# LLM Matrix

**None**

*None*

*None*

# Table of contents

# 1. LLM Matrix

LLM Matrix is a tool for running, evaluating, and comparing different language models across a matrix of hyperparameters.

## 1.1 Overview

LLM Matrix enables you to:

- Define test suites with multiple cases in YAML
- Run test cases across different models and parameters
- Evaluate model responses using specialized metrics
- Generate comprehensive reports to analyze performance
- Compare results across models and parameters

## 1.2 Key Concepts

### 1.2.1 Test Suites

A test suite is a YAML file that defines:

- Test cases (inputs and expected outputs)
- Templates for system and user prompts
- Models to test
- Hyperparameters to vary
- Metrics for evaluation

### 1.2.2 Matrix Execution

LLM Matrix runs each test case against a matrix of parameters, such as:

- Different LLM models (e.g., GPT-4, Claude)
- Various temperature settings
- Different prompt templates

### 1.2.3 Metrics

Custom metrics evaluate model outputs by comparing them to ideal answers:

- QA metrics for question-answering tasks
- Binary classification metrics (YES/NO questions)
- List matching metrics for enumeration tasks

### 1.2.4 Result Analysis

After running a test suite, you can analyze:

- Model performance by score
- Statistical summaries
- Per-case breakdowns
- Comparisons across models and parameters

## 1.3 Installation

```
pip install llm-matrix
```

For optional features:

```
# Excel report support
pip install "llm-matrix[excel]"

# MLflow integration
pip install "llm-matrix[mlflow]"

# LinkML mapping support
pip install "llm-matrix[map]"
```

## 1.4 Quick Links

- Tutorial - Get started with LLM Matrix
- Configuration - Learn how to configure test suites
- Metrics - Available metrics for evaluation
- CLI - Command-line interface reference
- API Reference - Python API documentation

# 2. Getting Started

## 2.1 Getting Started with LLM Matrix

This tutorial walks you through the basic usage of LLM Matrix, from installation to running your first evaluation.

### 2.1.1 Installation

Install LLM Matrix using pip:

```
pip install llm-matrix
```

### 2.1.2 Creating Your First Test Suite

Create a file named `first-test.yaml` with the following content:

```yaml
name: my-first-test
template: simple_qa
templates:
  simple_qa:
    system: Answer the following question accurately and concisely.
    prompt: "{input}"
    metrics:
      - qa_simple
matrix:
  hyperparameters:
    model: [gpt-3.5-turbo]
    temperature: [0.0]
cases:
  - input: What is the capital of France?
    ideal: Paris
    tags: [geography]
  - input: What is 2+2?
    ideal: "4"
    tags: [math]
```

**Understanding the YAML Structure**

- `name` : Identifier for your test suite
- `template` : Default template to use
- `templates` : Define different prompt templates
- `system` : System instructions to the model
- `prompt` : Template for user prompt (with variables in curly braces)
- `metrics` : Evaluation metrics to apply
- `matrix.hyperparameters` : Parameters to vary in the test
- `cases` : Individual test cases
- `input` : Query to send to the model
- `ideal` : Expected response (used for evaluation)
- `tags` : Optional categorization

### 2.1.3 Running the Test Suite

Execute your test suite with:

```
llm-matrix run first-test.yaml
```

This will: 1. Create a database file to cache results 2. Run each test case against the specified model(s) 3. Generate evaluation reports

## 2.1.4 Viewing Results

After running the test suite, you'll find a `first-test-output` directory with:

- `results.csv` : Raw results for all runs

- `summary.csv` : Statistical summary of model performance

- `by_model.csv` : Performance breakdown by model

- `grouped_by_input.tsv` : Detailed breakdown by test case

## 2.1.5 Next Steps

- Configuring Test Suites: Learn about advanced configuration options

- Understanding Metrics: Explore different evaluation metrics

- CLI Reference: Discover all command-line options

# 3. User Guide

## 3.1 Configuration Guide

This guide covers how to configure LLM Matrix test suites and runners.

### 3.1.1 Test Suite Configuration

Test suites are defined in YAML files with the following structure:

```yaml
name: suite-name
template: default-template-name
templates:
  template-name:
    system: System prompt text
    prompt: User prompt template with {variables}
    metrics:
      - metric_name1
      - metric_name2
matrix:
  hyperparameters:
    model: [model1, model2]
    temperature: [0.0, 0.7]
    # Other parameters as needed
cases:
  - input: Input text for the first case
    ideal: Expected output text
    tags: [tag1, tag2]
  - input: Input text for the second case
    ideal: Expected output text
    tags: [tag3]
    # Case-specific parameters to override defaults
    template: special-template
```

**Core Components**

**TEMPLATES**

Templates define how inputs are formatted for the LLM:

```yaml
templates:
  binary_qa:
    system: Answer the following question with YES or NO only.
    prompt: "Question: {input}"
    metrics:
      - binary_exact

  complex_qa:
    system: Provide a detailed answer to the following question.
    prompt: "Question: {input}\nPlease explain in detail."
    metrics:
      - qa_with_explanation
```

**MATRIX CONFIGURATION**

The matrix defines parameter combinations to test:

```yaml
matrix:
  hyperparameters:
    model: [gpt-4o, gpt-3.5-turbo, claude-3-opus]
    temperature: [0.0, 0.7]
    max_tokens: [100, 500]
```

LLM Matrix will run each test case with every combination of these parameters.

**TEST CASES**

Individual test cases define inputs and expected outputs:

```yaml
cases:
  - input: Is water wet?
    ideal: "YES"
    tags: [science, basic]
```

```
  - input: |
      What are the three primary colors?
    ideal: "The three primary colors are red, blue, and yellow."
    tags: [art]
    template: list_qa  # Override the default template
```

## 3.1.2 Runner Configuration

You can customize the runner behavior with a separate config file:

```
concurrency: 5  # Number of concurrent API calls
retries: 3      # Number of retry attempts for failed calls
timeout: 30     # Timeout in seconds for API calls
plugins:
  - citeseek    # Enable plugins
```

Pass this config to the CLI with:

```
llm-matrix run test-suite.yaml --runner-config runner-config.yaml
```

## 3.1.3 Advanced Configuration

### Using Variables

You can use variables in your prompts:

```
templates:
  complex_prompt:
    system: "You are an expert in {domain}."
    prompt: "Question about {domain}: {input}"
```

Then in your cases:

```
cases:
  - input: What is DNA?
    ideal: "DNA is a molecule that carries genetic information..."
    variables:
      domain: genetics
```

### Tags for Analysis

Tags help organize and filter your results:

```
cases:
  - input: What is photosynthesis?
    ideal: "Photosynthesis is a process used by plants..."
    tags: [biology, plants, difficulty:easy]
```

You can later filter or group results by these tags.

