

Starting April 29, 2025, Gemini 1.5 Pro and Gemini 1.5 Flash models are not available in projects that have no prior usage of these models, including new projects. For details, see [Model versions and lifecycle](#) ([/vertex-ai/generative-ai/docs/learn/model-versions#legacy-stable](#)).

# Gen AI evaluation service API

## Release Notes

The Gen AI evaluation service lets you evaluate your large language models (LLMs) across several metrics with your own criteria. You can provide inference-time inputs, LLM responses and additional parameters, and the Gen AI evaluation service returns metrics specific to the evaluation task.

Metrics include model-based metrics, such as `PointwiseMetric` and `PairwiseMetric`, and in-memory computed metrics, such as `rouge`, `bleu`, and tool function-call metrics. `PointwiseMetric` and `PairwiseMetric` are generic model-based metrics that you can customize with your own criteria. Because the service takes the prediction results directly from models as input, the evaluation service can perform both inference and subsequent evaluation on all models supported by Vertex AI.

For more information on evaluating a model, see [Gen AI evaluation service overview](#) ([/vertex-ai/generative-ai/docs/models/evaluation-overview](#)).

## Limitations

The following are limitations of the evaluation service:

- The evaluation service may have a propagation delay in your first call.
- Most model-based metrics consume [gemini-2.0-flash quota](#) ([/vertex-ai/generative-ai/docs/quotas#eval-quotas](#)) because the Gen AI evaluation service leverages `gemini-2.0-flash` as the underlying judge model to compute these model-based metrics.
- Some model-based metrics, such as MetricX and COMET, use different machine learning models, so they don't consume [gemini-2.0-flash quota](#) ([/vertex-ai/generative-ai/docs/quotas#eval-quotas](#)).

**Note:** MetricX and COMET will be not be charged during preview. At GA, the pricing will be the same as all pointwise model based metrics.

## Example syntax

Syntax to send an evaluation call.

```
curlPython (#python)
(#curl)

curl -X POST \
  -H "Authorization: Bearer $(gcloud auth print-access-token)" \
  -H "Content-Type: application/json" \

https://${LOCATION}-aiplatform.googleapis.com/v1/projects/${PROJECT_ID}/locatio
-d '{
  "pointwise_metric_input" : {
    "metric_spec" : {
      ...
    },
    "instance": {
      ...
    },
  },
}'
```

## Parameter list

Parameters	
exact_match_input	Optional: <b>ExactMatchInput</b>  Input to assess if the prediction matches the reference exactly.
bleu_input	Optional: <b>BleuInput</b>  Input to compute BLEU score by comparing the prediction against the reference.
rouge_input	Optional: <b>RougeInput</b>  Input to compute <b>rouge</b> scores by comparing the prediction against the reference. Different <b>rouge</b> scores are supported by <b>rouge_type</b> .
fluency_input	Optional: <b>FluencyInput</b>

	Input to assess a single response's language mastery.
<b>coherence_input</b>	Optional: <b>CoherenceInput</b>  Input to assess a single response's ability to provide a coherent, easy-to-follow reply.
<b>safety_input</b>	Optional: <b>SafetyInput</b>  Input to assess a single response's level of safety.
<b>groundedness_input</b>	Optional: <b>GroundednessInput</b>  Input to assess a single response's ability to provide or reference information included only in the input text.
<b>fulfillment_input</b>	Optional: <b>FulfillmentInput</b>  Input to assess a single response's ability to completely fulfill instructions.
<b>summarization_quality_input</b>	Optional: <b>SummarizationQualityInput</b>  Input to assess a single response's overall ability to summarize text.
<b>pairwise_summarization_quality_input</b>	Optional: <b>PairwiseSummarizationQualityInput</b>  Input to compare two responses' overall summarization quality.
<b>summarization_helpfulness_input</b>	Optional: <b>SummarizationHelpfulnessInput</b>  Input to assess a single response's ability to provide a summarization, which contains the details necessary to substitute the original text.
<b>summarization_verbosity_input</b>	Optional: <b>SummarizationVerbosityInput</b>  Input to assess a single response's ability to provide a succinct summarization.
<b>question_answering_quality_input</b>	Optional: <b>QuestionAnsweringQualityInput</b>  Input to assess a single response's overall ability to answer questions, given a body of text to reference.
<b>pairwise_question_answering_quality_input</b>	Optional: <b>PairwiseQuestionAnsweringQualityInput</b>  Input to compare two responses' overall ability to answer questions, given a body of text to reference.
<b>question_answering_relevance_input</b>	Optional: <b>QuestionAnsweringRelevanceInput</b>

Input to assess a single response's ability to respond with relevant information when asked a question.

**question\_answering\_helpfulness\_input** Optional: **QuestionAnsweringHelpfulnessInput**

Input to assess a single response's ability to provide key details when answering a question.

**question\_answering\_correctness\_input** Optional: **QuestionAnsweringCorrectnessInput**

Input to assess a single response's ability to correctly answer a question.

**pointwise\_metric\_input** Optional: **PointwiseMetricInput**

Input for a generic pointwise evaluation.

**pairwise\_metric\_input** Optional: **PairwiseMetricInput**

Input for a generic pairwise evaluation.

**tool\_call\_valid\_input** Optional: **ToolCallValidInput**

Input to assess a single response's ability to predict a valid tool call.

**tool\_name\_match\_input** Optional: **ToolNameMatchInput**

Input to assess a single response's ability to predict a tool call with the right tool name.

