

Introduction

This API reference describes the RESTful, streaming, and realtime APIs you can use to interact with the OpenAI platform. REST APIs are usable via HTTP in any environment that supports HTTP requests. Language-specific SDKs are listed [on the libraries page](#).

Authentication

The OpenAI API uses API keys for authentication. Create, manage, and learn more about API keys in your [organization settings](#).

Remember that your API key is a secret! Do not share it with others or expose it in any client-side code (browsers, apps). API keys should be securely loaded from an environment variable or key management service on the server.

API keys should be provided via [HTTP Bearer authentication](#).

Authorization: Bearer OPENAI_API_KEY

If you belong to multiple organizations or access projects through a legacy user API key, pass a header to specify which organization and project to use for an API request:

1 curl https://api.openai.com/v1/models \
2 -H "Authorization: Bearer \$OPENAI_API_KEY" \
3 -H "OpenAI-Organization: org-4VsgZaw4tQ5kPK1gFZlB3cI7" \
4 -H "OpenAI-Project: \$PROJECT_ID"

Usage from these API requests counts as usage for the specified organization and project. 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.

Debugging requests

In addition to [error codes](#) returned from API responses, you can inspect HTTP response headers containing the unique ID of a particular API request or information about rate limiting applied to your requests. Below is an incomplete list of HTTP headers returned with API responses:

API meta information

- `openai-organization` : The [organization](#) associated with the request
- `openai-processing-ms` : Time taken processing your API request
- `openai-version` : REST API version used for this request (currently `2020-10-01`)
- `x-request-id` : Unique identifier for this API request (used in troubleshooting)

Rate limiting information

- `x-ratelimit-limit-requests`
- `x-ratelimit-limit-tokens`
- `x-ratelimit-remaining-requests`
- `x-ratelimit-remaining-tokens`
- `x-ratelimit-reset-requests`

- `x-ratelimit-reset-tokens`

OpenAI recommends logging request IDs in production deployments for more efficient troubleshooting with our [support team](#), should the need arise. Our [official SDKs](#) provide a property on top-level response objects containing the value of the `x-request-id` header.

Backward compatibility

OpenAI is committed to providing stability to API users by avoiding breaking changes in major API versions whenever reasonably possible. This includes:

- The REST API (currently `v1`)
- Our first-party [SDKs](#) (released SDKs adhere to [semantic versioning](#))
- [Model](#) families (like `gpt-4o` or `o4-mini`)

Model prompting behavior between snapshots is subject to change. Model outputs are by their nature variable, so expect changes in prompting and model behavior between snapshots. For example, if you moved from `gpt-4o-2024-05-13` to `gpt-4o-2024-08-06` , the same `system` or `user` messages could function differently between versions. The best way to ensure consistent prompting behavior and model output is to use pinned model versions, and to implement [evals](#) for your applications.

Backwards-compatible API changes:

- Adding new resources (URLs) to the REST API and SDKs
- Adding new optional API parameters
- Adding new properties to JSON response objects or event data
- Changing the order of properties in a JSON response object
- Changing the length or format of opaque strings, like resource identifiers and UUIDs
- Adding new event types (in either streaming or the Realtime API)

See the [changelog](#) for a list of backwards-compatible changes and rare breaking changes.

Responses

OpenAI's most advanced interface for generating model responses. Supports text and image inputs, and text outputs. Create stateful interactions with the model, using the output of previous responses as input. Extend the model's capabilities with built-in tools for file search, web search, computer use, and more. Allow the model access to external systems and data using function calling.

Related guides:

- [Quickstart](#)
- [Text inputs and outputs](#)
- [Image inputs](#)
- [Structured Outputs](#)
- [Function calling](#)
- [Conversation state](#)
- [Extend the models with tools](#)

Create a model response

POST `https://api.openai.com/v1/responses`

Text inputImage inputWeb searchFile search

Creates a model response. Provide **text** or **image** inputs to generate **text** or **JSON** outputs. Have the model call your own **custom code** or use built-in **tools** like **web search** or **file search** to use your own data as input for the model's response.

Request body

input string or array **Required**
Text, image, or file inputs to the model, used to generate a response.

Learn more:

- [Text inputs and outputs](#)
- [Image inputs](#)
- [File inputs](#)
- [Conversation state](#)
- [Function calling](#)

▼ Show possible types

model string **Required**
Model ID used to generate the response, like `gpt-4o` or `o3`. OpenAI offers a wide range of models with different capabilities, performance characteristics, and price points. Refer to the [model guide](#) to browse and compare available models.

include array or null **Optional**
Specify additional output data to include in the model response. Currently supported values are:

- `file_search_call.results`: Include the search results of the file search tool call.
- `message.input_image.image_url`: Include image urls from the input message.
- `computer_call_output.output.image_url`: Include image urls from the computer call output.

instructions string or null **Optional**
Inserts a system (or developer) message as the first item in the model's context.

