Introduction &

You can interact with the API through HTTP requests from any language, via our official Python bindings, our official Node.js library, or a **community-maintained library**.

To install the official Python bindings, run the following command:

pip install openai	ð
To install the official Node.js library, run the following command in your Node.js project directory:	
npm install openai@^4.0.0	Ð

Authentication &

API keys∂

The OpenAI API uses API keys for authentication. You can create API keys at a user or service account level. Service accounts are tied to a "bot" individual and should be used to provision access for production systems. Each API key can be scoped to one of the following,

Project keys - Provides access to a single project (**preferred option**); access **Project API keys** by selecting the specific project you wish to generate keys against.

User keys - Our legacy keys. Provides access to all organizations and all projects that user has been added to; access **API Keys** to view your available keys. We highly advise transitioning to project keys for best security practices, although access via this method is currently still supported.

Remember that your API key is a secret! Do not share it with others or expose it in any client-side code (browsers, apps). Production requests must be routed through your own backend server where your API key can be securely loaded from an environment variable or key management service.

All API requests should include your API key in an Authorization HTTP header as follows:

Authorization: Bearer OPENAI_API_KEY	Ó

Organizations and projects (optional)

For users who belong to multiple organizations or are accessing their projects through their legacy user API key, you can pass a header to specify which organization and project is used for an API request. Usage from these API requests will count as usage for the specified organization and project.

To access the	Default project	in an organization,	leave out the	OpenAI-Project	header

Example curl command:

1	curl https://api.openai.com/v1/models \	6
2	-H "Authorization: Bearer \$OPENAI_API_KEY" \	
3	-H "OpenAI-Organization: org-R4EOK4IIOxoFL8jKhBVr9u0b" \	
4	-H "OpenAI-Project: \$PROJECT_ID"	
Exa	imple with the openai Python package:	

```
\Box
 1 from openai import OpenAI
 2
 3 client = OpenAI(
 4 organization='org-R4EOK4IIOxoFL8jKhBVr9u0b',
 5 project='$PROJECT_ID',
 6)
Example with the openai Node.js package:
                                                                                                               D
  1 import OpenAl from "openai";
 2
 3 const openai = new OpenAI({
      organization: "org-R4EOK4IIOxoFL8jKhBVr9u0b",
      project: "$PROJECT_ID",
 5
 6 });
```

Organization IDs can be found on your **Organization settings** page. Project IDs can be found on your **General settings** page by selecting the specific project.

Making requests ₽

You can paste the command below into your terminal to run your first API request. Make sure to replace \$OPENAI_API_KEY with your secret API key. If you are using a legacy user key and you have multiple projects, you will also need to **specify the Project Id**. For improved security, we recommend transitioning to project based keys instead.

```
curl https://api.openai.com/v1/chat/completions \
-H "Content-Type: application/json" \
-H "Authorization: Bearer $OPENAI_API_KEY" \
-d '{
    "model": "gpt-4o-mini",
    "messages": [{"role": "user", "content": "Say this is a test!"}],
    "temperature": 0.7
}
```

This request queries the <code>[gpt-4o-mini]</code> model (which under the hood points to a <code>[gpt-4o-mini]</code> model variant) to complete the text starting with a prompt of "Say this is a test". You should get a response back that resembles the following:

```
1
   {
2
      "id": "chatcmpl-abc123",
3
      "object": "chat.completion",
4
      "created": 1677858242,
      "model": "gpt-4o-mini",
5
6
      "usage": {
7
         "prompt_tokens": 13,
8
         "completion_tokens": 7,
9
         "total tokens": 20
10
11
      "choices": [
12
13
            "message": {
```

```
"role": "assistant",
14
               "content": "\n\nThis is a test!"
15
16
            },
17
            "logprobs": null,
            "finish reason": "stop",
18
            "index": 0
19
20
         }
21
       1
22 }
```

Now that you've generated your first chat completion, let's break down the **response object**. We can see the finish_reason is stop which means the API returned the full chat completion generated by the model without running in any limits. In the choices list, we only generated a single message but you can set the n parameter to generate multiple messages choices.

Streaming 2

The OpenAI API provides the ability to stream responses back to a client in order to allow partial results for certain requests. To achieve this, we follow the **Server-sent events** standard. Our official **Node** and **Python** libraries include helpers to make parsing these events simpler.

Streaming is supported for both the **Chat Completions API** and the **Assistants API**. This section focuses on how streaming works for Chat Completions. Learn more about how streaming works in the Assistants API **here**.

In Python, a streaming request looks like:

```
ð
    from openai import OpenAI
1
2
3
    client = OpenAI()
4
5
    stream = client.chat.completions.create(
6
      model="gpt-4o-mini",
7
      messages=[{"role": "user", "content": "Say this is a test"}],
8
      stream=True,
9
   )
10 for chunk in stream:
11
      if chunk.choices[0].delta.content is not None:
12
         print(chunk.choices[0].delta.content, end="")
```

In Node / Typescript, a streaming request looks like:

```
import OpenAI from "openai";
1
2
3
    const openai = new OpenAI();
4
5
    async function main() {
6
      const stream = await openai.chat.completions.create({
7
         model: "gpt-4o-mini",
8
         messages: [{ role: "user", content: "Say this is a test" }],
9
         stream: true,
10
      });
11
      for await (const chunk of stream) {
12
         process.stdout.write(chunk.choices[0]?.delta?.content II "");
```

```
13 }
14 }
15
16 main();
```

Parsing Server-sent events ℰ

Parsing Server-sent events is non-trivial and should be done with caution. Simple strategies like splitting by a new line may result in parsing errors. We recommend using **existing client libraries** when possible.

Audio*⊗*

Learn how to turn audio into text or text into audio.

Related guide: Speech to text

Create speech ₽

POST https://api.openai.com/v1/audio/speech

Generates audio from the input text.

Request body

```
model string Required
One of the available TTS models: tts-1 or tts-1-hd

input string Required
The text to generate audio for. The maximum length is 4096 characters.

voice string Required
The voice to use when generating the audio. Supported voices are alloy, echo, fable, onyx, nova, and shimmer.

Previews of the voices are available in the Text to speech guide.

response_format string Optional Defaults to mp3
```

The format to audio in. Supported formats are mp3, opus, aac

The speed of the generated audio. Select a value from 0.25 to

Returns

The audio file content.

, flac , wav , and pcm .

4.0 . 1.0 is the default.

Create transcription ₽

speed number Optional Defaults to 1

POST https://api.openai.com/v1/audio/transcriptions

Transcribes audio into the input language.

Request body

file file Required

Default Word timestamps Segment timestamps Example request python ✓ □

1 from openai import OpenAl

2 client = OpenAI()

3

API Reference - OpenAI API

The audio file object (not file name) to transcribe, in one of these formats: flac, mp3, mp4, mpeg, mpga, m4a, ogg, wav, or webm.

omiais. nac, mps, mp4, mpeg, mpga, m4a, ogg, wav

model string Required

ID of the model to use. Only whisper-1 (which is powered by our open source Whisper V2 model) is currently available.

language string Optional

The language of the input audio. Supplying the input language in **ISO-639-1** format will improve accuracy and latency.

prompt string Optional

An optional text to guide the model's style or continue a previous audio segment. The **prompt** should match the audio language.

response_format string Optional Defaults to json
The format of the transcript output, in one of these options: json,
text, srt, verbose_json, or vtt.

temperature number Optional Defaults to 0

The sampling temperature, between 0 and 1. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic. If set to 0, the model will use **log probability** to automatically increase the temperature until certain thresholds are hit.

timestamp_granularities[] array Optional Defaults to segment The timestamp granularities to populate for this transcription.

[response_format] must be set [verbose_json] to use timestamp granularities. Either or both of these options are supported: [word], or [segment]. Note: There is no additional latency for segment timestamps, but generating word timestamps incurs additional latency.

Returns

The transcription object or a verbose transcription object.

Create translation ₽

POST https://api.openai.com/v1/audio/translations

Translates audio into English.

Request body

file file Required

The audio file object (not file name) translate, in one of these formats: flac, mp3, mp4, mpeg, mpga, m4a, ogg, wav, or webm.

model string Required

ID of the model to use. Only whisper-1 (which is powered by our open source Whisper V2 model) is currently available.

```
4 audio_file = open("speech.mp3", "rb")
5 transcript = client.audio.transcriptions.create(
6 model="whisper-1",
7 file=audio_file
8 )

Response

1 {
2 "text": "Imagine the wildest idea that you've ever had, and you'n
3 }
```

```
仚
Example request
                                            python ~
1 from openai import OpenAl
2 client = OpenAI()
3
4 audio_file = open("speech.mp3", "rb")
5 transcript = client.audio.translations.create(
    model="whisper-1",
7
    file=audio_file
8)
                                                         币
Response
1 {
    "text": "Hello, my name is Wolfgang and I come from Germany.
```

prompt string Optional

An optional text to guide the model's style or continue a previous audio segment. The **prompt** should be in English.

response_format string Optional Defaults to json
The format of the transcript output, in one of these options: json,
text, srt, verbose_json, or vtt.

temperature number Optional Defaults to 0

The sampling temperature, between 0 and 1. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic. If set to 0, the model will use **log probability** to automatically increase the temperature until certain thresholds are hit.

Returns

The translated text.

The transcription object (JSON)∂

Represents a transcription response returned by model, based on the provided input.

text string

The transcribed text.

```
The transcription object (JSON)

1 {
2 "text": "Imagine the wildest idea that you've ever had, and you'uluge of the state of the state
```

The transcription object (Verbose JSON)∂

Represents a verbose json transcription response returned by model, based on the provided input.

language string

The language of the input audio.

duration string

The duration of the input audio.

text string

The transcribed text.

words array

Extracted words and their corresponding timestamps.

∨Show properties

segments array

Segments of the transcribed text and their corresponding details.

∨Show properties

```
The transcription object (Verbose JSON)
                                                          凸
1
    {
     "task": "transcribe",
3
     "language": "english",
     "duration": 8.470000267028809,
5
      "text": "The beach was a popular spot on a hot summer day. F
6
     "segments": [
7
      {
8
        "id": 0,
        "seek": 0,
10
        "start": 0.0,
        "end": 3.319999933242798,
11
12
        "text": " The beach was a popular spot on a hot summer da
13
         50364, 440, 7534, 390, 257, 3743, 4008, 322, 257, 2368,
14
15
16
        "temperature": 0.0,
17
        "avg_logprob": -0.2860786020755768,
18
        "compression_ratio": 1.2363636493682861,
19
        "no_speech_prob": 0.00985979475080967
20
      },
21
22
    ]
23 }
```


Given a list of messages comprising a conversation, the model will return a response.

Related guide: Chat Completions

Create chat completion ₽

POST https://api.openai.com/v1/chat/completions

Creates a model response for the given chat conversation.

Request body

messages array Required

A list of messages comprising the conversation so far. **Example Python code**.

∨Show possible types

model string Required

ID of the model to use. See the **model endpoint compatibility** table for details on which models work with the Chat API.

frequency_penalty number or null Optional Defaults to 0 Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency in the text so far, decreasing the model's likelihood to repeat the same line verbatim.

See more information about frequency and presence penalties.

logit_bias map Optional Defaults to null
Modify the likelihood of specified tokens appearing in the
completion.

Accepts a JSON object that maps tokens (specified by their token ID in the tokenizer) to an associated bias value from -100 to 100. Mathematically, the bias is added to the logits generated by the model prior to sampling. The exact effect will vary per model, but values between -1 and 1 should decrease or increase likelihood of selection; values like -100 or 100 should result in a ban or exclusive selection of the relevant token.

logprobs boolean or null Optional Defaults to false Whether to return log probabilities of the output tokens or not. If true, returns the log probabilities of each output token returned in the content of message.

top_logprobs integer or null Optional

An integer between 0 and 20 specifying the number of most likely tokens to return at each token position, each with an associated log probability. logprobs must be set to true if this parameter is used.

max_tokens integer or null Optional

The maximum number of **tokens** that can be generated in the chat completion.

```
Default
                                                Functions
             Image input
                               Streaming
Example request
                                 gpt-4o ∨ python ∨
                                                          币
    from openai import OpenAI
2
    client = OpenAI()
3
    completion = client.chat.completions.create(
4
5
     model="gpt-40",
6
     messages=[
      {"role": "system", "content": "You are a helpful assistant."},
7
      {"role": "user", "content": "Hello!"}
9
     1
10 )
11
12 print(completion.choices[0].message)
                                                          币
Response
1
```

```
2
     "id": "chatcmpl-123",
3
     "object": "chat.completion",
     "created": 1677652288,
4
     "model": "gpt-4o-mini",
6
     "system_fingerprint": "fp_44709d6fcb",
7
     "choices": [{
8
      "index": 0,
9
      "message": {
10
        "role": "assistant",
11
        "content": "\n\nHello there, how may I assist you today?"
12
      },
      "loaprobs": null.
13
```

The total length of input tokens and generated tokens is limited by the model's context length. **Example Python code** for counting tokens.

n integer or null Optional Defaults to 1
How many chat completion choices to generate for each input message. Note that you will be charged based on the number of generated tokens across all of the choices. Keep n as 1 to minimize costs.

presence_penalty number or null Optional Defaults to 0 Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far, increasing the model's likelihood to talk about new topics.

See more information about frequency and presence penalties.

response format object Optional

An object specifying the format that the model must output. Compatible with **GPT-4 Turbo** and all GPT-3.5 Turbo models newer than gpt-3.5-turbo-1106.

Setting to { "type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨ Show properties

seed integer or null Optional

This feature is in Beta. If specified, our system will make a best effort to sample deterministically, such that repeated requests with the same seed and parameters should return the same result. Determinism is not guaranteed, and you should refer to the system_fingerprint response parameter to monitor changes in the backend.

service_tier string or null Optional Defaults to null Specifies the latency tier to use for processing the request. This parameter is relevant for customers subscribed to the scale tier service:

If set to 'auto', the system will utilize scale tier credits until they are exhausted.

If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency quarentee.

When this parameter is set, the response body will include the service_tier utilized.

stop string / array / null Optional Defaults to null

Up to 4 sequences where the API will stop generating further tokens.

stream boolean or null Optional Defaults to false
If set, partial message deltas will be sent, like in ChatGPT. Tokens
will be sent as data-only server-sent events as they become
available, with the stream terminated by a data: [DONE]
message. Example Python code.

stream_options object or null Optional Defaults to null Options for streaming response. Only set this when you set stream: true.

∨Show properties

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

We generally recommend altering this or top_p but not both.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

tools array Optional

A list of tools the model may call. Currently, only functions are supported as a tool. Use this to provide a list of functions the model may generate JSON inputs for. A max of 128 functions are supported.

√ Show properties

tool_choice string or object Optional

Controls which (if any) tool is called by the model. none means the model will not call any tool and instead generates a message.

auto means the model can pick between generating a message or calling one or more tools. required means the model must call one or more tools. Specifying a particular tool via

{"type": "function", "function": {"name": "my_function"}} forces the model to call that tool.

none is the default when no tools are present. auto is the default if tools are present.

∨ Show possible types

parallel_tool_calls boolean Optional Defaults to true
Whether to enable parallel function calling during tool use.

user string Optional

A unique identifier representing your end-user, which can help OpenAl to monitor and detect abuse. **Learn more**.

function_call Deprecated string or object Optional

Deprecated in favor of tool_choice .

Controls which (if any) function is called by the model. none means the model will not call a function and instead generates a message. auto means the model can pick between generating a message or calling a function. Specifying a particular function via {"name": "my_function"} forces the model to call that function.

none is the default when no functions are present. auto is the default if functions are present.

∨Show possible types

functions Deprecated array Optional

Deprecated in favor of tools .

A list of functions the model may generate JSON inputs for.

∨ Show properties

Returns

Returns a **chat completion** object, or a streamed sequence of **chat completion chunk** objects if the request is streamed.

The chat completion object@

Represents a chat completion response returned by model, based on the provided input.

id string

A unique identifier for the chat completion.

choices array

A list of chat completion choices. Can be more than one if $\begin{bmatrix} n \end{bmatrix}$ is greater than 1.

√Show properties

created integer

The Unix timestamp (in seconds) of when the chat completion was created.

model string

The model used for the chat completion.

service_tier string or null

The service tier used for processing the request. This field is only included if the service_tier parameter is specified in the request.

system_fingerprint string

This fingerprint represents the backend configuration that the model runs with.

Can be used in conjunction with the seed request parameter to understand when backend changes have been made that might impact determinism.

object string

The chat completion object

```
D
```

```
{
2
     "id": "chatcmpl-123",
     "object": "chat.completion",
3
     "created": 1677652288,
     "model": "gpt-4o-mini",
     "system_fingerprint": "fp_44709d6fcb",
     "choices": [{
      "index": 0.
q
      "message": {
       "role": "assistant",
11
        "content": "\n\nHello there, how may I assist you today?",
12
13
       "logprobs": null,
14
      "finish_reason": "stop"
15
16
     "usage": {
17
      "prompt_tokens": 9,
18
      "completion_tokens": 12,
19
      "total_tokens": 21
20
21 }
```

The object type, which is always chat.completion.

usage object

Usage statistics for the completion request.

∨Show properties

The chat completion chunk object *₽*

Represents a streamed chunk of a chat completion response returned by model, based on the provided input.

id string

A unique identifier for the chat completion. Each chunk has the

choices array

A list of chat completion choices. Can contain more than one elements if n is greater than 1. Can also be empty for the last chunk if you set | stream_options: {"include_usage": true} |. ∨Show properties

created integer

The Unix timestamp (in seconds) of when the chat completion was created. Each chunk has the same timestamp.

model string

The model to generate the completion.

service tier string or null

The service tier used for processing the request. This field is only included if the service_tier parameter is specified in the request.

system_fingerprint string

This fingerprint represents the backend configuration that the model runs with. Can be used in conjunction with the seed request parameter to understand when backend changes have been made that might impact determinism.

object string

The object type, which is always chat.completion.chunk .

usage object

An optional field that will only be present when you set stream_options: {"include_usage": true} in your request. When present, it contains a null value except for the last chunk which contains the token usage statistics for the entire request. ∨Show properties

Creates an embedding vector representing the input text.

Embeddings *∂*

Get a vector representation of a given input that can be easily consumed by machine learning models and algorithms.

Related guide: Embeddings

Create embeddings *₽*

POST https://api.openai.com/v1/embeddings

python ~



The chat completion chunk object

1 {"id":"chatcmpl-123","object":"chat.completion.chunk","created": 2

凸

3 {"id":"chatcmpl-123","object":"chat.completion.chunk","created"::

5 6

7 {"id":"chatcmpl-123","object":"chat.completion.chunk","created":

Response

Request body

input string or array Required

Input text to embed, encoded as a string or array of tokens. To embed multiple inputs in a single request, pass an array of strings or array of token arrays. The input must not exceed the max input tokens for the model (8192 tokens for text-embedding-ada-002), cannot be an empty string, and any array must be 2048 dimensions or less. **Example Python code** for counting tokens. \sim Show possible types

model string Required

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them.

encoding_format string Optional Defaults to float
The format to return the embeddings in. Can be either float or
base64.

dimensions integer Optional

The number of dimensions the resulting output embeddings should have. Only supported in text-embedding-3 and later models.

user string Optional

A unique identifier representing your end-user, which can help OpenAl to monitor and detect abuse. **Learn more**.

Returns

A list of embedding objects.

The embedding object *₽*

Represents an embedding vector returned by embedding endpoint.

index integer

The index of the embedding in the list of embeddings.

embedding array

The embedding vector, which is a list of floats. The length of vector depends on the model as listed in the **embedding guide**.

object string

The object type, which is always "embedding".

Fine-tuning *₽*

Manage fine-tuning jobs to tailor a model to your specific training data.