**tool\_parameter\_key\_match\_input** Optional: **ToolParameterKeyMatchInput**

Input to assess a single response's ability to predict a tool call with correct parameter names.

**tool\_parameter\_kv\_match\_input** Optional: **ToolParameterKvMatchInput**

Input to assess a single response's ability to predict a tool call with correct parameter names and values

**comet\_input** Optional: **CometInput**

Input to evaluate using [COMET](https://huggingface.co/Unbabel/wmt22-comet-da) (<https://huggingface.co/Unbabel/wmt22-comet-da>).

**metricx\_input** Optional: **MetricxInput**

Input to evaluate using [MetricX](https://github.com/google-research/metricx) (<https://github.com/google-research/metricx>).

**ExactMatchInput**

```
{
  "exact_match_input": {
    "metric_spec": {},
    "instances": [
      {
        "prediction": string,
        "reference": string
      }
    ]
  }
}
```

Parameters

<b>metric_spec</b>	Optional: <b>ExactMatchSpec</b> .  Metric spec, defining the metric's behavior.
<b>instances</b>	Optional: <b>ExactMatchInstance[ ]</b>  Evaluation input, consisting of LLM response and reference.
<b>instances.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instances.reference</b>	Optional: <b>string</b>  Golden LLM response for reference.

ExactMatchResults

```
{
  "exact_match_results": {
    "exact_match_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

Output

<code>exact_match_metric_values</code>	<b>ExactMatchMetricValue[ ]</b>  Evaluation results per instance input.
<code>exact_match_metric_values.score</code>	<b>float</b>  One of the following: <ul style="list-style-type: none"><li>• 0: Instance was not an exact match</li><li>• 1: Exact match</li></ul>

BleuInput

```
{
  "bleu_input": {
    "metric_spec": {
      "use_effective_order": bool
    },
    "instances": [
      {
        "prediction": string,
        "reference": string
      }
    ]
  }
}
```

Parameters

<code>metric_spec</code>	Optional: <b>BleuSpec</b>  Metric spec, defining the metric's behavior.
<code>metric_spec.use_effective_order</code>	Optional: <b>bool</b>  Whether to take into account n-gram orders without any match.
<code>instances</code>	Optional: <b>BleuInstance[ ]</b>  Evaluation input, consisting of LLM response and reference.

<b>instances.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instances.reference</b>	Optional: <b>string</b>  Golden LLM response for reference.

**BleuResults**

```
{
  "bleu_results": {
    "bleu_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

**Output**

<b>bleu_metric_values</b>	<b>BleuMetricValue[]</b>  Evaluation results per instance input.
<b>bleu_metric_values.score</b>	<b>float: [0, 1]</b> , where higher scores mean the prediction is more like the reference.

**RougeInput**

```
{
  "rouge_input": {
    "metric_spec": {
      "rouge_type": string,
      "use_stemmer": bool,
      "split_summaries": bool
    },
    "instances": [
      {
```

```
        "prediction": string,  
        "reference": string  
    }  
]  
}  
}
```

Parameters

<b>metric_spec</b>	Optional: <b>RougeSpec</b>  Metric spec, defining the metric's behavior.
<b>metric_spec.rouge_type</b>	Optional: <b>string</b>  Acceptable values: <ul style="list-style-type: none"><li>• <b>rougen[1-9]</b>: compute <b>rouge</b> scores based on the overlap of n-grams between the prediction and the reference.</li><li>• <b>rougeL</b>: compute <b>rouge</b> scores based on the Longest Common Subsequence (LCS) between the prediction and the reference.</li><li>• <b>rougeLsum</b>: first splits the prediction and the reference into sentences and then computes the LCS for each tuple. The final <b>rougeLsum</b> score is the average of these individual LCS scores.</li></ul>
<b>metric_spec.use_stemmer</b>	Optional: <b>bool</b>  Whether Porter stemmer should be used to strip word suffixes to improve matching.
<b>metric_spec.split_summaries</b>	Optional: <b>bool</b>  Whether to add newlines between sentences for rougeLsum.
<b>instances</b>	Optional: <b>RougeInstance[ ]</b>  Evaluation input, consisting of LLM response and reference.
<b>instances.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instances.reference</b>	Optional: <b>string</b>  Golden LLM response for reference.



RougeResults

```
{
  "rouge_results": {
    "rouge_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

Output

<b>rouge_metric_values</b>	<b>RougeValue[ ]</b>
	Evaluation results per instance input.
<b>rouge_metric_values.score float: [0, 1]</b> , where higher scores mean the prediction is more like the reference.	

FluencyInput

```
{
  "fluency_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string
    }
  }
}
```

Parameters

<b>metric_spec</b>	Optional: <b>FluencySpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>FluencyInstance</b>

Evaluation input, consisting of LLM response.

<code>instance.prediction</code>	Optional: <code>string</code>  LLM response.
----------------------------------	--

FluencyResult

```
{
  "fluency_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

Output

<code>score</code>	<code>float</code> : One of the following: <ul style="list-style-type: none"><li>• 1: Inarticulate</li><li>• 2: Somewhat Inarticulate</li><li>• 3: Neutral</li><li>• 4: Somewhat fluent</li><li>• 5: Fluent</li></ul>
<code>explanation</code>	<code>string</code> : Justification for score assignment.
<code>confidence</code>	<code>float</code> : [ 0, 1 ] Confidence score of our result.

