

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).

# Generate content with the Gemini API in Vertex AI

## Release Notes

Use `generateContent` or `streamGenerateContent` to generate content with Gemini.

The Gemini model family includes models that work with multimodal prompt requests. The term multimodal indicates that you can use more than one modality, or type of input, in a prompt. Models that aren't multimodal accept prompts only with text. Modalities can include text, audio, video, and more.

## Create a Google Cloud account to get started

To start using the Gemini API in Vertex AI, [create a Google Cloud account](#)

(<https://console.cloud.google.com/freetrial?redirectPath=/marketplace/product/google/cloudaicompanion.googleapis.com>)

After creating your account, use this document to review the Gemini model [request body](#) (#request), [model parameters](#) (#parameters), [response body](#) (#response), and some sample [requests](#) (#sample-requests).

When you're ready, see the [Gemini API in Vertex AI quickstart](#)

(/vertex-ai/generative-ai/docs/start/quickstarts/quickstart-multimodal) to learn how to send a request to the Gemini API in Vertex AI using a programming language SDK or the REST API.

## Supported models

All Gemini models support content generation.

**Note:** Adding a lot of images to a request increases response latency.

## Parameter list

See [examples](#) (#sample-requests) for implementation details.

## Request body

```
{
  "cachedContent": string,
  "contents": [
    {
      "role": string,
      "parts": [
        {
          // Union field data can be only one of the following:
          "text": string,
          "inlineData": {
            "mimeType": string,
            "data": string
          },
          "fileData": {
            "mimeType": string,
            "fileUri": string
          },
          // End of list of possible types for union field data.

          "videoMetadata": {
            "startOffset": {
              "seconds": integer,
              "nanos": integer
            },
            "endOffset": {
              "seconds": integer,
              "nanos": integer
            }
          }
        }
      ]
    }
  ],
}
```

```

"systemInstruction": {
  "role": string,
  "parts": [
    {
      "text": string
    }
  ]
},
"tools": [
  {
    "functionDeclarations": [
      {
        "name": string,
        "description": string,
        "parameters": {
          object (OpenAPIObjectSchema) (https://spec.openapis.org/oas/v3.0.3#schema)
        }
      }
    ]
  }
],
"safetySettings": [
  {
    "category": enum (HarmCategory),
    "threshold": enum (HarmBlockThreshold)
  }
],
"generationConfig": {
  "temperature": number,
  "topP": number,
  "topK": number,
  "candidateCount": integer,
  "maxOutputTokens": integer,
  "presencePenalty": float,
  "frequencyPenalty": float,
  "stopSequences": [
    string
  ],
  "responseMimeType": string,
  "responseSchema": schema (/vertex-ai/docs/reference/rest/v1/projects.locations.cachedContents#Schema)
  "seed": integer,
  "responseLogprobs": boolean,
  "logprobs": integer,
  "audioTimestamp": boolean
},
"labels": {
  string: string
}

```

```
}  
}
```

The request body contains data with the following parameters:

Parameters	
cachedContent	<p>Optional: <b>string</b></p> <p>The name of the cached content used as context to serve the prediction. Format: <b>projects/{project}/locations/{location}/cachedContents/{cachedContent}</b></p>
contents	<p>Required: <b>Content</b></p> <p>The content of the current conversation with the model.</p> <p>For single-turn queries, this is a single instance. For multi-turn queries, this is a repeated field that contains conversation history and the latest request.</p>
systemInstruction	<p>Optional: <b>Content</b></p> <p>Available for <b>gemini-2.0-flash</b> and <b>gemini-2.0-flash-lite</b>.</p> <p>Instructions for the model to steer it toward better performance. For example, "Answer as concisely as possible" or "Don't use technical terms in your response".</p> <p>The <b>text</b> strings count toward the token limit.</p> <p>The <b>role</b> field of <b>systemInstruction</b> is ignored and doesn't affect the performance of the model.</p> <p>★ <b>Note:</b> Only <b>text</b> should be used in <b>parts</b> and content in each <b>part</b> should be in a separate paragraph.</p>
tools	<p>Optional. A piece of code that enables the system to interact with external systems to perform an action, or set of actions, outside of knowledge and scope of the model. See <a href="#">Function calling</a> (/vertex-ai/generative-ai/docs/model-reference/function-calling).</p>
toolConfig	<p>Optional. See <a href="#">Function calling</a> (/vertex-ai/generative-ai/docs/model-reference/function-calling).</p>
safetySettings	<p>Optional: <b>SafetySetting</b></p>

Per request settings for blocking unsafe content.