When using along with `previous_response_id`, the instructions from a previous response will not be carried over to the next response. This makes it simple to swap out system (or developer) messages in new responses.

max_output_tokens integer or null **Optional**
An upper bound for the number of tokens that can be generated for a response, including visible output tokens and [reasoning tokens](#).

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, and querying for objects via API or the dashboard.

Keys are strings with a maximum length of 64 characters. Values are strings with a maximum length of 512 characters.

parallel_tool_calls boolean or null **Optional** Defaults to true
Whether to allow the model to run tool calls in parallel.

previous_response_id string or null **Optional**
The unique ID of the previous response to the model. Use this to create multi-turn conversations. Learn more about [conversation state](#).

reasoning object or null **Optional**
o-series models only
Configuration options for [reasoning models](#).

Example request

javascript

```
1 import OpenAI from "openai";
2
3 const openai = new OpenAI();
4
5 const response = await openai.responses.create
6   model: "gpt-4.1",
7   input: "Tell me a three sentence bedtime s
8 });
9
10 console.log(response);
```

Response

```
1 {
2   "id": "resp_67ccd2bed1ec8190b14f964abc0542670",
3   "object": "response",
4   "created_at": 1741476542,
5   "status": "completed",
6   "error": null,
7   "incomplete_details": null,
8   "instructions": null,
9   "max_output_tokens": null,
10  "model": "gpt-4.1-2025-04-14",
11  "output": [
12    {
13      "type": "message",
14      "id": "msg_67ccd2bf17f0819081ff3bb2cf6508",
15      "status": "completed",
16      "role": "assistant",
17      "content": [
18        {
19          "type": "output_text",
20          "text": "In a peaceful grove beneath
21            "annotations": []
22          }
23        ]
24      }
25    ],
26    "parallel_tool_calls": true,
27    "previous_response_id": null,
28    "reasoning": {
29      "effort": null,
30      "summary": null
31    },
32    "store": true,
33    "temperature": 1.0,
34    "text": {
35      "format": {
36        "type": "text"
37      }
38    },
39    "tool_choice": "auto",
40    "tools": [],
41    "top_p": 1.0,
42    "truncation": "disabled",
43    "usage": {
44      "input_tokens": 36,
45      "input_tokens_details": {
46        "cached_tokens": 0
47      },
48      "output_tokens": 87,
49      "output_tokens_details": {
50        "reasoning_tokens": 0
51      },
52      "total_tokens": 123
53    },
54    "user": null,
```

▼ Show properties

service_tier string or null Optional Defaults to auto

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', and the Project is Scale tier enabled, the system will utilize scale tier credits until they are exhausted.
- If set to 'auto', and the Project is not Scale tier enabled, the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'flex', the request will be processed with the Flex Processing service tier. [Learn more.](#)
- When not set, the default behavior is 'auto'.

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

store boolean or null Optional Defaults to true

Whether to store the generated model response for later retrieval via API.

stream boolean or null Optional Defaults to false

If set to true, the model response data will be streamed to the client as it is generated using [server-sent events](#). See the [Streaming section below](#) for more information.

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.

text object Optional

Configuration options for a text response from the model. Can be plain text or structured JSON data. [Learn more:](#)

- [Text inputs and outputs](#)
- [Structured Outputs](#)

▼ Show properties

tool_choice string or object Optional

How the model should select which tool (or tools) to use when generating a response. See the `tools` parameter to see how to specify which tools the model can call.

▼ Show possible types

tools array Optional

An array of tools the model may call while generating a response. You can specify which tool to use by setting the `tool_choice` parameter.

The two categories of tools you can provide the model are:

- **Built-in tools:** Tools that are provided by OpenAI that extend the model's capabilities, like [web search](#) or [file search](#). [Learn more about built-in tools.](#)
- **Function calls (custom tools):** Functions that are defined by you, enabling the model to call your own code. [Learn more about function calling.](#)

▼ Show possible types

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.

truncation string or null Optional Defaults to disabled

The truncation strategy to use for the model response.

- `auto` : If the context of this response and previous ones exceeds the model's context window size, the model will truncate the response to fit the context window by dropping input items in the middle of the conversation.
- `disabled` (default): If a model response will exceed the context window size for a model, the request will fail with a 400 error.

user string Optional
A unique identifier representing your end-user, which can help OpenAI to monitor and detect abuse. [Learn more](#).

