# Evaluation Framework for Knowledge Base

# Introduction 🔗

To evaluate the Knowledge Base (KB), we will implement a simple evaluation framework, based on [DeepEval](#) framework and Python package. The Knowledge Base is an implementation of a Retrieval Augmented Generation (RAG) system, and we will use the specific metrics for a RAG.

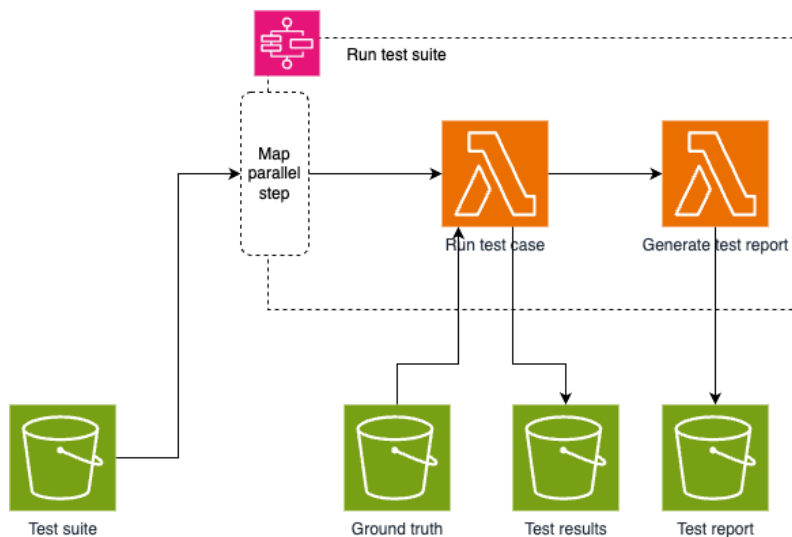Note: this document refers to the prototype for the evaluation framework, see: 📌 LGA-87: [TIMEBOX 2 days] Prototype RAG system evaluation framework `DONE`

The parameters we consider are the following:

- **input** - this is the user prompt or the initial query
- **actual_output** - what the KB returns as an answer to our query
- **expected_output** - what we expect that the KB will return as an answer to our query (ground truth)
- **retrieval_context** - the context retrieved from the KB vector database based on which the KB Agent will formulate the actual output (the answer).

These parameters are used to calculate the metrics (see below).

# Process 🔗

Note: the above approximative architecture diagram role is for understanding the process, should not be interpreted as prescriptive for the actual components to be included in the actual implementation.

## Metrics 🔗

The following table summarises the metrics, with their description, as well as the resources needed to generate each metric, and the parameters needed.

| Metric | Description | Resources needed | Parameters needed |
|---|---|---|---|
| **Rouge Metric** | **ROUGE**, which is short for **Recall-Oriented Understudy for Gisting Evaluation**, is a set of metrics used to evaluate how similar a generated text is to a reference text. It measures how much overlap there is between the model's output and a reference answer. | scorer | actual_output<br>expected_output |
| **Bert Score** | Measures a LLM. Leverages the pre-trained contextual embeddings from **BERT** and matches tokens in candidate and reference sentences by cosine similarity. It computes the precision, recall, and F1 score.<br><br>Note: We are only retaining and output F1 score in our implementation. | BERT model & tokenizer | actual_output<br>expected_output |
| **Faithfullness** | Measures the quality of **RAG retriever**/LLM by evaluating whether the actual_output factually aligns with retrieved_context - it is a self-explaining LLM evaluation | Claude model | input<br>actual_output<br>retrieval_context |
| **Answer Relevancy** | Measures the quality of **RAG retriever** by evaluating how relevant the actual_output of the application is | Claude model | input<br>actual_output |