Related guide: Fine-tune models

Create fine-tuning job ₽

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 client.embeddings.create(
5 model="text-embedding-ada-002",
6 input="The food was delicious and the waiter...",
7 encoding_format="float"
8 )
```

```
1
    {
2
     "object": "list",
     "data": [
4
5
        "object": "embedding",
6
       "embedding": [
7
         0.0023064255,
8
         -0.009327292,
9
         .... (1536 floats total for ada-002)
10
         -0.0028842222.
11
       ],
        "index": 0
12
13
      }
14
     1.
     "model": "text-embedding-ada-002",
15
16
     "usage": {
17
      "prompt_tokens": 8,
18
      "total_tokens": 8
19 }
```

凸

```
The embedding object
```

20 }

```
1 {
2  "object": "embedding",
3  "embedding": [
4     0.0023064255,
5     -0.009327292,
6     .... (1536 floats total for ada-002)
7     -0.0028842222,
8     ],
9     "index": 0
10 }
```

凸

POST https://api.openai.com/v1/fine_tuning/jobs

Creates a fine-tuning job which begins the process of creating a new model from a given dataset.

Response includes details of the enqueued job including job status and the name of the fine-tuned models once complete.

Learn more about fine-tuning

Request body

model string Required

The name of the model to fine-tune. You can select one of the **supported models**.

training_file string Required

The ID of an uploaded file that contains training data.

See upload file for how to upload a file.

Your dataset must be formatted as a JSONL file. Additionally, you must upload your file with the purpose fine-tune.

The contents of the file should differ depending on if the model uses the **chat** or **completions** format.

See the fine-tuning guide for more details.

hyperparameters object Optional

The hyperparameters used for the fine-tuning job.

∨Show properties

suffix string or null Optional Defaults to null

A string of up to 18 characters that will be added to your fine-tuned model name.

For example, a suffix of "custom-model-name" would produce a model name like

ft:gpt-3.5-turbo:openai:custom-model-name:7p4IURel .

validation_file string or null Optional

The ID of an uploaded file that contains validation data.

If you provide this file, the data is used to generate validation metrics periodically during fine-tuning. These metrics can be viewed in the fine-tuning results file. The same data should not be present in both train and validation files.

Your dataset must be formatted as a JSONL file. You must upload your file with the purpose fine-tune.

See the fine-tuning guide for more details.

integrations array or null Optional

A list of integrations to enable for your fine-tuning job.

∨Show properties

seed integer or null Optional

```
Default
             Epochs
                           Validation file
                                               W&B Integration
                                                          凸
Example request
                                             python ~
1 from openai import OpenAI
   client = OpenAI()
4 client.fine_tuning.jobs.create(
    training_file="file-abc123",
   model="gpt-3.5-turbo"
7)
                                                          币
Response
1
2
      "object": "fine_tuning.job",
3
     "id": "ftjob-abc123",
     "model": "gpt-3.5-turbo-0125",
4
5
     "created_at": 1614807352,
6
     "fine_tuned_model": null,
```

"organization_id": "org-123",

"training_file": "file-abc123",

"result_files": [],

"status": "queued",

"validation_file": null,

8

9

11

12 }

https://platform.openai.com/docs/api-reference/runs-v1

Example request

16 }

Response

The seed controls the reproducibility of the job. Passing in the same seed and job parameters should produce the same results, but may differ in rare cases. If a seed is not specified, one will be generated for you.

Returns

A fine-tuning.job object.

List fine-tuning jobs ₽

GET https://api.openai.com/v1/fine_tuning/jobs

List your organization's fine-tuning jobs

Query parameters

after string Optional

Identifier for the last job from the previous pagination request.

limit integer Optional Defaults to 20 Number of fine-tuning jobs to retrieve.

Returns

A list of paginated **fine-tuning job** objects.

List fine-tuning events ₽

GET https://api.openai.com/v1/fine_tuning/jobs/{fine_tuning_job_id}/e vents

Get status updates for a fine-tuning job.

Path parameters

fine_tuning_job_id string Required
The ID of the fine-tuning job to get events for.

Query parameters

after string Optional

Identifier for the last event from the previous pagination request.

limit integer Optional Defaults to 20

```
1 from openai import OpenAl
2 client = OpenAI()
4 client.fine_tuning.jobs.list()
                                                             币
Response
1
    {
      "object": "list",
      "data": [
3
4
5
        "object": "fine_tuning.job.event",
        "id": "ft-event-TjX0lMfOniCZX64t9PUQT5hn",
6
7
        "created_at": 1689813489,
        "level": "warn",
8
        "message": "Fine tuning process stopping due to job cance
10
        "data": null,
        "type": "message"
11
12
       },
13
       { ... },
14
      { ... }
15 ], "has_more": true
```

币

python ~

```
python ✓ ☐

from openai import OpenAI

client = OpenAI()

client.fine_tuning.jobs.list_events(

fine_tuning_job_id="ftjob-abc123",

limit=2

job_id="ftjob-abc123",

limit=2
```

```
1 {
2  "object": "list",
3  "data": [
4  {
5   "object": "fine_tuning.job.event",
```

币

Number of events to retrieve.

Returns

A list of fine-tuning event objects.

List fine-tuning checkpoints ∂

GET https://api.openai.com/v1/fine_tuning/jobs/{fine_tuning_job_id}/c heckpoints

List checkpoints for a fine-tuning job.

Path parameters

fine_tuning_job_id string Required

The ID of the fine-tuning job to get checkpoints for.

Query parameters

after string Optional

Identifier for the last checkpoint ID from the previous pagination request.

limit integer Optional Defaults to 10 Number of checkpoints to retrieve.

Returns

A list of fine-tuning checkpoint objects for a fine-tuning job.

Retrieve fine-tuning job@

GET https://api.openai.com/v1/fine_tuning/jobs/{fine_tuning_job_id}

Get info about a fine-tuning job.

Learn more about fine-tuning

Path parameters

```
API Reference - OpenAI API
```

```
"id": "ft-event-ddTJfwuMVpfLXseO0Am0Gqjm",
        "created_at": 1692407401,
8
        "level": "info",
9
        "message": "Fine tuning job successfully completed",
10
        "type": "message"
11
12
      },
13
14
        "object": "fine_tuning.job.event",
15
        "id": "ft-event-tyiGuB72evQncpH87xe505Sv",
        "created_at": 1692407400,
16
        "level": "info",
17
        "message": "New fine-tuned model created: ft:gpt-3.5-tur
18
19
        "tyna". "maccana"
```

Example request

curl∨ 🗇

- 1 curl https://api.openai.com/v1/fine_tuning/jobs/ftjob-abc123/chec
- 2 -H "Authorization: Bearer \$OPENAI_API_KEY"

```
Response
```



```
1
               {
                     "object": "list"
3
                     "data": [
4
                         {
                                "object": "fine_tuning.job.checkpoint",
5
                               "id": "ftckpt_zc4Q7MP6XxulcVzj4MZdwsAB",
6
7
                                "created_at": 1519129973,
                                "fine_tuned_model_checkpoint": "ft:gpt-3.5-turbo-0125:n
8
9
                                "metrics": {
10
                                    "full_valid_loss": 0.134,
11
                                    "full_valid_mean_token_accuracy": 0.874
12
                                "fine_tuning_job_id": "ftjob-abc123",
13
                                "step number": 2000,
14
15
                         },
16
                                 "object": "fine_tuning.job.checkpoint",
17
                                "id": "ftckpt_enQCFmOTGj3syEpYVhBRLTSy",
18
19
                                "created_at": 1519129833,
20
                                "fine_tuned_model_checkpoint": "ft:gpt-3.5-turbo-0125:m
21
                                "metrics": {
22
                                    "full_valid_loss": 0.167,
23
                                     "full_valid_mean_token_accuracy": 0.781
24
                                المراجع المستعدد المال الطالب المستعدد المستعدد
```

python ~



- 1 from openai import OpenAI
- 2 client = OpenAI()

3

4 client.fine_tuning.jobs.retrieve("ftjob-abc123")

fine_tuning_job_id string Required The ID of the fine-tuning job.

Returns

The fine-tuning object with the given ID.

Cancel fine-tuning *₽*

POST https://api.openai.com/v1/fine_tuning/jobs/{fine_tuning_job_i d}/cancel

Immediately cancel a fine-tune job.

Path parameters

fine_tuning_job_id string Required The ID of the fine-tuning job to cancel.

Returns

The cancelled fine-tuning object.

Training format for chat models@

The per-line training example of a fine-tuning input file for chat models

messages array

```
Response
```

```
1
2
     "object": "fine_tuning.job",
3
     "id": "ftjob-abc123",
     "model": "davinci-002",
5
     "created_at": 1692661014,
     "finished_at": 1692661190,
     "fine_tuned_model": "ft:davinci-002:my-org:custom_suffix:7
8
     "organization_id": "org-123",
     "result_files": [
        "file-abc123"
10
11
    1,
     "status": "succeeded",
12
13
     "validation_file": null,
     "training_file": "file-abc123",
14
15
     "hyperparameters": {
16
        "n_epochs": 4,
17
        "batch_size": 1,
18
        "learning_rate_multiplier": 1.0
19
     "trained_tokens": 5768,
20
21
     "integrations": [],
     "seed": 0,
22
23
     "estimated_finish": 0
```

```
Example request
```

```
python ∨ 🗇
```

币

```
    1 from openai import OpenAI
    2 client = OpenAI()
    3
    4 client.fine_tuning.jobs.cancel("ftjob-abc123")
```

Response

1

```
ð
```

```
2
     "object": "fine_tuning.job",
     "id": "ftjob-abc123",
     "model": "gpt-3.5-turbo-0125",
     "created_at": 1689376978,
     "fine_tuned_model": null,
     "organization_id": "org-123",
     "result_files": [],
     "hyperparameters": {
9
10
      "n_epochs": "auto"
11
     }.
     "status": "cancelled",
12
     "validation_file": "file-abc123",
13
     "training_file": "file-abc123"
15 }
```

```
Training format for chat models
```



```
1 {
2 "messages":[
```

```
API Reference - OpenAI API
               { "role": "user", "content": "What is the weather in San Fra
        4
        5
                "role": "assistant",
        6
                "tool_calls": [
        7
                   "id": "call id".
        8
        9
                   "type": "function",
        10
                   "function": {
        11
                    "name": "get_current_weather",
                    "arguments": "{\"location\": \"San Francisco, USA\", \'
        13
                 }
        15
                ]
        16
               }
        17
             1,
        18
              "parallel_tool_calls": false,
        19
              "tools": [
        20
        21
                 "type": "function",
                "function": {
        22
                  "name": "get_current_weather",
        23
                  "description": "Get the current weather",
        24
        25
                  "parameters": {
        26
                   "type": "object",
        27
                   "properties": {
        28
                    "location": {
                       "type": "string",
        29
        30
                       "description": "The city and country, eg. San France
        31
                    "format": { "type": "string", "enum": ["celsius", "fahren
        32
        33
                   "required": ["location", "format"]
        34
        35
                  }
        36
                }
        37
               }
        38
             ]
        39 }
```

Training format for completions models@

The per-line training example of a fine-tuning input file for completions models

prompt string

The input prompt for this training example.

completion string

The desired completion for this training example.

The fine-tuning job object@

The fine_tuning.job object represents a fine-tuning job that has been created through the API.

id string

The object identifier, which can be referenced in the API endpoints.

created_at integer

```
Training format for completions models

1 {
2  "prompt": "What is the answer to 2+2",
3  "completion": "4"
4 }
```

```
1  {
2    "object": "fine_tuning.job",
3    "id": "ftjob-abc123",
4    "model": "davinci-002",
5    "created_at": 1692661014,
6    "finished_at": 1692661190,
7    "fine_tuned_model": "ft:davinci-002:my-org:custom_suffix:7
```

The fine-tuning job object

币

The Unix timestamp (in seconds) for when the fine-tuning job was created.

error object or null

For fine-tuning jobs that have failed, this will contain more information on the cause of the failure.

√Show properties

fine_tuned_model string or null

The name of the fine-tuned model that is being created. The value will be null if the fine-tuning job is still running.

finished_at integer or null

The Unix timestamp (in seconds) for when the fine-tuning job was finished. The value will be null if the fine-tuning job is still running.

hyperparameters object

The hyperparameters used for the fine-tuning job. See the **fine-tuning guide** for more details.

√Show properties

model string

The base model that is being fine-tuned.

object string

The object type, which is always "fine_tuning.job".

organization id string

The organization that owns the fine-tuning job.

result_files array

The compiled results file ID(s) for the fine-tuning job. You can retrieve the results with the **Files API**.

status string

The current status of the fine-tuning job, which can be either validating_files, queued, running, succeeded, failed, or cancelled.

trained_tokens integer or null

The total number of billable tokens processed by this fine-tuning job. The value will be null if the fine-tuning job is still running.

training file string

The file ID used for training. You can retrieve the training data with the **Files API**.

validation_file string or null

The file ID used for validation. You can retrieve the validation results with the **Files API**.

integrations array or null

A list of integrations to enable for this fine-tuning job.

∨Show possible types

```
"organization_id": "org-123",
     "result_files": [
10
        "file-abc123"
11
     ],
     "status": "succeeded",
12
13
      "validation file": null.
     "training_file": "file-abc123",
14
15
     "hyperparameters": {
16
        "n_epochs": 4,
17
        "batch_size": 1,
        "learning_rate_multiplier": 1.0
18
19
20
     "trained_tokens": 5768,
21
      "integrations": [],
22
     "seed": 0,
23
     "estimated_finish": 0
24 }
```

seed integer

The seed used for the fine-tuning job.

estimated_finish integer or null

The Unix timestamp (in seconds) for when the fine-tuning job is estimated to finish. The value will be null if the fine-tuning job is not running.

The fine-tuning job event object ∂

Fine-tuning job event object

created_at integer
level string

id string

object string

message string

```
The fine-tuning job event object

1 {
2  "object": "fine_tuning.job.event",
3  "id": "ftevent-abc123"
4  "created_at": 1677610602,
5  "level": "info",
6  "message": "Created fine-tuning job"
7 }
```

The fine-tuning job checkpoint object *₽*

The [fine_tuning.job.checkpoint] object represents a model checkpoint for a fine-tuning job that is ready to use.

id string

The checkpoint identifier, which can be referenced in the API endpoints.

created_at integer

The Unix timestamp (in seconds) for when the checkpoint was created.

fine_tuned_model_checkpoint string

The name of the fine-tuned checkpoint model that is created.

step_number integer

The step number that the checkpoint was created at.

metrics object

Metrics at the step number during the fine-tuning job.

∨Show properties

fine_tuning_job_id string

The name of the fine-tuning job that this checkpoint was created from.

object string

The object type, which is always "fine_tuning.job.checkpoint".

Batch *₽*

Create large batches of API requests for asynchronous processing. The Batch API returns completions within 24 hours for a 50% discount.

The fine-tuning job checkpoint object

```
O
```

```
1
2
     "object": "fine_tuning.job.checkpoint",
3
     "id": "ftckpt_qtZ5Gyk4BLq1SfLFWp3RtO3P",
     "created_at": 1712211699,
     "fine_tuned_model_checkpoint": "ft:gpt-3.5-turbo-0125:my-org
5
6
     "fine_tuning_job_id": "ftjob-fpbNQ3H1GrMehXRf8cO97xTN",
     "metrics": {
8
      "step": 88,
      "train_loss": 0.478,
10
      "train_mean_token_accuracy": 0.924,
11
      "valid_loss": 10.112,
12
      "valid_mean_token_accuracy": 0.145,
      "full_valid_loss": 0.567,
      "full_valid_mean_token_accuracy": 0.944
14
15 },
     "step_number": 88
16
17 }
```

Related guide: Batch

POST https://api.openai.com/v1/batches

Creates and executes a batch from an uploaded file of requests

Request body

input_file_id string Required

The ID of an uploaded file that contains requests for the new batch.

See upload file for how to upload a file.

Your input file must be formatted as a JSONL file, and must be uploaded with the purpose | batch |. The file can contain up to 50,000 requests, and can be up to 100 MB in size.

endpoint string Required

The endpoint to be used for all requests in the batch. Currently /v1/chat/completions , /v1/embeddings , and /v1/completions are supported. Note that /v1/embeddings batches are also restricted to a maximum of 50,000 embedding inputs across all requests in the batch.

completion_window string Required

The time frame within which the batch should be processed. Currently only 24h is supported.

metadata object or null Optional

Optional custom metadata for the batch.

Returns

The created Batch object.

Retrieve batch@

GET https://api.openai.com/v1/batches/{batch_id}

Retrieves a batch.

Path parameters

batch id string Required The ID of the batch to retrieve.

Returns

The **Batch** object matching the specified ID.

Example request

1 from openai import OpenAI

2 client = OpenAI()

4 client.batches.create(

5 input_file_id="file-abc123",

endpoint="/v1/chat/completions", 6

7 completion_window="24h"

8)

Response

1 { 2 "id": "batch_abc123",

3 "object": "batch", "endpoint": "/v1/chat/completions", 4

5 "errors": null.

"input_file_id": "file-abc123", 6

7 "completion_window": "24h",

"status": "validating", 8

9 "output_file_id": null,

10 "error_file_id": null,

"created_at": 1711471533, 11

12 "in_progress_at": null,

13 "expires_at": null,

14 "finalizing_at": null,

15 "completed_at": null,

16 "failed_at": null,

"expired_at": null, 17

18 "cancelling_at": null,

19 "cancelled_at": null,

"request_counts": {

Example request

python ~

D

凸

python ~

币

币

1 from openai import OpenAI

2 client = OpenAI()

4 client.batches.retrieve("batch_abc123")

Response

1

2 "id": "batch_abc123",

"object": "batch", "endpoint": "/v1/completions", 4

5 "errors": null,

"input_file_id": "file-abc123", 6

"completion_window": "24h",

```
"status": "completed",
     "output_file_id": "file-cvaTdG",
10 "error_file_id": "file-HOWS94",
11
     "created_at": 1711471533,
12
     "in_progress_at": 1711471538,
13
     "expires_at": 1711557933,
14
     "finalizing_at": 1711493133,
15
     "completed_at": 1711493163,
16
     "failed_at": null,
17
     "expired_at": null,
     "cancelling_at": null,
18
     "cancelled_at": null,
19
20
     "request_counts": {
21
      "total": 100,
22
      "completed": 95,
23
      "failed": 5
```

python ~

仚

币

POST https://api.openai.com/v1/batches/{batch_id}/cancel

Cancels an in-progress batch. The batch will be in status cancelling for up to 10 minutes, before changing to cancelled, where it will have partial results (if any) available in the output file.

Path parameters

batch_id string Required The ID of the batch to cancel.

Returns

The **Batch** object matching the specified ID.

List batch@

GET https://api.openai.com/v1/batches

List your organization's batches.

Query parameters

```
Example request
```

1 from openai import OpenAI

2 client = OpenAI()

3

4 client.batches.cancel("batch_abc123")

Response

1 "id": "batch_abc123", 2 "object": "batch", 3

"endpoint": "/v1/chat/completions", 4

5 "errors": null,

"input_file_id": "file-abc123", 6

"completion_window": "24h",

"status": "cancelling", 8

"output_file_id": null,

10 "error_file_id": null,

"created_at": 1711471533, 11

"in_progress_at": 1711471538, 12

13 "expires_at": 1711557933,

"finalizing_at": null, 14

15 "completed_at": null,

16 "failed_at": null,

17 "expired_at": null,

"cancelling at": 1711475133, 18

19 "cancelled_at": null,

20 "request_counts": {

"total": 100, 21

22 "completed": 23,

23 "failed": 1

24 },

Example request

1 from openai import OpenAI

2 client = OpenAI()

4 client.batches.list()

凸

python ✓

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

Returns

A list of paginated Batch objects.

The batch object ₽

id string

object string

The object type, which is always batch.

endpoint string

The OpenAl API endpoint used by the batch.

errors object

∨Show properties

input_file_id string

The ID of the input file for the batch.

completion_window string

The time frame within which the batch should be processed.

status string

The current status of the batch.

output_file_id string

The ID of the file containing the outputs of successfully executed requests.

error_file_id string

The ID of the file containing the outputs of requests with errors.

created_at integer

```
Response
```

```
1
2
     "object": "list",
3
     "data": [
4
5
        "id": "batch_abc123",
6
        "object": "batch",
        "endpoint": "/v1/chat/completions",
8
        "errors": null,
9
        "input_file_id": "file-abc123",
        "completion_window": "24h",
10
11
        "status": "completed",
        "output_file_id": "file-cvaTdG",
12
        "error_file_id": "file-HOWS94",
13
14
        "created_at": 1711471533,
15
        "in_progress_at": 1711471538,
16
        "expires_at": 1711557933,
17
        "finalizing_at": 1711493133,
18
        "completed_at": 1711493163,
19
        "failed_at": null,
20
        "expired_at": null,
21
        "cancelling_at": null,
22
        "cancelled_at": null,
23
        "request_counts": {
```

仚

ð

The batch object

"total": 100,

24

29 }

```
1
    {
2
     "id": "batch_abc123",
3
     "object": "batch",
     "endpoint": "/v1/completions",
4
     "errors": null,
     "input_file_id": "file-abc123",
6
7
     "completion_window": "24h",
8
     "status": "completed",
     "output_file_id": "file-cvaTdG",
     "error_file_id": "file-HOWS94",
10
11
     "created_at": 1711471533,
12
     "in_progress_at": 1711471538,
13
     "expires_at": 1711557933,
14
     "finalizing_at": 1711493133,
15
     "completed_at": 1711493163,
16
     "failed at": null,
17
     "expired_at": null,
18
     "cancelling_at": null,
19
     "cancelled_at": null,
20
     "request_counts": {
21
       "total": 100,
22
       "completed": 95,
23
       "failed": 5
24
     },
25
     "metadata": {
26
       "customer_id": "user_123456789",
27
       "batch_description": "Nightly eval job",
28
   }
```

The Unix timestamp (in seconds) for when the batch was created.

in_progress_at integer

The Unix timestamp (in seconds) for when the batch started processing.

expires at integer

The Unix timestamp (in seconds) for when the batch will expire.

finalizing_at integer

The Unix timestamp (in seconds) for when the batch started finalizing.

completed_at integer

The Unix timestamp (in seconds) for when the batch was completed.

failed_at integer

The Unix timestamp (in seconds) for when the batch failed.

expired_at integer

The Unix timestamp (in seconds) for when the batch expired.

cancelling at integer

The Unix timestamp (in seconds) for when the batch started cancelling.

cancelled_at integer

The Unix timestamp (in seconds) for when the batch was cancelled.

request_counts object

The request counts for different statuses within the batch.

∨Show properties

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

The request input object@

The per-line object of the batch input file

custom_id string

A developer-provided per-request id that will be used to match outputs to inputs. Must be unique for each request in a batch.

method string

The HTTP method to be used for the request. Currently only POST is supported.

url string

The OpenAl API relative URL to be used for the request. Currently /v1/chat/completions , /v1/embeddings , and /v1/completions are

The request input object

D

{"custom_id": "request-1", "method": "POST", "url": "/v1/chat/comple

The request output object@

The per-line object of the batch output and error files

id string

custom_id string

A developer-provided per-request id that will be used to match outputs to inputs.

response object or null

∨Show properties

error object or null

For requests that failed with a non-HTTP error, this will contain more information on the cause of the failure.

√Show properties

Files *₽*

Files are used to upload documents that can be used with features like Assistants, Fine-tuning, and Batch API.

Upload file ₽

POST https://api.openai.com/v1/files

Upload a file that can be used across various endpoints. Individual files can be up to 512 MB, and the size of all files uploaded by one organization can be up to 100 GB.

The Assistants API supports files up to 2 million tokens and of specific file types. See the **Assistants Tools guide** for details.

The Fine-tuning API only supports .jsonl files. The input also has certain required formats for fine-tuning **chat** or **completions** models.

The Batch API only supports .jsonl files up to 100 MB in size. The input also has a specific required **format**.

Please **contact us** if you need to increase these storage limits.

Request body

file file Required

The File object (not file name) to be uploaded.

purpose string Required

The intended purpose of the uploaded file.

Use "assistants" for **Assistants** and **Message** files, "vision" for Assistants image file inputs, "batch" for **Batch API**, and "fine-tune" for **Fine-tuning**.