Enforced on `GenerateContentResponse.candidates`.

<b>generationConfig</b>	Optional: <b>GenerationConfig</b>  Generation configuration settings.
<b>labels</b>	Optional: <b>string</b>  Metadata that you can add to the API call in the format of key-value pairs.

**contents**

The base structured data type containing multi-part content of a message.

This class consists of two main properties: `role` and `parts`. The `role` property denotes the individual producing the content, while the `parts` property contains multiple elements, each representing a segment of data within a message.

**Parameters**

<b>role</b>	<b>string</b>  The identity of the entity that creates the message. The following values are supported: <ul style="list-style-type: none"><li>• <b>user</b>: This indicates that the message is sent by a real person, typically a user-generated message.</li><li>• <b>model</b>: This indicates that the message is generated by the model.</li></ul> The <b>model</b> value is used to insert messages from the model into the conversation during multi-turn conversations.
<b>parts</b>	<b>Part</b>  A list of ordered parts that make up a single message. Different parts may have different <a href="https://www.iana.org/assignments/media-types/media-types.xml">IANA MIME types</a> ( <a href="https://www.iana.org/assignments/media-types/media-types.xml">https://www.iana.org/assignments/media-types/media-types.xml</a> ).  For limits on the inputs, such as the maximum number of tokens or the number of images, see the model specifications on the <a href="/vertex-ai/generative-ai/docs/learn/models">Google models</a> ( <a href="/vertex-ai/generative-ai/docs/learn/models">/vertex-ai/generative-ai/docs/learn/models</a> ) page.  To compute the number of tokens in your request, see <a href="/vertex-ai/generative-ai/docs/multimodal/get-token-count">Get token count</a> ( <a href="/vertex-ai/generative-ai/docs/multimodal/get-token-count">/vertex-ai/generative-ai/docs/multimodal/get-token-count</a> ).

parts

A data type containing media that is part of a multi-part Content message.

Parameters

text	<p>Optional: <b>string</b></p> <p>A text prompt or code snippet.</p>
inlineData	<p>Optional: <b>Blob</b></p> <p>Inline data in raw bytes.</p> <p>For <b>gemini-2.0-flash-lite</b> and <b>gemini-2.0-flash</b>, you can specify up to 3000 images by using <b>inlineData</b>.</p>
fileData	<p>Optional: <b>fileData</b></p> <p>Data stored in a file.</p>
functionCall	<p>Optional: <b>FunctionCall</b>.</p> <p>It contains a string representing the <b>FunctionDeclaration.name</b> field and a structured JSON object containing any parameters for the function call predicted by the model.</p> <p>See <a href="#">Function calling</a> (/vertex-ai/generative-ai/docs/model-reference/function-calling).</p>
functionResponse	<p>Optional: <b>FunctionResponse</b>.</p> <p>The result output of a <b>FunctionCall</b> that contains a string representing the <b>FunctionDeclaration.name</b> field and a structured JSON object containing any output from the function call. It is used as context to the model.</p> <p>See <a href="#">Function calling</a> (/vertex-ai/generative-ai/docs/model-reference/function-calling).</p>
videoMetadata	<p>Optional: <b>VideoMetadata</b></p> <p>For video input, the start and end offset of the video in <a href="#">Duration</a> (<a href="https://protobuf.dev/reference/protobuf/google.protobuf/#duration">https://protobuf.dev/reference/protobuf/google.protobuf/#duration</a>) format. For example, to specify a 10 second clip starting at 1:00, set "<b>startOffset</b>": { "seconds": 60 } and "<b>endOffset</b>": { "seconds": 70 }.</p> <p>The metadata should only be specified while the video data is presented in <b>inlineData</b> or <b>fileData</b>.</p>

Content blob. If possible send as text rather than raw bytes.