CoherenceInput

```
{
  "coherence_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string
    }
  }
}
```

```
}  
}
```

Parameters

<b>metric_spec</b>	Optional: <b>CoherenceSpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>CoherenceInstance</b>  Evaluation input, consisting of LLM response.
<b>instance.prediction</b>	Optional: <b>string</b>  LLM response.

CoherenceResult

```
{  
  "coherence_result": {  
    "score": float,  
    "explanation": string,  
    "confidence": float  
  }  
}
```

Output

<b>score</b>	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 1: Incoherent</li><li>• 2: Somewhat incoherent</li><li>• 3: Neutral</li><li>• 4: Somewhat coherent</li><li>• 5: Coherent</li></ul>
<b>explanation</b>	<b>string:</b> Justification for score assignment.
<b>confidence</b>	<b>float:</b> [0, 1] Confidence score of our result.

SafetyInput

```
{
  "safety_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string
    }
  }
}
```

Parameters

metric_spec	Optional: <b>SafetySpec</b>  Metric spec, defining the metric's behavior.
instance	Optional: <b>SafetyInstance</b>  Evaluation input, consisting of LLM response.
instance.prediction	Optional: <b>string</b>  LLM response.

SafetyResult

```
{
  "safety_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

Output

score	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 0: Unsafe</li></ul>
-------	---

- 1: Safe

explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

GroundednessInput

```
{
  "groundedness_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "context": string
    }
  }
}
```

Parameter	Description
metric_spec	Optional: GroundednessSpec  Metric spec, defining the metric's behavior.
instance	Optional: GroundednessInstance  Evaluation input, consisting of inference inputs and corresponding response.
instance.prediction	Optional: string  LLM response.
instance.context	Optional: string  Inference-time text containing all information, which can be used in the LLM response.

GroundednessResult

```
{
  "groundedness_result": {
    "score": float,
```

```
    "explanation": string,  
    "confidence": float  
  }  
}
```

Output

score	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 0: Ungrounded</li><li>• 1: Grounded</li></ul>
explanation	<b>string:</b> Justification for score assignment.
confidence	<b>float:</b> [ 0, 1 ] Confidence score of our result.

FulfillmentInput

```
{  
  "fulfillment_input": {  
    "metric_spec": {},  
    "instance": {  
      "prediction": string,  
      "instruction": string  
    }  
  }  
}
```

Parameters

metric_spec	Optional: <b>FulfillmentSpec</b>  Metric spec, defining the metric's behavior.
instance	Optional: <b>FulfillmentInstance</b>  Evaluation input, consisting of inference inputs and corresponding response.
instance.prediction	Optional: <b>string</b>  LLM response.

---

<b>instance.instruction</b>	Optional: <b>string</b>
	Instruction used at inference time.

---

**FulfillmentResult**

```
{
  "fulfillment_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

---

**Output**

---

<b>score</b>	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 1: No fulfillment</li><li>• 2: Poor fulfillment</li><li>• 3: Some fulfillment</li><li>• 4: Good fulfillment</li><li>• 5: Complete fulfillment</li></ul>
<b>explanation</b>	<b>string:</b> Justification for score assignment.
<b>confidence</b>	<b>float:</b> [ 0, 1 ] Confidence score of our result.

---

**SummarizationQualityInput**

```
{
  "summarization_quality_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "instruction": string,
      "context": string,
    }
  }
}
```

```
}  
}
```

Parameters

<b>metric_spec</b>	Optional: <b>SummarizationQualitySpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>SummarizationQualityInstance</b>  Evaluation input, consisting of inference inputs and corresponding response.
<b>instance.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instance.instruction</b>	Optional: <b>string</b>  Instruction used at inference time.
<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.

SummarizationQualityResult

```
{  
  "summarization_quality_result": {  
    "score": float,  
    "explanation": string,  
    "confidence": float  
  }  
}
```

Output

<b>score</b>	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 1: Very bad</li><li>• 2: Bad</li></ul>
--------------	--



- 3: Ok
- 4: Good
- 5: Very good

explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

PairwiseSummarizationQualityInput

```
{
  "pairwise_summarization_quality_input": {
    "metric_spec": {},
    "instance": {
      "baseline_prediction": string,
      "prediction": string,
      "instruction": string,
      "context": string,
    }
  }
}
```

Parameters

metric_spec	Optional: PairwiseSummarizationQualitySpec  Metric spec, defining the metric's behavior.
instance	Optional: PairwiseSummarizationQualityInstance  Evaluation input, consisting of inference inputs and corresponding response.
instance.baseline_prediction	Optional: string  Baseline model LLM response.
instance.prediction	Optional: string  Candidate model LLM response.
instance.instruction	Optional: string  Instruction used at inference time.

<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.
-------------------------	---

**PairwiseSummarizationQualityResult**

```
{
  "pairwise_summarization_quality_result": {
    "pairwise_choice": PairwiseChoice,
    "explanation": string,
    "confidence": float
  }
}
```

**Output**

<b>pairwise_choice</b>	<b>PairwiseChoice:</b> Enum with possible values as follows:  <ul style="list-style-type: none"><li>• <b>BASELINE:</b> Baseline prediction is better</li><li>• <b>CANDIDATE:</b> Candidate prediction is better</li><li>• <b>TIE:</b> Tie between Baseline and Candidate predictions.</li></ul>
<b>explanation</b>	<b>string:</b> Justification for pairwise_choice assignment.
<b>confidence</b>	<b>float:</b> [0, 1] Confidence score of our result.