Returns

Returns a [Response](#) object.

Get a model response

GET `https://api.openai.com/v1/responses/{response_id}`

Retrieves a model response with the given ID.

Path parameters

response_id string **Required**
The ID of the response to retrieve.

Query parameters

include array Optional
Additional fields to include in the response. See the `include` parameter for Response creation above for more information.

Returns

The [Response](#) object matching the specified ID.

Example request javascript

```
1 import OpenAI from "openai";
2 const client = new OpenAI();
3
4 const response = await client.responses.retrieve(response_id);
5 console.log(response);
```

Response

```
1 {
2   "id": "resp_67cb71b351908190a308f3859487620d",
3   "object": "response",
4   "created_at": 1741386163,
5   "status": "completed",
6   "error": null,
7   "incomplete_details": null,
8   "instructions": null,
9   "max_output_tokens": null,
10  "model": "gpt-4o-2024-08-06",
11  "output": [
12    {
13      "type": "message",
14      "id": "msg_67cb71b3c2b0819084d481baaaf14",
15      "status": "completed",
16      "role": "assistant",
17      "content": [
18        {
19          "type": "output_text",
20          "text": "Silent circuits hum, \nThou",
21          "annotations": []
22        }
23      ]
24    }
25  ],
26  "parallel_tool_calls": true,
27  "previous_response_id": null,
28  "reasoning": {
29    "effort": null,
30    "summary": null
31  },
32  "store": true,
33  "temperature": 1.0,
34  "text": {
35    "format": {
```

```
37     "type": "text"
38   },
39 },
40 "tool_choice": "auto",
41 "tools": [],
42 "top_p": 1.0,
43 "truncation": "disabled",
44 "usage": {
45   "input_tokens": 32,
46   "input_tokens_details": {
47     "cached_tokens": 0
48   },
49   "output_tokens": 18,
50   "output_tokens_details": {
51     "reasoning_tokens": 0
52   },
53 },
54 "total_tokens": 50
55 },
56 "user": null,
57 "metadata": {}
58 }
```

Delete a model response

DELETE https://api.openai.com/v1/responses/{response_id}

Deletes a model response with the given ID.

Path parameters

response_id string **Required**
The ID of the response to delete.

Returns

A success message.

Example request

javascript ↕ 📄

```
1 import OpenAI from "openai";
2 const client = new OpenAI();
3
4 const response = await client.responses.del("response_id");
5 console.log(response);
```

Response

📄

```
1 {
2   "id": "resp_6786a1bec27481909a17d673315b29f6"
3   "object": "response",
4   "deleted": true
5 }
```

List input items

GET https://api.openai.com/v1/responses/{response_id}/input_items

Returns a list of input items for a given response.

Path parameters

response_id string **Required**
The ID of the response to retrieve input items for.

Query parameters

after string Optional
An item ID to list items after, used in pagination.

before string Optional
An item ID to list items before, used in pagination.

Example request

javascript ↕ 📄

```
1 import OpenAI from "openai";
2 const client = new OpenAI();
3
4 const response = await client.responses.inputItems({response_id});
5 console.log(response.data);
```

Response

📄

```
1 {
2   "object": "list",
3   "data": [
4     {
5       "id": "msg_abc123",
6       "type": "message",
7       "role": "user",
8       "content": [
9         {
10          "type": "input_text",
```


include array Optional
Additional fields to include in the response. See the `include` parameter for Response creation above for more information.

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
The order to return the input items in. Default is `asc`.

- `asc`: Return the input items in ascending order.
- `desc`: Return the input items in descending order.

Returns

A list of input item objects.

The response object

created_at number
Unix timestamp (in seconds) of when this Response was created.

error object or null
An error object returned when the model fails to generate a Response.

▼ Show properties

id string
Unique identifier for this Response.

incomplete_details object or null
Details about why the response is incomplete.

▼ Show properties

instructions string or null
Inserts a system (or developer) message as the first item in the model's context.

When using along with `previous_response_id`, the instructions from a previous response will not be carried over to the next response. This makes it simple to swap out system (or developer) messages in new responses.

max_output_tokens integer or null
An upper bound for the number of tokens that can be generated for a response, including visible output tokens and [reasoning tokens](#).

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, and querying for objects via API or the dashboard.