**mimeType**

string

 Click to expand MIME types

- For **gemini-2.0-flash-lite** and **gemini-2.0-flash**, the maximum length of an audio file is 8.4 hours and the maximum length of a video file (without audio) is one hour. For more information, see Gemini [audio](/vertex-ai/generative-ai/docs/multimodal/audio-understanding#audio-requirements) (/vertex-ai/generative-ai/docs/multimodal/audio-understanding#audio-requirements) and [video](/vertex-ai/generative-ai/docs/multimodal/video-understanding#video-requirements) (/vertex-ai/generative-ai/docs/multimodal/video-understanding#video-requirements) requirements.

Text files must be UTF-8 encoded. The contents of the text file count toward the token limit.

There is no limit on image resolution.

data	bytes
	The <a href="#">base64 encoding</a> (/vertex-ai/generative-ai/docs/image/base64-encode) of the image, PDF, or video to include inline in the prompt. When including media inline, you must also specify the media type ( <b>mimeType</b> ) of the data.
	Size limit: 20MB

FileData

URI or web-URL data.

Parameters	
mimeType	string
	<a href="#">IANA MIME type</a> ( <a href="https://www.iana.org/assignments/media-types/media-types.xml">https://www.iana.org/assignments/media-types/media-types.xml</a> ) of the data.
fileUri	string
	<p>The URI or URL of the file to include in the prompt. Acceptable values include the following:</p> <ul style="list-style-type: none"><li>• <b>Cloud Storage bucket URI:</b> The object must either be publicly readable or reside in the same Google Cloud project that's sending the request. For <b>gemini-2.0-flash</b> and <b>gemini-2.0-flash-lite</b>, the size limit is 2 GB.</li><li>• <b>HTTP URL:</b> The file URL must be publicly readable. You can specify one video file, one audio file, and up to 10 image files per request. Audio files, video files, and documents can't exceed 15 MB.</li><li>• <b>YouTube video URL:</b>The YouTube video must be either owned by the account that you used to sign in to the Google Cloud console or is public. Only one YouTube video URL is supported per request.</li></ul> <p>When specifying a <b>fileURI</b>, you must also specify the media type (<b>mimeType</b>) of the file. If VPC Service Controls is enabled, specifying a media file URL for <b>fileURI</b> is not supported.</p>

functionCall



A predicted `functionCall` returned from the model that contains a string representing the `functionDeclaration.name` and a structured JSON object containing the parameters and their values.

Parameters	
name	<div>string</div> <div>The name of the function to call.</div>
args	<div>Struct</div> <div>The function parameters and values in JSON object format.</div> <div>See <a href="#">Function calling</a> (/vertex-ai/generative-ai/docs/model-reference/function-calling) for parameter details.</div>

**functionResponse**

The resulting output from a `FunctionCall` that contains a string representing the `FunctionDeclaration.name`. Also contains a structured JSON object with the output from the function (and uses it as context for the model). This should contain the result of a `FunctionCall` made based on model prediction.

Parameters	
name	<div>string</div> <div>The name of the function to call.</div>
response	<div>Struct</div> <div>The function response in JSON object format.</div>

**videoMetadata**

Metadata describing the input video content.

Parameters	
startOffset	<div>Optional: <code>google.protobuf.Duration</code></div> <div>The start offset of the video.</div>

endOffset

Optional: `google.protobuf.Duration`

The end offset of the video.

safetySetting

Safety settings.

Parameters

category

Optional: `HarmCategory`

The safety category to configure a threshold for. Acceptable values include the following:

Click to expand safety categories

- `HARM_CATEGORY_SEXUALLY_EXPLICIT`
- `HARM_CATEGORY_HATE_SPEECH`
- `HARM_CATEGORY_HARASSMENT`
- `HARM_CATEGORY_DANGEROUS_CONTENT`

threshold

Optional: `HarmBlockThreshold`

The threshold for blocking responses that could belong to the specified safety category based on probability.