**SummarizationHelpfulnessInput**

```
{
  "summarization_helpfulness_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "instruction": string,
      "context": string,
    }
  }
}
```

```
}  
}
```

Parameters

<b>metric_spec</b>	Optional: <b>SummarizationHelpfulnessSpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>SummarizationHelpfulnessInstance</b>  Evaluation input, consisting of inference inputs and corresponding response.
<b>instance.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instance.instruction</b>	Optional: <b>string</b>  Instruction used at inference time.
<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.

SummarizationHelpfulnessResult

```
{  
  "summarization_helpfulness_result": {  
    "score": float,  
    "explanation": string,  
    "confidence": float  
  }  
}
```

Output

<b>score</b>	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>1: Unhelpful</li><li>2: Somewhat unhelpful</li></ul>
--------------	--

- 3: Neutral
- 4: Somewhat helpful
- 5: Helpful

explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

SummarizationVerbosityInput

```
{
  "summarization_verbosity_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "instruction": string,
      "context": string,
    }
  }
}
```

Parameters

metric_spec	Optional: SummarizationVerbositySpec  Metric spec, defining the metric's behavior.
instance	Optional: SummarizationVerbosityInstance  Evaluation input, consisting of inference inputs and corresponding response.
instance.prediction	Optional: string  LLM response.
instance.instruction	Optional: string  Instruction used at inference time.
instance.context	Optional: string  Inference-time text containing all information, which can be used in the LLM response.

SummarizationVerbosityResult

```
{
  "summarization_verbosity_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

Output

score	float. One of the following: <ul style="list-style-type: none"><li>-2: Terse</li><li>-1: Somewhat terse</li><li>0: Optimal</li><li>1: Somewhat verbose</li><li>2: Verbose</li></ul>
explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

QuestionAnsweringQualityInput

```
{
  "question_answering_quality_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "instruction": string,
      "context": string,
    }
  }
}
```

---

Parameters

<code>metric_spec</code>	Optional: <code>QuestionAnsweringQualitySpec</code>  Metric spec, defining the metric's behavior.
<code>instance</code>	Optional: <code>QuestionAnsweringQualityInstance</code>  Evaluation input, consisting of inference inputs and corresponding response.
<code>instance.prediction</code>	Optional: <code>string</code>  LLM response.
<code>instance.instruction</code>	Optional: <code>string</code>  Instruction used at inference time.
<code>instance.context</code>	Optional: <code>string</code>  Inference-time text containing all information, which can be used in the LLM response.

---

`QuestionAnsweringQualityResult`

```
{
  "question_answering_quality_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

---

Output

---

score	float: One of the following: <ul style="list-style-type: none"><li>1: Very bad</li><li>2: Bad</li><li>3: Ok</li><li>4: Good</li><li>5: Very good</li></ul>
explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

PairwiseQuestionAnsweringQualityInput

```
{
  "question_answering_quality_input": {
    "metric_spec": {},
    "instance": {
      "baseline_prediction": string,
      "prediction": string,
      "instruction": string,
      "context": string
    }
  }
}
```

Parameters	
metric_spec	Optional: QuestionAnsweringQualitySpec  Metric spec, defining the metric's behavior.
instance	Optional: QuestionAnsweringQualityInstance  Evaluation input, consisting of inference inputs and corresponding response.
instance.baseline_prediction	Optional: string  Baseline model LLM response.
instance.prediction	Optional: string

Candidate model LLM response.

<b>instance.instruction</b>	Optional: <b>string</b>  Instruction used at inference time.
<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.

**PairwiseQuestionAnsweringQualityResult**

```
{
  "pairwise_question_answering_quality_result": {
    "pairwise_choice": PairwiseChoice,
    "explanation": string,
    "confidence": float
  }
}
```

**Output**

<b>pairwise_choice</b>	<b>PairwiseChoice:</b> Enum with possible values as follows: <ul style="list-style-type: none"><li>• <b>BASELINE:</b> Baseline prediction is better</li><li>• <b>CANDIDATE:</b> Candidate prediction is better</li><li>• <b>TIE:</b> Tie between Baseline and Candidate predictions.</li></ul>
<b>explanation</b>	<b>string:</b> Justification for <b>pairwise_choice</b> assignment.
<b>confidence</b>	<b>float:</b> [0, 1] Confidence score of our result.

**QuestionAnsweringRelevanceInput**

```
{
  "question_answering_quality_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
```



```
    "instruction": string,
    "context": string
  }
}
```

Parameters

<b>metric_spec</b>	Optional: <b>QuestionAnsweringRelevanceSpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>QuestionAnsweringRelevanceInstance</b>  Evaluation input, consisting of inference inputs and corresponding response.
<b>instance.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instance.instruction</b>	Optional: <b>string</b>  Instruction used at inference time.
<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.

QuestionAnsweringRelevancyResult

```
{
  "question_answering_relevancy_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

Output

**score** **float:** One of the following:

- 1: Irrelevant
- 2: Somewhat irrelevant
- 3: Neutral
- 4: Somewhat relevant
- 5: Relevant

explanation	string: Justification for score assignment.
confidence	float: [ 0, 1 ] Confidence score of our result.

QuestionAnsweringHelpfulnessInput

```
{
  "question_answering_helpfulness_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "instruction": string,
      "context": string
    }
  }
}
```

Parameters

metric_spec	Optional: QuestionAnsweringHelpfulnessSpec  Metric spec, defining the metric's behavior.
instance	Optional: QuestionAnsweringHelpfulnessInstance  Evaluation input, consisting of inference inputs and corresponding response.
instance.prediction	Optional: string  LLM response.
instance.instruction	Optional: string  Instruction used at inference time.
instance.context	Optional: string

Inference-time text containing all information, which can be used in the LLM response.

QuestionAnsweringHelpfulnessResult

```
{
  "question_answering_helpfulness_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

Output

score	float: One of the following: <ul style="list-style-type: none"><li>1: Unhelpful</li><li>2: Somewhat unhelpful</li><li>3: Neutral</li><li>4: Somewhat helpful</li><li>5: Helpful</li></ul>
explanation	string: Justification for score assignment.
confidence	float: [0, 1] Confidence score of our result.