```
The request output object

{"id": "batch_req_wnaDys", "custom_id": "request-2", "response": {"
```

```
Example request python ✓ ☐

1 from openai import OpenAI

2 client = OpenAI()

3

4 client.files.create(

5 file=open("mydata.jsonl", "rb"),

6 purpose="fine-tune"

7 )

Response

1 {
2 "id": "file-abc123",

3 "object": "file",

4 "bytes": 120000,
```

"created_at": 1677610602,

"filename": "mydata.jsonl",

"purpose": "fine-tune",

8 }

Response

Returns

The uploaded File object.

List files@

GET https://api.openai.com/v1/files

Returns a list of files that belong to the user's organization.

Query parameters

purpose string Optional

Only return files with the given purpose.

Returns

A list of File objects.

Retrieve file

GET https://api.openai.com/v1/files/{file_id}

Returns information about a specific file.

Path parameters

file_id string Required

The ID of the file to use for this request.

Returns

The File object matching the specified ID.

```
Example request python ✓ ☐

1 from openai import OpenAl
2 client = OpenAl()
3
4 client.files.list()
```

O

D

凸

python ~

```
1
2
     "data": [
3
      {
4
        "id": "file-abc123",
        "object": "file",
5
6
        "bytes": 175,
7
        "created_at": 1613677385,
        "filename": "salesOverview.pdf",
8
9
        "purpose": "assistants",
10
      },
11
      {
12
        "id": "file-abc123",
13
        "object": "file",
        "bytes": 140,
14
        "created_at": 1613779121,
15
16
        "filename": "puppy.jsonl",
        "purpose": "fine-tune",
17
18
      }
19
20
     "object": "list"
21 }
```

```
1 from openai import OpenAI
2 client = OpenAI()
3
```

4 client.files.retrieve("file-abc123")

Example request

Response

```
1 {
2  "id": "file-abc123",
3  "object": "file",
4  "bytes": 120000,
5  "created_at": 1677610602,
6  "filename": "mydata.jsonl",
7  "purpose": "fine-tune",
8 }
```

Delete file∂

DELETE https://api.openai.com/v1/files/{file_id}

Delete a file.

Path parameters

file_id string Required

The ID of the file to use for this request.

Returns

Deletion status.

Retrieve file contents

GET https://api.openai.com/v1/files/{file_id}/content

Returns the contents of the specified file.

Path parameters

file_id string Required

The ID of the file to use for this request.

Returns

The file content.

The file object ∂

The File object represents a document that has been uploaded to OpenAI.

id string

The file identifier, which can be referenced in the API endpoints.

bytes integer

The size of the file, in bytes.

created_at integer

The Unix timestamp (in seconds) for when the file was created.

filename string

The name of the file.

object string

The object type, which is always file .

purpose string

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 client.files.delete("file-abc123")

Response

1 {
2 "id": "file-abc123",
3 "object": "file",
4 "deleted": true
5 }
```

1 from openai import OpenAI
2 client = OpenAI()
3

4 content = client.files.content("file-abc123")

 \Box

币

python ~

The file object

Example request

```
1 {
2    "id": "file-abc123",
3    "object": "file",
4    "bytes": 120000,
5    "created_at": 1677610602,
6    "filename": "salesOverview.pdf",
7    "purpose": "assistants",
8 }
```

```
The intended purpose of the file. Supported values are assistants, assistants_output, batch, batch_output, fine-tune, fine-tune-results and vision.

status Deprecated string
Deprecated. The current status of the file, which can be either uploaded, processed, or error.

status_details Deprecated string
Deprecated. For details on why a fine-tuning training file failed validation, see the error field on fine_tuning.job.
```

Uploads ∂

Allows you to upload large files in multiple parts.

Create upload ∂

POST https://api.openai.com/v1/uploads

Creates an intermediate **Upload** object that you can add **Parts** to. Currently, an Upload can accept at most 8 GB in total and expires after an hour after you create it.

Once you complete the Upload, we will create a **File** object that contains all the parts you uploaded. This File is usable in the rest of our platform as a regular File object.

For certain purpose s, the correct mime_type must be specified. Please refer to documentation for the supported MIME types for your use case:

Assistants

For guidance on the proper filename extensions for each purpose, please follow the documentation on **creating a File**.

Request body

filename string Required The name of the file to upload.

purpose string Required

The intended purpose of the uploaded file.

See the documentation on File purposes.

bytes integer Required

The number of bytes in the file you are uploading.

mime_type string Required
The MIME type of the file.

This must fall within the supported MIME types for your file purpose. See the supported MIME types for assistants and vision.

```
Example request
                                                curl ∨
                                                          币
1 curl https://api.openai.com/v1/uploads \
    -H "Authorization: Bearer $OPENAI_API_KEY" \
3
   -d '{
     "purpose": "fine-tune",
     "filename": "training_examples.jsonl",
     "bytes": 2147483648,
     "mime_type": "text/jsonl"
7
                                                           凸
Response
1
2
     "id": "upload_abc123",
     "object": "upload",
3
     "bytes": 2147483648,
     "created_at": 1719184911,
     "filename": "training_examples.jsonl",
     "purpose": "fine-tune",
     "status": "pending",
     "expires_at": 1719127296
10 }
```

Returns

The **Upload** object with status pending.

Add upload part@

POST https://api.openai.com/v1/uploads/{upload_id}/parts

Adds a **Part** to an **Upload** object. A Part represents a chunk of bytes from the file you are trying to upload.

Each Part can be at most 64 MB, and you can add Parts until you hit the Upload maximum of 8 GB.

It is possible to add multiple Parts in parallel. You can decide the intended order of the Parts when you **complete the Upload**.

Path parameters

upload_id string Required
The ID of the Upload.

Request body

data file Required

The chunk of bytes for this Part.

Returns

The upload Part object.

Complete upload ₽

POST https://api.openai.com/v1/uploads/{upload_id}/complete

Completes the Upload.

Within the returned Upload object, there is a nested **File** object that is ready to use in the rest of the platform.

You can specify the order of the Parts by passing in an ordered list of the Part IDs.

The number of bytes uploaded upon completion must match the number of bytes initially specified when creating the Upload object. No Parts may be added after an Upload is completed.

Path parameters

upload_id string Required

Example request 1 curl https://api.openai.com/v1/uploads/upload_abc123/parts 2 -F data="aHR0cHM6Ly9hcGkub3BlbmFpLmNvbS92MS91cGx Response

"object": "upload.part",
"created_at": 1719185911,

"upload_id": "upload_abc123"

5

6 }

Response

```
Example request curl Curl Curl Curl Curl Curl Curl https://api.openai.com/v1/uploads/upload_abc123/complete
2 -d '{
3 "part_ids": ["part_def456", "part_ghi789"]
4 }'
```

```
1 {
2  "id": "upload_abc123",
3  "object": "upload",
4  "bytes": 2147483648,
5  "created_at": 1719184911,
6  "filename": "training_examples.jsonl",
7  "purpose": "fine-tune",
8  "status": "completed",
9  "expires_at": 1719127296,
10  "file": {
```

币

The ID of the Upload.

Request body

part_ids array Required
The ordered list of Part IDs.

md5 string Optional

The optional md5 checksum for the file contents to verify if the bytes uploaded matches what you expect.

Returns

The **Upload** object with status completed with an additional file property containing the created usable File object.

Cancel upload ∂

POST https://api.openai.com/v1/uploads/{upload_id}/cancel

Cancels the Upload. No Parts may be added after an Upload is cancelled.

Path parameters

upload_id string RequiredThe ID of the Upload.

Returns

The **Upload** object with status cancelled.

The upload object

The Upload object can accept byte chunks in the form of Parts.

id string

The Upload unique identifier, which can be referenced in API endpoints.

created at integer

The Unix timestamp (in seconds) for when the Upload was created.

filename string

The name of the file to be uploaded.

bytes integer

The intended number of bytes to be uploaded.

```
Example request
```

API Reference - OpenAI API

12

14

15

16

17 }

18 }

"id": "file-xyz321",

"bytes": 2147483648,

"purpose": "fine-tune",

"created_at": 1719186911,

"filename": "training_examples.jsonl",

"object": "file",

curl ~

凸

curl https://api.openai.com/v1/uploads/upload_abc123/cancel

Response



The upload object



```
1
2
     "id": "upload_abc123",
     "object": "upload",
3
     "bytes": 2147483648,
5
     "created_at": 1719184911,
     "filename": "training_examples.jsonl",
7
     "purpose": "fine-tune",
8
     "status": "completed",
9
     "expires_at": 1719127296,
10
11
      "id": "file-xyz321",
12
      "object": "file",
13
      "bytes": 2147483648,
14
      "created_at": 1719186911,
15
      "filename": "training_examples.jsonl",
      "purpose": "fine-tune",
```

18 }

purpose string

The intended purpose of the file. **Please refer here** for acceptable values.

status string

The status of the Upload.

expires_at integer

The Unix timestamp (in seconds) for when the Upload was created.

object string

The object type, which is always "upload".

file

The File object represents a document that has been uploaded to OpenAI.

The upload part object@

The upload Part represents a chunk of bytes we can add to an Upload object.

id string

The upload Part unique identifier, which can be referenced in API endpoints.

created_at integer

The Unix timestamp (in seconds) for when the Part was created.

upload_id string

The ID of the Upload object that this Part was added to.

object string

The object type, which is always upload.part .

Images∂

Given a prompt and/or an input image, the model will generate a new image.

Related guide: Image generation

Create image *₽*

POST https://api.openai.com/v1/images/generations

Creates an image given a prompt.

Request body

prompt string Required

A text description of the desired image(s). The maximum length is 1000 characters for dall-e-2 and 4000 characters for dall-e-3.

model string Optional Defaults to dall-e-2

The model to use for image generation.

The upload part object

1 {
2 "id": "part_def456",
3 "object": "upload.part",
4 "created_at": 1719186911,
5 "upload_id": "upload_abc123"
6 }

Example request

python ~



O

1 from openai import OpenAI
2 client = OpenAI()
3
4 client.images.generate(
5 model="dall-e-3",
6 prompt="A cute baby sea otter",
7 n=1,
8 size="1024x1024"
9)

Response



```
\begin{array}{ll} \textbf{n} & \text{integer or null} & \text{Optional} & \text{Defaults to 1} \\ \\ \text{The number of images to generate. Must be between 1 and 10.} \\ \\ \text{For } \left( \begin{array}{ll} \text{dall-e-3} \end{array} \right), \text{ only } \left( \begin{array}{ll} \text{n=1} \end{array} \right) \text{ is supported.} \end{array}
```

quality string Optional Defaults to standard

The quality of the image that will be generated. Indicreates images with finer details and greater consistency across the image. This param is only supported for Idall-e-3.

response_format string or null Optional Defaults to url
The format in which the generated images are returned. Must be
one of url or b64_json. URLs are only valid for 60 minutes after
the image has been generated.

```
size string or null Optional Defaults to 1024x1024

The size of the generated images. Must be one of 256x256,

512x512, or 1024x1024 for dall-e-2. Must be one of 1024x1024, 1792x1024, or 1024x1792 for dall-e-3 models.
```

style string or null Optional Defaults to vivid

The style of the generated images. Must be one of vivid or natural. Vivid causes the model to lean towards generating hyper-real and dramatic images. Natural causes the model to produce more natural, less hyper-real looking images. This param is only supported for dall-e-3.

user string Optional

A unique identifier representing your end-user, which can help OpenAl to monitor and detect abuse. **Learn more**.

Returns

Returns a list of image objects.

Create image edit ∂

POST https://api.openai.com/v1/images/edits

Creates an edited or extended image given an original image and a prompt.

Request body

image file Required

The image to edit. Must be a valid PNG file, less than 4MB, and square. If mask is not provided, image must have transparency, which will be used as the mask.

prompt string Required

A text description of the desired image(s). The maximum length is 1000 characters.

```
mask file Optional
```

```
1
      "created": 1589478378,
2
3
     "data": [
4
        "url": "https://..."
5
6
       },
7
       {
        "url": "https://..."
8
10
    ]
11 }
```

```
python ~
Example request
                                                         币
    from openai import OpenAI
2
    client = OpenAI()
3
4
   client.images.edit(
     image=open("otter.png", "rb"),
5
     mask=open("mask.png", "rb"),
     prompt="A cute baby sea otter wearing a beret",
9
     size="1024x1024"
10 )
```

Response

凸

An additional image whose fully transparent areas (e.g. where alpha is zero) indicate where image should be edited. Must be a valid PNG file, less than 4MB, and have the same dimensions as image.

model string Optional Defaults to dall-e-2
The model to use for image generation. Only dall-e-2 is supported at this time.

n integer or null Optional Defaults to 1The number of images to generate. Must be between 1 and 10.

size string or null Optional Defaults to 1024x1024

The size of the generated images. Must be one of 256x256, 512x512, or 1024x1024.

response_format string or null Optional Defaults to url
The format in which the generated images are returned. Must be
one of url or b64_json. URLs are only valid for 60 minutes after
the image has been generated.

user string Optional

A unique identifier representing your end-user, which can help OpenAl to monitor and detect abuse. **Learn more**.

Returns

Returns a list of image objects.

Create image variation@

POST https://api.openai.com/v1/images/variations

Creates a variation of a given image.

Request body

image file Required

The image to use as the basis for the variation(s). Must be a valid PNG file, less than 4MB, and square.

model string Optional Defaults to dall-e-2
The model to use for image generation. Only dall-e-2 is supported at this time.

n integer or null Optional Defaults to 1
 The number of images to generate. Must be between 1 and 10.
 For dall-e-3, only n=1 is supported.

response_format string or null Optional Defaults to url
The format in which the generated images are returned. Must be
one of url or b64_json. URLs are only valid for 60 minutes after
the image has been generated.

```
4 {
5 "url": "https://..."
6 },
7 {
8 "url": "https://..."
9 }
10 ]
11 }
```

API Reference - OpenAI API

```
Example request
                                             python ~
                                                           O
1 from openai import OpenAl
2 client = OpenAI()
3
   response = client.images.create_variation(
    image=open("image_edit_original.png", "rb"),
    n=2.
    size="1024x1024"
8)
                                                           凸
Response
2
      "created": 1589478378.
3
     "data": [
5
        "url": "https://..."
      }.
7
        "url": "https://..."
9
      }
10
11 }
```

```
size string or null Optional Defaults to 1024x1024

The size of the generated images. Must be one of 256x256, 512x512, or 1024x1024.
```

user string Optional

A unique identifier representing your end-user, which can help OpenAI to monitor and detect abuse. **Learn more**.

Returns

Returns a list of image objects.

The image object *₽*

Represents the url or the content of an image generated by the OpenAl API.

b64_json string

```
The base64-encoded JSON of the generated image, if response_format is b64_json.
```

url string

The URL of the generated image, if response_format is url (default).

revised_prompt string

The prompt that was used to generate the image, if there was any revision to the prompt.

Models *₽*

List and describe the various models available in the API. You can refer to the **Models** documentation to understand what models are available and the differences between them.

List models@

GET https://api.openai.com/v1/models

Lists the currently available models, and provides basic information about each one such as the owner and availability.

Returns

A list of model objects.

```
The image object

1 {
2 "url": "...",
3 "revised_prompt": "..."
4 }
```

```
python > 

1 from openai import OpenAI

2 client = OpenAI()

3

4 client.models.list()
```

```
币
Response
      "object": "list",
2
3
     "data": [
        "id": "model-id-0",
5
        "object": "model",
7
        "created": 1686935002,
8
        "owned_by": "organization-owner"
9
10
        "id": "model-id-1",
```

```
12
        "object": "model",
        "created": 1686935002,
13
        "owned_by": "organization-owner",
14
15
      },
16
      {
        "id": "model-id-2".
17
       "object": "model",
18
        "created": 1686935002,
19
        "owned_by": "openai"
20
21
22 ],
23
     "object": "list"
```

Retrieve model@

GET https://api.openai.com/v1/models/{model}

Retrieves a model instance, providing basic information about the model such as the owner and permissioning.

Path parameters

model string Required

The ID of the model to use for this request

Returns

The model object matching the specified ID.

Delete a fine-tuned model

DELETE https://api.openai.com/v1/models/{model}

Delete a fine-tuned model. You must have the Owner role in your organization to delete a model.

Path parameters

model string Required The model to delete

Returns

Deletion status.

The model object@

Describes an OpenAI model offering that can be used with the API.

id string

The model identifier, which can be referenced in the API endpoints.

created integer

```
Example req... gpt-3.5-turbo-instruct ✓ python ✓
```

凸

ð

D

币

币

```
1 from openai import OpenAI
```

- 2 client = OpenAI()
- 3

Response

4 client.models.retrieve("gpt-3.5-turbo-instruct")

```
1 {
2 "id": "gpt-3.5-turbo-instruct",
3 "object": "model",
4 "created": 1686935002,
5 "owned_by": "openai"
6 }
```

gpt-3.5-turbo-instruct ∨

Example request python ∨

- 1 from openai import OpenAI
- 2 client = OpenAI()
- 3

Response

5 }

4 client.models.delete("ft:gpt-3.5-turbo:acemeco:suffix:abc123")

```
1 {
2 "id": "ft:gpt-3.5-turbo:acemeco:suffix:abc123",
3 "object": "model",
4 "deleted": true
```

The model object

```
1 {
2  "id": "davinci",
3  "object": "model",
4  "created": 1686935002,
5  "owned_by": "openai"
6 }
```

The Unix timestamp (in seconds) when the model was created.

object string

The object type, which is always "model".

owned_by string

The organization that owns the model.

Moderations *₽*

Given some input text, outputs if the model classifies it as potentially harmful across several categories.

Related guide: Moderations

Create moderation &

POST https://api.openai.com/v1/moderations

Classifies if text is potentially harmful.

Request body

input string or array Required The input text to classify

model string Optional Defaults to text-moderation-latest Two content moderations models are available:

text-moderation-stable and text-moderation-latest.

The default is text-moderation-latest which will be automatically upgraded over time. This ensures you are always using our most accurate model. If you use text-moderation-stable, we will provide advanced notice before updating the model. Accuracy of text-moderation-stable may be slightly lower than for text-moderation-latest .

Returns

A moderation object.

