# Agora: A Comprehensive Framework for Evaluating and Generating Synthetic Data with Language Models

OpenHands AI Assistant
Technical Report
December 2023

*Abstract*—**This technical report presents an in-depth analysis of Agora, an innovative open-source framework designed for generating and evaluating synthetic data using Large Language Models (LLMs). We explore the architecture, key components, and methodological approaches implemented in the framework. Agora provides a standardized environment for assessing LLMs' capabilities in data generation across multiple domains including mathematics, general instruction-following, and code generation. The framework introduces AgoraBench, a comprehensive benchmark that enables systematic comparison of different LLMs as data generators. We discuss the framework's modular design, evaluation metrics, and practical applications in enhancing LLM training and evaluation. Our analysis reveals that Agora represents a significant advancement in synthetic data generation and evaluation, offering valuable insights into the capabilities and limitations of modern language models.**

## I. Introduction

The increasing importance of high-quality training data in machine learning has led to growing interest in synthetic data generation using Large Language Models (LLMs). The Agora framework, inspired by the ancient Athenian marketplace where knowledge was freely exchanged, provides a systematic approach to generating and evaluating synthetic data. This framework addresses a critical gap in the field by offering standardized methods for comparing different LLMs' capabilities as data generators.

### A. Background and Motivation

The development of large language models has revolutionized natural language processing, but their effectiveness is heavily dependent on the quality and quantity of training data. Traditional data collection methods often face challenges such as:

- Limited availability of domain-specific data
- High costs associated with manual data annotation
- Privacy concerns and data protection regulations
- Bias and quality control issues in collected datasets

Synthetic data generation using LLMs offers a promising solution to these challenges, but it introduces new questions about the quality and reliability of generated data. The Agora framework aims to address these concerns by providing:

- Standardized evaluation metrics for synthetic data quality
- Reproducible benchmarking methodology
- Domain-specific validation techniques
- Comprehensive analysis tools for generated datasets

### B. Related Work

Previous research in synthetic data generation has primarily focused on specific domains or applications. Notable approaches include:

- Template-based generation methods
- Rule-based synthetic data creation
- GAN-based data synthesis
- Few-shot learning approaches

Agora builds upon these foundations while introducing novel approaches to evaluation and quality control. The framework's design incorporates lessons learned from previous work while addressing their limitations.

## II. Framework Architecture

Agora is built with a modular architecture that facilitates customization and extension. The framework's design emphasizes flexibility, scalability, and ease of use, making it suitable for both research and production environments.

### A. Core Components

The framework consists of several key components, each designed to handle specific aspects of the data generation and evaluation process:

*1) Prompt Loader:* The Prompt Loader component is responsible for:

- Template management and customization
- Context window optimization
- Few-shot example selection
- Dynamic prompt construction

It implements sophisticated algorithms for selecting and arranging in-context examples, ensuring optimal prompt construction for different tasks and domains. The component supports various placeholder formats:

```
placeholder_formats = {
    "demonstration_input_placeholder": "<
        input@>",
    "demonstration_output_placeholder": "<
        output@>",
    "test_input_placeholder": "<input>",
    "test_output_placeholder": "<o>",
    "test_input_trigger": "INPUT:",
    "test_output_trigger": "OUTPUT:",
    "stop_phrase": "[END]"
}
```

*2) Parser:* The Parser component handles:

- Output format standardization
- Structured data extraction
- Error handling and recovery
- Format validation and normalization

The parser is designed to be robust against variations in model outputs while maintaining strict format compliance. Here's an example of a custom parser implementation:

```
1  class InstanceGenerationParser(Parser):
2      def parse(self, prompt,
           teacher_model_output,
3                placeholder_formats):
4          instruction = (
5              teacher_model_output.split(
6                  placeholder_formats["
                      test_input_trigger"]
7              )[-1]
8              .split(
9                  placeholder_formats["
                      test_output_trigger"]
10             )[0]
11             .strip()
12         )
13         response = (
14             teacher_model_output.split(
15                 placeholder_formats["
                      test_output_trigger"]
16             )[-1]
17             .split(
18                 placeholder_formats["
                      stop_phrase"]
19             )[0]
20             .strip()
21         )
22         return {
23             "instruction": instruction,
24             "response": response
25         }
```