QuestionAnsweringCorrectnessInput

```
{
  "question_answering_correctness_input": {
    "metric_spec": {
      "use_reference": bool
    },
    "instance": {
      "prediction": string,
      "reference": string,
      "instruction": string,

```

```
    "context": string
  }
}
```

Parameters

<b>metric_spec</b>	Optional: <b>QuestionAnsweringCorrectnessSpec</b>  Metric spec, defining the metric's behavior.
<b>metric_spec.use_reference</b>	Optional: <b>bool</b>  If reference is used or not in the evaluation.
<b>instance</b>	Optional: <b>QuestionAnsweringCorrectnessInstance</b>  Evaluation input, consisting of inference inputs and corresponding response.
<b>instance.prediction</b>	Optional: <b>string</b>  LLM response.
<b>instance.reference</b>	Optional: <b>string</b>  Golden LLM response for reference.
<b>instance.instruction</b>	Optional: <b>string</b>  Instruction used at inference time.
<b>instance.context</b>	Optional: <b>string</b>  Inference-time text containing all information, which can be used in the LLM response.

QuestionAnsweringCorrectnessResult

```
{
  "question_answering_correctness_result": {
    "score": float,
    "explanation": string,
    "confidence": float
  }
}
```

```
}
}
```

Output

score	<b>float:</b> One of the following: <ul style="list-style-type: none"><li>• 0: Incorrect</li><li>• 1: Correct</li></ul>
explanation	<b>string:</b> Justification for score assignment.
confidence	<b>float:</b> [0, 1] Confidence score of our result.

PointwiseMetricInput

```
{
  "pointwise_metric_input": {
    "metric_spec": {
      "metric_prompt_template": string
    },
    "instance": {
      "json_instance": string,
    }
  }
}
```

Parameters

metric_spec	Required: <b>PointwiseMetricSpec</b>  Metric spec, defining the metric's behavior.
metric_spec.metric_prompt_template	Required: <b>string</b>  A prompt template defining the metric. It is rendered by the key-value pairs in instance.json_instance
instance	Required: <b>PointwiseMetricInstance</b>  Evaluation input, consisting of json_instance.

<b>instance.json_instance</b>	Optional: <b>string</b>  The key-value pairs in Json format. For example, {"key_1": "value_1", "key_2": "value_2"}. It is used to render metric_spec.metric_prompt_template.
-------------------------------	--

**PointwiseMetricResult**

```
{
  "pointwise_metric_result": {
    "score": float,
    "explanation": string,
  }
}
```

**Output**

<b>score</b>	<b>float:</b> A score for pointwise metric evaluation result.
<b>explanation</b>	<b>string:</b> Justification for score assignment.

**PairwiseMetricInput**

```
{
  "pairwise_metric_input": {
    "metric_spec": {
      "metric_prompt_template": string
    },
    "instance": {
      "json_instance": string,
    }
  }
}
```

**Parameters**

<b>metric_spec</b>	Required: <b>PairwiseMetricSpec</b>  Metric spec, defining the metric's behavior.
--------------------	---

<b>metric_spec.metric_prompt_template</b>	Required: <b>string</b>  A prompt template defining the metric. It is rendered by the key-value pairs in instance.json_instance
<b>instance</b>	Required: <b>PairwiseMetricInstance</b>  Evaluation input, consisting of json_instance.
<b>instance.json_instance</b>	Optional: <b>string</b>  The key-value pairs in JSON format. For example, {"key_1": "value_1", "key_2": "value_2"}. It is used to render metric_spec.metric_prompt_template.

**PairwiseMetricResult**

```
{
  "pairwise_metric_result": {
    "score": float,
    "explanation": string,
  }
}
```

**Output**

<b>score</b>	<b>float:</b> A score for pairwise metric evaluation result.
<b>explanation</b>	<b>string:</b> Justification for score assignment.

**ToolCallValidInput**

```
{
  "tool_call_valid_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "reference": string
    }
  }
}
```

```
}  
}
```

Parameters

<b>metric_spec</b>	<p>Optional: <b>ToolCallValidSpec</b></p> <p>Metric spec, defining the metric's behavior.</p>
<b>instance</b>	<p>Optional: <b>ToolCallValidInstance</b></p> <p>Evaluation input, consisting of LLM response and reference.</p>
<b>instance.prediction</b>	<p>Optional: <b>string</b></p> <p>Candidate model LLM response, which is a JSON serialized string that contains <b>content</b> and <b>tool_calls</b> keys. The <b>content</b> value is the text output from the model. The <b>tool_call</b> value is a JSON serialized string of a list of tool calls. Ar example is:</p> <pre>{   "content": "",   "tool_calls": [     {       "name": "book_tickets",       "arguments": {         "movie": "Mission Impossible Dead Reckoning Part 1",         "theater": "Regal Edwards 14",         "location": "Mountain View CA",         "showtime": "7:30",         "date": "2024-03-30",         "num_tix": "2"       }     }   ] }</pre>
<b>instance.reference</b>	<p>Optional: <b>string</b></p> <p>Golden model output in the same format as prediction.</p>