The moderation object@

Represents if a given text input is potentially harmful.

id string

The unique identifier for the moderation request.

model string

The model used to generate the moderation results.

results array

```
Example request python > 

1 from openai import OpenAI
2 client = OpenAI()
3
4 moderation = client.moderations.create(input="I want to kill them
5 print(moderation)
```

```
凸
Response
1
    {
      "id": "modr-XXXXX",
      "model": "text-moderation-005",
3
4
      "results": [
5
       {
6
         "flagged": true,
7
         "categories": {
          "sexual": false,
8
9
          "hate": false,
10
          "harassment": false,
          "self-harm": false,
11
12
          "sexual/minors": false,
13
          "hate/threatening": false,
14
          "violence/graphic": false,
15
          "self-harm/intent": false,
          "self-harm/instructions": false,
16
17
          "harassment/threatening": true,
18
          "violence": true,
19
20
         "category_scores": {
21
          "sexual": 1.2282071e-06,
22
          "hate": 0.010696256,
```

```
The moderation object

1 {
2 "id": "modr-XXXXX",
3 "model": "text-moderation-005",
4 "results": [
5 {
6 "flagged": true,
7 "categories": {
```

API Reference - OpenAI API

```
A list of moderation objects
```

```
∨ Show properties
```

```
"sexual": false,
         "hate": false,
9
10
         "harassment": false,
11
         "self-harm": false,
12
         "sexual/minors": false,
13
         "hate/threatening": false,
14
         "violence/graphic": false,
15
         "self-harm/intent": false,
16
         "self-harm/instructions": false,
17
         "harassment/threatening": true,
         "violence": true,
18
19
20
        "category_scores": {
21
         "sexual": 1.2282071e-06,
22
         "hate": 0.010696256,
23
         "harassment": 0.29842457,
24
         "self-harm": 1.5236925e-08,
25
         "sexual/minors": 5.7246268e-08,
26
         "hate/threatening": 0.0060676364,
27
         "violence/graphic": 4.435014e-06,
         "self-harm/intent": 8.098441e-10,
28
29
         "self-harm/instructions": 2.8498655e-11,
30
         "harassment/threatening": 0.63055265,
31
         "violence": 0.99011886,
32
33
      }
34
35 }
```

Assistants Beta &

Build assistants that can call models and use tools to perform tasks.

Get started with the Assistants API

Create assistant Beta &

POST https://api.openai.com/v1/assistants

Create an assistant with a model and instructions.

Request body

model string Required

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them.

name string or null Optional

The name of the assistant. The maximum length is 256 characters.

description string or null Optional

The description of the assistant. The maximum length is 512 characters.

instructions string or null Optional

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array Optional Defaults to []

```
Example request
```

python ∨ 🗇

Code Interpreter

Files

```
from openal import OpenAl
client = OpenAl()

my_assistant = client.beta.assistants.create(
instructions="You are a personal math tutor. When asked a cname="Math Tutor",
tools=[{"type": "code_interpreter"}],
model="gpt-4-turbo",
)
print(my_assistant)
```

Response



```
1 {
2  "id": "asst_abc123",
3  "object": "assistant",
4  "created_at": 1698984975,
5  "name": "Math Tutor",
6  "description": null,
7  "model": "gpt-4-turbo",
```

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, file_search, or function.

tool_resources object or null Optional

A set of resources that are used by the assistant's tools. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

> Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

response_format string or object Optional Specifies the format that the model must output. Compatible with GPT-4o, GPT-4 Turbo, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to { "type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

An assistant object.

List assistants Beta &

GET https://api.openai.com/v1/assistants

Returns a list of assistants.

```
8  "instructions": "You are a personal math tutor. When asked
9  "tools": [
10     {
11         "type": "code_interpreter"
12     }
13     ],
14  "metadata": {},
```

Example request

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

Returns

A list of assistant objects.

Retrieve assistant Beta &

GET https://api.openai.com/v1/assistants/{assistant_id}

Retrieves an assistant.

Path parameters

assistant_id string Required

The ID of the assistant to retrieve.

Returns

The assistant object matching the specified ID.

```
1 from openai import OpenAl
2 client = OpenAI()
3
4 my_assistants = client.beta.assistants.list(
5
     order="desc",
6
     limit="20",
7)
8 print(my_assistants.data)
                                                              合
Response
1
    {
2
      "object": "list",
      "data": [
4
5
        "id": "asst_abc123",
6
        "object": "assistant",
7
        "created_at": 1698982736,
8
        "name": "Coding Tutor",
        "description": null,
9
10
        "model": "gpt-4-turbo",
        "instructions": "You are a helpful assistant designed to m
11
12
        "tools": [],
        "tool_resources": {},
13
14
        "metadata": {},
15
        "top_p": 1.0,
16
        "temperature": 1.0,
```

Example request

},

{

"response_format": "auto"

17

18

19

python ~

币

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 my_assistant = client.beta.assistants.retrieve("asst_abc123")
5 print(my_assistant)
```

```
Response
```

凸

```
1
   {
2
     "id": "asst_abc123",
3
     "object": "assistant",
     "created_at": 1699009709,
4
     "name": "HR Helper",
5
6
     "description": null,
     "model": "gpt-4-turbo",
8
     "instructions": "You are an HR bot, and you have access to
9
     "tools": [
10
        "type": "file_search"
11
12
      }
13
     "metadata": {},
14
15
     "top_p": 1.0,
```

```
16 "temperature": 1.0,17 "response_format": "auto"18 }
```

Modify assistant Beta €

POST https://api.openai.com/v1/assistants/{assistant_id}

Modifies an assistant.

Path parameters

assistant_id string Required
The ID of the assistant to modify.

Request body

model Optional

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them.

name string or null Optional

The name of the assistant. The maximum length is 256 characters.

description string or null Optional

The description of the assistant. The maximum length is 512 characters.

instructions string or null Optional

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array Optional Defaults to []

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, file_search, or function.

∨Show possible types

tool_resources object or null Optional

A set of resources that are used by the assistant's tools. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

∨Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower

Example request

python ~



```
from openai import OpenAl
client = OpenAl()

my_updated_assistant = client.beta.assistants.update(
    "asst_abc123",
    instructions="You are an HR bot, and you have access to files
    name="HR Helper",
    tools=[{"type": "file_search"}],
    model="gpt-4-turbo"
)

print(my_updated_assistant)
```

Response



```
1
   {
2
     "id": "asst_123",
3
     "object": "assistant",
     "created_at": 1699009709,
     "name": "HR Helper",
6
     "description": null,
     "model": "gpt-4-turbo",
7
8
     "instructions": "You are an HR bot, and you have access to
9
     "tools": [
10
     {
11
        "type": "file_search"
12
      }
13
14
     "tool_resources": {
```

values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

response_format string or object Optional Specifies the format that the model must output. Compatible with GPT-4o, GPT-4 Turbo, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to { "type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

The modified assistant object.

Delete assistant Beta &

DELETE https://api.openai.com/v1/assistants/{assistant_id}

Delete an assistant.

Path parameters

assistant_id string Required
The ID of the assistant to delete.

Returns

Deletion status

The assistant object Beta ₽

Represents an assistant that can call the model and use tools.

id string

```
Example request python > 

1 from openai import OpenAI
2 client = OpenAI()
3
4 response = client.beta.assistants.delete("asst_abc123")
5 print(response)

Response

1 {
2 "id": "asst_abc123",
3 "object": "assistant.deleted",
4 "deleted": true
5 }
```

The assistant object

"id": "asst_abc123",

"object": "assistant",

1

币

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always assistant .

created at integer

The Unix timestamp (in seconds) for when the assistant was created.

name string or null

The name of the assistant. The maximum length is 256 characters.

description string or null

The description of the assistant. The maximum length is 512 characters.

model string

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them.

instructions string or null

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, file_search, or function.

∨Show possible types

tool_resources object or null

A set of resources that are used by the assistant's tools. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

metadata map

∨ Show properties

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null

What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null

An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

```
"created_at": 1698984975,
5
     "name": "Math Tutor",
6
     "description": null,
7
     "model": "gpt-4-turbo",
8
     "instructions": "You are a personal math tutor. When asked
9
     "tools": I
10
     {
11
        "type": "code_interpreter"
12
      }
13
     1.
     "metadata": {},
14
     "top_p": 1.0
15
     "temperature": 1.0,
16
     "response_format": "auto"
17
18 }
```

response_format string or object

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to { "type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Threads Beta &

Create threads that assistants can interact with.

Related guide: Assistants

Create thread Beta &

POST https://api.openai.com/v1/threads

Create a thread.

Request body

messages array Optional

A list of messages to start the thread with.

∨Show properties

tool_resources object or null Optional

A set of resources that are made available to the assistant's tools in this thread. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

∨ Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

A thread object.

Retrieve thread Beta &

GET https://api.openai.com/v1/threads/{thread_id}

Retrieves a thread.

```
币
Example request
                                           python ✓
1 from openai import OpenAI
2 client = OpenAI()
4 empty_thread = client.beta.threads.create()
5 print(empty_thread)
                                                         O
Response
1 {
   "id": "thread_abc123",
    "object": "thread",
   "created_at": 1699012949,
   "metadata": {},
    "tool_resources": {}
7 }
```

Example request

Empty

Messages

D

python ~

Path parameters

thread_id string Required
The ID of the thread to retrieve.

Returns

The thread object matching the specified ID.

Modify thread Beta €

POST https://api.openai.com/v1/threads/{thread_id}

Modifies a thread.

Path parameters

thread_id string Required

The ID of the thread to modify. Only the metadata can be modified.

Request body

tool_resources object or null Optional

A set of resources that are made available to the assistant's tools in this thread. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

> Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

The modified **thread** object matching the specified ID.

Delete thread Beta &

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 my_thread = client.beta.threads.retrieve("thread_abc123")
5 print(my_thread)
```

```
Response
                                                           ð
1
    {
2
     "id": "thread_abc123",
     "object": "thread",
     "created_at": 1699014083,
5
     "metadata": {},
6
     "tool_resources": {
7
      "code_interpreter": {
8
        "file_ids": []
9
10
    }
11 }
```

```
Example request
```

python V



币

```
from openai import OpenAI
2
   client = OpenAI()
3
   my_updated_thread = client.beta.threads.update(
5
    "thread_abc123",
6
    metadata={
7
     "modified": "true",
      "user": "abc123"
8
9
10 )
11 print(my_updated_thread)
```

```
Response
```

10 }

```
1  {
2    "id": "thread_abc123",
3    "object": "thread",
4    "created_at": 1699014083,
5    "metadata": {
6        "modified": "true",
7        "user": "abc123"
8    },
```

"tool_resources": {}

DELETE https://api.openai.com/v1/threads/{thread_id}

Delete a thread.

Path parameters

thread_id string Required The ID of the thread to delete.

Returns

Deletion status

The thread object

Represents a thread that contains messages.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.

created_at integer

The Unix timestamp (in seconds) for when the thread was created.

tool_resources object or null

A set of resources that are made available to the assistant's tools in this thread. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

∨ Show properties

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

Messages

Create messages within threads

Related guide: Assistants

Beta 🔗 Create message

POST https://api.openai.com/v1/threads/{thread_id}/messages

Create a message.

Path parameters

thread id string Required

- 1 from openai import OpenAI
- 2 client = OpenAI()

Example request

- 4 response = client.beta.threads.delete("thread_abc123")
- 5 print(response)

Response

币

凸

币

python ~

```
1 {
   "id": "thread_abc123",
```

- "object": "thread.deleted",
- "deleted": true
- 5 }

6 }

The thread object

1 { "id": "thread_abc123", "object": "thread",

"created_at": 1698107661, "metadata": {}

Example request

python ✓

O

1 from openai import OpenAl

- 2 client = OpenAI()
- 4 thread_message = client.beta.threads.messages.create(
- "thread_abc123",
- role="user",

The ID of the thread to create a message for.

Request body

role string Required

The role of the entity that is creating the message. Allowed values include:

user: Indicates the message is sent by an actual user and should be used in most cases to represent user-generated messages.

assistant: Indicates the message is generated by the assistant. Use this value to insert messages from the assistant into the conversation.

content string or array Required

∨Show possible types

attachments array or null Optional

A list of files attached to the message, and the tools they should be added to.

∨ Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

A message object.

List messages Beta &

GET https://api.openai.com/v1/threads/{thread_id}/messages

Returns a list of messages for a given thread.

Path parameters

thread_id string Required

The ID of the thread the messages belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

```
7 content="How does AI work? Explain it in simple terms.",
8 )
9 print(thread_message)
```

```
凸
Response
1
2
      "id": "msg_abc123",
3
      "object": "thread.message",
      "created_at": 1713226573,
4
      "assistant_id": null,
5
6
      "thread_id": "thread_abc123",
      "run_id": null,
8
      "role": "user",
9
      "content": [
10
         "type": "text",
11
12
         "text": {
          "value": "How does Al work? Explain it in simple terms."
13
          "annotations": []
14
16
       }
     ],
```

```
Example request
```

python ✓



- 1 from openai import OpenAI
- 2 client = OpenAI()
- 3
- 4 thread_messages = client.beta.threads.messages.list("thread_a
- 5 print(thread_messages.data)

Response



```
{
2
     "object": "list",
3
     "data": [
4
        "id": "msg_abc123",
5
6
        "object": "thread.message",
7
        "created_at": 1699016383,
8
        "assistant_id": null,
9
        "thread_id": "thread_abc123",
```

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

run_id string Optional

Filter messages by the run ID that generated them.

Returns

A list of message objects.

Retrieve message Beta €

GET https://api.openai.com/v1/threads/{thread_id}/messages/{message_id}

Retrieve a message.

Path parameters

thread_id string Required

The ID of the thread to which this message belongs.

message_id string Required

The ID of the message to retrieve.

Returns

The message object matching the specified ID.

```
10
        "run_id": null,
        "role": "user",
11
12
        "content": [
13
           "type": "text",
14
15
           "text": {
            "value": "How does AI work? Explain it in simple term
16
17
            "annotations": []
18
19
         }
20
        ],
        "attachments": [],
21
```

```
Example request
```

python ✓

```
ð
```

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 message = client.beta.threads.messages.retrieve(
5 message_id="msg_abc123",
6 thread_id="thread_abc123",
7 )
8 print(message)
```

```
Response
```

```
ð
```

```
1
2
     "id": "msg_abc123",
     "object": "thread.message",
     "created_at": 1699017614,
5
     "assistant_id": null,
6
     "thread_id": "thread_abc123",
7
     "run_id": null,
8
     "role": "user",
9
     "content": [
10
        "type": "text",
11
12
13
         "value": "How does AI work? Explain it in simple terms."
14
         "annotations": []
15
16
      }
17
18
     "attachments": [],
```

Modify message

POST https://api.openai.com/v1/threads/{thread_id}/messages/{mess

Modifies a message.

Path parameters

thread_id string Required

The ID of the thread to which this message belongs.

message id string Required The ID of the message to modify.

Request body

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

Returns

The modified message object.

Delete message

DELETE https://api.openai.com/v1/threads/{thread_id}/messages/{m essage_id}

Deletes a message.

Path parameters

thread_id string Required

The ID of the thread to which this message belongs.

message_id string Required

The ID of the message to delete.

Returns

Deletion status

```
Example request
```

python ~

```
\Box
```

```
from openai import OpenAI
2
   client = OpenAI()
3
4
    message = client.beta.threads.messages.update(
5
     message_id="msg_abc12",
     thread_id="thread_abc123",
7
     metadata={
8
      "modified": "true",
9
      "user": "abc123",
10 },
11 )
12 print(message)
```

Response



```
1
     "id": "msg_abc123",
2
3
     "object": "thread.message",
4
     "created_at": 1699017614,
5
     "assistant_id": null,
6
     "thread_id": "thread_abc123",
     "run_id": null,
7
     "role": "user",
8
9
     "content": [
10
      {
        "type": "text",
11
12
         "value": "How does Al work? Explain it in simple terms."
13
14
         "annotations": []
15
```

Example request

python ~



```
1 from openai import OpenAI
2 client = OpenAI()
4 deleted_message = client.beta.threads.messages.delete(
   message_id="msg_abc12",
   thread_id="thread_abc123",
6
7)
8 print(deleted_message)
```

Response



```
"id": "msg_abc123",
    "object": "thread.message.deleted",
    "deleted": true
5 }
```

The message object Beta &

Represents a message within a thread.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.message .

created_at integer

The Unix timestamp (in seconds) for when the message was created.

thread_id string

The thread ID that this message belongs to.

status string

The status of the message, which can be either in_progress, incomplete, or completed.

incomplete_details object or null

On an incomplete message, details about why the message is incomplete.

∨Show properties

completed_at integer or null

The Unix timestamp (in seconds) for when the message was completed.

incomplete_at integer or null

The Unix timestamp (in seconds) for when the message was marked as incomplete.

role string

The entity that produced the message. One of user or assistant .

content array

The content of the message in array of text and/or images.

∨Show possible types

assistant_id string or null

If applicable, the ID of the assistant that authored this message.

run_id string or null

The ID of the **run** associated with the creation of this message. Value is null when messages are created manually using the create message or create thread endpoints.

attachments array or null

A list of files attached to the message, and the tools they were added to.

√Show properties

metadata map

The message object

```
2
     "id": "msg_abc123",
3
     "object": "thread.message",
4
     "created_at": 1698983503,
5
     "thread id": "thread abc123",
     "role": "assistant",
     "content": [
8
      {
9
        "type": "text",
10
        "text": {
11
         "value": "Hi! How can I help you today?",
         "annotations": []
12
13
14
      }
15
     "assistant_id": "asst_abc123",
16
17
     "run_id": "run_abc123",
     "attachments": [],
18
19
     "metadata": {}
20 }
```

6

7

8

10

11

12

13

14

16

"thread_id": "thread_abc123",

"completed_at": 1699063291,

"started_at": 1699063290,

"status": "queued",

"expires_at": null,

"failed_at": null,

"last_error": null,

"instructions": null,
"incomplete_details": null,

"tools": [

"model": "gpt-4-turbo",

"cancelled_at": null,

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Runs Beta &

Represents an execution run on a thread.

Related guide: Assistants

Create run Beta &

POST https://api.openai.com/v1/threads/{thread_id}/runs

Create a run.

Path parameters

thread_id string Required
The ID of the thread to run.

Request body

assistant_id string Required

The ID of the assistant to use to execute this run.

model string Optional

The ID of the **Model** to be used to execute this run. If a value is provided here, it will override the model associated with the assistant. If not, the model associated with the assistant will be used.

instructions string or null Optional

Overrides the **instructions** of the assistant. This is useful for modifying the behavior on a per-run basis.

additional_instructions string or null Optional

Appends additional instructions at the end of the instructions for the run. This is useful for modifying the behavior on a per-run basis without overriding other instructions.

additional_messages array or null Optional

Adds additional messages to the thread before creating the run.

∨Show properties

tools array or null Optional

Override the tools the assistant can use for this run. This is useful for modifying the behavior on a per-run basis.

∨Show possible types

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

```
Default
                   Streaming
                                   Streaming with Functions
Example request
                                            python ~
                                                         凸
1 from openai import OpenAl
2 client = OpenAI()
3
4 run = client.beta.threads.runs.create(
  thread_id="thread_abc123",
   assistant_id="asst_abc123"
7)
9 print(run)
Response
                                                         ð
1
2
     "id": "run_abc123",
3
     "object": "thread.run",
     "created_at": 1699063290,
4
5
     "assistant_id": "asst_abc123",
```

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1

An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

max_prompt_tokens integer or null Optional

The maximum number of prompt tokens that may be used over the course of the run. The run will make a best effort to use only the number of prompt tokens specified, across multiple turns of the run. If the run exceeds the number of prompt tokens specified, the run will end with status incomplete. See incomplete_details for more info.

max_completion_tokens integer or null Optional

The maximum number of completion tokens that may be used over the course of the run. The run will make a best effort to use only the number of completion tokens specified, across multiple turns of the run. If the run exceeds the number of completion tokens specified, the run will end with status incomplete. See incomplete_details for more info.

truncation_strategy object Optional

Controls for how a thread will be truncated prior to the run. Use this to control the intial context window of the run.

∨Show properties

tool_choice string or object Optional

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling one or more tools. required means the model must call one or more tools before responding to the user. Specifying a particular tool like {"type": "file_search"} or {"type": "function", "function": {"name": "my_function"}} forces the model to call that tool.

∨Show possible types

parallel_tool_calls boolean Optional Defaults to true
Whether to enable parallel function calling during tool use.

response_format string or object Optional

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to {"type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

A run object.

Create thread and run Beta &

POST https://api.openai.com/v1/threads/runs

Create a thread and run it in one request.

Request body

assistant id string Required

The ID of the assistant to use to execute this run.

thread object Optional

∨Show properties

model string Optional

The ID of the **Model** to be used to execute this run. If a value is provided here, it will override the model associated with the assistant. If not, the model associated with the assistant will be used.

instructions string or null Optional

Override the default system message of the assistant. This is useful for modifying the behavior on a per-run basis.

tools array or null Optional

Override the tools the assistant can use for this run. This is useful for modifying the behavior on a per-run basis.

tool_resources object or null Optional

A set of resources that are used by the assistant's tools. The resources are specific to the type of tool. For example, the code_interpreter tool requires a list of file IDs, while the file_search tool requires a list of vector store IDs.

∨Show properties

Default Streaming Streaming with Functions

Example request

python ~



```
from openai import OpenAI
2
    client = OpenAI()
3
4
   run = client.beta.threads.create_and_run(
5
     assistant_id="asst_abc123",
6
     thread={
7
      "messages": [
8
       {"role": "user", "content": "Explain deep learning to a 5 ye
9
10
    }
11 )
13 print(run)
```