- `OFF`
- `BLOCK_NONE`
- `BLOCK_LOW_AND_ABOVE`
- `BLOCK_MEDIUM_AND_ABOVE`
- `BLOCK_ONLY_HIGH`

method

Optional: `HarmBlockMethod`

Specify if the threshold is used for probability or severity score. If not specified, the threshold is used for probability score.

harmCategory

Harm categories that block content.

Parameters	
HARM_CATEGORY_UNSPECIFIED	The harm category is unspecified.
HARM_CATEGORY_HATE_SPEECH	The harm category is hate speech.
HARM_CATEGORY_DANGEROUS_CONTENT	The harm category is dangerous content.
HARM_CATEGORY_HARASSMENT	The harm category is harassment.
HARM_CATEGORY_SEXUALLY_EXPLICIT	The harm category is sexually explicit content.

harmBlockThreshold

Probability thresholds levels used to block a response.

Parameters	
HARM_BLOCK_THRESHOLD_UNSPECIFIED	Unspecified harm block threshold.
BLOCK_LOW_AND_ABOVE	Block low threshold and higher (i.e. block more).
BLOCK_MEDIUM_AND_ABOVE	Block medium threshold and higher.
BLOCK_ONLY_HIGH	Block only high threshold (i.e. block less).
BLOCK_NONE	Block none.
OFF	Switches off safety if all categories are turned OFF

harmBlockMethod

A probability threshold that blocks a response based on a combination of probability and severity.

Parameters	
HARM_BLOCK_METHOD_UNSPECIFIED	The harm block method is unspecified.
SEVERITY	The harm block method uses both probability and severity scores.
PROBABILITY	The harm block method uses the probability score.

# generationConfig

Configuration settings used when generating the prompt.

## Parameters

temperature	<p>Optional: <code>float</code></p> <p>The temperature is used for sampling during response generation, which occurs when <code>topP</code> and <code>topK</code> are applied. Temperature controls the degree of randomness in token selection. Lower temperatures are good for prompts that require a less open-ended or creative response, while higher temperatures can lead to more diverse or creative results. A temperature of <code>0</code> means that the highest probability tokens are always selected. In this case, responses for a given prompt are mostly deterministic, but a small amount of variation is still possible.</p> <p>If the model returns a response that's too generic, too short, or the model gives a fallback response, try increasing the temperature.</p> <ul style="list-style-type: none"><li>• Range for <code>gemini-2.0-flash-lite</code>: <code>0.0</code> - <code>2.0</code> (default: <code>1.0</code>)</li><li>• Range for <code>gemini-2.0-flash</code>: <code>0.0</code> - <code>2.0</code> (default: <code>1.0</code>)</li></ul> <p>For more information, see <a href="#">Content generation parameters (/vertex-ai/generative-ai/docs/multimodal/content-generation-parameters#temperature)</a>.</p>
topP	<p>Optional: <code>float</code></p> <p>If specified, nucleus sampling is used.</p> <p><u>Top-P</u> (<a href="#">/vertex-ai/generative-ai/docs/multimodal/content-generation-parameters#top-p</a>)</p> <p>changes how the model selects tokens for output. Tokens are selected from the most (see top-K) to least probable until the sum of their probabilities equals the top-P value. For example, if tokens A, B, and C have a probability of 0.3, 0.2, and 0.1 and the top-P value is <code>0.5</code>, then the model will select either A or B as the next token by using temperature and excludes C as a candidate.</p> <p>Specify a lower value for less random responses and a higher value for more random responses.</p> <ul style="list-style-type: none"><li>• Range: <code>0.0</code> - <code>1.0</code></li><li>• Default for <code>gemini-2.0-flash-lite</code>: <code>0.95</code></li><li>• Default for <code>gemini-2.0-flash</code>: <code>0.95</code></li></ul>

---

**candidateCount**Optional: **int**

The number of response variations to return. For each request, you're charged for the output tokens of all candidates, but are only charged once for the input tokens.

Specifying multiple candidates is a Preview feature that works with **generateContent** (**streamGenerateContent** is not supported). The following models are supported:

- **gemini-2.0-flash-lite**: 1-8, default: 1
- **gemini-2.0-flash**: 1-8, default: 1

---

**maxOutputTokens**Optional: **int**

Maximum number of tokens that can be generated in the response. A token is approximately four characters. 100 tokens correspond to roughly 60-80 words.