| | | | |
|---|---|---|---|
| | compared with the provided input - it is a self-explaining LLM evaluation | | |
| **Contextual Precision** | Measures the **RAG retriever** by evaluating whether nodes in retrieval_context that are relevant to the input are ranked higher than irrelevant ones - it is a self-explaining LLM evaluation | Claude model | input<br><br>actual_output<br><br>expected_output<br><br>retrieval_context |
| **Contextual Recall** | Measures the quality of **RAG retriever** by evaluating the extent of which the retrieval_context aligns with the expected_output - it is a self-explaining LLM evaluation | Claude model | input<br><br>actual_output<br><br>expected_output<br><br>retrieval_context |
| **Contextual Relevancy** | Measures the relevance of **RAG retriever** by evaluating the overall relevance of the information presented in your retrieval_context for a given input - it is a self-explaining LLM evaluation | Claude model | input<br><br>actual_output<br><br>retrieval_context |
| **Bias** | Determines whether LLM output contains gender, racial, or political bias. This can occur after fine-tuning a custom model from any RLHF or optimization. It is a referenceless metric | Claude model | input<br><br>actual_output |
| **Toxicity** | Evaluates toxicity in LLM output. Particularly useful for a fine-tuning case.  It is a referenceless metric. | Claude model | input<br><br>actual_output |

# Run tests 🔗

## Run a test case 🔗

To run a test case, will run the metrics ensemble against the definition of a test_case.

```
1   # Test case definition
2   test_case = LLMTestCase(
3       input= input_data,
4       actual_output=actual_output_data,
5       expected_output=expected_output_data,
6       retrieval_context=retrieval_context_data
7   )
8
9   # Metrics evaluation
10  metrics = [
11      RougeMetric(),
12      BERTScoreMetric(),
13      FaithfulnessMetric(model=aws_bedrock),
14      AnswerRelevancyMetric(model=aws_bedrock),
```

```
15        ContextualPrecisionMetric(model=aws_bedrock),
16        ContextualRecallMetric(model=aws_bedrock),
17        ContextualRelevancyMetric(model=aws_bedrock),
18        BiasMetric(model=aws_bedrock),
19        ToxicityMetric(model=aws_bedrock)
20    ]
21    # Run evaluation on the test case
22    for metric in metrics:
23        metric.measure(test_case)
```

## Run a test suite 🔗

To run a test suite, we run repeatedly (in parallel) the lambda that run a test case. Each run will save the metric measure / test case on S3.

## Testing report 🔗

The testing report contains:

- the result for each test case, with this structure:

```
1    {
2        "metric": metric.__name__,
3        "score": metric.score,
4        "success": metric.success,
5        "reason": metric.reason
6    }
```

- the aggregated report for entire test suite, here is an example:

|   | name | Average Score | Pass rate |
|---|---|---|---|
| 0 | Rouge Metric | 0.461070 | 0.4 |
| 1 | BERT Score | 0.571103 | 0.0 |
| 2 | Faithfulness | 1.000000 | 1.0 |
| 3 | Answer Relevancy | 1.000000 | 1.0 |
| 4 | Contextual Precision | 0.916667 | 1.0 |
| 5 | Contextual Recall | 1.000000 | 1.0 |
| 6 | Contextual Relevancy | 0.800000 | 1.0 |
| 7 | Bias | 0.000000 | 1.0 |
| 8 | Toxicity | 0.000000 | 1.0 |

## Discussion 🔗

The impact of poorly formulated expected_output will be primarily on:

- Rouge Metric, BERT Score - very small score, 0 or very small pass rate.
- Contextual recall, Contextual precision - reduced score

# Executing Tests 🔗

## Create a Test Suite 🔗

1. Use the following template to create your test suite, placing your list of questions and expected answers in the relevant columns. 📗 test-suite-template.xlsx
2. Upload your file to the TestRunnerBucket in the test-suites folder.

## Execute a Test Suite 🔗

The test runner is a step function which can run a test suite against either another step function, a lambda or a specific bedrock agent alias.