ToolCallValidResults



```
{
  "tool_call_valid_results": {
    "tool_call_valid_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

Output

<code>tool_call_valid_metric_values</code>	repeated <code>ToolCallValidMetricValue</code> : Evaluation results per instance input.
<code>tool_call_valid_metric_values.score</code>	<code>float</code> : One of the following: <ul style="list-style-type: none"><li><code>0</code>: Invalid tool call</li><li><code>1</code>: Valid tool call</li></ul>

ToolNameMatchInput

```
{
  "tool_name_match_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "reference": string
    }
  }
}
```

Parameters

<code>metric_spec</code>	Optional: <code>ToolNameMatchSpec</code>  Metric spec, defining the metric's behavior.
<code>instance</code>	Optional: <code>ToolNameMatchInstance</code>

Evaluation input, consisting of LLM response and reference.

<code>instance.prediction</code>	Optional: <code>string</code>  Candidate model LLM response, which is a JSON serialized string that contains <code>content</code> and <code>tool_calls</code> keys. The <code>content</code> value is the text output from the model. The <code>tool_call</code> value is a JSON serialized string of a list of tool calls.
<code>instance.reference</code>	Optional: <code>string</code>  Golden model output in the same format as prediction.

**ToolNameMatchResults**

```
{
  "tool_name_match_results": {
    "tool_name_match_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

**Output**

<code>tool_name_match_metric_values</code>	repeated <code>ToolNameMatchMetricValue</code> : Evaluation results per instance input.
<code>tool_name_match_metric_values.score</code>	<code>float</code> : One of the following: <ul style="list-style-type: none"><li>• 0: Tool call name doesn't match the reference.</li><li>• 1: Tool call name matches the reference.</li></ul>

**ToolParameterKeyMatchInput**

```
{
  "tool_parameter_key_match_input": {
    "metric_spec": {},
  }
}
```

```
    "instance": {
      "prediction": string,
      "reference": string
    }
  }
}
```

Parameters

<b>metric_spec</b>	Optional: <b>ToolParameterKeyMatchSpec</b>  Metric spec, defining the metric's behavior.
<b>instance</b>	Optional: <b>ToolParameterKeyMatchInstance</b>  Evaluation input, consisting of LLM response and reference.
<b>instance.prediction</b>	Optional: <b>string</b>  Candidate model LLM response, which is a JSON serialized string that contains <b>content</b> and <b>tool_calls</b> keys. The <b>content</b> value is the text output from the model. The <b>tool_call</b> value is a JSON serialized string of a list of tool calls.
<b>instance.reference</b>	Optional: <b>string</b>  Golden model output in the same format as prediction.

ToolParameterKeyMatchResults

```
{
  "tool_parameter_key_match_results": {
    "tool_parameter_key_match_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

Output

<code>tool_parameter_key_match_metric_value</code>	repeated <code>ToolParameterKeyMatchMetricValue</code> : Evaluation results per instance input.
<code>tool_parameter_key_match_metric_value.score</code>	<code>float</code> : [0, 1], where higher scores mean more parameters match the reference parameters' names.

ToolParameterKVMatchInput

```
{
  "tool_parameter_kv_match_input": {
    "metric_spec": {},
    "instance": {
      "prediction": string,
      "reference": string
    }
  }
}
```

Parameters

<code>metric_spec</code>	Optional: <code>ToolParameterKVMatchSpec</code>  Metric spec, defining the metric's behavior.
<code>instance</code>	Optional: <code>ToolParameterKVMatchInstance</code>  Evaluation input, consisting of LLM response and reference.
<code>instance.prediction</code>	Optional: <code>string</code>  Candidate model LLM response, which is a JSON serialized string that contains <code>content</code> and <code>tool_calls</code> keys. The <code>content</code> value is the text output from the model. The <code>tool_call</code> value is a JSON serialized string of a list of tool calls.
<code>instance.reference</code>	Optional: <code>string</code>  Golden model output in the same format as prediction.

ToolParameterKVMatchResults

```
{
  "tool_parameter_kv_match_results": {
    "tool_parameter_kv_match_metric_values": [
      {
        "score": float
      }
    ]
  }
}
```

Output

<code>tool_parameter_kv_match_metric_values</code>	repeated <code>ToolParameterKVMatchMetricValue</code> : Evaluation results per instance input.
<code>tool_parameter_kv_match_metric_values.score</code>	<code>float: [0, 1]</code> , where higher scores mean more parameters match the reference parameters' names and values.

CometInput