*3) Validator:* The Validator ensures data quality through:

- Semantic consistency checking
- Domain-specific rule validation
- Format compliance verification
- Quality metrics computation

Example validator implementation:

```
1  class CustomValidator(Validator):
2      def validate(self, instruction, response):
3          # Check minimum length requirements
4          if len(instruction) < 10 or len(
              response) < 10:
5              return False
6
7          # Verify format compliance
8          if not self._check_format(instruction,
              response):
9              return False
10
11         # Domain-specific validation
12         if not self._validate_domain_rules(
              instruction, response):
13             return False
14
15
16         return True
```

*4) LLM Interface:* The LLM Interface provides:

- Unified API for multiple model providers
- Batch processing capabilities
- Error handling and retry logic
- Resource optimization

Example usage of the LLM interface:

```
llm = OpenAILLM(
    model_name="gpt-4",
    api_key="YOUR_API_KEY"
)

sampling_params = {
    "max_tokens": 4096,
    "temperature": 1.0,
    "top_p": 0.9,
    "stop": "[END]"
}

result = llm.generate(
    prompt=prompt,
    **sampling_params
)
```

## B. System Integration

The components are integrated through a well-defined pipeline that ensures:

- Efficient data flow
- Error isolation and handling
- Scalability and parallelization
- Monitoring and logging

Example pipeline configuration:

```
1  agora = Agora(
2      llm=llm,
3      placeholder_formats=placeholder_formats,
4      prompt_loader=prompt_loader,
5      parser=parser,
6      validator=validator,
7      sampling_params=sampling_params
8  )
9
10 result = agora.run(
11     num_instances=10000,
12     num_threads=16,
13     output_file="./results/final_result.json"
14 )
15
16 \section{AgoraBench}
17 AgoraBench serves as a standardized evaluation
        framework for assessing LLMs' data
        generation capabilities. It provides
        comprehensive metrics and evaluation
        methodologies across multiple domains.
18
19 \subsection{Evaluation Domains}
20 The benchmark covers three primary domains:
21
22 \subsubsection{Mathematics}
23 Mathematics evaluation includes:
24 \begin{itemize}
25     \item Problem complexity assessment
26     \item Solution correctness verification
```

```
27        \item Step-by-step explanation quality
28        \item Mathematical notation accuracy
29   \end{itemize}
30
31   Specific datasets used include:
32   \begin{itemize}
33        \item GSM8K for grade school mathematics
34        \item MATH for advanced mathematical
             reasoning
35   \end{itemize}
36
37   Example GSM8K problem generation:
38   \begin{lstlisting}
39   # Prompt template for math problem generation
40   template = """Generate a grade school math
41   problem that tests {concept} understanding.
42   The problem should be clear and solvable.
43
44   Format:
45   Problem: [Your problem here]
46   Solution: [Step-by-step solution]
47
48   Examples:
49   {examples}
50
51   Now generate a new problem:"""
52
53   # Example usage
54   problem = generate_math_problem(
55        template,
56        concept="percentage calculation",
57        examples=gsm8k_examples[:3]
58   )
```

*1) General Instruction Following:* This domain evaluates:

- Task comprehension accuracy
- Response relevance and completeness
- Instruction adherence
- Output format compliance

Evaluation tools include:

- AlpacaEval 2.0 for instruction following
- Arena-Hard for complex task handling

Example instruction generation:

```
1    # Generate complex instruction
2    instruction = generate_instruction(
3        domain="task_planning",
4        complexity_level="high",
5        required_steps=5,
6        context_length="medium"
7    )
8
9    # Validate instruction quality
10   quality_score = evaluate_instruction(
11        instruction,
12        metrics=["clarity", "complexity", "
             feasibility"]
13   )
```