## 3.1.4 Examples

See the examples directory for complete examples of test suite configurations.

## 3.2 Evaluation Metrics

LLM Matrix provides various metrics to evaluate model responses against expected outputs. These metrics are specified in the template configuration.

### 3.2.1 Available Metrics

**Basic Metrics**

`binary_exact`

Evaluates binary (YES/NO) responses with exact matching.

```
templates:
  binary_qa:
    system: Answer the following question with YES or NO only.
    prompt: "Question: {input}"
    metrics:
      - binary_exact
```

• Score: 1.0 for exact match, 0.0 otherwise
• Best for: Simple YES/NO questions where only the answer matters, not explanation

`qa_simple`

Simple string matching for question answering.

```
metrics:
  - qa_simple
```

• Score: 1.0 for exact match, partial scores for close matches
• Best for: Factual questions with specific answers

**Explanation Metrics**

`qa_with_explanation`

Evaluates both the answer and explanation.

```
metrics:
  - qa_with_explanation
```

• Score: Combined score for answer correctness and explanation quality
• Best for: Questions requiring both correct answers and explanations

`binary_with_explanation`

For YES/NO questions with explanations.

```
metrics:
  - binary_with_explanation
```

• Score: 1.0 for correct binary answer with good explanation, lower for partial matches
• Best for: Binary questions where explanation is important

**List Metrics**

`list_comparison`

Compares lists of items.

```yaml
metrics:
  - list_comparison
```

- Score: Based on overlap between expected and actual lists

- Best for: Enumeration tasks (e.g., "List the planets in the solar system")

## 3.2.2 Creating Custom Metrics

You can create custom metrics by implementing the `Metric` class:

```python
from llm_matrix.metrics import Metric
from llm_matrix.schema import EvalResult, TestCase, LLMResponse

class MyCustomMetric(Metric):
    def evaluate(self, case: TestCase, response: LLMResponse) -> EvalResult:
        # Implement your evaluation logic
        score = calculate_score(case.ideal, response.text)

        return EvalResult(
            score=score,
            explanation="Reason for this score"
        )
```

Register your custom metric:

```python
from llm_matrix.metrics import register_metric

register_metric("my_custom_metric", MyCustomMetric())
```

Then use it in your templates:

```yaml
templates:
  my_template:
    system: Custom prompt
    prompt: "{input}"
    metrics:
      - my_custom_metric
```

## 3.2.3 Best Practices

- Choose metrics appropriate for your task type

- Consider using multiple metrics for complex tasks

- For binary tasks, prefer specialized binary metrics over general ones

- When using list-based metrics, ensure your ideal answer is formatted as expected

- For most accurate scoring, define clear evaluation criteria in your prompts

## 3.2.4 Metrics API Reference

For detailed API documentation of all metrics, see the Metrics API Reference.

## 3.3 Command Line Interface

LLM Matrix provides a command-line interface for running evaluations and managing results.

### 3.3.1 Basic Usage

```
llm-matrix run my-suite.yaml
```

### 3.3.2 Global Options

```
--verbose, -v     Increase verbosity (can be used multiple times)
--help            Show help message and exit
```

### 3.3.3 Commands

`run`

Run a test suite against specified models.

```
llm-matrix run <suite-path> [options]
```

**ARGUMENTS**

- `suite-path` : Path to the evaluation suite YAML file

**OPTIONS**

- `--store-path, -s <path>` : Path to the cache store (defaults to same directory as suite with .db extension)
- `--runner-config, -C <path>` : Path to the runner config file
- `--output-file, -o <path>` : Path to save output file
- `--output-dir, -D <path>` : Directory to save output files (defaults to suite-name-output)
- `--output-format, -F <format>` : Output format (csv, tsv, excel, jsonl, json, yaml)

**EXAMPLES**

```
# Basic usage
llm-matrix run my-suite.yaml

# Specify custom output location
llm-matrix run my-suite.yaml -D ./results

# Use custom runner configuration
llm-matrix run my-suite.yaml -C runner-config.yaml

# Increase verbosity for debugging
llm-matrix run my-suite.yaml -vv

# Output as Excel file
llm-matrix run my-suite.yaml -o results.xlsx -F excel
```

`convert`

Convert between different file formats.

```
llm-matrix convert <input-files> [options]
```

**ARGUMENTS**

- `input-files` : Paths to files to be converted

**OPTIONS**

- `--source, -s <field>` : Source field

- `--target, -t <field>` : Target field

### 3.3.4 Output Formats

LLM Matrix supports the following output formats:

- `csv` : Comma-separated values
- `tsv` : Tab-separated values
- `excel` : Microsoft Excel format (requires the `excel` extra)
- `jsonl` : JSON Lines format (one JSON object per line)
- `json` : Single JSON array with all results
- `yaml` : YAML format

### 3.3.5 Output Directory Structure

When using the default output directory, LLM Matrix creates:

```
suite-name-output/
├── by_model.csv            # Performance summary by model
├── by_model_ideal.csv      # Performance by model and ideal answer
├── grouped_by_input.tsv    # Detailed view of each test case
├── grouped_by_input.xlsx   # Excel version of the above
├── results.csv             # Raw results
├── results.html            # HTML view of raw results
├── summary.csv             # Statistical summary
└── summary.html            # HTML view of summary
```

# 4. API Reference

This section provides detailed documentation of the LLM Matrix Python API.

## 4.1 Core Modules

### 4.1.1 LLMRunner

The main class for running evaluations.

```python
from llm_matrix import LLMRunner

runner = LLMRunner(store_path="results.db")
results = runner.run(suite)
```