Keys are strings with a maximum length of 64 characters. Values are strings with a maximum length of 512 characters.

model string
Model ID used to generate the response, like `gpt-4o` or `o3`. OpenAI offers a wide range of models with different capabilities, performance characteristics, and price points. Refer to the [model guide](#) to browse and compare available models.

```
11         "text": "Tell me a three sentence be
12     }
13 }
14 }
15 ],
16 "first_id": "msg_abc123",
17 "last_id": "msg_abc123",
18 "has_more": false
19 }
```

OBJECT The response object



```
1 {
2   "id": "resp_67ccd3a9da748190baa7f1570fe91ac60",
3   "object": "response",
4   "created_at": 1741476777,
5   "status": "completed",
6   "error": null,
7   "incomplete_details": null,
8   "instructions": null,
9   "max_output_tokens": null,
10  "model": "gpt-4o-2024-08-06",
11  "output": [
12    {
13      "type": "message",
14      "id": "msg_67ccd3acc8d48190a77525dc6de641",
15      "status": "completed",
16      "role": "assistant",
17      "content": [
18        {
19          "type": "output_text",
20          "text": "The image depicts a scenic
21          "annotations": []
22        }
23      ]
24    }
25  ],
26  "parallel_tool_calls": true,
27  "previous_response_id": null,
28  "reasoning": {
29    "effort": null,
30    "summary": null
31  },
32  "store": true,
33  "temperature": 1.0,
34  "text": {
35    "format": {
36      "type": "text"
37    }
38  },
39  "tool_choice": "auto",
40  "tools": [],
41 }
```

object string

The object type of this resource - always set to `response` .

output array

An array of content items generated by the model.

- The length and order of items in the `output` array is dependent on the model's response.
- Rather than accessing the first item in the `output` array and assuming it's an `assistant` message with the content generated by the model, you might consider using the `output_text` property where supported in SDKs.

▼ Show possible types

output_text string or null SDK Only

SDK-only convenience property that contains the aggregated text output from all `output_text` items in the `output` array, if any are present. Supported in the Python and JavaScript SDKs.

parallel_tool_calls boolean

Whether to allow the model to run tool calls in parallel.

previous_response_id string or null

The unique ID of the previous response to the model. Use this to create multi-turn conversations. Learn more about [conversation state](#).

reasoning object or null

o-series models only

Configuration options for [reasoning models](#).

▼ Show properties

service_tier string or 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', and the Project is Scale tier enabled, the system will utilize scale tier credits until they are exhausted.
- If set to 'auto', and the Project is not Scale tier enabled, the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'flex', the request will be processed with the Flex Processing service tier. [Learn more](#).
- When not set, the default behavior is 'auto'.

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

status string

The status of the response generation. One of `completed` , `failed` , `in_progress` , or `incomplete` .

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. We generally recommend altering this or `top_p` but not both.

text object

Configuration options for a text response from the model. Can be plain text or structured JSON data. Learn more:

- [Text inputs and outputs](#)
- [Structured Outputs](#)

```
43  "top_p": 1.0,
44  "truncation": "disabled",
45  "usage": {
46    "input_tokens": 328,
47    "input_tokens_details": {
48      "cached_tokens": 0
49    },
50    "output_tokens": 52,
51    "output_tokens_details": {
52      "reasoning_tokens": 0
53    },
54    "total_tokens": 380
55  },
56  "user": null,
  "metadata": {}
}
```


▼ Show properties

tool_choice string or object

How the model should select which tool (or tools) to use when generating a response. See the `tools` parameter to see how to specify which tools the model can call.

▼ Show possible types

tools array

An array of tools the model may call while generating a response. You can specify which tool to use by setting the `tool_choice` parameter.

The two categories of tools you can provide the model are:

- **Built-in tools:** Tools that are provided by OpenAI that extend the model's capabilities, like [web search](#) or [file search](#). Learn more about [built-in tools](#).
- **Function calls (custom tools):** Functions that are defined by you, enabling the model to call your own code. Learn more about [function calling](#).

▼ Show possible types

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.

truncation string or null

The truncation strategy to use for the model response.

- `auto` : If the context of this response and previous ones exceeds the model's context window size, the model will truncate the response to fit the context window by dropping input items in the middle of the conversation.
- `disabled` (default): If a model response will exceed the context window size for a model, the request will fail with a 400 error.

usage object

Represents token usage details including input tokens, output tokens, a breakdown of output tokens, and the total tokens used.

▼ Show properties

user string

A unique identifier representing your end-user, which can help OpenAI to monitor and detect abuse. [Learn more](#).

The input item list

A list of Response items.

data array

A list of items used to generate this response.