```
Response
```



```
1
   {
2
     "id": "run_abc123",
3
     "object": "thread.run",
     "created_at": 1699076792,
5
     "assistant_id": "asst_abc123",
6
     "thread id": "thread abc123",
7
     "status": "queued",
8
     "started_at": null,
9
     "expires_at": 1699077392,
10
     "cancelled_at": null,
     "failed_at": null,
11
12
     "completed_at": null,
13
     "required action": null.
```

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

max prompt tokens integer or null Optional

The maximum number of prompt tokens that may be used over the course of the run. The run will make a best effort to use only the number of prompt tokens specified, across multiple turns of the run. If the run exceeds the number of prompt tokens specified, the run will end with status <code>[incomplete]</code>. See <code>[incomplete_details]</code> for more info.

max_completion_tokens integer or null Optional

The maximum number of completion tokens that may be used over the course of the run. The run will make a best effort to use only the number of completion tokens specified, across multiple turns of the run. If the run exceeds the number of completion tokens specified, the run will end with status incomplete. See incomplete_details for more info.

truncation_strategy object Optional

Controls for how a thread will be truncated prior to the run. Use this to control the intial context window of the run.

∨Show properties

tool_choice string or object Optional

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling one or more tools. required means the model must call one or more tools before responding to the user. Specifying a particular tool like {"type": "file_search"} or {"type": "function", "function": {"name": "my_function"}} forces the

model to call that tool.

∨Show possible types

parallel_tool_calls boolean Optional Defaults to true
Whether to enable parallel function calling during tool use.

response_format string or object Optional Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to [{"type": "json_object" }] enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

A run object.

List runs Beta &

GET https://api.openai.com/v1/threads/{thread_id}/runs

Returns a list of runs belonging to a thread.

Path parameters

thread_id string Required
The ID of the thread the run belongs to.

Query parameters

limit integer Optional Defaults to 20
A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

```
python > □

from openai import OpenAl
client = OpenAl()

runs = client.beta.threads.runs.list(
"thread_abc123"
)

python > □

python > □

python > □

print(runs)
```

```
Response
```

```
1
     {
       "object": "list",
2
3
       "data": [
4
5
         "id": "run_abc123",
         "object": "thread.run",
6
7
         "created_at": 1699075072,
8
         "assistant_id": "asst_abc123",
9
         "thread_id": "thread_abc123",
10
         "status": "completed",
         "started_at": 1699075072,
11
12
         "expires at": null,
13
         "cancelled_at": null,
14
         "failed_at": null,
         "completed_at": 1699075073,
15
         "last_error": null,
16
17
         "model": "gpt-4-turbo",
18
         "instructions": null,
```

币

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

```
"incomplete_details": null,
20
          "tools": [
```

Returns

A list of run objects.

Retrieve run Beta 🔗

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}

Retrieves a run.

Path parameters

thread_id string Required The ID of the thread that was run.

run_id string Required The ID of the run to retrieve.

Returns

The run object matching the specified ID.

Modify run Beta 🔗

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}

Modifies a run.

Path parameters

thread id string Required The ID of the thread that was run.

run_id string Required

```
Example request
```

1 from openai import OpenAI

2 client = OpenAI()

3

4 run = client.beta.threads.runs.retrieve(

python ✓

ð

凸

thread_id="thread_abc123",

run_id="run_abc123" 7)

9 print(run)

Response

```
1
   {
2
     "id": "run_abc123",
     "object": "thread.run",
     "created_at": 1699075072,
5
     "assistant_id": "asst_abc123",
     "thread_id": "thread_abc123",
6
7
     "status": "completed",
8
     "started_at": 1699075072,
     "expires_at": null,
10
     "cancelled_at": null,
11
     "failed_at": null,
12
     "completed_at": 1699075073,
     "last_error": null,
13
```

"instructions": null, 15

16 "incomplete_details": null,

"model": "gpt-4-turbo",

17 "tools": [

18

14

19 "type": "code_interpreter"

Example request

```
from openai import OpenAI
1
2
   client = OpenAI()
3
4
    run = client.beta.threads.runs.update(
5
```

thread_id="thread_abc123",

6 run_id="run_abc123",

7 metadata={"user_id": "user_abc123"},

8

54/106

python ✓

ð

The ID of the run to modify.

Request body

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

Returns

The modified **run** object matching the specified ID.

Submit tool outputs to run

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/sub mit_tool_outputs

When a run has the status:

"requires_action" and required_action.type is submit_tool_outputs, this endpoint can be used to submit the outputs from the tool calls once they're all completed. All outputs must be submitted in a single request.

Path parameters

thread_id string Required

The ID of the thread to which this run belongs.

run_id string Required

The ID of the run that requires the tool output submission.

Request body

tool_outputs array Required

A list of tools for which the outputs are being submitted.

∨ Show properties

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

Returns

```
10 print(run)
Response
1
    {
2
      "id": "run abc123",
3
      "object": "thread.run",
      "created_at": 1699075072,
4
5
      "assistant_id": "asst_abc123",
```

10 "cancelled_at": null, 11 "failed_at": null, "completed_at": 1699075073, 12

"thread_id": "thread_abc123",

"status": "completed",

"expires_at": null,

"started_at": 1699075072,

13 "last_error": null, "model": "gpt-4-turbo", 14

"instructions": null, 15 "incomplete_details": null, 16

17 "tools": [18 {

7

8

9

Default Streaming

Example request

python ~



```
from openai import OpenAI
2
    client = OpenAI()
3
    run = client.beta.threads.runs.submit tool outputs(
5
     thread_id="thread_123",
6
     run_id="run_123",
     tool_outputs=[
7
8
        "tool_call_id": "call_001",
9
10
       "output": "70 degrees and sunny."
11
12
    ]
13
```

Response



```
1
2
     "id": "run 123",
3
     "object": "thread.run",
     "created_at": 1699075592,
     "assistant_id": "asst_123",
5
6
     "thread_id": "thread_123",
7
     "status": "queued",
8
     "started_at": 1699075592,
      "expires_at": 1699076192,
10
     "cancelled_at": null,
11
     "failed at": null,
     "completed_at": null,
12
     "last_error": null
```

The modified run object matching the specified ID.

Cancel a run Beta &

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/can cel

Cancels a run that is in_progress.

Path parameters

thread_id string Required

The ID of the thread to which this run belongs.

run id string Required

The ID of the run to cancel.

Returns

The modified run object matching the specified ID.

The run object Beta €

Represents an execution run on a thread.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.run .

created_at integer

The Unix timestamp (in seconds) for when the run was created.

thread_id string

The ID of the thread that was executed on as a part of this run.

assistant_id string

The ID of the assistant used for execution of this run.

```
status string
```

```
https://platform.openai.com/docs/api-reference/runs-v1
```

```
Example request python ✓ ☐

1 from openai import OpenAI

2 client = OpenAI()

3

4 run = client.beta.threads.runs.cancel(
```

```
5 thread_id="thread_abc123",
6 run_id="run_abc123"
7 )
8
9 print(run)
```

Response



```
1
2
     "id": "run_abc123",
     "object": "thread.run",
3
     "created_at": 1699076126,
     "assistant_id": "asst_abc123",
     "thread_id": "thread_abc123",
     "status": "cancelling",
7
8
     "started_at": 1699076126,
9
     "expires_at": 1699076726,
10
    "cancelled_at": null,
     "failed_at": null,
11
12
     "completed_at": null,
13
     "last_error": null,
     "model": "gpt-4-turbo",
14
     "instructions": "You summarize books.",
15
     "tools": [
16
17
      {
        "type": "file_search"
18
19
```

The run object



```
1
   {
     "id": "run_abc123",
2
3
     "object": "thread.run",
     "created_at": 1698107661,
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
7
     "status": "completed",
8
     "started_at": 1699073476,
     "expires_at": null,
10
     "cancelled_at": null,
11
     "failed_at": null,
12
     "completed_at": 1699073498,
     "last_error": null,
13
14
     "model": "gpt-4-turbo",
15
     "instructions": null,
     "tools": [{"type": "file_search"}, {"type": "code_interpreter"}],
     "metadata": {},
17
     "incomplete_details": null,
```

```
The status of the run, which can be either | queued |, | in_progress |,
 requires_action, cancelling, cancelled, failed, completed
 incomplete, or expired.
required_action object or null
Details on the action required to continue the run. Will be |null | if
no action is required.
∨Show properties
last_error object or null
The last error associated with this run. Will be | null | if there are no
```

errors.

∨ Show properties

expires_at integer or null

The Unix timestamp (in seconds) for when the run will expire.

started_at integer or null

The Unix timestamp (in seconds) for when the run was started.

cancelled at integer or null

The Unix timestamp (in seconds) for when the run was cancelled.

failed_at integer or null

The Unix timestamp (in seconds) for when the run failed.

completed_at integer or null

The Unix timestamp (in seconds) for when the run was completed.

incomplete_details object or null

Details on why the run is incomplete. Will be | null | if the run is not incomplete.

∨Show properties

model string

The model that the assistant used for this run.

instructions string

The instructions that the assistant used for this run.

tools array

The list of tools that the assistant used for this run.

∨ Show possible types

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

usage object or null

Usage statistics related to the run. This value will be | null | if the run is not in a terminal state (i.e. in_progress , queued , etc.).

∨Show properties

```
19
     "usage": {
20
      "prompt_tokens": 123,
21
      "completion_tokens": 456,
22
      "total_tokens": 579
23
24
     "temperature": 1.0,
25
     "top_p": 1.0,
26
     "max_prompt_tokens": 1000,
27
     "max_completion_tokens": 1000,
28
     "truncation_strategy": {
29
      "type": "auto",
      "last_messages": null
31
    },
32
     "response_format": "auto",
     "tool_choice": "auto",
33
34
    "parallel_tool_calls": true
35 }
```

temperature number or null

The sampling temperature used for this run. If not set, defaults to 1.

top_p number or null

The nucleus sampling value used for this run. If not set, defaults to

max_prompt_tokens integer or null

The maximum number of prompt tokens specified to have been used over the course of the run.

max_completion_tokens integer or null

The maximum number of completion tokens specified to have been used over the course of the run.

truncation_strategy object

Controls for how a thread will be truncated prior to the run. Use this to control the intial context window of the run.

√ Show properties

tool_choice string or object

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling one or more tools. required means the model must call one or more tools before responding to the user. Specifying a particular tool like {"type": "file_search"} or {"type": "function", "function": {"name": "my_function"}} forces the model to call that tool.

∨Show possible types

parallel_tool_calls boolean

Whether to enable parallel function calling during tool use.

response_format string or object

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to [{"type": "json_object" }] enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Run Steps Beta &

Represents the steps (model and tool calls) taken during the run.

Related guide: Assistants

List run steps Beta &

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/steps

Returns a list of run steps belonging to a run.

Path parameters

thread_id string Required

The ID of the thread the run and run steps belong to.

run id string Required

The ID of the run the run steps belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

Returns

A list of run step objects.

Retrieve run step Beta &

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/step s/{step_id}

Retrieves a run step.

Path parameters

thread_id string Required

The ID of the thread to which the run and run step belongs.

```
Example request
```

python ~

```
<sub>(1)</sub>
```

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 run_steps = client.beta.threads.runs.steps.list(
5 thread_id="thread_abc123",
6 run_id="run_abc123"
7 )
8
9 print(run_steps)
```

Response



```
{
2
     "object": "list",
3
     "data": [
4
5
        "id": "step_abc123",
        "object": "thread.run.step",
6
7
        "created_at": 1699063291,
        "run_id": "run_abc123",
9
        "assistant_id": "asst_abc123",
10
        "thread_id": "thread_abc123",
11
        "type": "message_creation",
12
        "status": "completed",
13
        "cancelled_at": null,
14
        "completed_at": 1699063291,
15
        "expired_at": null,
16
        "failed_at": null,
17
        "last_error": null,
18
        "step_details": {
19
         "type": "message_creation",
```

python ✓



```
from openai import OpenAl
client = OpenAl()

run_step = client.beta.threads.runs.steps.retrieve(
thread_id="thread_abc123",
run_id="run_abc123",
step_id="step_abc123"
```

run_id string Required

The ID of the run to which the run step belongs.

step_id string Required

The ID of the run step to retrieve.

Returns

The run step object matching the specified ID.

The run step object Beta €

Represents a step in execution of a run.

id string

The identifier of the run step, which can be referenced in API endpoints.

object string

The object type, which is always thread.run.step .

created_at integer

The Unix timestamp (in seconds) for when the run step was created.

assistant_id string

The ID of the assistant associated with the run step.

thread_id string

The ID of the thread that was run.

run_id string

The ID of the run that this run step is a part of.

type string

The type of run step, which can be either message_creation or tool_calls .

status string

The status of the run step, which can be either in_progress, cancelled, failed, completed, or expired.

step_details object

```
8 )910 print(run_step)
```

Response

```
1
   {
2
     "id": "step_abc123",
3
     "object": "thread.run.step",
     "created_at": 1699063291,
     "run_id": "run_abc123",
     "assistant_id": "asst_abc123",
     "thread id": "thread abc123",
Я
     "type": "message_creation",
     "status": "completed",
     "cancelled_at": null,
10
     "completed_at": 1699063291,
12
     "expired_at": null,
     "failed_at": null,
14
     "last_error": null,
15
     "step_details": {
16
      "type": "message_creation",
17
      "message_creation": {
       "message_id": "msg_abc123"
```

币

币

The run step object

```
1
     "id": "step_abc123",
     "object": "thread.run.step",
3
4
     "created_at": 1699063291,
     "run id": "run abc123",
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
8
     "type": "message_creation",
     "status": "completed",
10
     "cancelled_at": null,
     "completed_at": 1699063291,
     "expired_at": null,
12
13
     "failed_at": null,
14
     "last_error": null,
15
     "step_details": {
16
      "type": "message_creation",
17
      "message_creation": {
18
       "message_id": "msg_abc123"
19
      }
20
     },
21
     "usage": {
22
      "prompt_tokens": 123,
23
      "completion_tokens": 456,
      "total_tokens": 579
25
26 }
```

The details of the run step.

∨Show possible types

last error object or null

The last error associated with this run step. Will be |null | if there are no errors.

∨Show properties

expired_at integer or null

The Unix timestamp (in seconds) for when the run step expired. A step is considered expired if the parent run is expired.

cancelled_at integer or null

The Unix timestamp (in seconds) for when the run step was cancelled.

failed_at integer or null

The Unix timestamp (in seconds) for when the run step failed.

completed_at integer or null

The Unix timestamp (in seconds) for when the run step completed.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

usage object or null

Usage statistics related to the run step. This value will be null while the run step's status is in_progress

∨Show properties

Vector Stores Beta @

Vector stores are used to store files for use by the file_search tool.

Related guide: File Search

Create vector store

POST https://api.openai.com/v1/vector_stores

Create a vector store.

Request body

file_ids array Optional

A list of File IDs that the vector store should use. Useful for tools like file_search that can access files.

name string Optional

The name of the vector store.

expires after object Optional

The expiration policy for a vector store.

∨Show properties

Example request

python ~

仚

币

```
1 from openai import OpenAl
2 client = OpenAI()
3
  vector_store = client.beta.vector_stores.create(
   name="Support FAQ"
6)
7 print(vector_store)
```

Response

```
1
     "id": "vs_abc123",
     "object": "vector_store",
3
4
    "created_at": 1699061776,
5
     "name": "Support FAQ",
```

"bytes": 139920,

chunking strategy object Optional

The chunking strategy used to chunk the file(s). If not set, will use the auto strategy. Only applicable if file_ids is non-empty.

∨Show possible types

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

```
7 "file_counts": {
8     "in_progress": 0,
9     "completed": 3,
10     "failed": 0,
11     "cancelled": 0,
12     "total": 3
13     }
14 }
```

Returns

A vector store object.

List vector stores Beta &

GET https://api.openai.com/v1/vector_stores

Returns a list of vector stores.

Query parameters

limit integer Optional Defaults to 20
A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list

Returns

A list of vector store objects.

Retrieve vector store Beta &

GET https://api.openai.com/v1/vector_stores/{vector_store_id}

Retrieves a vector store.

```
Example request python ✓ ☐

1 from openai import OpenAI

2 client = OpenAI()

3

4 vector_stores = client.beta.vector_stores.list()

5 print(vector_stores)

Response
```

```
1
    {
2
     "object": "list",
3
     "data": [
4
        "id": "vs_abc123",
5
        "object": "vector_store",
6
7
        "created_at": 1699061776,
8
        "name": "Support FAQ",
9
        "bytes": 139920,
        "file_counts": {
10
11
         "in_progress": 0,
12
         "completed": 3,
         "failed": 0,
13
14
         "cancelled": 0.
         "total": 3
15
16
        }
17
      },
18
      {
        "id": "vs_abc456",
19
20
        "object": "vector_store",
21
        "created_at": 1699061776,
22
        "name": "Support FAQ v2",
23
        "bytes": 139920,
```

Example request

python ~



Path parameters

vector_store_id string Required
The ID of the vector store to retrieve.

Returns

The vector store object matching the specified ID.

Modify vector store Beta €

POST https://api.openai.com/v1/vector_stores/{vector_store_id}

Modifies a vector store.

Path parameters

vector_store_id string Required
The ID of the vector store to modify.

Request body

name string or null Optional The name of the vector store.

expires_after object Optional

The expiration policy for a vector store.

√Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

The modified vector store object.

Delete vector store Beta &

DELETE https://api.openai.com/v1/vector_stores/{vector_store_id}

Delete a vector store.

Path parameters

```
1 from openal import OpenAl
2 client = OpenAl()
3
4 vector_store = client.beta.vector_stores.retrieve(
5 vector_store_id="vs_abc123"
6 )
7 print(vector_store)
```

Response

1 {
2 "id": "vs_abc123",
3 "object": "vector_store",
4 "created_at": 1699061776
5 }

Example request

python ✓ 🗇

币

合

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 vector_store = client.beta.vector_stores.update(
5 vector_store_id="vs_abc123",
6 name="Support FAQ"
7 )
8 print(vector_store)
```

Response

1 { "id": "vs_abc123", 3 "object": "vector_store", "created_at": 1699061776, "name": "Support FAQ", 6 "bytes": 139920, 7 "file_counts": { "in_progress": 0, 8 "completed": 3, 10 "failed": 0, 11 "cancelled": 0, "total": 3 12 13 }

Example request

- 1 from openai import OpenAI
- 2 client = OpenAI()
- 3

14 }

4 deleted_vector_store = client.beta.vector_stores.delete(

币

python ~

6)

vector_store_id string Required

The ID of the vector store to delete.

Returns

Deletion status

The vector store object Beta ∂

A vector store is a collection of processed files can be used by the file_search tool.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always vector_store.

created_at integer

The Unix timestamp (in seconds) for when the vector store was created.

name string

The name of the vector store.

usage_bytes integer

The total number of bytes used by the files in the vector store.

file_counts object

∨Show properties

status strino

The status of the vector store, which can be either expired, in_progress, or completed. A status of completed indicates that the vector store is ready for use.

expires_after object

The expiration policy for a vector store.

∨Show properties

expires at integer or null

The Unix timestamp (in seconds) for when the vector store will expire.

last active at integer or null

The Unix timestamp (in seconds) for when the vector store was last active.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in

```
Response
```

```
1 {
2  id: "vs_abc123",
3  object: "vector_store.deleted",
4  deleted: true
5 }
```

5 vector_store_id="vs_abc123"

7 print(deleted_vector_store)

The vector store object

```
1
   {
2
     "id": "vs_123",
     "object": "vector_store",
     "created_at": 1698107661,
     "usage_bytes": 123456,
6
     "last_active_at": 1698107661,
7
     "name": "my_vector_store",
8
     "status": "completed",
     "file_counts": {
      "in_progress": 0,
10
11
      "completed": 100,
12
      "cancelled": 0.
      "failed": 0,
13
      "total": 100
15 },
     "metadata": {},
     "last_used_at": 1698107661
17
18 }
```

币

币

a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

Vector Store Files Beta ∂

Vector store files represent files inside a vector store.

Related guide: File Search

Create vector store file Beta ∂

POST https://api.openai.com/v1/vector_stores/{vector_store_id}/files

Create a vector store file by attaching a File to a vector store.

Path parameters

vector_store_id string Required

The ID of the vector store for which to create a File.

Request body

file_id string Required

A File ID that the vector store should use. Useful for tools like file_search that can access files.

chunking strategy object Optional

The chunking strategy used to chunk the file(s). If not set, will use the auto strategy.

∨Show possible types

Returns

A vector store file object.

List vector store files

GET https://api.openai.com/v1/vector_stores/{vector_store_id}/files

Returns a list of vector store files.

Path parameters

vector_store_id string Required

The ID of the vector store that the files belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

```
Example request
```

python ~

币

```
1 from openai import OpenAl
```

```
2 client = OpenAI()
```

3

4 vector_store_file = client.beta.vector_stores.files.create(

```
vector_store_id="vs_abc123",
```

6 file_id="file-abc123"

7)

8 print(vector_store_file)

Response

9 }

币

```
1 {
2
    "id": "file-abc123",
   "object": "vector_store.file",
   "created_at": 1699061776,
    "usage_bytes": 1234,
    "vector_store_id": "vs_abcd",
   "status": "completed",
  "last_error": null
```

Example request

python ✓

O

```
1 from openai import OpenAI
```

```
2 client = OpenAI()
```

4 vector_store_files = client.beta.vector_stores.files.list(

vector_store_id="vs_abc123"

6)

7 print(vector_store_files)

Response

币

```
"object": "list",
     "data": [
3
4
         "id": "file-abc123",
```

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

filter string Optional

Filter by file status. One of in_progress , completed , failed , cancelled .

Returns

A list of **vector store file** objects.

Retrieve vector store file

GET https://api.openai.com/v1/vector_stores/{vector_store_id}/files/{files/files/ le_id}

Retrieves a vector store file.

Path parameters

vector store id string Required

The ID of the vector store that the file belongs to.

file_id string Required

The ID of the file being retrieved.

Returns

The vector store file object.

Delete vector store file Beta ∂

DELETE https://api.openai.com/v1/vector_stores/{vector_store_id}/fil es/{file_id}

Delete a vector store file. This will remove the file from the vector store but the file itself will not be deleted. To delete

```
API Reference - OpenAI API
```

```
"object": "vector_store.file",
7
        "created_at": 1699061776,
8
        "vector_store_id": "vs_abc123"
9
      },
10
      {
        "id": "file-abc456".
11
12
        "object": "vector_store.file",
13
        "created_at": 1699061776,
14
        "vector_store_id": "vs_abc123"
15
16
    ],
     "first_id": "file-abc123",
17
     "last_id": "file-abc456",
18
     "has_more": false
```

Example request

凸 python ~

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 vector_store_file = client.beta.vector_stores.files.retrieve(
   vector store id="vs abc123",
   file_id="file-abc123"
7)
8 print(vector_store_file)
```

Response

8 }

1 { "id": "file-abc123", "object": "vector_store.file", 4 "created_at": 1699061776, "vector_store_id": "vs_abcd", "status": "completed", "last_error": null

Example request

python ~



ð

1 from openai import OpenAI

2 client = OpenAI()

3

Response

the file, use the delete file endpoint.