Specify a lower value for shorter responses and a higher value for potentially longer responses.

For more information, see [Content generation parameters](#) (/vertex-ai/generative-ai/docs/multimodal/content-generation-parameters#max-output-tokens)

.

---

**stopSequences**Optional: **List[string]**

Specifies a list of strings that tells the model to stop generating text if one of the strings is encountered in the response. If a string appears multiple times in the response, then the response truncates where it's first encountered. The strings are case-sensitive.

For example, if the following is the returned response when **stopSequences** isn't specified:

```
public static string reverse(string myString)
```

Then the returned response with **stopSequences** set to [ "Str", "reverse" ] is:

```
public static string
```

Maximum 5 items in the list.

For more information, see [Content generation parameters](#) (/vertex-ai/generative-ai/docs/multimodal/content-generation-parameters#stop-sequences)

<b>presencePenalty</b>	<p>Optional: <b>float</b></p> <p>Positive penalties.</p> <p>Positive values penalize tokens that already appear in the generated text, increasing the probability of generating more diverse content.</p> <p>The maximum value for <b>presencePenalty</b> is up to, but not including, <b>2.0</b>. Its minimum value is <b>-2.0</b>.</p> <p>Supported by Gemini 2.0 Flash-Lite and Gemini 2.0 Flash.</p>
<b>frequencyPenalty</b>	<p>Optional: <b>float</b></p> <p>Positive values penalize tokens that repeatedly appear in the generated text, decreasing the probability of repeating content.</p> <p>This maximum value for <b>frequencyPenalty</b> is up to, but not including, <b>2.0</b>. Its minimum value is <b>-2.0</b>.</p> <p>Supported by Gemini 2.0 Flash-Lite and Gemini 2.0 Flash.</p>
<b>responseMimeType</b>	<p>Optional: <b>string (enum)</b></p> <p>Available for the following models:</p> <ul style="list-style-type: none"><li>• Gemini 2.0 Flash-Lite</li><li>• Gemini 2.0 Flash</li></ul> <p>The output response MIME type of the generated candidate text.</p> <p>The following MIME types are supported:</p> <ul style="list-style-type: none"><li>• <b>application/json</b>: JSON response in the candidates.</li><li>• <b>text/plain</b> (default): Plain text output.</li><li>• <b>text/x.enum</b>: For classification tasks, output an enum value as defined in the response schema.</li></ul> <p>Specify the appropriate response type to avoid unintended behaviors. For example, if you require a JSON-formatted response, specify <b>application/json</b> and not <b>text/plain</b>.</p>
<b>responseSchema</b>	<p>Optional: <u><b>schema</b></u> (/vertex-ai/docs/reference/rest/v1/projects.locations.cachedContents#Schema)</p> <p>The schema that generated candidate text must follow. For more information, see <u><b>Control generated output</b></u></p>

(/vertex-ai/generative-ai/docs/multimodal/control-generated-output).

You must specify the `responseMimeType` parameter to use this parameter.

Available for the following models:

- Gemini 2.0 Flash-Lite
- Gemini 2.0 Flash

seed	<p>Optional: <code>int</code></p> <p>When <code>seed</code> is fixed to a specific value, the model makes a best effort to provide the same response for repeated requests. Deterministic output isn't guaranteed. Also, changing the model or parameter settings, such as the temperature, can cause variations in the response even when you use the same seed value. By default, a random seed value is used.</p> <p>Available for the following models:</p> <ul style="list-style-type: none"><li>• Gemini 2.5 Flash</li><li>• Gemini 2.5 Pro</li><li>• Gemini 2.0 Flash-Lite</li><li>• Gemini 2.0 Flash</li></ul>
responseLogprobs	<p>Optional: <code>boolean</code></p> <p>If true, returns the log probabilities of the tokens that were chosen by the model at each step. By default, this parameter is set to <code>false</code>. The daily limit for requests using <code>responseLogprobs</code> is 1.</p> <p>Available for the following models:</p> <ul style="list-style-type: none"><li>• Gemini 2.0 Flash-Lite</li><li>• Gemini 2.0 Flash</li></ul> <p>This is a preview feature.</p>
logprobs	<p>Optional: <code>int</code></p> <p>Returns the log probabilities of the top candidate tokens at each generation step. The model's chosen token might not be the same as the top candidate token at each step. Specify the number of candidates to return by using an integer value in the range of 1-5.</p> <p>You must enable <code>responseLogprobs</code> (<code>#responseLogprobs</code>) to use this parameter. The daily limit for requests using <code>logprobs</code> is 1.</p> <p>This is a preview feature.</p>