▼ Show possible types

first_id string

The ID of the first item in the list.

has_more boolean

Whether there are more items available.

last_id string

The ID of the last item in the list.

OBJECT The input item list



```
1 {
2   "object": "list",
3   "data": [
4     {
5       "id": "msg_abc123",
6       "type": "message",
7       "role": "user",
8       "content": [
9         {
10          "type": "input_text",
11          "text": "Tell me a three sentence be
12        }
13      ]
14    }
15  ]
}
```

object string

The type of object returned, must be `list` .

```
16  ],
17  "first_id": "msg_abc123",
18  "last_id": "msg_abc123",
19  "has_more": false
}
```

Streaming

When you [create a Response](#) with `stream` set to `true` , the server will emit server-sent events to the client as the Response is generated. This section contains the events that are emitted by the server.

[Learn more about streaming responses.](#)

response.created

An event that is emitted when a response is created.

response object

The response that was created.

▼ Show properties

type string

The type of the event. Always `response.created` .

OBJECT response.created

```
1  {
2    "type": "response.created",
3    "response": {
4      "id": "resp_67ccfcdd16748190a91872c75d3853",
5      "object": "response",
6      "created_at": 1741487325,
7      "status": "in_progress",
8      "error": null,
9      "incomplete_details": null,
10     "instructions": null,
11     "max_output_tokens": null,
12     "model": "gpt-4o-2024-08-06",
13     "output": [],
14     "parallel_tool_calls": true,
15     "previous_response_id": null,
16     "reasoning": {
17       "effort": null,
18       "summary": null
19     },
20   },
21   "store": true,
22   "temperature": 1,
23   "text": {
24     "format": {
25       "type": "text"
26     }
27   },
28   "tool_choice": "auto",
29   "tools": [],
30   "top_p": 1,
31   "truncation": "disabled",
32   "usage": null,
33   "user": null,
34   "metadata": {}
35 }
```

response.in_progress

Emitted when the response is in progress.

response object

The response that is in progress.

▼ Show properties

OBJECT response.in_progress

```
1  {
2    "type": "response.in_progress",
3    "response": {
4      "id": "resp_67ccfcdd16748190a91872c75d3853",
5
```

type string
The type of the event. Always `response.in_progress` .

```
6  "object": "response",
7  "created_at": 1741487325,
8  "status": "in_progress",
9  "error": null,
10 "incomplete_details": null,
11 "instructions": null,
12 "max_output_tokens": null,
13 "model": "gpt-4o-2024-08-06",
14 "output": [],
15 "parallel_tool_calls": true,
16 "previous_response_id": null,
17 "reasoning": {
18   "effort": null,
19   "summary": null
20 },
21 "store": true,
22 "temperature": 1,
23 "text": {
24   "format": {
25     "type": "text"
26   }
27 },
28 "tool_choice": "auto",
29 "tools": [],
30 "top_p": 1,
31 "truncation": "disabled",
32 "usage": null,
33 "user": null,
34 "metadata": {}
35 }
```

response.completed

Emitted when the model response is complete.

response object
Properties of the completed response.
▼ Show properties

type string
The type of the event. Always `response.completed` .

OBJECT response.completed

```
1  {
2    "type": "response.completed",
3    "response": {
4      "id": "resp_123",
5      "object": "response",
6      "created_at": 1740855869,
7      "status": "completed",
8      "error": null,
9      "incomplete_details": null,
10     "input": [],
11     "instructions": null,
12     "max_output_tokens": null,
13     "model": "gpt-4o-mini-2024-07-18",
14     "output": [
15       {
16         "id": "msg_123",
17         "type": "message",
18         "role": "assistant",
19         "content": [
20           {
21             "type": "output_text",
22             "text": "In a shimmering forest un",
23             "annotations": []
24           }
25         ]
26       }
27     ],
28     "previous_response_id": null,
29     "reasoning_effort": null,
30     "store": false,
31     "temperature": 1,
32     "text": {
```

```
35     "format": {
36       "type": "text"
37     }
38   },
39   "tool_choice": "auto",
40   "tools": [],
41   "top_p": 1,
42   "truncation": "disabled",
43   "usage": {
44     "input_tokens": 0,
45     "output_tokens": 0,
46     "output_tokens_details": {
47       "reasoning_tokens": 0
48     },
49     "total_tokens": 0
50   },
51   "user": null,
52   "metadata": {}
53 }
54 }
```

response.failed

An event that is emitted when a response fails.

response object
The response that failed.
Show properties

type string
The type of the event. Always `response.failed`.