**Source code in** `src/llm_matrix/runner.py` ⌄

```python
19    @dataclass
20    class LLMRunner:
21
22        store_path: Optional[Path] = None
23        _aimodels: Optional[Dict[tuple, AIModel]] = None
24        _store: Optional[Store] = None
25        config: Optional[LLMRunnerConfig] = None
26
27        def run(self, suite: Suite) -> List[TestCaseResult]:
28            """
29            Run the suite of cases
30
31            :param suite:
32            :return:
33            """
34            return list(self.run_iter(suite))
35
36        def run_iter(self, suite: Suite) -> Iterator[TestCaseResult]:
37            """
38            Run the suite of cases iterating over the results.
39
40            :param suite:
41            :return:
42            """
43            logger.info(f"Running suite {suite.name}")
44            for params in iter_hyperparameters(suite.matrix):
45                logger.info(f"Running {len(suite.cases)} with {params}")
46                for n, case in enumerate(suite.cases):
47                    logger.debug(f"Running case {n+1}/{len(suite.cases)}")
48                    result = self.run_case(case, params, suite)
49                    yield result
50
51        def run_case(self, case: TestCase, params: Dict[str, Any], suite: Suite) -> TestCaseResult:
52            logger.info(f"Running case {case.input} with {params}")
53            store = self._get_store()
54            cached = store.get_result(suite, case, params)
55            if cached:
56                return cached
57            actual_params = copy(params)
58            if self.config and "model" in params:
59                model_logical_name = params["model"]
60                actual_params["model"] = self.config.model_name_map.get(model_logical_name, model_logical_name)
61                logger.info(f"Mapping model {model_logical_name} to {actual_params['model']}")
62            model = self.get_aimodel(actual_params, suite=suite)
63            template = self.get_template(case, suite)
64            response = model.prompt(case.input, template=template, case=case)
65            result = TestCaseResult(
66                case=case,
67                response=response,
68                hyperparameters=params,
69                metrics=template.metrics if template else None,
70            )
71            from llm_matrix.metrics import evaluate_result
72            evaluate_result(result, runner=self)
73            store.add_result(suite, result)
74            return result
75
76        def get_template(self, case, suite: Suite) -> Optional[Template]:
77            tn = case.template
78            if not tn:
79                tn = suite.template
80            if not tn:
81                return None
82            return suite.templates[tn]
83
84        def get_aimodel(self, params: Dict[str, Any], suite: Suite=None) -> AIModel:
85            if not self._aimodels:
86                self._aimodels = {}
87            key = tuple(sorted(params.items()))
88            if key not in self._aimodels:
89                if suite and suite.models:
90                    bespoke_models = suite.models
91                    model_name = params.get("model")
92                    if model_name and model_name in bespoke_models:
93                        model_info = bespoke_models[model_name]
94                        params = {**params, **model_info.parameters}
95                        plugins = model_info.plugins or []
96                        # TODO: proper mechanism for plugins
97                        if "citeseek" in plugins:
98                            from llm_matrix.plugins.citeseek_plugin import CiteseekPlugin
99                            self._aimodels[key] = CiteseekPlugin(parameters=params)
100               if key not in self._aimodels:
101                   self._aimodels[key] = AIModel(parameters=params)
102           m = self._aimodels[key]
103           return self._aimodels[key]
104
105       def _get_store(self) -> Store:
106           if not self._store:
107               if not self.store_path:
108                   self.store_path = Path("results") / "cache.db"
109               self.store_path.parent.mkdir(parents=True, exist_ok=True)
110               self._store = Store(self.store_path)
111           return self._store
```

## 4.2 `run(suite)`

Run the suite of cases

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| suite | Suite | | *required* |

**Returns:**

| Type | Description |
|------|-------------|
| List[TestCaseResult] | |

**"** Source code in `src/llm_matrix/runner.py` ⌄

```
27    def run(self, suite: Suite) -> List[TestCaseResult]:
28        """
29        Run the suite of cases
30
31        :param suite:
32        :return:
33        """
34        return list(self.run_iter(suite))
```

## 4.3 `run_iter(suite)`

Run the suite of cases iterating over the results.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| suite | Suite | | *required* |

**Returns:**

| Type | Description |
|------|-------------|
| Iterator[TestCaseResult] | |

**"** Source code in `src/llm_matrix/runner.py` ⌄

```
36    def run_iter(self, suite: Suite) -> Iterator[TestCaseResult]:
37        """
38        Run the suite of cases iterating over the results.
39
40        :param suite:
41        :return:
42        """
43        logger.info(f"Running suite {suite.name}")
44        for params in iter_hyperparameters(suite.matrix):
45            logger.info(f"Running {len(suite.cases)} with {params}")
46            for n, case in enumerate(suite.cases):
47                logger.debug(f"Running case {n+1}/{len(suite.cases)}")
48                result = self.run_case(case, params, suite)
49                yield result
```

### 4.3.1 Schema

Data models for test cases, templates, responses, and results.

## 4.4 `StrictBaseModel`

Bases: `BaseModel`

Base class for Pydantic models that forbids extra fields.

> **Source code in `src/llm_matrix/schema.py`** ˅
>
> ```python
> 18    class StrictBaseModel(BaseModel):
> 19        """
> 20        Base class for Pydantic models that forbids extra fields.
> 21        """
> 22        model_config = ConfigDict(extra="forbid")
> 23
> 24        def as_flat_dict(self, simple=True, prefix = None) -> Dict[str, Any]:
> 25            """
> 26            Convert the model to a flat dictionary.
> 27
> 28            Suitable for conversion to a DataFrame.
> 29
> 30            :param simple:
> 31            :param prefix:
> 32            :return:
> 33            """
> 34            def is_complex(v):
> 35                if isinstance(v, dict):
> 36                    return True
> 37                if isinstance(v, list) and any(isinstance(i, dict) for i in v):
> 38                    return True
> 39                return False
> 40            def mk_key(k):
> 41                if prefix:
> 42                    return f"{prefix}_{k}"
> 43                return k
> 44            d = {mk_key(k): v for k, v in self.model_dump().items() if not is_complex(v)}
> 45            return d
> ```

### 4.4.1 `as_flat_dict(simple=True, prefix=None)`

Convert the model to a flat dictionary.

Suitable for conversion to a DataFrame.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| `simple` | | | `True` |
| `prefix` | | | `None` |

**Returns:**

| Type | Description |
|------|-------------|
| `Dict[str, Any]` | |

**Source code in** `src/llm_matrix/schema.py`  ⌄

```python
24    def as_flat_dict(self, simple=True, prefix = None) -> Dict[str, Any]:
25        """
26        Convert the model to a flat dictionary.
27
28        Suitable for conversion to a DataFrame.
29
30        :param simple:
31        :param prefix:
32        :return:
33        """
34        def is_complex(v):
35            if isinstance(v, dict):
36                return True
37            if isinstance(v, list) and any(isinstance(i, dict) for i in v):
38                return True
39            return False
40        def mk_key(k):
41            if prefix:
42                return f"{prefix}_{k}"
43            return k
44        d = {mk_key(k): v for k, v in self.model_dump().items() if not is_complex(v)}
45        return d
```

## 4.5  `MetricEnum`

Bases: `Enum`

Enum for metrics to be evaluated.

Designed to be extensible.

**Source code in** `src/llm_matrix/schema.py`  ⌄

```python
47    class MetricEnum(Enum):
48        """
49        Enum for metrics to be evaluated.
50
51        Designed to be extensible.
52        """
53        QA_WITH_EXPLANATION = "qa_with_explanation"
54        SIMPLE_QUESTION = "simple_question"
55        LIST_MEMBERSHIP = "list_membership"
56        RANKED_LIST = "ranked_list"
57        NARRATIVE_SIMILARITY = "narrative_similarity"
58        REVIEW = "review"
59        BLEU = "bleu"
60        ROUGE = "rouge"
61        METEOR = "meteor"
```

## 4.6  `Response`

Bases: `StrictBaseModel`

Response from the AI model.