Path parameters

vector_store_id string Required

The ID of the vector store that the file belongs to.

file id string Required

The ID of the file to delete.

Returns

Deletion status

The vector store file object Beta &

A list of files attached to a vector store.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always vector_store.file .

usage bytes integer

The total vector store usage in bytes. Note that this may be different from the original file size.

created at integer

The Unix timestamp (in seconds) for when the vector store file was created.

vector store id string

The ID of the vector store that the File is attached to.

status string

The status of the vector store file, which can be either in_progress, completed, cancelled, or failed. The status completed indicates that the vector store file is ready for use.

last_error object or null

The last error associated with this vector store file. Will be null if there are no errors.

∨Show properties

chunking_strategy object

The strategy used to chunk the file.

∨Show possible types

Vector Store File Batches Beta &

Vector store file batches represent operations to add multiple files to a vector store.

Related guide: File Search

```
4 deleted_vector_store_file = client.beta.vector_stores.files.dele
5    vector_store_id="vs_abc123",
6    file_id="file-abc123"
7 )
8    print(deleted_vector_store_file)
```

```
1 {
2 id: "file-abc123",
3 object: "vector_store.file.deleted",
4 deleted: true
5 }
```

O

ð

```
1
    {
2
     "id": "file-abc123",
     "object": "vector_store.file",
3
     "usage_bytes": 1234,
     "created_at": 1698107661,
6
     "vector_store_id": "vs_abc123",
7
     "status": "completed",
8
     "last_error": null,
     "chunking_strategy": {
      "type": "static",
10
11
      "static": {
12
        "max_chunk_size_tokens": 800,
13
        "chunk_overlap_tokens": 400
      }
15 }
16 }
```

The vector store file object

Create vector store file batch Beta &

POST https://api.openai.com/v1/vector_stores/{vector_store_id}/file_ batches

Create a vector store file batch.

Path parameters

vector_store_id string Required

The ID of the vector store for which to create a File Batch.

Request body

file_ids array Required

A list of **File** IDs that the vector store should use. Useful for tools like file_search that can access files.

chunking_strategy object Optional

The chunking strategy used to chunk the file(s). If not set, will use the auto strategy.

∨Show possible types

Returns

A vector store file batch object.

Retrieve vector store file batch Beta &

GET https://api.openai.com/v1/vector_stores/{vector_store_id}/file_b atches/{batch_id}

Retrieves a vector store file batch.

Path parameters

vector_store_id string Required

The ID of the vector store that the file batch belongs to.

batch_id string Required

The ID of the file batch being retrieved.

Returns

The vector store file batch object.

```
Example request python ✓ ☐

1 from openai import OpenAI

2 client = OpenAI()

3

4 vector_store_file_batch = client.beta.vector_stores.file_batches.

5 vector_store_id="vs_abc123",

6 file_ids=["file-abc123", "file-abc456"]

7 )

8 print(vector_store_file_batch)

Response ☐
```

```
1
   {
2
     "id": "vsfb_abc123",
3
     "object": "vector_store.file_batch",
4
     "created_at": 1699061776,
     "vector_store_id": "vs_abc123",
5
6
     "status": "in_progress",
     "file_counts": {
8
      "in_progress": 1,
9
       "completed": 1,
10
      "failed": 0.
      "cancelled": 0,
11
12
      "total": 0,
13 }
14 }
```

```
Example request python ✓ □
```

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 vector_store_file_batch = client.beta.vector_stores.file_batches.
5 vector_store_id="vs_abc123",
6 batch_id="vsfb_abc123"
7 )
8 print(vector_store_file_batch)
```

```
Response
```

```
1
2
     "id": "vsfb abc123",
     "object": "vector_store.file_batch",
     "created_at": 1699061776,
5
     "vector_store_id": "vs_abc123",
6
     "status": "in_progress",
7
     "file_counts": {
8
      "in_progress": 1,
      "completed": 1,
10
      "failed": 0,
11
      "cancelled": 0.
      "total": 0,
12
```

仚

13 }

Cancel vector store file batch Beta &

POST https://api.openai.com/v1/vector_stores/{vector_store_id}/file_batches/{batch_id}/cancel

Cancel a vector store file batch. This attempts to cancel the processing of files in this batch as soon as possible.

Path parameters

vector_store_id string RequiredThe ID of the vector store that the file batch belongs to.

batch_id string Required
The ID of the file batch to cancel.

Returns

The modified vector store file batch object.

List vector store files in a batch Beta d

GET https://api.openai.com/v1/vector_stores/{vector_store_id}/file_b atches/{batch_id}/files

Returns a list of vector store files in a batch.

Path parameters

vector_store_id string Required

The ID of the vector store that the files belong to.

batch id string Required

The ID of the file batch that the files belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

```
Example request
```

 $\text{python}\,\checkmark$

```
1 from openai import OpenAl
2 client = OpenAl()
3
4 deleted_vector_store_file_batch = client.beta.vector_stores.file_
5    vector_store_id="vs_abc123",
6    file_batch_id="vsfb_abc123"
7 )
8 print(deleted_vector_store_file_batch)
```

```
Response
```



仚

```
1
2
     "id": "vsfb_abc123",
     "object": "vector_store.file_batch",
     "created_at": 1699061776,
5
     "vector_store_id": "vs_abc123",
     "status": "cancelling",
7
     "file counts": {
8
      "in_progress": 12,
      "completed": 3,
10
      "failed": 0,
      "cancelled": 0,
11
      "total": 15,
12
13 }
14 }
```

Example request

python ~



```
1 from openai import OpenAl
2 client = OpenAl()
3
4 vector_store_files = client.beta.vector_stores.file_batches.list_fil
5 vector_store_id="vs_abc123",
6 batch_id="vsfb_abc123"
7 )
8 print(vector_store_files)
```

Response



```
1
     "object": "list",
2
3
     "data": [
4
      {
        "id": "file-abc123",
5
6
        "object": "vector_store.file",
7
        "created_at": 1699061776,
8
        "vector_store_id": "vs_abc123"
9
      },
10
        "id": "file-abc456",
11
```

```
after string Optional
```

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list

```
filter string Optional

Filter by file status. One of in_progress, completed, failed, cancelled.
```

12 "object": "vector_store.file", 13 "created_at": 1699061776, 14 "vector_store_id": "vs_abc123" 15 } 16], 17 "first_id": "file-abc123", 18 "last_id": "file-abc456", 19 "has_more": false

Returns

A list of vector store file objects.

The vector store files batch object Beta &

A batch of files attached to a vector store.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always vector_store.file_batch .

created at integer

The Unix timestamp (in seconds) for when the vector store files batch was created.

vector_store_id string

The ID of the vector store that the File is attached to.

status string

The status of the vector store files batch, which can be either in_progress, completed, cancelled or failed.

file counts object

∨Show properties

Streaming Beta &

Stream the result of executing a Run or resuming a Run after submitting tool outputs.

You can stream events from the **Create Thread and Run**, **Create Run**, and **Submit Tool Outputs** endpoints by passing "stream": true. The response will be a **Server-Sent events** stream.

Our Node and Python SDKs provide helpful utilities to make streaming easy. Reference the **Assistants API** quickstart to learn more.

```
1
   {
     "id": "vsfb_123",
2
     "object": "vector_store.files_batch",
3
     "created_at": 1698107661,
     "vector_store_id": "vs_abc123",
5
     "status": "completed",
7
     "file_counts": {
8
      "in_progress": 0,
9
      "completed": 100,
10
      "failed": 0,
      "cancelled": 0.
11
12
      "total": 100
13 }
14 }
```

The vector store files batch object

仚

The message delta object Beta &

Represents a message delta i.e. any changed fields on a message during streaming.

id string

The identifier of the message, which can be referenced in API endpoints.

object string

The object type, which is always thread.message.delta .

delta object

The delta containing the fields that have changed on the Message. ~Show properties

The run step delta object Beta ∂

Represents a run step delta i.e. any changed fields on a run step during streaming.

id string

The identifier of the run step, which can be referenced in API endpoints.

object string

The object type, which is always thread.run.step.delta .

delta object

The delta containing the fields that have changed on the run step. VShow properties

The message delta object

```
1
2
     "id": "msg_123",
3
     "object": "thread.message.delta",
4
     "delta": {
5
       "content": [
6
        {
          "index": 0.
8
         "type": "text",
9
          "text": { "value": "Hello", "annotations": [] }
10
11
      ]
12
    }
13 }
```

币

凸

The run step delta object

```
1
    {
2
      "id": "step_123",
3
      "object": "thread.run.step.delta",
4
      "delta": {
       "step_details": {
5
6
        "type": "tool_calls",
7
         "tool_calls": [
8
           "index": 0,
           "id": "call_123",
10
11
           "type": "code_interpreter",
           "code_interpreter": { "input": "", "outputs": [] }
12
13
         }
14
        1
15
       }
16 }
17 }
```

Assistant stream events Beta &

Represents an event emitted when streaming a Run.

Each event in a server-sent events stream has an event and data property:

```
event: thread.created
data: {"id": "thread_123", "object": "thread", ...}
```

We emit events whenever a new object is created, transitions to a new state, or is being streamed in parts (deltas). For example, we emit thread.run.created when a new run is created, thread.run.completed when a run completes, and so on. When an Assistant chooses to

create a message during a run, we emit a thread.message.created event, a thread.message.in_progress event, many thread.message.delta events, and finally a thread.message.completed event.

We may add additional events over time, so we recommend handling unknown events gracefully in your code. See the **Assistants API quickstart** to learn how to integrate the Assistants API with streaming.

thread.created data is a thread

Occurs when a new thread is created.

thread.run.created data is a run

Occurs when a new run is created.

thread.run.queued data is a run

Occurs when a **run** moves to a queued status.

thread.run.in_progress data is a run

Occurs when a **run** moves to an in_progress status.

thread.run.requires_action data is a run

Occurs when a **run** moves to a requires_action status.

thread.run.completed data is a run

Occurs when a run is completed.

thread.run.incomplete data is a run

Occurs when a **run** ends with status incomplete.

thread.run.failed data is a run

Occurs when a run fails.

thread.run.cancelling data is a run

Occurs when a **run** moves to a cancelling status.

thread.run.cancelled data is a run

Occurs when a **run** is cancelled.

thread.run.expired data is a run

Occurs when a run expires.

thread.run.step.created data is a run step

Occurs when a run step is created.

thread.run.step.in_progress | data | is a run step

Occurs when a **run step** moves to an in_progress state.

thread.run.step.delta | data | is a run step delta Occurs when parts of a run step are being streamed. thread.run.step.completed | data | is a run step Occurs when a run step is completed. thread.run.step.failed data is a run step Occurs when a run step fails. thread.run.step.cancelled data is a run step Occurs when a run step is cancelled. thread.run.step.expired data is a run step Occurs when a run step expires. thread.message.created | data | is a message Occurs when a message is created. thread.message.in_progress | data | is a message Occurs when a **message** moves to an in_progress state. thread.message.delta data is a message delta Occurs when parts of a Message are being streamed. thread.message.completed | data | is a message Occurs when a message is completed. thread.message.incomplete data is a message Occurs when a message ends before it is completed. error data is an error Occurs when an error occurs. This can happen due to an internal server error or a timeout. done data is [DONE] Occurs when a stream ends.

Completions Legacy ∂

Given a prompt, the model will return one or more predicted completions along with the probabilities of alternative tokens at each position. Most developer should use our **Chat Completions API** to leverage our best and newest models.

Create completion Legacy €

POST https://api.openai.com/v1/completions

Creates a completion for the provided prompt and parameters.

Request body

Example req... gpt-3.5-turbo-instruct ✓ python ✓

No streaming

- 1 from openai import OpenAI
- 2 client = OpenAI()
- 3
- 4 client.completions.create(

Streaming

model string Required

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them.

prompt string or array Required

The prompt(s) to generate completions for, encoded as a string, array of strings, array of tokens, or array of token arrays.

Note that <lendoftextl> is the document separator that the model sees during training, so if a prompt is not specified the model will generate as if from the beginning of a new document.

best_of integer or null Optional Defaults to 1

Generates best_of completions server-side and returns the "best" (the one with the highest log probability per token). Results cannot be streamed.

When used with $\begin{bmatrix} n \end{bmatrix}$, $\begin{bmatrix} best_of \end{bmatrix}$ controls the number of candidate completions and $\begin{bmatrix} n \end{bmatrix}$ specifies how many to return $\begin{bmatrix} best_of \end{bmatrix}$ must be greater than $\begin{bmatrix} n \end{bmatrix}$.

Note: Because this parameter generates many completions, it can quickly consume your token quota. Use carefully and ensure that you have reasonable settings for max_tokens and stop.

echo boolean or null Optional Defaults to false Echo back the prompt in addition to the completion

frequency_penalty number or null Optional Defaults to 0 Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency in the text so far, decreasing the model's likelihood to repeat the same line verbatim.

See more information about frequency and presence penalties.

logit_bias map Optional Defaults to null
Modify the likelihood of specified tokens appearing in the
completion.

Accepts a JSON object that maps tokens (specified by their token ID in the GPT tokenizer) to an associated bias value from -100 to 100. You can use this **tokenizer tool** to convert text to token IDs. Mathematically, the bias is added to the logits generated by the model prior to sampling. The exact effect will vary per model, but values between -1 and 1 should decrease or increase likelihood of selection; values like -100 or 100 should result in a ban or exclusive selection of the relevant token.

As an example, you can pass {"50256": -100} to prevent the <lendoftextl> token from being generated.

logprobs integer or null Optional Defaults to null Include the log probabilities on the logprobs most likely output tokens, as well the chosen tokens. For example, if logprobs is 5, the API will return a list of the 5 most likely tokens. The API will always return the logprob of the sampled token, so there may be up to logprobs+1 elements in the response.

```
5 model="gpt-3.5-turbo-instruct",
6 prompt="Say this is a test",
7 max_tokens=7,
8 temperature=0
9 )
```

仚

```
Response gpt-3.5-turbo-instruct ∨
```

```
1
   {
2
     "id": "cmpl-uqkvlQyYK7bGYrRHQ0eXlWi7",
     "object": "text_completion",
3
4
     "created": 1589478378,
5
     "model": "gpt-3.5-turbo-instruct",
6
     "system_fingerprint": "fp_44709d6fcb",
7
     "choices": [
Я
      {
9
        "text": "\n\nThis is indeed a test",
        "index": 0,
10
11
        "logprobs": null,
        "finish_reason": "length"
12
13
      }
14
    ],
15
     "usage": {
16
      "prompt_tokens": 5,
17
      "completion_tokens": 7,
```

The maximum value for logprobs is 5.

max_tokens integer or null Optional Defaults to 16
The maximum number of tokens that can be generated in the completion.

The token count of your prompt plus max_tokens cannot exceed the model's context length. **Example Python code** for counting tokens.

n integer or null Optional Defaults to 1How many completions to generate for each prompt.

Note: Because this parameter generates many completions, it can quickly consume your token quota. Use carefully and ensure that you have reasonable settings for max_tokens and stop.

presence_penalty number or null Optional Defaults to 0 Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far, increasing the model's likelihood to talk about new topics.

See more information about frequency and presence penalties.

seed integer or null Optional

If specified, our system will make a best effort to sample deterministically, such that repeated requests with the same seed and parameters should return the same result.

Determinism is not guaranteed, and you should refer to the system_fingerprint response parameter to monitor changes in the backend.

stop string / array / null Optional Defaults to null Up to 4 sequences where the API will stop generating further tokens. The returned text will not contain the stop sequence.

stream boolean or null Optional Defaults to false
Whether to stream back partial progress. If set, tokens will be sent as data-only **server-sent events** as they become available, with the stream terminated by a data: [DONE] message. **Example Python code**.

suffix string or null Optional Defaults to null
The suffix that comes after a completion of inserted text.

This parameter is only supported for gpt-3.5-turbo-instruct.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

We generally recommend altering this or top_p but not both.

top_p number or null Optional Defaults to 1

An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

user string Optional

A unique identifier representing your end-user, which can help OpenAl to monitor and detect abuse. **Learn more**.

Returns

Returns a **completion** object, or a sequence of completion objects if the request is streamed.

The completion object Legacy €

Represents a completion response from the API. Note: both the streamed and non-streamed response objects share the same shape (unlike the chat endpoint).

id string

A unique identifier for the completion.

choices array

The list of completion choices the model generated for the input prompt.

√Show properties

created integer

The Unix timestamp (in seconds) of when the completion was created.

model string

The model used for completion.

system_fingerprint string

This fingerprint represents the backend configuration that the model runs with.

Can be used in conjunction with the seed request parameter to understand when backend changes have been made that might impact determinism.

object string

The object type, which is always "text_completion"

usage object

Usage statistics for the completion request.

∨ Show properties

The completion object

O

```
1
2
     "id": "cmpl-uqkvlQyYK7bGYrRHQ0eXlWi7",
3
     "object": "text_completion",
     "created": 1589478378,
4
     "model": "gpt-4-turbo",
6
     "choices": [
7
      {
8
        "text": "\n\nThis is indeed a test",
q
       "index": 0,
10
       "logprobs": null,
        "finish_reason": "length"
11
12
13 ],
14
     "usage": {
15
      "prompt_tokens": 5,
      "completion_tokens": 7,
      "total_tokens": 12
17
18 }
19 }
```

Example request

Response

Assistants (v1) Legacy &

Build assistants that can call models and use tools to perform tasks.

Get started with the Assistants API

Create assistant (v1) Legacy ∂

POST https://api.openai.com/v1/assistants

Create an assistant with a model and instructions.

Request body

model string Required

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them. type: string

name string or null Optional

The name of the assistant. The maximum length is 256 characters.

description string or null Optional

The description of the assistant. The maximum length is 512 characters.

instructions string or null Optional

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array Optional Defaults to []

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, retrieval, or function.

∨Show possible types

file ids array Optional Defaults to []

A list of **file** IDs attached to this assistant. There can be a maximum of 20 files attached to the assistant. Files are ordered by their creation date in ascending order.

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1 An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the toke

sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

```
Code Interpreter Files
```

python ~

O

币

```
from openai import OpenAl
client = OpenAl()

my_assistant = client.beta.assistants.create(
    instructions="You are a personal math tutor. When asked a coname="Math Tutor",
    tools=[{"type": "code_interpreter"}],
    model="gpt-4-turbo",
    )

print(my_assistant)
```

```
1
2
     "id": "asst_abc123",
3
     "object": "assistant",
     "created_at": 1698984975,
5
     "name": "Math Tutor",
     "description": null,
     "model": "gpt-4-turbo",
7
     "instructions": "You are a personal math tutor. When asked
9
     "tools": [
10
         "type": "code_interpreter"
11
12
13
     ],
     "file_ids": [],
14
     "metadata": {}.
15
```

We generally recommend altering this or temperature but not both.

response format string or object Optional Specifies the format that the model must output. Compatible with GPT-4o, GPT-4 Turbo, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106 .

Setting to { "type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if | finish_reason="length" |, which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

An assistant object.

Create assistant file (v1) Legacy ∂

POST https://api.openai.com/v1/assistants/{assistant_id}/files

Create an assistant file by attaching a File to an assistant.

Path parameters

assistant_id string Required The ID of the assistant for which to create a File.

Request body

```
file_id string Required
A File ID (with purpose="assistants" ) that the assistant should
use. Useful for tools like retrieval and code_interpreter that can
access files.
```

Returns

An assistant file object.

List assistants (v1) Legacy &

GET https://api.openai.com/v1/assistants

Returns a list of assistants.

Query parameters