OBJECT response.failed

```
1  {
2    "type": "response.failed",
3    "response": {
4      "id": "resp_123",
5      "object": "response",
6      "created_at": 1740855869,
7      "status": "failed",
8      "error": {
9        "code": "server_error",
10       "message": "The model failed to generate
11     },
12     "incomplete_details": null,
13     "instructions": null,
14     "max_output_tokens": null,
15     "model": "gpt-4o-mini-2024-07-18",
16     "output": [],
17     "previous_response_id": null,
18     "reasoning_effort": null,
19     "store": false,
20     "temperature": 1,
21     "text": {
22       "format": {
23         "type": "text"
24       }
25     },
26   },
27   "tool_choice": "auto",
28   "tools": [],
29   "top_p": 1,
30   "truncation": "disabled",
31   "usage": null,
32   "user": null,
33   "metadata": {}
34 }
35 }
```

response.incomplete

An event that is emitted when a response finishes as incomplete.

OBJECT response.incomplete

response object
The response that was incomplete.
▼ Show properties

type string
The type of the event. Always `response.incomplete` .

```
1 {
2   "type": "response.incomplete",
3   "response": {
4     "id": "resp_123",
5     "object": "response",
6     "created_at": 1740855869,
7     "status": "incomplete",
8     "error": null,
9     "incomplete_details": {
10      "reason": "max_tokens"
11    },
12    "instructions": null,
13    "max_output_tokens": null,
14    "model": "gpt-4o-mini-2024-07-18",
15    "output": [],
16    "previous_response_id": null,
17    "reasoning_effort": null,
18    "store": false,
19    "temperature": 1,
20    "text": {
21      "format": {
22        "type": "text"
23      }
24    },
25  },
26  "tool_choice": "auto",
27  "tools": [],
28  "top_p": 1,
29  "truncation": "disabled",
30  "usage": null,
31  "user": null,
32  "metadata": {}
33 }
```

response.output_item.added

Emitted when a new output item is added.

item object
The output item that was added.
▼ Show possible types

output_index integer
The index of the output item that was added.

type string
The type of the event. Always `response.output_item.added` .

OBJECT response.output_item.added

```
1 {
2   "type": "response.output_item.added",
3   "output_index": 0,
4   "item": {
5     "id": "msg_123",
6     "status": "in_progress",
7     "type": "message",
8     "role": "assistant",
9     "content": []
10  }
11 }
```

response.output_item.done

Emitted when an output item is marked done.

item object
The output item that was marked done.
▼ Show possible types

output_index integer
The index of the output item that was marked done.

type string

OBJECT response.output_item.done

```
1 {
2   "type": "response.output_item.done",
3   "output_index": 0,
4   "item": {
5     "id": "msg_123",
6     "status": "completed",
7     "type": "message",
8     "role": "assistant",
9     "content": [
10     {
11     }
```

The type of the event. Always `response.output_item.done` .

```
12         "type": "output_text",
13         "text": "In a shimmering forest under a sky",
14         "annotations": []
15     }
16 ]
17 }
}
```

response.content_part.added

Emitted when a new content part is added.

content_index integer

The index of the content part that was added.

item_id string

The ID of the output item that the content part was added to.

output_index integer

The index of the output item that the content part was added to.

part object

The content part that was added.

⌵ Show possible types

type string

The type of the event. Always `response.content_part.added` .

OBJECT response.content_part.added

```
1 {
2   "type": "response.content_part.added",
3   "item_id": "msg_123",
4   "output_index": 0,
5   "content_index": 0,
6   "part": {
7     "type": "output_text",
8     "text": "",
9     "annotations": []
10  }
11 }
```

response.content_part.done

Emitted when a content part is done.

content_index integer

The index of the content part that is done.

item_id string

The ID of the output item that the content part was added to.

output_index integer

The index of the output item that the content part was added to.

part object

The content part that is done.

⌵ Show possible types

type string

The type of the event. Always `response.content_part.done` .

OBJECT response.content_part.done

```
1 {
2   "type": "response.content_part.done",
3   "item_id": "msg_123",
4   "output_index": 0,
5   "content_index": 0,
6   "part": {
7     "type": "output_text",
8     "text": "In a shimmering forest under a sky",
9     "annotations": []
10  }
11 }
```

response.output_text.delta

Emitted when there is an additional text delta.

content_index integer

The index of the content part that the text delta was added to.

delta string

OBJECT response.output_text.delta

```
1 {
2   "type": "response.output_text.delta",
3   "item_id": "msg_123",
4   "output_index": 0,
5   "content_index": 0,
```


The text delta that was added.

item_id string

The ID of the output item that the text delta was added to.