**Source code in** `src/llm_matrix/schema.py`  ⌄

```python
63    class Response(StrictBaseModel):
64        """
65        Response from the AI model.
66        """
67        text : str = Field(..., description="The text of the response from the AI model")
68        prompt: Optional[str] = Field(None, description="The prompt used to generate the response")
69        system: Optional[str] = Field(None, description="The system prompt used to generate the response")
```

## 4.7 `Template`

Bases: `StrictBaseModel`

Template for generating prompts to the AI model.

> **Source code in `src/llm_matrix/schema.py`** ∨
>
> ```python
> 72   class Template(StrictBaseModel):
> 73       """
> 74       Template for generating prompts to the AI model.
> 75       """
> 76       system: Optional[FormatString] = Field(None, description="System prompt")
> 77       prompt: Optional[FormatString] = Field(None, description="Prompt format string to the model")
> 78       metrics: Optional[List[Metric]] = Field(None, description="Metrics to be evaluated")
> ```

## 4.8 `Matrix`

Bases: `StrictBaseModel`

Specifies a combination of hyperparameters to be evaluated.

> **Source code in `src/llm_matrix/schema.py`** ∨
>
> ```python
> 80   class Matrix(StrictBaseModel):
> 81       """
> 82       Specifies a combination of hyperparameters to be evaluated.
> 83       """
> 84       hyperparameters: Dict[Hyperparameter, List[Any]] = Field(..., description="Hyperparameters to be evaluated")
> 85       exclude: Optional["Matrix"] = Field(None, description="Hyperparameters to be excluded")
> ```

## 4.9 `TestCase`

Bases: `StrictBaseModel`

Test case for the AI model.

> **Source code in `src/llm_matrix/schema.py`** ∨
>
> ```python
> 87    class TestCase(StrictBaseModel):
> 88        """
> 89        Test case for the AI model.
> 90        """
> 91        input: str = Field(..., description="Input to the model")
> 92        original_input: Optional[Dict[str, Any]] = Field(
> 93            None,
> 94            description="""Original input to the model, prior to transformation into text.
> 95            An example is a structured data record.
> 96            """
> 97        )
> 98        #output: Optional[str] = Field(None, description="Provided output from the model")
> 99        ideal: Optional[str] = Field(None, description="Ideal output from the model")
> 100       template: Optional[TemplateName] = Field(None, description="Template for the test case")
> 101       tags: Optional[List[str]] = Field(None, description="Tags for the test case")
> 102       comments: Optional[List[str]] = Field(None, description="Comments for the test case")
> 103
> 104       def as_flat_dict(self, **kwargs) -> Dict[str, Any]:
> 105           top_dict = super().as_flat_dict(**kwargs)
> 106           orig = self.original_input or {}
> 107           return {**top_dict, **orig}
> ```

## 4.10 `TestCaseResult`

Bases: `StrictBaseModel`

Result of running an individual case on a model with hyperparameters.

**Source code in** `src/llm_matrix/schema.py` ⌄

```python
110   class TestCaseResult(StrictBaseModel):
111       """
112       Result of running an individual case on a model with hyperparameters.
113       """
114       case: TestCase = Field(..., description="Test case")
115       response: Response = Field(..., description="Response from the model")
116       hyperparameters: Dict[Hyperparameter, Any] = Field(..., description="Hyperparameters used for the test case")
117       metrics: Optional[List[Metric]] = Field(None, description="Metrics to be evaluated")
118       score: Optional[float] = Field(None, description="Score for the test case", ge=0, le=1)
119       evaluation_message: Optional[str] = Field(None, description="Message from the evaluation")
120
121       def as_flat_dict(self, **kwargs) -> Dict[str, Any]:
122           top_dict = super().as_flat_dict(**kwargs)
123           case_dict = self.case.as_flat_dict(prefix="case")
124           response_dict = self.response.as_flat_dict(prefix="response")
125           hyperparameters_dict = {k: str(v) for k, v in self.hyperparameters.items()}
126           hyperparameters_str = "_".join(f"{k}={v}" for k, v in hyperparameters_dict.items())
127           return {"hyperparameters": hyperparameters_str, **top_dict, **case_dict, **response_dict, **hyperparameters_dict}
```

## 4.10.1 Metrics

Evaluation metrics for comparing model outputs to expected answers.

## 4.11 `MetricEvaluator`

Bases: `ABC`

Base class for metric evaluators.

**Source code in** `src/llm_matrix/metrics.py` ⌄

```python
13   class MetricEvaluator(ABC):
14       """Base class for metric evaluators."""
15
16       @abstractmethod
17       def evaluate(
18           self,
19           actual_output: str,
20           expected_output: str,
21           runner: Optional[LLMRunner] = None,
22           result: Optional[TestCaseResult] = None
23       ) -> float:
24           """
25           Evaluate the result and return a score between 0 and 1.
26
27           :param actual_output: The actual output from the model
28           :param expected_output: The expected output
29           :param runner: The LLMRunner instance
30           :param result: The TestCaseResult instance, for storing evaluation messages
31           :return: A score between 0 and 1
32           """
33           pass
```

### 4.11.1 `evaluate(actual_output, expected_output, runner=None, result=None)` `abstractmethod`

Evaluate the result and return a score between 0 and 1.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| actual_output | str | The actual output from the model | *required* |
| expected_output | str | The expected output | *required* |
| runner | Optional[LLMRunner] | The LLMRunner instance | None |
| result | Optional[TestCaseResult] | The TestCaseResult instance, for storing evaluation messages | None |

**Returns:**

| Type | Description |
|------|-------------|
| float | A score between 0 and 1 |

**Source code in** `src/llm_matrix/metrics.py` ⌄

```
16    @abstractmethod
17    def evaluate(
18        self,
19        actual_output: str,
20        expected_output: str,
21        runner: Optional[LLMRunner] = None,
22        result: Optional[TestCaseResult] = None
23    ) -> float:
24        """
25        Evaluate the result and return a score between 0 and 1.
26
27        :param actual_output: The actual output from the model
28        :param expected_output: The expected output
29        :param runner: The LLMRunner instance
30        :param result: The TestCaseResult instance, for storing evaluation messages
31        :return: A score between 0 and 1
32        """
33        pass
```

## 4.12 `QAWithExplanationEvaluator`

Bases: `MetricEvaluator`

Evaluator for QA with explanation metrics.

```python
36    class QAWithExplanationEvaluator(MetricEvaluator):
37        """Evaluator for QA with explanation metrics."""
38
39        def evaluate(
40            self,
41            actual_output: str,
42            expected_output: str,
43            runner: Optional[LLMRunner] = None,
44            result: Optional[TestCaseResult] = None
45        ) -> float:
46            """Evaluate QA with explanation result."""
47            # First token regex
48            pattern = re.compile(r"^(\w+)")
49            match = pattern.match(actual_output)
50            actual_answer = match.group(1).upper() if match else "OTHER"
51
52            expected_match = pattern.match(expected_output)
53            if not expected_match:
54                logger.warning(f"Could not extract first token from expected output: {expected_output}")
55                return 0.0
56
57            expected_answer = expected_match.group(1).upper()
58
59            if actual_answer == expected_answer:
60                score = 1.0
61            elif actual_answer == "OTHER" or expected_answer == "OTHER":
62                score = 0.5
63            else:
64                score = 0.0
65
66            return score
```

## 4.12.1 `evaluate(actual_output, expected_output, runner=None, result=None)`

Evaluate QA with explanation result.