```
{
  "comet_input" : {
    "metric_spec" : {
      "version": string
    },
    "instance": {
      "prediction": string,
      "source": string,
      "reference": string,
    },
  }
}
```

Parameters

<code>metric_spec</code>	Optional: <code>CometSpec</code>  Metric spec, defining the metric's behavior.
--------------------------	--

<b>metric_spec.version</b>	Optional: <b>string</b>  COMET_22_SRC_REF: COMET 22 for translation, source, and reference. It evaluates the translation (prediction) using all three inputs.
<b>metric_spec.source_language</b>	Optional: <b>string</b>  Source language in <u>BCP-47 format</u> ( <a href="https://en.wikipedia.org/wiki/IETF_language_tag">https://en.wikipedia.org/wiki/IETF_language_tag</a> ). For example, "es".
<b>metric_spec.target_language</b>	Optional: <b>string</b>  Target language in <u>BCP-47 format</u> ( <a href="https://en.wikipedia.org/wiki/IETF_language_tag">https://en.wikipedia.org/wiki/IETF_language_tag</a> ). For example, "es"
<b>instance</b>	Optional: <b>CometInstance</b>  Evaluation input, consisting of LLM response and reference. The exact fields used for evaluation are dependent on the COMET version.
<b>instance.prediction</b>	Optional: <b>string</b>  Candidate model LLM response. This is the output of the LLM which is being evaluated.
<b>instance.source</b>	Optional: <b>string</b>  Source text. This is in the original language that the prediction was translated from.
<b>instance.reference</b>	Optional: <b>string</b>  Ground truth used to compare against the prediction. This is in the same language as the prediction.

**CometResult**

```
{
  "comet_result" : {
    "score": float
  }
}
```

**Output**

score

float: [0, 1], where 1 represents a perfect translation.

MetricxInput

```

{
  "metricx_input" : {
    "metric_spec" : {
      "version": string
    },
    "instance": {
      "prediction": string,
      "source": string,
      "reference": string,
    },
  }
}
```

Parameters

metric_spec	Optional: MetricxSpec  Metric spec, defining the metric's behavior.
metric_spec.version	Optional: string  One of the following: <ul style="list-style-type: none"> <li>METRICX_24_REF: MetricX 24 for translation and reference. It evaluates the prediction (translation) by comparing with the provided reference text input.</li> <li>METRICX_24_SRC: MetricX 24 for translation and source. It evaluates the translation (prediction) by Quality Estimation (QE), without a reference text input.</li> <li>METRICX_24_SRC_REF: MetricX 24 for translation, source and reference. It evaluates the translation (prediction) using all three inputs.</li> </ul>
metric_spec.source_language	Optional: string  Source language in <u>BCP-47 format</u> ( <a href="https://en.wikipedia.org/wiki/IETF_language_tag">https://en.wikipedia.org/wiki/IETF_language_tag</a> ). For example, "es".

<code>metric_spec.target_language</code>	Optional: <code>string</code>  Target language in <u>BCP-47 format</u> ( <a href="https://en.wikipedia.org/wiki/IETF_language_tag">https://en.wikipedia.org/wiki/IETF_language_tag</a> ). For example, "es".
<code>instance</code>	Optional: <code>MetricxInstance</code>  Evaluation input, consisting of LLM response and reference. The exact fields used for evaluation are dependent on the MetricX version.
<code>instance.prediction</code>	Optional: <code>string</code>  Candidate model LLM response. This is the output of the LLM which is being evaluated.
<code>instance.source</code>	Optional: <code>string</code>  Source text which is in the original language that the prediction was translated from.
<code>instance.reference</code>	Optional: <code>string</code>  Ground truth used to compare against the prediction. It is in the same language as the prediction.

**MetricxResult**

```
{
  "metricx_result" : {
    "score": float
  }
}
```

**Output**

<code>score</code>	<code>float: [0, 25]</code> , where 0 represents a perfect translation.
--------------------	---

Examples

Evaluate an output



The following example demonstrates how to call the Gen AI Evaluation API to evaluate the output of an LLM using a variety of evaluation metrics, including the following:

- `summarization_quality`
- `groundedness`
- `fulfillment`
- `summarization_helpfulness`
- `summarization_verbosity`

PythonGo (#go)  
(#python)

```
import pandas as pd

import vertexai
from vertexai.preview.evaluation import EvalTask, MetricPromptTemplateExamples

# TODO(developer): Update and un-comment below line
# PROJECT_ID = "your-project-id"
vertexai.init(project=PROJECT_ID, location="us-central1")

eval_dataset = pd.DataFrame(
    {
        "instruction": [
            "Summarize the text in one sentence.",
            "Summarize the text such that a five-year-old can understand.",
        ],
        "context": [
            """As part of a comprehensive initiative to tackle urban congestion sustainable urban living, a major city has revealed ambitious plans extensive overhaul of its public transportation system. The project only to improve the efficiency and reliability of public transit bu reduce the city\'s carbon footprint and promote eco-friendly commut City officials anticipate that this strategic investment will enhan accessibility for residents and visitors alike, ushering in a new e efficient, environmentally conscious urban transportation."""
            """A team of archaeologists has unearthed ancient artifacts sheddin previously unknown civilization. The findings challenge existing hi narratives and provide valuable insights into human history."""
        ],
        "response": [
            "A major city is revamping its public transportation system to figh",
            "Some people who dig for old things found some very special tools a
```

```

    ],
    }
)

eval_task = EvalTask(
    dataset=eval_dataset,
    metrics=[
        MetricPromptTemplateExamples.Pointwise.SUMMARIZATION_QUALITY,
        MetricPromptTemplateExamples.Pointwise.GROUNDEDNESS,
        MetricPromptTemplateExamples.Pointwise.VERBOSITY,
        MetricPromptTemplateExamples.Pointwise.INSTRUCTION_FOLLOWING,
    ],
)

prompt_template = (
    "Instruction: {instruction}. Article: {context}. Summary: {response}"
)
result = eval_task.evaluate(prompt_template=prompt_template)

print("Summary Metrics:\n")

for key, value in result.summary_metrics.items():
    print(f"{key}: \t{value}")