output_index integer

The index of the output item that the text delta was added to.

type string

The type of the event. Always `response.output_text.delta` .

response.output_text.annotation.added

Emitted when a text annotation is added.

annotation object

▼ Show possible types

annotation_index integer

The index of the annotation that was added.

content_index integer

The index of the content part that the text annotation was added to.

item_id string

The ID of the output item that the text annotation was added to.

output_index integer

The index of the output item that the text annotation was added to.

type string

The type of the event. Always `response.output_text.annotation.added` .

```
6  "delta": "In"
7  }
```

OBJECT response.output_text.annotation.added

```
1  {
2    "type": "response.output_text.annotation.added",
3    "item_id": "msg_abc123",
4    "output_index": 1,
5    "content_index": 0,
6    "annotation_index": 0,
7    "annotation": {
8      "type": "file_citation",
9      "index": 390,
10     "file_id": "file-4wDz5b167pAf72nx1h9eiN",
11     "filename": "dragons.pdf"
12   }
13 }
```

response.output_text.done

Emitted when text content is finalized.

content_index integer

The index of the content part that the text content is finalized.

item_id string

The ID of the output item that the text content is finalized.

output_index integer

The index of the output item that the text content is finalized.

text string

The text content that is finalized.

type string

The type of the event. Always `response.output_text.done` .

OBJECT response.output_text.done

```
1  {
2    "type": "response.output_text.done",
3    "item_id": "msg_123",
4    "output_index": 0,
5    "content_index": 0,
6    "text": "In a shimmering forest under a sky fr
7  }
```

response.refusal.delta

Emitted when there is a partial refusal text.

content_index integer
The index of the content part that the refusal text is added to.

delta string
The refusal text that is added.

item_id string
The ID of the output item that the refusal text is added to.

output_index integer
The index of the output item that the refusal text is added to.

type string
The type of the event. Always `response.refusal.delta`.

OBJECT response.refusal.delta

```
1 {
2   "type": "response.refusal.delta",
3   "item_id": "msg_123",
4   "output_index": 0,
5   "content_index": 0,
6   "delta": "refusal text so far"
7 }
```

response.refusal.done

Emitted when refusal text is finalized.

content_index integer
The index of the content part that the refusal text is finalized.

item_id string
The ID of the output item that the refusal text is finalized.

output_index integer
The index of the output item that the refusal text is finalized.

refusal string
The refusal text that is finalized.

type string
The type of the event. Always `response.refusal.done`.

OBJECT response.refusal.done

```
1 {
2   "type": "response.refusal.done",
3   "item_id": "item-abc",
4   "output_index": 1,
5   "content_index": 2,
6   "refusal": "final refusal text"
7 }
```

response.function_call_arguments.delta

Emitted when there is a partial function-call arguments delta.

delta string
The function-call arguments delta that is added.

item_id string
The ID of the output item that the function-call arguments delta is added to.

output_index integer
The index of the output item that the function-call arguments delta is added to.

type string
The type of the event. Always `response.function_call_arguments.delta`.

OBJECT response.function_call_arguments.delta

```
1 {
2   "type": "response.function_call_arguments.delta",
3   "item_id": "item-abc",
4   "output_index": 0,
5   "delta": "{ \"arg\": \""
6 }
```

response.function_call_arguments.done

Emitted when function-call arguments are finalized.

arguments string
The function-call arguments.

item_id string
The ID of the item.

output_index integer
The index of the output item.

type string

OBJECT response.function_call_arguments.done

```
1 {
2   "type": "response.function_call_arguments.done",
3   "item_id": "item-abc",
4   "output_index": 1,
5   "arguments": "{ \"arg\": 123 }"
6 }
```

response.file_search_call.in_progress

Emitted when a file search call is initiated.

item_id string
The ID of the output item that the file search call is initiated.

output_index integer
The index of the output item that the file search call is initiated.

type string
The type of the event. Always response.file_search_call.in_progress .

OBJECT response.file_search_call.in_progress

```
1 {
2   "type": "response.file_search_call.in_progress",
3   "output_index": 0,
4   "item_id": "fs_123",
5 }
```

response.file_search_call.searching

Emitted when a file search is currently searching.

item_id string
The ID of the output item that the file search call is initiated.

output_index integer
The index of the output item that the file search call is searching.

type string
The type of the event. Always response.file_search_call.searching .

OBJECT response.file_search_call.searching

```
1 {
2   "type": "response.file_search_call.searching",
3   "output_index": 0,
4   "item_id": "fs_123",
5 }
```

response.file_search_call.completed

Emitted when a file search call is completed (results found).

item_id string
The ID of the output item that the file search call is initiated.

output_index integer
The index of the output item that the file search call is initiated.