*2) Code Generation:* Code generation assessment covers:

- Functional correctness
- Code quality metrics
- Documentation quality
- Error handling implementation

Benchmark suites include:

- MBPP for basic programming tasks
- HumanEval for realistic coding scenarios

Example code generation evaluation:

```
1    def evaluate_code_solution(
2        problem_spec,
3        generated_code,
4        test_cases
5    ):
6        # Compile and run tests
7        results = run_test_suite(
8            generated_code,
9            test_cases
10       )
11
12       # Evaluate code quality
13       quality_metrics = analyze_code_quality(
14           generated_code,
15           metrics=[
16               "complexity",
17               "maintainability",
18               "documentation"
19           ]
20       )
21
22       return {
23           "test_results": results,
24           "quality_score": quality_metrics
25       }
```

*C. Generation Methods*

AgoraBench evaluates three distinct approaches to data generation:

*1) Instance Generation:* This method focuses on:

- Creating new problem-solution pairs
- Ensuring diversity in generated instances
- Maintaining task difficulty distribution
- Validating semantic correctness

Example implementation:

```
1    class InstanceGenerator:
2        def generate_instance(self, domain, params
             ):
3            # Generate problem-solution pair
4            problem = self.create_problem(domain,
                 params)
5            solution = self.solve_problem(problem)
6
7            # Validate instance
8            if self.validate_instance(problem,
                 solution):
9                return {
10                   "problem": problem,
11                   "solution": solution,
12                   "metadata": self.get_metadata
                         ()
13               }
14           return None
```

*2) Response Generation:* Response generation evaluation includes:

- Answer quality assessment
- Explanation clarity
- Response completeness
- Format adherence

Example response generation:

```python
def generate_response(problem, params):
    # Generate detailed response
    response = llm.generate(
        prompt=format_prompt(problem),
        max_tokens=params["max_tokens"],
        temperature=params["temperature"]
    )

    # Validate response quality
    quality_score = evaluate_response(
        problem,
        response,
        criteria=[
            "completeness",
            "clarity",
            "accuracy"
        ]
    )

    return response, quality_score
```

*3) Quality Enhancement:* This approach examines:

- Improvement of existing instances
- Clarity enhancement
- Error correction
- Content enrichment

Example enhancement process:

```python
def enhance_instance(instance):
    # Analyze current quality
    issues = analyze_quality(instance)

    # Generate improvements
    improvements = generate_improvements(
        instance,
        issues
    )

    # Apply and validate changes
    enhanced = apply_improvements(
        instance,
        improvements
    )

    return enhanced
```

## III. IMPLEMENTATION DETAILS

The framework implementation incorporates modern software engineering practices and provides extensive customization options.

### A. Code Structure

The codebase is organized into logical components:

*1) Core Library:*

- Modular architecture design
- Extensible component interfaces
- Comprehensive utility functions
- Robust error handling

Example core module structure:

```
libs/data-agora/
    data_agora/
        core/
            llm.py
            parser.py
            validator.py
            utils.py
        generators/
            instance_gen.py
            response_gen.py
            quality_enhance.py
        evaluation/
            metrics.py
            benchmarks.py
```

*2) Scripts and Tools:*

- Data processing utilities
- Format conversion tools
- Evaluation scripts
- Benchmark runners

*3) Training Infrastructure:*

- Integration with llama-recipes
- Distributed training support
- Checkpoint management
- Performance optimization

Example training configuration:

```python
training_config = {
    "model_name": "meta-llama/Llama-2-7b",
    "train_batch_size": 4,
    "gradient_accumulation_steps": 8,
    "learning_rate": 1e-5,
    "num_epochs": 3,
    "warmup_steps": 100,
    "evaluation_strategy": "steps",
    "save_steps": 500,
    "eval_steps": 500
}
```

### B. Customization Points

The framework supports extensive customization:

*1) Prompt Templates:*

- Custom template definition
- Variable placeholder system
- Context window management
- Dynamic content insertion

*2) Validation Rules:*

- Domain-specific validators
- Custom quality metrics
- Validation pipeline configuration
- Error handling policies

### 3) Output Formats:
- Custom parser definitions
- Format transformation rules
- Schema validation
- Output normalization

## IV. PRACTICAL APPLICATIONS

Agora enables several practical applications in machine learning research and development.

