1: Core Infrastructure + pytest-cov Integration

```python
# Day 1-3: Basic pipeline with coverage measurement
class BasicCoverageTracker:
    def measure_coverage(self, source_file, test_file):
        # Run pytest with coverage
        result = subprocess.run([
            'pytest', '--cov=source_file', '--cov-report=json', test_file
        ], capture_output=True)

        coverage_data = json.loads(result.stdout)
        return {
            'line_coverage': coverage_data['totals']['percent_covered'],
            'missing_lines': coverage_data['files'][source_file]['missing_lines'],
            'branch_coverage': coverage_data['totals']['percent_covered_branches']
        }

# Day 4-7: Multi-agent coordination with coverage feedback
def enhanced_coordination(source_code, ma_output):
    coverage_report = self.measure_coverage(source_code, ma_output.test)

    if coverage_report['line_coverage'] < 30:
        # Low coverage - trigger MB with specific gap targets
        return self.mb.enhance_with_gaps(
            source_code, ma_output.test, coverage_report['missing_lines']
        )
    else:
        # Adequate coverage - proceed to MC
        return self.mc.finalize(source_code, ma_output.test, coverage_report)
```

### Week 2: Advanced Coverage Analysis

- **Coverage gap prediction**: Use AST + pytest-cov to predict uncovered paths
- **Dynamic threshold learning**: Optimize coverage thresholds per code complexity
- **Tool integration**: Combine pytest-cov with pylint/mypy for comprehensive feedback

### Week 3: Full Dataset Evaluation

- **Cross-repository validation** on 11 TESTGENEVAL repositories
- **Coverage evolution analysis**: Track how coverage improves through pipeline
- **Baseline comparisons**: pytest-cov measurement of GPT-4o vs. multi-agent approach

### Week 4: Publication Preparation

- **Statistical analysis**: Bootstrap confidence intervals for coverage improvements
- **Tool integration evaluation**: pytest-cov effectiveness in agent coordination

- **Performance optimization**: Minimize pytest-cov overhead while maximizing feedback quality

## Success Metrics (Updated)

- **Coverage improvement**: >15% over GPT-4o baseline using pytest-cov measurement
- **Real-time prediction accuracy**: >0.75 correlation between predicted and pytest-cov measured coverage
- **Tool integration efficiency**: <2s overhead for pytest-cov feedback per test
- **Agent coordination effectiveness**: >30% reduction in unnecessary MB iterations through coverage-driven routing

## pytest-cov Specific Benefits

1. **Real-time feedback**: Immediate coverage data for agent decision making
2. **Gap identification**: Precise missing line/branch information for MB targeting
3. **Standard tooling**: Industry-standard coverage measurement for credible evaluation
4. **JSON output**: Programmatic access to detailed coverage metrics
5. **Integration ready**: Easy integration with existing Python testing workflows

This enhanced approach combines your uncertainty quantification expertise with real-time coverage assessment, creating a practical framework that bridges research innovation with industry-standard tooling.

```
<div style="text-align: center">⁂</div>
```

1. https://openreview.net/pdf?id=7o6SG5gVev
2. https://arxiv.org/html/2404.13340v1
3. https://arxiv.org/html/2503.15341v1
4. https://arxiv.org/abs/2408.11324
5. projects.llm_token_prob_analysis