---

**audioTimestamp**Optional: **boolean**

Available for the following models:

- Gemini 2.0 Flash-Lite
- Gemini 2.0 Flash

Enables timestamp understanding for audio-only files.

This is a preview feature.

---

## Response body

```
{
  "candidates": [
    {
      "content": {
        "parts": [
          {
            "text": string
          }
        ]
      },
      "finishReason": enum (FinishReason),
      "safetyRatings": [
        {
          "category": enum (HarmCategory),
          "probability": enum (HarmProbability),
          "blocked": boolean
        }
      ],
      "citationMetadata": {
        "citations": [
          {
            "startIndex": integer,
            "endIndex": integer,
            "uri": string,
            "title": string,
            "license": string,
            "publicationDate": {
              "year": integer,
              "month": integer,
              "day": integer
            }
          }
        ]
      }
    }
  ]
}
```



```
    }
  ]
},
"avgLogprobs": double,
"logprobsResult": {
  "topCandidates": [
    {
      "candidates": [
        {
          "token": string,
          "logProbability": float
        }
      ]
    }
  ],
  "chosenCandidates": [
    {
      "token": string,
      "logProbability": float
    }
  ]
}
}
},
"usageMetadata": {
  "promptTokenCount": integer,
  "candidatesTokenCount": integer,
  "totalTokenCount": integer
},
"modelVersion": string
}
```

Response element	Description
modelVersion	The model and version used for generation. For example: <b>gemini-2.0-flash-lite-001</b> .
text	The generated text.
finishReason	<p>The reason why the model stopped generating tokens. If empty, the model has not stopped generating the tokens. Because the response uses the prompt for context, it's not possible to change the behavior of how the model stops generating tokens.</p> <ul style="list-style-type: none"><li><b>FINISH_REASON_STOP</b>: Natural stop point of the model or provided stop sequence.</li><li><b>FINISH_REASON_MAX_TOKENS</b>: The maximum number of tokens as specified in the request was reached.</li></ul>

- **FINISH\_REASON\_SAFETY:** Token generation was stopped because the response was flagged for safety reasons. Note that `Candidate.content` is empty if content filters block the output.
- **FINISH\_REASON\_RECITATION:** The token generation was stopped because the response was flagged for unauthorized citations.
- **FINISH\_REASON\_BLOCKLIST:** Token generation was stopped because the response includes blocked terms.
- **FINISH\_REASON\_PROHIBITED\_CONTENT:** Token generation was stopped because the response was flagged for prohibited content, such as child sexual abuse material (CSAM).
- **FINISH\_REASON\_SPII:** Token generation was stopped because the response was flagged for sensitive personally identifiable information (SPII).
- **FINISH\_REASON\_MALFORMED\_FUNCTION\_CALL:** Candidates were blocked because of malformed and unparsable function call.
- **FINISH\_REASON\_OTHER:** All other reasons that stopped the token
- **FINISH\_REASON\_UNSPECIFIED:** The finish reason is unspecified.

category	<p>The safety category to configure a threshold for. Acceptable values include the following:</p> <div><div><div>+</div></div>Click to expand safety categories</div> <ul style="list-style-type: none"><li>• HARM_CATEGORY_SEXUALLY_EXPLICIT</li><li>• HARM_CATEGORY_HATE_SPEECH</li><li>• HARM_CATEGORY_HARASSMENT</li><li>• HARM_CATEGORY_DANGEROUS_CONTENT</li></ul>
probability	<p>The harm probability levels in the content.</p> <ul style="list-style-type: none"><li>• HARM_PROBABILITY_UNSPECIFIED</li><li>• NEGLIGIBLE</li><li>• LOW</li><li>• MEDIUM</li><li>• HIGH</li></ul>
blocked	<p>A boolean flag associated with a safety attribute that indicates if the model's input or output was blocked.</p>
startIndex	<p>An integer that specifies where a citation starts in the <code>content</code>.</p>