Source code in `src/llm_matrix/metrics.py` ⌄

```python
39    def evaluate(
40        self,
41        actual_output: str,
42        expected_output: str,
43        runner: Optional[LLMRunner] = None,
44        result: Optional[TestCaseResult] = None
45    ) -> float:
46        """Evaluate QA with explanation result."""
47        # First token regex
48        pattern = re.compile(r"^(\w+)")
49        match = pattern.match(actual_output)
50        actual_answer = match.group(1).upper() if match else "OTHER"
51
52        expected_match = pattern.match(expected_output)
53        if not expected_match:
54            logger.warning(f"Could not extract first token from expected output: {expected_output}")
55            return 0.0
56
57        expected_answer = expected_match.group(1).upper()
58
59        if actual_answer == expected_answer:
60            score = 1.0
61        elif actual_answer == "OTHER" or expected_answer == "OTHER":
62            score = 0.5
63        else:
64            score = 0.0
65
66        return score
```

## 4.13 `LLMBasedEvaluator`

Bases: `MetricEvaluator`

Base class for evaluators that use LLMs for scoring.

**Source code in** `src/llm_matrix/metrics.py` ⌄

```python
69    class LLMBasedEvaluator(MetricEvaluator):
70        """Base class for evaluators that use LLMs for scoring."""
71
72        def __init__(self, system_prompt: str, user_input_template: str):
73            self.system_prompt = system_prompt
74            self.user_input_template = user_input_template
75
76        def evaluate(
77            self,
78            actual_output: str,
79            expected_output: str,
80            runner: Optional[LLMRunner] = None,
81            result: Optional[TestCaseResult] = None
82        ) -> float:
83            """Evaluate using an LLM."""
84            if not runner:
85                logger.error("Runner is required for LLM-based evaluation")
86                return 0.0
87
88            eval_model = self._get_eval_model(runner)
89            if not eval_model:
90                logger.error("Could not get evaluation model")
91                return 0.0
92
93            eval_response = self._prompt_model(eval_model, actual_output, expected_output)
94            eval_response_text = eval_response.text.strip()
95
96            if result:
97                result.evaluation_message = eval_response_text
98
99            score = self._extract_score(eval_response_text)
100           return score
101
102       def _get_eval_model(self, runner: LLMRunner):
103           """Get the evaluation model."""
104           if runner.config and runner.config.evaluation_model_name:
105               eval_model_name = runner.config.evaluation_model_name
106           else:
107               eval_model_name = DEFAULT_EVALUATION_MODEL_NAME
108
109           return runner.get_aimodel({"model": eval_model_name})
110
111       def _prompt_model(self, model, actual_output: str, expected_output: str):
112           """Prompt the evaluation model."""
113           user_input = self._format_user_input(actual_output, expected_output)
114
115           return model.prompt(
116               system_prompt=self.system_prompt,
117               user_input=user_input
118           )
119
120       def _format_user_input(self, actual_output: str, expected_output: str) -> str:
121           """Format the user input for the evaluation model."""
122           return self.user_input_template.format(
123               actual_output=actual_output,
124               expected_output=expected_output
125           )
126
127       def _extract_score(self, response_text: str) -> float:
128           """Extract a score from the LLM response."""
129           pattern = re.compile(r"(\d+(\.\d+)?)")
130           matches = pattern.match(response_text)
131
132           if matches:
133               return float(matches.group(1))
134           else:
135               logger.error(f"Could not parse score from {response_text}")
136               raise ValueError(f"Could not parse score from {response_text}")
```

### 4.13.1 `evaluate(actual_output, expected_output, runner=None, result=None)`

Evaluate using an LLM.

Source code in `src/llm_matrix/metrics.py` ⌄

```
 76    def evaluate(
 77        self,
 78        actual_output: str,
 79        expected_output: str,
 80        runner: Optional[LLMRunner] = None,
 81        result: Optional[TestCaseResult] = None
 82    ) -> float:
 83        """Evaluate using an LLM."""
 84        if not runner:
 85            logger.error("Runner is required for LLM-based evaluation")
 86            return 0.0
 87
 88        eval_model = self._get_eval_model(runner)
 89        if not eval_model:
 90            logger.error("Could not get evaluation model")
 91            return 0.0
 92
 93        eval_response = self._prompt_model(eval_model, actual_output, expected_output)
 94        eval_response_text = eval_response.text.strip()
 95
 96        if result:
 97            result.evaluation_message = eval_response_text
 98
 99        score = self._extract_score(eval_response_text)
100        return score
```

### 4.13.2 _get_eval_model(runner)

Get the evaluation model.

Source code in `src/llm_matrix/metrics.py` ⌄

```
102    def _get_eval_model(self, runner: LLMRunner):
103        """Get the evaluation model."""
104        if runner.config and runner.config.evaluation_model_name:
105            eval_model_name = runner.config.evaluation_model_name
106        else:
107            eval_model_name = DEFAULT_EVALUATION_MODEL_NAME
108
109        return runner.get_aimodel({"model": eval_model_name})
```

### 4.13.3 _prompt_model(model, actual_output, expected_output)

Prompt the evaluation model.

Source code in `src/llm_matrix/metrics.py` ⌄

```
111    def _prompt_model(self, model, actual_output: str, expected_output: str):
112        """Prompt the evaluation model."""
113        user_input = self._format_user_input(actual_output, expected_output)
114
115        return model.prompt(
116            system_prompt=self.system_prompt,
117            user_input=user_input
118        )
```

### 4.13.4 _format_user_input(actual_output, expected_output)

Format the user input for the evaluation model.

> **Source code in** `src/llm_matrix/metrics.py` ⌄

```python
120   def _format_user_input(self, actual_output: str, expected_output: str) -> str:
121       """Format the user input for the evaluation model."""
122       return self.user_input_template.format(
123           actual_output=actual_output,
124           expected_output=expected_output
125       )
```

## 4.13.5 `_extract_score(response_text)`

Extract a score from the LLM response.