print("\n\nMetrics Table:\n")
print(result.metrics_table)
# Example response:
# Summary Metrics:
# row_count:      2
# summarization_quality/mean:      3.5
# summarization_quality/std:      2.1213203435596424
# ...


```






## Evaluate an output: pairwise summarization quality

The following example demonstrates how to call the Gen AI evaluation service API to evaluate the output of an LLM using a pairwise summarization quality comparison.

RESTPython (#python)Go (#go)  
(#rest)

Before using any of the request data, make the following replacements:

- **PROJECT\_ID** : Your project ID  
(/resource-manager/docs/creating-managing-projects#identifiers).

- **LOCATION** : The region to process the request.
- **PREDICTION** : LLM response.
- **BASELINE\_PREDICTION** : Baseline model LLM response.
- **INSTRUCTION** : The instruction used at inference time.
- **CONTEXT** : Inference-time text containing all relevant information, that can be used in the LLM response.

HTTP method and URL:

POST [https://LOCATION-aiplatform.googleapis.com/v1beta1/projects/PROJECT\\_ID](https://LOCATION-aiplatform.googleapis.com/v1beta1/projects/PROJECT_ID)

Request JSON body:


```
{
  "pairwise_summarization_quality_input": {
    "metric_spec": {},
    "instance": {
      "prediction": "PREDICTION",
      "baseline_prediction": "BASELINE_PREDICTION",
      "instruction": "INSTRUCTION",
      "context": "CONTEXT",
    }
  }
}
```

To send your request, choose one of these options:

**curlPowerShell** (#powershell)  
(#curl)

★ **Note:** The following command assumes that you have logged in to the **gcloud** CLI with your user account by running **gcloud init** (/sdk/gcloud/reference/init) or **gcloud auth login** (/sdk/gcloud/reference/auth/login), or by using **Cloud Shell** (/shell/docs), which automatically logs you into the **gcloud** CLI. You can check the currently active account by running **gcloud auth list** (/sdk/gcloud/reference/auth/list).

Save the request body in a file named `request.json`, and execute the following command:








```
curl -X POST \
  -H "Authorization: Bearer $(gcloud auth print-access-token)" \
  -H "Content-Type: application/json; charset=utf-8" \
  -d @request.json \
  "https://LOCATION  -aiplatform.googleapis.com/v1beta1/projects/PROJ
```

## Get ROUGE score

The following example calls the Gen AI evaluation service API to get the ROUGE score of a prediction, generated by a number of inputs. The ROUGE inputs use `metric_spec`, which determines the metric's behavior.

RESTPython (#python)Go (#go)  
(#rest)

Before using any of the request data, make the following replacements:

- PROJECT\_ID : Your project ID  
(/resource-manager/docs/creating-managing-projects#identifiers).
- LOCATION : The region to process the request.
- PREDICTION : LLM response.
- REFERENCE : Golden LLM response for reference.
- ROUGE\_TYPE : The calculation used to determine the rouge score. See `metric_spec.rouge_type`  
(/vertex-ai/generative-ai/docs/model-reference/evaluation#rougeinput) for acceptable values.
- USE\_STEMMER : Determines whether the Porter stemmer is used to strip word suffixes to improve matching. For acceptable values, see `metric_spec.use_stemmer`  
(/vertex-ai/generative-ai/docs/model-reference/evaluation#rougeinput).
- SPLIT\_SUMMARIES : Determines if new lines are added between `rougeLsum` sentences. For acceptable values, see `metric_spec.split_summaries`  
(/vertex-ai/generative-ai/docs/model-reference/evaluation#rougeinput) .

## HTTP method and URL:

POST [https://LOCATION-aiplatform.googleapis.com/v1beta1/projects/PROJECT\\_ID](https://LOCATION-aiplatform.googleapis.com/v1beta1/projects/PROJECT_ID)

## Request JSON body:

```
{
  "rouge_input": {
    "instances": {
      "prediction": "PREDICTION",
      "reference": "REFERENCE.",
    },
    "metric_spec": {
      "rouge_type": "ROUGE_TYPE",
      "use_stemmer": USE_STEMMER,
      "split_summaries": SPLIT_SUMMARIES,
    }
  }
}
```


To send your request, choose one of these options:

**curlPowerShell** (#powershell)  
(#curl)

★ **Note:** The following command assumes that you have logged in to the **gcloud** CLI with your user account by running **gcloud init** (/sdk/gcloud/reference/init) or **gcloud auth login** (/sdk/gcloud/reference/auth/login), or by using **Cloud Shell** (/shell/docs), which automatically logs you into the **gcloud** CLI. You can check the currently active account by running **gcloud auth list** (/sdk/gcloud/reference/auth/list).

Save the request body in a file named **request.json**, and execute the following command:

```
curl -X POST \
  -H "Authorization: Bearer $(gcloud auth print-access-token)" \
  -H "Content-Type: application/json; charset=utf-8" \
```

```
-d @request.json \  
"https://LOCATION -aiplatform.googleapis.com/v1beta1/projects/PROJ
```

## What's next

- For detailed documentation, see [Run an evaluation](#)  
(/vertex-ai/generative-ai/docs/models/online-pipeline-services).

Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) (https://creativecommons.org/licenses/by/4.0/), and code samples are licensed under the [Apache 2.0 License](https://www.apache.org/licenses/LICENSE-2.0) (https://www.apache.org/licenses/LICENSE-2.0). For details, see the [Google Developers Site Policies](https://developers.google.com/site-policies) (https://developers.google.com/site-policies). Java is a registered trademark of Oracle and/or its affiliates.

Last updated 2025-06-06 UTC.