endIndex	An integer that specifies where a citation ends in the <b>content</b> .
url	The URL of a citation source. Examples of a URL source might be a news website or a GitHub repository.
title	The title of a citation source. Examples of source titles might be that of a news article or a book.
license	The license associated with a citation.
publicationDate	The date a citation was published. Its valid formats are YYYY, YYYY-MM, and YYYY-MM-DD.
avgLogprobs	Average log probability of the candidate.
logprobsResult	Returns the top candidate tokens ( <b>topCandidates</b> ) and the actual chosen tokens ( <b>chosenCandidates</b> ) at each step.
token	Generative AI models break down text data into tokens for processing, which can be characters, words, or phrases.
logProbability	A log probability value that indicates the model's confidence for a particular token.
promptTokenCount	Number of tokens in the request.
candidatesTokenCount	Number of tokens in the response(s).
totalTokenCount	Number of tokens in the request and response(s).

## Examples

### Text Generation

Generate a text response from a text input.

Gen AI SDK for PythonPython (OpenAI)...

Go (#go)

```
from google import genai
from google.genai.types import HttpOptions

client = genai.Client(http_options=HttpOptions(api_version="v1"))
response = client.models.generate_content(
    model="gemini-2.5-flash-preview-05-20",
    contents="How does AI work?",
```

```

)
print(response.text)
# Example response:
# Okay, let's break down how AI works. It's a broad field, so I'll focus on the
#
# Here's a simplified overview:
# ...

```

## Using multimodal prompt

Generate a text response from a multimodal input, such as text and an image.

<a href="#">Gen AI SDK for Python</a> Python (OpenAI)... (#gen-ai-sdk-for-python)	<a href="#">Go</a> (#go)
<pre> from google import genai from google.genai.types import HttpOptions, Part  client = genai.Client(http_options=HttpOptions(api_version="v1")) response = client.models.generate_content(     model="gemini-2.5-flash-preview-05-20",     contents=[         "What is shown in this image?",         Part.from_uri(             file_uri="gs://cloud-samples-data/generative-ai/image/scones.jpg",             mime_type="image/jpeg",         ),     ], ) print(response.text) # Example response: # The image shows a flat lay of blueberry scones arranged on parchment paper. T </pre>	

## Streaming text response

Generate a streaming model response from a text input.

<a href="#">Gen AI SDK for Python</a> Python (OpenAI)... (#gen-ai-sdk-for-python)	<a href="#">Go</a> (#go)

```
from google import genai
from google.genai.types import HttpOptions

client = genai.Client(http_options=HttpOptions(api_version="v1"))

for chunk in client.models.generate_content_stream(
    model="gemini-2.5-flash-preview-05-20",
    contents="Why is the sky blue?",
):
    print(chunk.text, end="")
# Example response:
# The
# sky appears blue due to a phenomenon called Rayleigh scattering. Here's
# a breakdown of why:
# ...
```

## Model versions

To use the [auto-updated version](#) (/vertex-ai/generative-ai/docs/learn/model-versioning#auto-updated-version), specify the model name without the trailing version number, for example **gemini-2.0-flash** instead of **gemini-2.0-flash-001**.

For more information, see [Gemini model versions and lifecycle](#) (/vertex-ai/generative-ai/docs/learn/model-versioning#gemini-model-versions).

## What's next

- Learn more about the [Gemini API in Vertex AI](#) (/vertex-ai/generative-ai/docs/model-reference/gemini).
- Learn more about [Function calling](#) (/vertex-ai/generative-ai/docs/multimodal/function-calling).
- Learn more about [Grounding responses for Gemini models](#) (/vertex-ai/generative-ai/docs/multimodal/ground-gemini).

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/>), and code samples are licensed under the [Apache 2.0 License](#) (<https://www.apache.org/licenses/LICENSE-2.0>). For details, see the [Google Developers 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.