```
O
Example request
                                            python ~
1 from openai import OpenAI
2 client = OpenAI()
3
4 assistant_file = client.beta.assistants.files.create(
   assistant_id="asst_abc123",
   file_id="file-abc123"
6
7)
8 print(assistant_file)
                                                         币
Response
```

```
1 {
   "id": "file-abc123",
2
   "object": "assistant.file",
```

"created_at": 1699055364, "assistant_id": "asst_abc123"

Example request

python ~



1 from openai import OpenAI

2 client = OpenAI()

6 }

Response

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

Returns

A list of assistant objects.

List assistant files (v1) Legacy &

GET https://api.openai.com/v1/assistants/{assistant_id}/files

Returns a list of assistant files.

Path parameters

assistant_id string Required

The ID of the assistant the file belongs to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

```
4 my_assistants = client.beta.assistants.list(
5 order="desc",
6 limit="20",
7 )
8 print(my_assistants.data)
```

币

python ~

币

币

```
1
    {
2
      "object": "list",
3
     "data": [
4
        "id": "asst_abc123",
6
        "object": "assistant",
7
        "created_at": 1698982736,
R
        "name": "Coding Tutor",
        "description": null,
10
        "model": "gpt-4-turbo",
11
        "instructions": "You are a helpful assistant designed to m
12
        "tools": [],
13
        "file_ids": [],
        "metadata": {},
14
15
        "top_p": 1.0,
16
        "temperature": 1.0,
17
        "response_format": "auto"
18
      }.
19
      {
```

```
Example request
```

1 from openai import OpenAl2 client = OpenAl()

3

4 assistant_files = client.beta.assistants.files.list(

5 assistant_id="asst_abc123"

6)

7 print(assistant_files)

Response

16],

```
1
     "object": "list",
2
3
      "data": [
4
      {
5
        "id": "file-abc123",
6
        "object": "assistant.file",
7
        "created_at": 1699060412,
8
        "assistant id": "asst abc123"
q
      },
10
      {
        "id": "file-abc456",
11
12
        "object": "assistant.file",
13
        "created_at": 1699060412,
14
        "assistant_id": "asst_abc123"
15
      }
```

Example request

18 }

2

3

7)

Example request

1 from openai import OpenAI

file_id="file-abc123"

8 print(assistant_file)

assistant_id="asst_abc123",

assistant_file = client.beta.assistants.files.retrieve(

client = OpenAI()

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

17 "first_id": "file-abc123", "last_id": "file-abc456", "has_more": false 20 }

Returns

A list of assistant file objects.

Retrieve assistant (v1)

GET https://api.openai.com/v1/assistants/{assistant_id}

Retrieves an assistant.

Path parameters

assistant_id string Required The ID of the assistant to retrieve.

Returns

The assistant object matching the specified ID.

GET https://api.openai.com/v1/assistants/{assistant_id}/files/{file_id}

Retrieves an AssistantFile.

Path parameters

assistant_id string Required The ID of the assistant who the file belongs to.

Retrieve assistant file (v1)

file_id string Required

```
1 from openai import OpenAI
2 client = OpenAI()
4 my_assistant = client.beta.assistants.retrieve("asst_abc123")
5 print(my_assistant)
                                                             凸
Response
1
    {
2
      "id": "asst_abc123",
      "object": "assistant",
      "created_at": 1699009709,
     "name": "HR Helper",
      "description": null,
6
7
      "model": "gpt-4-turbo",
8
      "instructions": "You are an HR bot, and you have access to file
9
      "tools": [
10
      {
         "type": "retrieval"
11
12
13
      "file_ids": [
14
       "file-abc123"
15
16
     1,
17
      "metadata": {}
```

D

python ~

币

python ✓

The ID of the file we're getting.

Returns

The assistant file object matching the specified ID.

Modify assistant (v1) Legacy ∂

POST https://api.openai.com/v1/assistants/{assistant_id}

Modifies an assistant.

Path parameters

assistant_id string Required
The ID of the assistant to modify.

Request body

model Optional

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them. type: string

name string or null Optional

The name of the assistant. The maximum length is 256 characters.

description string or null Optional

The description of the assistant. The maximum length is 512 characters.

instructions string or null Optional

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array Optional Defaults to []

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, retrieval, or function.

>Show possible types

file_ids array Optional Defaults to []

A list of **File** IDs attached to this assistant. There can be a maximum of 20 files attached to the assistant. Files are ordered by their creation date in ascending order. If a file was previously attached to the list but does not show up in the list, it will be deleted from the assistant.

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in

```
Response
```

```
1 {
2  "id": "file-abc123",
3  "object": "assistant.file",
4  "created_at": 1699055364,
5  "assistant_id": "asst_abc123"
6 }
```

Example request

python ✓



币

```
from openai import OpenAI
1
2
   client = OpenAI()
    my_updated_assistant = client.beta.assistants.update(
     "asst abc123".
6
     instructions="You are an HR bot, and you have access to files
7
     name="HR Helper",
     tools=[{"type": "retrieval"}],
8
     model="gpt-4-turbo",
10
    file_ids=["file-abc123", "file-abc456"],
11 )
12
13 print(my_updated_assistant)
```

Response



```
1
   {
2
     "id": "asst_abc123",
     "object": "assistant",
     "created_at": 1699009709,
     "name": "HR Helper",
6
     "description": null,
7
     "model": "gpt-4-turbo",
8
     "instructions": "You are an HR bot, and you have access to
9
     "tools": [
10
      {
        "type": "retrieval"
11
12
13
    ],
     nena nalani r
```

a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

response_format string or object Optional Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to [{"type": "json_object"}] enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

The modified assistant object.

Delete assistant (v1) Legacy ∂

DELETE https://api.openai.com/v1/assistants/{assistant_id}

Delete an assistant.

Path parameters

assistant_id string Required
The ID of the assistant to delete.

Returns

Deletion status

```
Example request python > 

1 from openai import OpenAI
2 client = OpenAI()
3
4 response = client.beta.assistants.delete("asst_abc123")
5 print(response)

Response

1 {
2 "id": "asst_abc123",
3 "object": "assistant.deleted",
4 "deleted": true
```

5 }

1 {

5 }

id: "file-abc123",

deleted: true

object: "assistant.file.deleted",

Delete assistant file (v1) Legacy ∂

DELETE https://api.openai.com/v1/assistants/{assistant_id}/files/{file id}

Delete an assistant file.

Path parameters

assistant_id string Required

The ID of the assistant that the file belongs to.

file_id string Required

The ID of the file to delete.

Returns

Deletion status

The assistant object (v1) Legacy €

Represents an assistant that can call the model and use tools.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always assistant.

created at integer

The Unix timestamp (in seconds) for when the assistant was created.

name string or null

The name of the assistant. The maximum length is 256 characters.

description string or null

The description of the assistant. The maximum length is 512 characters.

model

ID of the model to use. You can use the **List models** API to see all of your available models, or see our **Model overview** for descriptions of them. type: string

instructions string or null

The system instructions that the assistant uses. The maximum length is 256,000 characters.

tools array

A list of tool enabled on the assistant. There can be a maximum of 128 tools per assistant. Tools can be of types code_interpreter, retrieval, or function.

```
ð
The assistant object (v1)
1
2
      "id": "asst_abc123",
3
      "object": "assistant",
      "created_at": 1698984975,
5
      "name": "Math Tutor",
6
      "description": null,
      "model": "gpt-4-turbo",
8
     "instructions": "You are a personal math tutor. When asked a
      "tools": [
10
11
         "type": "code_interpreter"
12
       }
13
     ],
14
     "file_ids": [],
      "metadata": {},
      "top_p": 1.0,
16
     "temperature": 1.0,
     "response_format": "auto"
18
19 }
```

∨Show possible types

file_ids array

A list of **file** IDs attached to this assistant. There can be a maximum of 20 files attached to the assistant. Files are ordered by their creation date in ascending order.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null

What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top p number or null

An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

response_format string or object

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to {"type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

The assistant file object (v1) Legacy ∂

A list of Files attached to an assistant.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always assistant.file.

created_at integer

The Unix timestamp (in seconds) for when the assistant file was created.

assistant_id string

The assistant file object (v1)

```
1 {
2 "id": "file-abc123",
3 "object": "assistant.file",
4 "created_at": 1699055364,
5 "assistant_id": "asst_abc123"
6 }
```

84/106

The assistant ID that the file is attached to.

Threads (v1) Legacy ∂

Create threads that assistants can interact with.

Related guide: Assistants

Create thread (v1) Legacy &

POST https://api.openai.com/v1/threads

Create a thread.

Request body

messages array Optional

A list of messages to start the thread with.

√Show properties

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

A thread object.

Retrieve thread (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}

Retrieves a thread.

Path parameters

thread_id string Required
The ID of the thread to retrieve.

Returns

The thread object matching the specified ID.

Modify thread (v1) Legacy ∂

POST https://api.openai.com/v1/threads/{thread_id}

Modifies a thread.

Path parameters

Example request

python ✓ ਿ

Messages

Empty

```
1 from openai import OpenAI
```

- 2 client = OpenAI()
- 3
- 4 empty_thread = client.beta.threads.create()
- 5 print(empty_thread)

Response

6 }



```
1 {
2 "id": "thread_abc123",
3 "object": "thread",
4 "created_at": 1699012949,
5 "metadata": {}
```

Example request

python ~



```
1 from openai import OpenAI
```

- 2 client = OpenAI()
- 3
- 4 my_thread = client.beta.threads.retrieve("thread_abc123")
- 5 print(my_thread)

Response



```
1 {
2  "id": "thread_abc123",
3  "object": "thread",
4  "created_at": 1699014083,
5  "metadata": {}
6 }
```

Example request

python ~



1 from openai import OpenAI

- 2 client = OpenAI()
- 3
 - 4 my_updated_thread = client.beta.threads.update(

Example request

Response

6 }

thread_id string Required

The ID of the thread to modify. Only the metadata can be modified.

Request body

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

The modified thread object matching the specified ID.

Delete thread (v1) Legacy &

DELETE https://api.openai.com/v1/threads/{thread_id}

Delete a thread.

Path parameters

thread_id string Required
The ID of the thread to delete.

Returns

Deletion status

The thread object (v1) Legacy &

Represents a thread that contains messages.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.

created_at integer

The Unix timestamp (in seconds) for when the thread was created.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

```
5  "thread_abc123",
6  metadata={
7    "modified": "true",
8    "user": "abc123"
9  }
10 )
11  print(my_updated_thread)
```

```
Response

1 {
2  "id": "thread_abc123",
3  "object": "thread",
4  "created_at": 1699014083,
5  "metadata": {
6   "modified": "true",
7   "user": "abc123"
8  }
9 }
```

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 response = client.beta.threads.delete("thread_abc123")
5 print(response)
```

```
1 {
2 "id": "thread_abc123",
3 "object": "thread.deleted",
4 "deleted": true
5 }
```

```
The thread object (v1)

1 {
2  "id": "thread_abc123",
3  "object": "thread",
4  "created_at": 1698107661,
5  "metadata": {}
```

D

凸

python ~

Response

Messages (v1) Legacy ∂

Create messages within threads

Related guide: Assistants

Create message (v1) Legacy &

POST https://api.openai.com/v1/threads/{thread_id}/messages

Create a message.

Path parameters

thread_id string Required

The ID of the thread to create a message for.

Request body

role string Required

The role of the entity that is creating the message. Allowed values include:

user: Indicates the message is sent by an actual user and should be used in most cases to represent user-generated messages.

assistant: Indicates the message is generated by the assistant. Use this value to insert messages from the assistant into the conversation.

content string Required

The content of the message.

file_ids array Optional Defaults to []

A list of **File** IDs that the message should use. There can be a maximum of 10 files attached to a message. Useful for tools like retrieval and code_interpreter that can access and use files.

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

A message object.

List messages (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/messages

Returns a list of messages for a given thread.

Path parameters

```
from openai import OpenAl
client = OpenAl()
thread_message = client.beta.threads.messages.create(
thread_abc123",
role="user",
content="How does Al work? Explain it in simple terms.",
print(thread_message)
```

```
1
    {
     "id": "msg_abc123",
2
     "object": "thread.message",
     "created_at": 1699017614,
4
     "thread_id": "thread_abc123",
6
     "role": "user",
7
     "content": [
8
      {
9
        "type": "text",
10
        "text": {
         "value": "How does AI work? Explain it in simple terms."
11
         "annotations": []
12
13
        }
14
      }
15
16
     "file_ids": [],
     "assistant_id": null,
     11...... (1111)
```

python ~



凸

1 from openai import OpenAl

2 client = OpenAI()
3

Response

thread_id string Required

The ID of the **thread** the messages belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

run_id string Optional

Filter messages by the run ID that generated them.

Returns

A list of **message** objects.

List message files (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/messages/{message_id}/files

Returns a list of message files.

Path parameters

thread_id string Required

The ID of the thread that the message and files belong to.

message_id string Required

The ID of the message that the files belongs to.

Query parameters

limit integer Optional Defaults to 20

```
4 thread_messages = client.beta.threads.messages.list("thread
```

```
1
    {
2
      "object": "list",
      "data": [
4
5
        "id": "msg_abc123",
6
        "object": "thread.message",
7
        "created_at": 1699016383,
8
        "thread_id": "thread_abc123",
9
        "role": "user",
10
        "content": [
11
12
           "type": "text",
13
           "text": {
             "value": "How does AI work? Explain it in simple term
14
15
            "annotations": []
16
           }
17
         }
18
19
        "file_ids": [],
20
        "assistant_id": null,
21
        "run_id": null,
```

```
Example request
```

python ~



凸

```
from openai import OpenAl
client = OpenAl()
message_files = client.beta.threads.messages.files.list(
thread_id="thread_abc123",
message_id="msg_abc123"

print(message_files)
```

```
Response
```



A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

```
order string Optional Defaults to desc
```

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list

Returns

A list of message file objects.

Retrieve message (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/messages/{message_id}

Retrieve a message.

Path parameters

thread id string Required

The ID of the thread to which this message belongs.

message_id string Required

The ID of the message to retrieve.

Returns

The message object matching the specified ID.

```
"object": "thread.message.file",
        "created_at": 1699061776,
8
        "message_id": "msg_abc123"
9
      },
10
      {
        "id": "file-abc123".
11
        "object": "thread.message.file",
12
13
        "created_at": 1699061776,
14
        "message_id": "msg_abc123"
15
16
     ],
     "first_id": "file-abc123",
17
     "last_id": "file-abc123",
18
     "has_more": false
```

```
Example request
```

python ✓



```
1 from openai import OpenAl
2 client = OpenAl()
3
4 message = client.beta.threads.messages.retrieve(
5 message_id="msg_abc123",
6 thread_id="thread_abc123",
7 )
8 print(message)
```

Response



```
1
2
     "id": "msq_abc123",
3
     "object": "thread.message",
     "created_at": 1699017614,
     "thread_id": "thread_abc123",
5
6
     "role": "user",
7
     "content": [
8
        "type": "text".
9
10
        "text": {
11
         "value": "How does AI work? Explain it in simple terms."
         "annotations": []
12
13
14
      }
15
     ],
     "file_ids": [],
16
     "assistant_id": null,
```

18 "run_id": null,

Retrieve message file (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/messages/{message_id}/file_id}

Retrieves a message file.

Path parameters

thread_id string Required

The ID of the thread to which the message and File belong.

message_id string Required

The ID of the message the file belongs to.

file_id string Required

The ID of the file being retrieved.

Returns

The message file object.

Modify message (v1) Legacy ∂

POST https://api.openai.com/v1/threads/{thread_id}/messages/{message_id}

Modifies a message.

Path parameters

thread_id string Required

The ID of the thread to which this message belongs.

message_id string Required

The ID of the message to modify.

Request body

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

The modified message object.

```
Example request
```

python ~

```
· fi
```

```
from openai import OpenAl
client = OpenAl()
message_files = client.beta.threads.messages.files.retrieve(
thread_id="thread_abc123",
message_id="msg_abc123",
file_id="file-abc123"
)
print(message_files)
```

Response

ப

```
1 {
2  "id": "file-abc123",
3  "object": "thread.message.file",
4  "created_at": 1699061776,
5  "message_id": "msg_abc123"
6 }
```

Example request

python ~



```
from openai import OpenAI
2
   client = OpenAI()
3
4
   message = client.beta.threads.messages.update(
5
     message_id="msg_abc12",
6
     thread_id="thread_abc123",
7
     metadata={
8
      "modified": "true",
9
      "user": "abc123",
10 },
11 )
12 print(message)
```

Response

Ð

```
1
     "id": "msg_abc123",
3
     "object": "thread.message",
4
     "created_at": 1699017614,
     "thread_id": "thread_abc123",
     "role": "user",
6
7
     "content": [
R
9
        "type": "text",
10
11
         "value": "How does AI work? Explain it in simple terms."
12
         "annotations": []
13
       }
14
      }
```

The message object (v1) Legacy &

Represents a message within a thread.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.message .

created at integer

The Unix timestamp (in seconds) for when the message was created.

thread id string

The thread ID that this message belongs to.

status string

The status of the message, which can be either in_progress, incomplete, or completed.

incomplete_details object or null

On an incomplete message, details about why the message is incomplete.

∨Show properties

completed_at integer or null

The Unix timestamp (in seconds) for when the message was completed.

incomplete_at integer or null

The Unix timestamp (in seconds) for when the message was marked as incomplete.

role string

The entity that produced the message. One of user or assistant.

content array

The content of the message in array of text and/or images.

∨Show possible types

assistant id string or null

If applicable, the ID of the assistant that authored this message.

run id string or null

The ID of the **run** associated with the creation of this message. Value is null when messages are created manually using the create message or create thread endpoints.

file ids array

A list of **file** IDs that the assistant should use. Useful for tools like retrieval and code_interpreter that can access files. A maximum of 10 files can be attached to a message.

```
The message object (v1)
```

```
1
   {
2
     "id": "msg_abc123",
3
     "object": "thread.message",
     "created_at": 1698983503,
     "thread_id": "thread_abc123",
     "role": "assistant",
6
7
     "content": [
R
        "type": "text",
9
10
        "text": {
11
         "value": "Hi! How can I help you today?",
12
         "annotations": []
13
14
      }
15
16
     "file_ids": [],
     "assistant_id": "asst_abc123",
17
     "run_id": "run_abc123",
19
     "metadata": {}
20 }
```



metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

The message file object (v1) Legacy ∂

A list of files attached to a message.

id string

The identifier, which can be referenced in API endpoints.

object string

The object type, which is always thread.message.file .

created_at integer

The Unix timestamp (in seconds) for when the message file was created.

message_id string

The ID of the **message** that the **File** is attached to.

Runs (v1) Legacy &

Represents an execution run on a thread.

Related guide: Assistants

Create run (v1) Legacy &

POST https://api.openai.com/v1/threads/{thread_id}/runs

Create a run.

Path parameters

thread_id string Required
The ID of the thread to run.

Request body

assistant id string Required

The ID of the assistant to use to execute this run.

model string Optional

The ID of the **Model** to be used to execute this run. If a value is provided here, it will override the model associated with the assistant. If not, the model associated with the assistant will be used.

instructions string or null Optional

Overrides the **instructions** of the assistant. This is useful for modifying the behavior on a per-run basis.

additional_instructions string or null Optional

```
The message file object (v1)

1 {
2 "id": "file-abc123",
3 "object": "thread.message.file",
4 "created_at": 1698107661,
5 "message_id": "message_QLoltBbqwyAJEzlTy4y9kOMM",
6 "file_id": "file-abc123"
7 }
```

Default Streaming Streaming with Functions

Example request

python ~

ð

```
from openai import OpenAl
client = OpenAl()

run = client.beta.threads.runs.create(
thread_id="thread_abc123",
assistant_id="asst_abc123"

print(run)
```

Response



```
{
1
     "id": "run_abc123",
2
3
     "object": "thread.run",
     "created_at": 1699063290,
5
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
     "status": "queued",
8
     "started_at": 1699063290,
     "expires_at": null,
10
     "cancelled_at": null,
11
     "failed at": null,
     "completed_at": 1699063291,
```

Appends additional instructions at the end of the instructions for the run. This is useful for modifying the behavior on a per-run basis without overriding other instructions.

additional_messages array or null Optional

Adds additional messages to the thread before creating the run.

∨Show properties

tools array or null Optional

Override the tools the assistant can use for this run. This is useful for modifying the behavior on a per-run basis.

∨Show possible types

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maxium of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1 An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

max_prompt_tokens integer or null Optional

The maximum number of prompt tokens that may be used over the course of the run. The run will make a best effort to use only the number of prompt tokens specified, across multiple turns of the run. If the run exceeds the number of prompt tokens specified, the run will end with status complete . See incomplete_details for more info.

max completion tokens integer or null Optional

The maximum number of completion tokens that may be used over the course of the run. The run will make a best effort to use only the number of completion tokens specified, across multiple turns of the run. If the run exceeds the number of completion tokens specified, the run will end with status | complete |. See incomplete_details for more info.

truncation_strategy object Optional

∨ Show properties

tool_choice string or object Optional

- "last_error": null, 13
- 14 "model": "gpt-4-turbo",

"incomplete_details": null,

- 15 "instructions": null,
- 17 "toolo": [

16

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling a tool. Specifying a particular tool like {"type": "TOOL_TYPE"} or {"type": "function", "function": {"name": "my_function"}} forces the model to call that tool.