1. Create a new execution of the TestRunnerStateMachine
2. provide an event.json which follows the format

```
1  {
2    "key": "test-suites/YOUR_TEST_SUITE_FILE.XLSX",
3    "target": {
4    // either
5      "agentId":"AGENT_ID_HERE",
6      "agentAliasId":"AGENT_ALIAS_ID",
7
8    // or
9      "stepFunctionArn":"STEP_FUNCTION_ARN_HERE",
10
11   // or
12     "lambdaArn": "LAMBDA_ARN_HERE"
13   }
14 }
15
16 /*
17  if you update the target object to include a "metric": "NAME OF METRIC" it will evaluate that metric
   specifically
18
19  here is the list of names to use (exact match)
20  Rouge Metric
21  BERT Score
22  Bias
23  Toxicity
24  Faithfulness
25  Answer Relevancy
26  Contextual Precision
27  Contextual Recall
28  Contextual Relevancy
29 */
30
31 // examples
32
33 {
34   "key": "test-suites/dp-test-suite.xlsx",
35   "target": {
36     "agentId": "E8NYLN6J5L",
37     "agentAliasId": "IL55MJT2ZB",
38   }
39 }
```

```
40
41  {
42    "key": "test-suites/dp-test-suite.xlsx",
43    "target": {
44      "stepFunctionArn": "arn:aws:states:eu-west-
      2:420498525515:stateMachine:CombinedRagFusionChatRagFusionChatStateMachine7FADB45B-TcTnLbpEm5Ns"
45    }
46  }
47
48  {
49    "key": "test-suites/dp-test-suite.xlsx",
50    "target": {
51      "lambdaArn": "arn:aws:lambda:eu-west-2:420498525515:function:lhasa-genai-development-s-
      RagFusionChatBaseAgent9B-d9bh3SWThLdv"
52    }
53  }
54
55  {
56    "key": "test-suites/dp-test-suite-2.xlsx",
57    "target": {
58      "stepFunctionArn": "arn:aws:states:eu-west-
      2:420498525515:stateMachine:CombinedRagFusionChatRagFusionChatStateMachine7FADB45B-TcTnLbpEm5Ns",
59      "metric": "BERT Score"
60    }
61  }
62
```

    a. when using an agent the question will be passed to the agent.

    b. when using a lambda or step function arn the target will be provided with a json object with a `question` property `{
"question" : "your test suite question here"}`

3. The runner will execute the test suite, generate a report and provide the result to the step function in json. Additionally it creates two csv reports the first `details.csv` containing the details and evaluation metrics of every question and a second `result.csv` containing the aggregated result of the full test suite execution. these are stored in the [TestRunnerBucket](#) using the following using the target and timestamp of the execution in the path format:

    a. **Agents** - /test-reports/**agentId-agentAliasId/yymmdd-hhmmss**-[details/result].csv

    b. **Step Function**- /test-reports/**stepFunctionName/yymmdd-hhmmss**-[details/result].csv

    c. **Lambda** - /test-reports/**lambdaName/yymmdd-hhmmss**-[details/result].csv

If the target provides a response the following metrics will be evaluated:

- Rouge
- BERTScore
- Bias
- Toxicity

If the target responds with the citations used, it will also evaluate:

- Faithfulness
- Answer Relevancy
- Contextual Precision
- Contextual Recall
- Contextual Relevancy

There is an opportunity to include the following metrics once multi-turn history is implemented

- Conversation Relevancy

- Conversation Completeness

when targeting a lambdaArn you may encounter **An error occurred (AccessDeniedException) when calling the Invoke operation**. To resolve this you must add a resource based policy statement to the target lambda function allowing invocation by the test runner state machine. here is an example of the policy



Lastly, the test runner internally parses the responses of lambdas and stepfunctions differently and therefore we should standardise response jsons from these services so that the test runner can be simplified.

## Update 1: 🔗

Now you can now specify a target metric so that only that metric is evaluated. Use one of the following metric names:

- Rouge Metric
- BERT Score
- Bias
- Toxicity
- Faithfulness
- Answer Relevancy
- Contextual Precision
- Contextual Recall
- Contextual Relevancy

like so:

```
{
  "key": "test-suites/dp-test-suite-2.xlsx",
  "target": {
    "stepFunctionArn": "arn:aws:states:eu-west-2:420498525515:stateMachine:CombinedRagFusionChatRagFusionChatStateMachine7FADB45B-TcTnLbpEm5Ns",
    "metric": "BERT Score"
  }
}
```

## Update 2: 🔗

Now you can specify the number of iterations to have the test runner run the test suite. The test runner will run for each iteration, the full response generation and evaluation flow for each question in the test suite

Specify the number of iterations like so:

```
1  {
2    "key": "test-suites/dp-test-suite-2.xlsx",
3    "iterations": 2,
4    "target": {
5      "stepFunctionArn": "arn:aws:states:eu-west-
   2:420498525515:stateMachine:CombinedRagFusionChatRagFusionChatStateMachine7FADB45B-TcTnLbpEm5Ns",
6      "metric": "BERT Score"
7    }
8  }
```

**IMPORTANT:**

the execution input "iterations" property ...

1. must contain a positive integer value

2. is always required and does not have a default value

3. the execution will fail if missing from the input