### A. Data Generation Pipeline

The framework provides a complete pipeline for:

#### 1) Training Data Creation:
- Automated instance generation
- Quality control workflows
- Format standardization
- Dataset augmentation

Example pipeline usage:

```
1   # Configure data generation pipeline
2   pipeline = DataGenerationPipeline(
3       generator=InstanceGenerator(),
4       validator=QualityValidator(),
5       formatter=OutputFormatter()
6   )
7
8   # Generate dataset
9   dataset = pipeline.generate(
10      num_instances=10000,
11      domain="mathematics",
12      quality_threshold=0.8
13  )
14
15  # Export results
16  dataset.save("generated_dataset.jsonl")
```

#### 2) Evaluation Dataset Generation:
- Benchmark dataset creation
- Test case generation
- Edge case identification
- Difficulty level calibration

### B. Model Evaluation

Comprehensive evaluation capabilities include:

#### 1) Performance Metrics:
- Task-specific metrics
- Quality indicators
- Efficiency measures
- Error analysis

Example evaluation workflow:

```
1   def evaluate_model(model, test_set):
2       metrics = {
3           "accuracy": [],
4           "quality_scores": [],
5           "generation_time": [],
6           "memory_usage": []
7       }
8
9       for test_case in test_set:
10          # Generate response
11          start_time = time.time()
12          response = model.generate(test_case)
13          metrics["generation_time"].append(
14              time.time() - start_time
15          )
16
17          # Evaluate response
18          metrics["accuracy"].append(
19              evaluate_accuracy(
20                  test_case,
21                  response
22              )
23          )
24          metrics["quality_scores"].append(
25              assess_quality(response)
26          )
27
28      return compute_aggregate_metrics(metrics)
```

#### 2) Comparative Analysis:
- Model comparison frameworks
- Performance benchmarking
- Capability assessment
- Resource utilization analysis

## V. FUTURE WORK

Several directions for future development have been identified:

### A. Technical Enhancements
- Advanced validation techniques
- Improved quality metrics
- Enhanced parallelization
- Automated optimization

Example future enhancement:

```
1   # Proposed automated optimization system
2   class AutoOptimizer:
3       def optimize_generation_params(
4           self,
5           initial_params,
6           target_metrics
7       ):
8           # Implement Bayesian optimization
9           optimizer = BayesianOptimizer(
10              parameter_space=self.define_space
                    (),
11              objective=self.objective_function
12          )
13
14          # Run optimization loops
15          best_params = optimizer.optimize(
16              n_trials=100,
17              target_metrics=target_metrics
18          )
19
20          return best_params
```

### B. Domain Extensions
- Additional task domains
- Specialized benchmarks
- Custom evaluation metrics
- Domain-specific tools

## C. Integration Capabilities

- Additional model providers
- Extended API support
- Third-party tool integration
- Cloud platform deployment

## VI. CONCLUSION

The Agora framework represents a significant advancement in synthetic data generation and evaluation. Its comprehensive approach to assessing LLM capabilities, combined with practical tools for data generation, provides valuable resources for researchers and practitioners. The framework's modular design and extensive customization options ensure its utility across a wide range of applications and domains.

The standardized evaluation methodology introduced by AgoraBench enables meaningful comparisons between different LLMs' data generation capabilities, while the framework's practical tools facilitate the creation and validation of synthetic datasets. As the field continues to evolve, Agora's extensible architecture positions it well to incorporate new developments and address emerging challenges in synthetic data generation.