Show possible types

response_format string or object Optional Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to {"type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

A run object.

Create thread and run (v1) Legacy &

POST https://api.openai.com/v1/threads/runs

Create a thread and run it in one request.

Request body

assistant_id string Required

The ID of the assistant to use to execute this run.

thread object Optional

∨Show properties

model string Optional

The ID of the **Model** to be used to execute this run. If a value is provided here, it will override the model associated with the assistant. If not, the model associated with the assistant will be used.

instructions string or null Optional

Override the default system message of the assistant. This is useful for modifying the behavior on a per-run basis.

tools array or null Optional

```
Default Streaming Streaming with Functions
```

```
Example request
                                             python ~
                                                          仚
    from openai import OpenAI
2
    client = OpenAI()
3
   run = client.beta.threads.create_and_run(
     assistant_id="asst_abc123",
5
6
     thread={
7
      "messages": [
8
        {"role": "user", "content": "Explain deep learning to a 5 ye
9
10
     }
11 )
12
13 print(run)
```

```
Response 

1 {
2  "id": "run_abc123",
3  "object": "thread.run",
4  "created_at": 1699076792,
5  "assistant_id": "asst_abc123",
```

Override the tools the assistant can use for this run. This is useful for modifying the behavior on a per-run basis.

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

temperature number or null Optional Defaults to 1 What sampling temperature to use, between 0 and 2. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic.

top_p number or null Optional Defaults to 1
An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with top_p probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered.

We generally recommend altering this or temperature but not both.

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

max_prompt_tokens integer or null Optional

The maximum number of prompt tokens that may be used over the course of the run. The run will make a best effort to use only the number of prompt tokens specified, across multiple turns of the run. If the run exceeds the number of prompt tokens specified, the run will end with status complete. See incomplete_details for more info.

max completion tokens integer or null Optional

The maximum number of completion tokens that may be used over the course of the run. The run will make a best effort to use only the number of completion tokens specified, across multiple turns of the run. If the run exceeds the number of completion tokens specified, the run will end with status complete. See incomplete_details for more info.

truncation_strategy object Optional

∨Show properties

tool_choice string or object Optional

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling a tool. Specifying a particular tool like {"type": "TOOL_TYPE"} or {"type": "function", "function": {"name": "my_function"}} forces the model to call that tool.

∨Show possible types

response_format string or object Optional

- 6 "thread_id": "thread_abc123",
- 7 "status": "queued",
- 8 "started_at": null,
- 9 "expires_at": 1699077392,
- 10 "cancelled_at": null,
- 11 "failed_at": null,
- 12 "completed_at": null,
- 13 "last_error": null,

Response

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to {"type": "json_object" } enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length", which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

Returns

A run object.

List runs (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/runs

Returns a list of runs belonging to a thread.

Path parameters

thread_id string Required
The ID of the thread the run belongs to.

Query parameters

limit integer Optional Defaults to 20
A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

```
python ✓ ☐

from openai import OpenAl

client = OpenAl()

runs = client.beta.threads.runs.list(

"thread_abc123"

h)

print(runs)
```

```
1
   {
2
     "object": "list",
     "data": [
4
      {
5
        "id": "run_abc123",
6
        "object": "thread.run",
7
        "created_at": 1699075072,
8
        "assistant_id": "asst_abc123",
9
        "thread_id": "thread_abc123",
10
        "status": "completed",
11
        "started_at": 1699075072,
12
        "expires at": null,
13
        "cancelled_at": null,
14
        "failed_at": null,
15
        "completed_at": 1699075073,
16
        "last_error": null,
17
        "model": "gpt-4-turbo",
18
        "instructions": null,
19
        "incomplete_details": null,
20
        "tools": [
```

币

Returns

A list of run objects.

List run steps (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/step

Returns a list of run steps belonging to a run.

Path parameters

thread_id string Required

The ID of the thread the run and run steps belong to.

run_id string Required

The ID of the run the run steps belong to.

Query parameters

limit integer Optional Defaults to 20

A limit on the number of objects to be returned. Limit can range between 1 and 100, and the default is 20.

order string Optional Defaults to desc

Sort order by the created_at timestamp of the objects. asc for ascending order and desc for descending order.

after string Optional

A cursor for use in pagination. after is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include after=obj_foo in order to fetch the next page of the list.

before string Optional

A cursor for use in pagination. before is an object ID that defines your place in the list. For instance, if you make a list request and receive 100 objects, ending with obj_foo, your subsequent call can include before=obj_foo in order to fetch the previous page of the list.

Returns

A list of run step objects.

Retrieve run (v1) Legacy &

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}

Retrieves a run.

Path parameters

```
Example request
```

python ~

 \Box

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 run_steps = client.beta.threads.runs.steps.list(
     thread_id="thread_abc123",
     run_id="run_abc123"
7)
9 print(run_steps)
```

Response



```
1
2
     "object": "list",
     "data": [
4
5
        "id": "step_abc123",
6
        "object": "thread.run.step",
        "created_at": 1699063291,
        "run_id": "run_abc123",
9
        "assistant_id": "asst_abc123",
        "thread_id": "thread_abc123",
11
        "type": "message_creation",
12
        "status": "completed",
13
        "cancelled_at": null,
14
        "completed_at": 1699063291,
15
        "expired_at": null,
16
        "failed_at": null,
17
        "last_error": null,
18
        "step_details": {
19
         "type": "message_creation",
```

Example request

python ~



1 from openai import OpenAI

2 client = OpenAI()

3

https://platform.openai.com/docs/api-reference/runs-v1

thread_id string Required

The ID of the **thread** that was run.

run_id string Required

The ID of the run to retrieve.

Returns

The run object matching the specified ID.

Retrieve run step (v1) Legacy ∂

GET https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/steps/{step_id}

Retrieves a run step.

Path parameters

thread_id string Required

The ID of the thread to which the run and run step belongs.

run_id string Required

The ID of the run to which the run step belongs.

step_id string Required

The ID of the run step to retrieve.

Returns

The run step object matching the specified ID.

```
4 run = client.beta.threads.runs.retrieve(
5 thread_id="thread_abc123",
6 run_id="run_abc123"
7 )
8
9 print(run)
```

ð

```
Response
```

```
1
   {
2
     "id": "run_abc123",
3
     "object": "thread.run",
4
     "created_at": 1699075072,
5
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
7
     "status": "completed",
8
     "started_at": 1699075072,
     "expires_at": null,
9
10
     "cancelled_at": null,
11
     "failed_at": null,
12
     "completed_at": 1699075073,
     "last_error": null,
13
14
     "model": "gpt-4-turbo",
     "instructions": null,
15
     "incomplete_details": null,
     "tools": [
17
18
      {
19
        "type": "code_interpreter"
```

Example request

python ~

 \Box

币

```
from openai import OpenAl
client = OpenAl()

run_step = client.beta.threads.runs.steps.retrieve(
thread_id="thread_abc123",
 run_id="run_abc123",
step_id="step_abc123"

print(run_step)
```

Response

14

15

"last_error": null,
"step_details": {

```
1
     "id": "step_abc123",
2
     "object": "thread.run.step",
3
     "created_at": 1699063291,
4
5
     "run_id": "run_abc123",
     "assistant id": "asst abc123",
6
7
     "thread_id": "thread_abc123",
     "type": "message_creation",
9
     "status": "completed",
10
    "cancelled_at": null,
     "completed_at": 1699063291,
11
12
     "expired_at": null,
13
     "failed_at": null,
```

Example request

Modify run (v1) Legacy ∂

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}

Modifies a run.

Path parameters

thread_id string Required
The ID of the thread that was run.

run_id string RequiredThe ID of the run to modify.

Request body

metadata map Optional

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

Returns

The modified run object matching the specified ID.

Submit tool outputs to run (v1) Legacy &

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/sub_mit_tool_outputs

When a run has the status:

"requires_action" and required_action.type is submit_tool_outputs, this endpoint can be used to submit the outputs from the tool calls once they're all completed. All outputs must be submitted in a single request.

Path parameters

thread id string Required

The ID of the thread to which this run belongs.

run id string Required

The ID of the run that requires the tool output submission.

```
"type": "message_creation","message_creation": {
```

```
4 from annual immed Ones Al
```

```
from openai import OpenAl
client = OpenAl()

run = client.beta.threads.runs.update(
thread_id="thread_abc123",
 run_id="run_abc123",
 metadata={"user_id": "user_abc123"},

print(run)
```

Response

12

13 14

15

16

17

```
1
   {
2
     "id": "run_abc123",
     "object": "thread.run",
     "created_at": 1699075072,
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
7
     "status": "completed",
8
     "started_at": 1699075072,
9
     "expires_at": null,
10
     "cancelled at": null,
11
     "failed_at": null,
```

Default Streaming

币

币

python ~

Example request

"completed_at": 1699075073,

"last_error": null,

"instructions": null,

"tools": [{

"model": "gpt-4-turbo",

"incomplete_details": null,

```
python ~
```

ð

```
from openai import OpenAI
1
    client = OpenAI()
3
4
    run = client.beta.threads.runs.submit_tool_outputs(
5
     thread_id="thread_123",
6
     run_id="run_123",
7
     tool_outputs=[
8
        "tool_call_id": "call_001",
9
10
       "output": "70 degrees and sunny."
11
12
    ]
13
```

Response



Request body

tool_outputs array Required

A list of tools for which the outputs are being submitted.

∨Show properties

stream boolean or null Optional

If true, returns a stream of events that happen during the Run as server-sent events, terminating when the Run enters a terminal state with a data: [DONE] message.

Returns

The modified run object matching the specified ID.

Cancel a run (v1) Legacy &

POST https://api.openai.com/v1/threads/{thread_id}/runs/{run_id}/can cel

Cancels a run that is in_progress.

Path parameters

thread_id string Required

The ID of the thread to which this run belongs.

run_id string Required

The ID of the run to cancel.

Returns

The modified run object matching the specified ID.

```
1
      "id": "run 123",
2
     "object": "thread.run",
3
     "created_at": 1699075592,
5
     "assistant_id": "asst_123",
     "thread_id": "thread_123",
7
     "status": "queued",
8
     "started_at": 1699075592,
      "expires_at": 1699076192,
9
10
     "cancelled_at": null,
11
     "failed at": null,
12
     "completed_at": null,
     "last_error": null
```

```
Example request
```

python ~

O

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 run = client.beta.threads.runs.cancel(
5 thread_id="thread_abc123",
6 run_id="run_abc123"
7 )
8
9 print(run)
```

Response

```
币
```

```
1
   {
     "id": "run_abc123",
2
     "object": "thread.run",
     "created_at": 1699076126,
     "assistant_id": "asst_abc123",
6
     "thread_id": "thread_abc123",
7
     "status": "cancelling",
     "started_at": 1699076126,
8
9
     "expires_at": 1699076726,
10
     "cancelled_at": null,
11
     "failed_at": null,
12
     "completed at": null,
     "last_error": null,
13
14
     "model": "gpt-4-turbo",
     "instructions": "You summarize books.",
15
16
     "tools": [
17
18
        "type": "retrieval"
19
```

The run object (v1) Legacy ∂

Represents an execution run on a thread.

id string

The identifier, which can be referenced in API endpoints.

```
1 {
2 "id": "run_abc123",
3 "object": "thread.run",
```

The run object (v1)

凸

object string

The object type, which is always thread.run.

created_at integer

The Unix timestamp (in seconds) for when the run was created.

thread id string

The ID of the thread that was executed on as a part of this run.

assistant_id string

The ID of the assistant used for execution of this run.

status string

The status of the run, which can be either queued, in_progress, requires_action, cancelling, cancelled, failed, completed, or expired.

required_action object or null

Details on the action required to continue the run. Will be null if no action is required.

∨Show properties

last_error object or null

The last error associated with this run. Will be null if there are no errors.

∨Show properties

expires_at integer or null

The Unix timestamp (in seconds) for when the run will expire.

started at integer or null

The Unix timestamp (in seconds) for when the run was started.

cancelled_at integer or null

The Unix timestamp (in seconds) for when the run was cancelled.

failed at integer or null

The Unix timestamp (in seconds) for when the run failed.

completed_at integer or null

The Unix timestamp (in seconds) for when the run was completed.

incomplete_details object or null

Details on why the run is incomplete. Will be <u>null</u> if the run is not incomplete.

∨Show properties

model string

The model that the assistant used for this run.

instructions string

The instructions that the assistant used for this run.

tools array

```
API Reference - OpenAI API
```

35 }

```
"created_at": 1698107661,
     "assistant_id": "asst_abc123",
5
6
     "thread_id": "thread_abc123",
7
     "status": "completed",
8
     "started at": 1699073476,
     "expires_at": null,
10
     "cancelled_at": null,
11
     "failed at": null,
12
     "completed_at": 1699073498,
13
     "last_error": null,
     "model": "gpt-4-turbo",
14
     "instructions": null,
15
     "tools": [{"type": "retrieval"}, {"type": "code_interpreter"}],
16
17
18
     "metadata": {},
19
     "incomplete_details": null,
20
     "usage": {
21
       "prompt_tokens": 123,
22
       "completion_tokens": 456,
23
       "total_tokens": 579
24
     },
25
     "temperature": 1.0,
26
     "top_p": 1.0,
27
     "max_prompt_tokens": 1000,
28
     "max_completion_tokens": 1000,
29
     "truncation_strategy": {
30
      "type": "auto",
31
       "last_messages": null
32 },
     "response_format": "auto",
34
     "tool_choice": "auto"
```

The list of tools that the assistant used for this run.

∨Show possible types

file ids array

The list of File IDs the assistant used for this run.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

usage

temperature number or null

The sampling temperature used for this run. If not set, defaults to

top_p number or null

The nucleus sampling value used for this run. If not set, defaults to

max_prompt_tokens integer or null

The maximum number of prompt tokens specified to have been used over the course of the run.

max completion tokens integer or null

The maximum number of completion tokens specified to have been used over the course of the run.

truncation strategy object

∨Show properties

tool_choice string or object

Controls which (if any) tool is called by the model. none means the model will not call any tools and instead generates a message. auto is the default value and means the model can pick between generating a message or calling a tool. Specifying a particular tool like {"type": "TOOL_TYPE"} or

["type": "function", "function": {"name": "my_function"}}] forces the model to call that tool.

∨Show possible types

response_format string or object

Specifies the format that the model must output. Compatible with **GPT-4o**, **GPT-4 Turbo**, and all GPT-3.5 Turbo models since gpt-3.5-turbo-1106.

Setting to [{"type": "json_object"}] enables JSON mode, which guarantees the message the model generates is valid JSON.

Important: when using JSON mode, you must also instruct the model to produce JSON yourself via a system or user message. Without this, the model may generate an unending stream of whitespace until the generation reaches the token limit, resulting in a long-running and seemingly "stuck" request. Also note that the message content may be partially cut off if finish_reason="length",

which indicates the generation exceeded max_tokens or the conversation exceeded the max context length.

∨Show possible types

The run step object (v1) Legacy ∂

Represents a step in execution of a run.

id string

The identifier of the run step, which can be referenced in API endpoints.

object string

The object type, which is always thread.run.step

created_at integer

The Unix timestamp (in seconds) for when the run step was created.

assistant_id string

The ID of the assistant associated with the run step.

thread id string

The ID of the thread that was run.

run_id string

The ID of the run that this run step is a part of.

type string

The type of run step, which can be either message_creation or tool_calls .

status string

The status of the run step, which can be either in_progress, cancelled, failed, completed, or expired.

step_details object

The details of the run step.

∨Show possible types

last_error object or null

The last error associated with this run step. Will be null if there are no errors.

∨Show properties

expired_at integer or null

The Unix timestamp (in seconds) for when the run step expired. A step is considered expired if the parent run is expired.

cancelled_at integer or null

The Unix timestamp (in seconds) for when the run step was cancelled.

failed at integer or null

The Unix timestamp (in seconds) for when the run step failed.

```
The run step object (v1)
```

```
1
   {
2
     "id": "step_abc123",
     "object": "thread.run.step",
3
4
     "created_at": 1699063291,
5
     "run_id": "run_abc123",
     "assistant_id": "asst_abc123",
     "thread_id": "thread_abc123",
     "type": "message_creation",
     "status": "completed",
     "cancelled_at": null,
10
11
     "completed_at": 1699063291,
     "expired_at": null,
12
13
     "failed_at": null,
14
     "last_error": null,
     "step_details": {
15
16
      "type": "message_creation",
17
      "message_creation": {
18
        "message_id": "msg_abc123"
19
      }
20
     },
21
     "usage": {
22
      "prompt_tokens": 123,
      "completion_tokens": 456,
23
      "total_tokens": 579
25 }
26 }
```

completed_at integer or null

The Unix timestamp (in seconds) for when the run step completed.

metadata map

Set of 16 key-value pairs that can be attached to an object. This can be useful for storing additional information about the object in a structured format. Keys can be a maximum of 64 characters long and values can be a maximum of 512 characters long.

usage

Streaming (v1) Legacy &

Stream the result of executing a Run or resuming a Run after submitting tool outputs.

You can stream events from the **Create Thread and Run**, **Create Run**, and **Submit Tool Outputs** endpoints by passing "stream": true. The response will be a **Server-Sent events** stream.

Our Node and Python SDKs provide helpful utilities to make streaming easy. Reference the **Assistants API quickstart** to learn more.

The message delta object (v1) Legacy ℰ

Represents a message delta i.e. any changed fields on a message during streaming.

id string

The identifier of the message, which can be referenced in API endpoints.

object string

The object type, which is always thread.message.delta .

delta object

The delta containing the fields that have changed on the Message.

Show properties

```
币
The message delta object (v1)
1
      "id": "msg_123",
2
3
      "object": "thread.message.delta",
4
      "delta": {
       "content": [
5
6
7
          "index": 0.
8
          "type": "text",
          "text": { "value": "Hello", "annotations": [] }
10
11
       1
12 }
13 }
```

The run step delta object (v1) Legacy ℰ

Represents a run step delta i.e. any changed fields on a run step during streaming.

id string

The identifier of the run step, which can be referenced in API endpoints.

object string

The object type, which is always thread.run.step.delta .

delta object

The delta containing the fields that have changed on the run step. VShow properties

```
1
   {
2
     "id": "step_123",
3
     "object": "thread.run.step.delta",
4
      "delta": {
       "step_details": {
        "type": "tool_calls",
6
7
         "tool calls": [
R
         {
           "index": 0,
10
           "id": "call_123",
11
           "type": "code_interpreter",
           "code_interpreter": { "input": "", "outputs": [] }
12
13
         }
14
        ]
15
       }
```

The run step delta object (v1)

凸

16 } 17 }

Assistant stream events (v1) Legacy &

Represents an event emitted when streaming a Run.

Each event in a server-sent events stream has an event and data property:



We emit events whenever a new object is created, transitions to a new state, or is being streamed in parts (deltas). For example, we emit thread.run.created when a new run is created, thread.run.completed when a run completes, and so on. When an Assistant chooses to create a message during a run, we emit a thread.message.created event, a thread.message.in_progress event, many thread.message.delta events, and finally a thread.message.completed event.

We may add additional events over time, so we recommend handling unknown events gracefully in your code. See the **Assistants API quickstart** to learn how to integrate the Assistants API with streaming.

thread.created data is a thread
Occurs when a new thread is created.
thread.run.created data is a run
Occurs when a new run is created.
thread.run.queued data is a run
Occurs when a run moves to a queued status.
thread.run.in_progress data is a run
Occurs when a run moves to an in_progress status.
thread.run.requires_action data is a run
Occurs when a run moves to a requires_action status.
thread.run.completed data is a run
Occurs when a run is completed.
thread.run.failed data is a run
Occurs when a run fails.
thread.run.cancelling data is a run
Occurs when a run moves to a cancelling status.
thread.run.cancelled data is a run
Occurs when a run is cancelled.

thread.run.expired data is a run

Occurs when a run expires.

thread.run.step.created data is a run step Occurs when a run step is created. thread.run.step.in_progress | data | is a run step Occurs when a **run step** moves to an in_progress state. thread.run.step.delta data is a run step delta Occurs when parts of a run step are being streamed. thread.run.step.completed | data | is a run step Occurs when a run step is completed. thread.run.step.failed data is a run step Occurs when a run step fails. thread.run.step.cancelled data is a run step Occurs when a run step is cancelled. thread.run.step.expired data is a run step Occurs when a run step expires. thread.message.created data is a message Occurs when a message is created. thread.message.in_progress | data | is a message Occurs when a **message** moves to an in_progress state. thread.message.delta data is a message delta Occurs when parts of a Message are being streamed. thread.message.completed | data | is a message Occurs when a message is completed. thread.message.incomplete | data | is a message Occurs when a message ends before it is completed.

Occurs when an error occurs. This can happen due to an internal server error or a timeout

https://platform.openai.com/docs/api-reference/runs-v1

error data is an error