> **Source code in** `src/llm_matrix/metrics.py` ⌄

```python
127   def _extract_score(self, response_text: str) -> float:
128       """Extract a score from the LLM response."""
129       pattern = re.compile(r"(\d+(\.\d+)?)")
130       matches = pattern.match(response_text)
131
132       if matches:
133           return float(matches.group(1))
134       else:
135           logger.error(f"Could not parse score from {response_text}")
136           raise ValueError(f"Could not parse score from {response_text}")
```

## 4.14 `ListMembershipEvaluator`

Bases: `LLMBasedEvaluator`

Evaluator for list membership metrics.

> **Source code in** `src/llm_matrix/metrics.py` ⌄

```python
139   class ListMembershipEvaluator(LLMBasedEvaluator):
140       """Evaluator for list membership metrics."""
141
142       def __init__(self):
143           super().__init__(
144               system_prompt=(
145                   "Check if all the expected list items are present in the text. "
146                   "Your response should be an overlap score between 0 and 1, where 1 is a perfect "
147                   "match (all members match) and 0 is the worst possible match (no members match). "
148                   "Your response should be the score followed by any explanatory text. "
149                   "For example, '0.5 Only half of the items matched'. "
150                   "Do NOT put ANY text before the score. ALWAYS start with the score. "
151                   "Note the text you are evaluating may have additional verbiage, do not "
152                   "penalize this. Your task is just to determine if the list is presented clearly "
153                   "and if the items match"
154               ),
155               user_input_template=(
156                   "The expected list is: {expected_output}. "
157                   "The text: {actual_output}. "
158               )
159           )
```

## 4.15 `ReviewEvaluator`

Bases: `LLMBasedEvaluator`

Evaluator for review metrics.

> **Source code in** `src/llm_matrix/metrics.py` ∨
>
> ```
> 162    class ReviewEvaluator(LLMBasedEvaluator):
> 163        """Evaluator for review metrics."""
> 164
> 165        def __init__(self):
> 166            super().__init__(
> 167                system_prompt=(
> 168                    "Review the output for correctness, completeness, and clarity. "
> 169                    "The response should be a score between 0 (worst) and 1 (best). "
> 170                    "Your response should be the score followed by any explanatory text. "
> 171                    "For example, '0.3 The response has many inaccuracies'. "
> 172                ),
> 173                user_input_template="The output to score is: {actual_output}. "
> 174            )
> ```

## 4.16 `RankedListEvaluator`

Bases: `LLMBasedEvaluator`

Evaluator for ranked list metrics.

> **Source code in** `src/llm_matrix/metrics.py` ∨
>
> ```
> 177    class RankedListEvaluator(LLMBasedEvaluator):
> 178        """Evaluator for ranked list metrics."""
> 179
> 180        def __init__(self):
> 181            super().__init__(
> 182                system_prompt=(
> 183                    "Compare the ranked list to the expected output. "
> 184                    "The response should be a score between 0 and 1. "
> 185                    "If the item ranked first is equal to the expected item, score is 1. "
> 186                    "If there is no overlap between the ranked list and the expected list, score is 0. "
> 187                    "Otherwise score according to rank, with 0.5 for 2nd, 0.25 for 3rd, and so on."
> 188                ),
> 189                user_input_template=(
> 190                    "The expected answer is: {expected_output}. "
> 191                    "The output to score is: {actual_output}. "
> 192                )
> 193            )
> ```

## 4.17 `SimpleQuestionEvaluator`

Bases: `LLMBasedEvaluator`

Evaluator for simple question metrics.

> **Source code in** `src/llm_matrix/metrics.py` ∨
>
> ```
> 196    class SimpleQuestionEvaluator(LLMBasedEvaluator):
> 197        """Evaluator for simple question metrics."""
> 198
> 199        def __init__(self):
> 200            super().__init__(
> 201                system_prompt=(
> 202                    "Compare the answer given to the expected output. "
> 203                    "The response should be a score between 0 and 1. "
> 204                    "The answer should be provided first, explanations may follow "
> 205                    "A precise correct answer is 1, a wrong answer is 0."
> 206                    "You can use values in between for imprecise answers"
> 207                ),
> 208                user_input_template=(
> 209                    "The expected answer is: {expected_output}. "
> 210                    "The output to score is: {actual_output}. "
> 211                )
> 212            )
> ```

## 4.18 `register_metric_evaluator(metric_name, evaluator)`

Register a custom metric evaluator.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| metric_name | str | The name of the metric | *required* |
| evaluator | MetricEvaluator | The evaluator instance | *required* |

**Source code in `src/llm_matrix/metrics.py`** ⌄

```python
225    def register_metric_evaluator(metric_name: str, evaluator: MetricEvaluator) -> None:
226        """
227        Register a custom metric evaluator.
228
229        :param metric_name: The name of the metric
230        :param evaluator: The evaluator instance
231        """
232        METRIC_REGISTRY[metric_name] = evaluator
```

## 4.19 `evaluate_result(result, runner=None)`

Evaluate the result of a test case.

Example:

```
>>> result = TestCaseResult(
...     case=TestCase(input="What is II+IV?", ideal="VI. Blah"),
...     response=Response(text="VI"),
...     hyperparameters={"model": "gpt-4o"},
...     metrics=["qa_with_explanation"],
... )
>>> evaluate_result(result)
>>> result.score
1.0
>>> result.response.text = "VII"
>>> evaluate_result(result)
>>> result.score
0.0
>>> result.response.text = "Other"
>>> evaluate_result(result)
>>> result.score
0.5
```

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| result | TestCaseResult | The test case result to evaluate | *required* |
| runner | Optional[LLMRunner] | The LLMRunner instance | None |

**Source code in** `src/llm_matrix/metrics.py` ⌄

```
235    def evaluate_result(result: TestCaseResult, runner: Optional[LLMRunner] = None):
236        """
237        Evaluate the result of a test case.
238
239        Example:
240
241            >>> result = TestCaseResult(
242            ...     case=TestCase(input="What is II+IV?", ideal="VI. Blah"),
243            ...     response=Response(text="VI"),
244            ...     hyperparameters={"model": "gpt-4o"},
245            ...     metrics=["qa_with_explanation"],
246            ... )
247            >>> evaluate_result(result)
248            >>> result.score
249            1.0
250            >>> result.response.text = "VII"
251            >>> evaluate_result(result)
252            >>> result.score
253            0.0
254            >>> result.response.text = "Other"
255            >>> evaluate_result(result)
256            >>> result.score
257            0.5
258
259        :param result: The test case result to evaluate
260        :param runner: The LLMRunner instance
261        """
262        actual_output = result.response.text
263        expected_output = result.case.ideal
264
265        scores = []
266
267        for metric_name in result.metrics or []:
268            if metric_name not in METRIC_REGISTRY:
269                raise NotImplementedError(f"Metric {metric_name} not implemented")
270
271            evaluator = METRIC_REGISTRY[metric_name]
272            try:
273                score = evaluator.evaluate(
274                    actual_output=actual_output,
275                    expected_output=expected_output,
276                    runner=runner,
277                    result=result
278                )
279                scores.append(score)
280            except Exception as e:
281                logger.error(f"Error evaluating metric {metric_name}: {e}")
282                # Optionally fall back to a default score
283                # scores.append(0.0)
284                # Or re-raise the exception
285                raise
286
287        if scores:
288            result.score = sum(scores) / len(scores)
```

## 4.19.1 Store

Cache for storing results and avoiding redundant API calls.

## 4.20 `Store` `dataclass`

A persistent store using DuckDB to cache test results with JSON support for Pydantic models.

Example:

```
>>> store = Store("test-cache.db")
>>> case = TestCase(input="1+1", ideal="2")
>>> suite = Suite(name="test", cases=[case], matrix={"hyperparameters": {}})
>>> response = Response(text="2")
>>> result = TestCaseResult(case=case, response=response, hyperparameters={"model": "gpt-4"})
>>> store.add_result(suite, result)
>>> cached = store.get_result(suite, case, {"model": "gpt-4"})
>>> assert cached.response == response
```

To use an in-memory database, pass `None` as the `db_path`:

```
>>> store = Store(None)
```

> **Source code in** `src/llm_matrix/store.py` ⌄

```
23    @dataclass
24    class Store:
25        """A persistent store using DuckDB to cache test results with JSON support for Pydantic models.
26
27        Example:
28
29            >>> store = Store("test-cache.db")
30            >>> case = TestCase(input="1+1", ideal="2")
31            >>> suite = Suite(name="test", cases=[case], matrix={"hyperparameters": {}})
32            >>> response = Response(text="2")
33            >>> result = TestCaseResult(case=case, response=response, hyperparameters={"model": "gpt-4"})
34            >>> store.add_result(suite, result)
35            >>> cached = store.get_result(suite, case, {"model": "gpt-4"})
36            >>> assert cached.response == response
37
38        To use an in-memory database, pass `None` as the `db_path`:
39
40            >>> store = Store(None)
41
42        """
43        db_path: Optional[str] = None
44        _conn: Optional[duckdb.DuckDBPyConnection] = None
45
46        def __post_init__(self):
47            """Initialize the database connection and create the table if it doesn't exist."""
48            self._conn = duckdb.connect(str(self.db_path) if self.db_path else ":memory:")
49            # Using JSON type for storing Pydantic models and hyperparameters
50            self._conn.execute("""
51                CREATE TABLE IF NOT EXISTS results (
52                    suite_name VARCHAR,
53                    test_case VARCHAR,
54                    ideal VARCHAR,
55                    hyperparameters JSON,
56                    result JSON,
57                    PRIMARY KEY (suite_name, test_case, ideal, hyperparameters)
58                )
59            """)
60
61        def add_result(self, suite: Suite, result: TestCaseResult):
62            """Add a result to the store."""
63            self._conn.execute("""
64                INSERT OR REPLACE INTO results
65                (suite_name, test_case, ideal, hyperparameters, result)
66                VALUES (?, ?, ?, ?, ?)
67            """, (
68                *unique_key(suite, result.case, result.hyperparameters),
69                result.model_dump_json(exclude_unset=True),
70            ))
71            logger.debug(f"Added result for {suite.name} {result.case} {result.hyperparameters}")
72            self._conn.commit()
73
74        def get_result(self, suite: Suite, case: TestCase, hyperparameters: dict) -> Optional[TestCaseResult]:
75            """Get a result from the store."""
76            result = self._conn.execute("""
77                SELECT result
78                FROM results
79                WHERE suite_name = ?
80                AND test_case = ?
81                AND ideal = ?
82                AND hyperparameters = ?
83            """, (
84                *unique_key(suite, case, hyperparameters),
85            )).fetchone()
86
87            logger.debug(f"Present: {result is not None} when looking up {suite.name} {case} {hyperparameters}")
88
89            if result:
90                return TestCaseResult.model_validate_json(result[0])
91            return None
92
93        @property
94        def size(self) -> int:
95            """Get the number of results in the store."""
96            return self._conn.execute("SELECT COUNT(*) FROM results").fetchone()[0]
97
98        def __del__(self):
99            """Close the database connection when the object is destroyed."""
100           if self._conn:
101               self._conn.close()
```

### 4.20.1 `size` `property`

Get the number of results in the store.

### 4.20.2 __post_init__()

Initialize the database connection and create the table if it doesn't exist.

> **Source code in** `src/llm_matrix/store.py` ⌄
>
> ```python
> 46   def __post_init__(self):
> 47       """Initialize the database connection and create the table if it doesn't exist."""
> 48       self._conn = duckdb.connect(str(self.db_path) if self.db_path else ":memory:")
> 49       # Using JSON type for storing Pydantic models and hyperparameters
> 50       self._conn.execute("""
> 51           CREATE TABLE IF NOT EXISTS results (
> 52               suite_name VARCHAR,
> 53               test_case VARCHAR,
> 54               ideal VARCHAR,
> 55               hyperparameters JSON,
> 56               result JSON,
> 57               PRIMARY KEY (suite_name, test_case, ideal, hyperparameters)
> 58           )
> 59       """)
> ```

### 4.20.3 add_result(suite, result)

Add a result to the store.

> **Source code in** `src/llm_matrix/store.py` ⌄
>
> ```python
> 61   def add_result(self, suite: Suite, result: TestCaseResult):
> 62       """Add a result to the store."""
> 63       self._conn.execute("""
> 64           INSERT OR REPLACE INTO results
> 65           (suite_name, test_case, ideal, hyperparameters, result)
> 66           VALUES (?, ?, ?, ?, ?)
> 67       """, (
> 68           *unique_key(suite, result.case, result.hyperparameters),
> 69           result.model_dump_json(exclude_unset=True),
> 70       ))
> 71       logger.debug(f"Added result for {suite.name} {result.case} {result.hyperparameters}")
> 72       self._conn.commit()
> ```

### 4.20.4 get_result(suite, case, hyperparameters)

Get a result from the store.

> **Source code in** `src/llm_matrix/store.py` ⌄
>
> ```python
> 74   def get_result(self, suite: Suite, case: TestCase, hyperparameters: dict) -> Optional[TestCaseResult]:
> 75       """Get a result from the store."""
> 76       result = self._conn.execute("""
> 77           SELECT result
> 78           FROM results
> 79           WHERE suite_name = ?
> 80           AND test_case = ?
> 81           AND ideal = ?
> 82           AND hyperparameters = ?
> 83       """, (
> 84           *unique_key(suite, case, hyperparameters),
> 85       )).fetchone()
> 86
> 87       logger.debug(f"Present: {result is not None} when looking up {suite.name} {case} {hyperparameters}")
> 88
> 89       if result:
> 90           return TestCaseResult.model_validate_json(result[0])
> 91       return None
> ```

### 4.20.5 __del__()

Close the database connection when the object is destroyed.

**Source code in** `src/llm_matrix/store.py` ⌄

```
98    def __del__(self):
99        """Close the database connection when the object is destroyed."""
100       if self._conn:
101           self._conn.close()
```

## 4.21 `unique_key(suite, case, hyperparameters)`

Generate a unique key for a test result.

**Source code in** `src/llm_matrix/store.py` ⌄

```
11    def unique_key(suite: Suite, case: TestCase, hyperparameters: dict) -> tuple:
12        """Generate a unique key for a test result."""
13        suite_name = suite.name
14        if suite.version:
15            suite_name += f"--{suite.version}"
16
17        def empty(v):
18            return v is None or (isinstance(v, (str, list, dict)) and not v)
19        return suite_name, case.input, case.ideal or "", {k: v for k, v in hyperparameters.items() if not empty(v)}
```

# 4.22 Using the API

## 4.22.1 Basic Example

```python
from llm_matrix import LLMRunner
from llm_matrix.schema import load_suite

# Load a test suite from YAML
suite = load_suite("my-suite.yaml")

# Create a runner
runner = LLMRunner(store_path="results.db")

# Run the suite
results = runner.run(suite)

# Process results
for result in results:
    print(f"Case: {result.case.input}")
    print(f"Score: {result.score}")
    print(f"Response: {result.response.text}")
```

## 4.22.2 Custom Configuration

```python
from llm_matrix import LLMRunner
from llm_matrix.runner import LLMRunnerConfig

# Create custom configuration
config = LLMRunnerConfig(
    concurrency=5,
    retries=3,
    timeout=30
)

# Initialize runner with configuration
runner = LLMRunner(
    store_path="results.db",
    config=config
)
```

## 4.22.3 Working with Results

```python
from llm_matrix.schema import results_to_dataframe

# Convert results to pandas DataFrame
df = results_to_dataframe(results)
```

```
# Calculate statistics
print(df.describe())

# Group by model
model_performance = df.groupby("model").agg({
    "score": ["mean", "std", "count"]
})
```

## 4.22.4 Creating Test Suites Programmatically

```python
from llm_matrix.schema import TestSuite, TestCase, Template

# Create templates
templates = {
    "qa_template": Template(
        system="Answer the following question accurately.",
        prompt="{input}",
        metrics=["qa_simple"]
    )
}

# Create test cases
cases = [
    TestCase(
        input="What is the capital of France?",
        ideal="Paris",
        tags=["geography"]
    ),
    TestCase(
        input="What is 2+2?",
        ideal="4",
        tags=["math"]
    )
]

# Create test suite
suite = TestSuite(
    name="programmatic-suite",
    template="qa_template",
    templates=templates,
    cases=cases,
    matrix={
        "hyperparameters": {
            "model": ["gpt-3.5-turbo", "gpt-4o"],
            "temperature": [0.0]
        }
    }
)

# Run the suite
runner = LLMRunner(store_path="results.db")
results = runner.run(suite)
```

# 5. Development

## 5.1 Contributing to LLM Matrix

We welcome contributions to LLM Matrix! This document provides guidelines and instructions for contributing.

### 5.1.1 Development Environment

**Prerequisites**

- Python 3.11 or higher
- Poetry for dependency management

**Setup**

1. Clone the repository:

```
git clone https://github.com/monarch-initiative/llm-matrix.git
cd llm-matrix
```

1. Install dependencies with Poetry:

```
poetry install
```

1. Activate the virtual environment:

```
poetry shell
```

### 5.1.2 Development Workflow

**Code Style**

We use the following tools to maintain code quality:

- **Black**: Code formatter with 120 character line length
- **Ruff**: Linter with various rules (flake8, isort, etc.)
- **MyPy**: Static type checking

To check code quality:

```
# Run type checking
make mypy

# Run linters
tox -e lint

# Auto-fix linting issues where possible
tox -e lint-fix

# Check spelling
tox -e codespell
```

**Running Tests**

```
# Run all tests, doctest, mypy, and codespell
make test

# Run only pytest tests
make pytest

# Run a specific test
poetry run pytest tests/test_file.py::TestClass::test_function
```

### 5.1.3 Making Changes

1. Create a new branch for your changes:

```
git checkout -b feature/your-feature-name
```

1. Make your changes following our code style guidelines.

2. Add tests for new functionality.

3. Ensure all tests pass:

```
make test
```

1. Update documentation as needed.

### 5.1.4 Submitting Changes

1. Push your changes to your fork:

```
git push origin feature/your-feature-name
```

1. Create a pull request on GitHub.

2. Ensure the PR description clearly describes the problem and solution.

### 5.1.5 Adding New Metrics

To add a new evaluation metric:

1. Create a new class in `src/llm_matrix/metrics.py` that inherits from `Metric`:

```python
class MyNewMetric(Metric):
    def evaluate(self, case: TestCase, response: LLMResponse) -> EvalResult:
        # Implement your evaluation logic
        score = calculate_score(case.ideal, response.text)

        return EvalResult(
            score=score,
            explanation="Reason for this score"
        )
```

1. Register your metric:

```
register_metric("my_new_metric", MyNewMetric())
```

1. Add tests for your metric in `tests/test_metrics.py`.

2. Update documentation in `docs/metrics/index.md`.

### 5.1.6 Creating Plugins

To create a new plugin:

1. Create a new file in `src/llm_matrix/plugins/` (e.g., `my_plugin.py`).

2. Implement the plugin interface:

```python
from llm_matrix.plugins.plugin import Plugin

class MyPlugin(Plugin):
    def name(self) -> str:
        return "my_plugin"
```

```
# Implement required methods
```

1. Register your plugin in `src/llm_matrix/plugins/__init__.py`.
2. Add tests for your plugin in `tests/test_plugins/`.
3. Update documentation as needed.

## 5.1.7 Documentation

We use MkDocs with the Material theme for documentation:

1. Update documentation as needed in the `docs/` directory.
2. Preview documentation locally:

```
mkdocs serve
```

1. Ensure all links and references are correct.

## 5.1.8 Release Process

1. Update version number in `pyproject.toml`.
2. Update CHANGELOG.md with changes in the release.
3. Create a tag with the new version number.
4. A GitHub Action will automatically build and publish to PyPI.

## 5.1.9 Questions?

If you have any questions, please open an issue on GitHub or reach out to the maintainers.