OBJECT response.file_search_call.completed

```
1 {
2   "type": "response.file_search_call.completed",
3   "output_index": 0,
4   "item_id": "fs_123",
5 }
```

type string
The type of the event. Always `response.file_search_call.completed` .

response.web_search_call.in_progress

Emitted when a web search call is initiated.

item_id string
Unique ID for the output item associated with the web search call.

output_index integer
The index of the output item that the web search call is associated with.

type string
The type of the event. Always `response.web_search_call.in_progress` .

OBJECT response.web_search_call.in_progress

```
1 {
2   "type": "response.web_search_call.in_progress",
3   "output_index": 0,
4   "item_id": "ws_123",
5 }
```

response.web_search_call.searching

Emitted when a web search call is executing.

item_id string
Unique ID for the output item associated with the web search call.

output_index integer
The index of the output item that the web search call is associated with.

type string
The type of the event. Always `response.web_search_call.searching` .

OBJECT response.web_search_call.searching

```
1 {
2   "type": "response.web_search_call.searching",
3   "output_index": 0,
4   "item_id": "ws_123",
5 }
```

response.web_search_call.completed

Emitted when a web search call is completed.

item_id string
Unique ID for the output item associated with the web search call.

output_index integer
The index of the output item that the web search call is associated with.

type string
The type of the event. Always `response.web_search_call.completed` .

OBJECT response.web_search_call.completed

```
1 {
2   "type": "response.web_search_call.completed",
3   "output_index": 0,
4   "item_id": "ws_123",
5 }
```

error

Emitted when an error occurs.

code string or null
The error code.

OBJECT error

```
1 {
2   "type": "error",
3   "code": "ERR_SOMETHING",
```

message string
The error message.

param string or null
The error parameter.

type string
The type of the event. Always `error`.

```
4  "message": "Something went wrong",
5  "param": null
6  }
```

Chat Completions

The Chat Completions API endpoint will generate a model response from a list of messages comprising a conversation.

Related guides:

- Quickstart
- Text inputs and outputs
- Image inputs
- Audio inputs and outputs
- Structured Outputs
- Function calling
- Conversation state

Starting a new project? We recommend trying [Responses](#) to take advantage of the latest OpenAI platform features. Compare [Chat Completions with Responses](#).

Create chat completion

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

Starting a new project? We recommend trying [Responses](#) to take advantage of the latest OpenAI platform features. Compare [Chat Completions with Responses](#).

Creates a model response for the given chat conversation. Learn more in the [text generation](#), [vision](#), and [audio](#) guides. Parameter support can differ depending on the model used to generate the response, particularly for newer reasoning models. Parameters that are only supported for reasoning models are noted below. For the current state of unsupported parameters in reasoning models, [refer to the reasoning guide](#).

Request body

messages array **Required**
A list of messages comprising the conversation so far. Depending on the [model](#) you use, different message types (modalities) are supported, like [text](#), [images](#), and [audio](#).
Show possible types

model string **Required**

Default Image input Streaming Functions Logp

Example requestgpt-4.1curl

```
1  curl https://api.openai.com/v1/chat/completion:
2    -H "Content-Type: application/json" \
3    -H "Authorization: Bearer $OPENAI_API_KEY" \
4    -d '{
5      "model": "gpt-4.1",
6      "messages": [
7        {
8          "role": "developer",
9          "content": "You are a helpful assistant"
10       },
11       {
12         "role": "user",
13         "content": "Hello!"
14       }
15     ]
16   }'
```

Response

```
1  {
2    "id": "chatcpl-B9MBs8Cjcv0U2jLn4n570S5qMJKc",
3    "object": "chat.completion",
4    "created": 1741569952,
5    "model": "gpt-4.1-2025-04-14",
```

Model ID used to generate the response, like `gpt-4o` or `o3`. OpenAI offers a wide range of models with different capabilities, performance characteristics, and price points. Refer to the [model guide](#) to browse and compare available models.

audio object or null Optional

Parameters for audio output. Required when audio output is requested with `modalities: ["audio"]`. [Learn more](#).

▼ Show properties

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.

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

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`.

max_completion_tokens integer or null Optional

An upper bound for the number of tokens that can be generated for a completion, including visible output tokens and [reasoning tokens](#).

max_tokens Deprecated integer or null Optional

The maximum number of [tokens](#) that can be generated in the chat completion. This value can be used to control [costs](#) for text generated via API.

This value is now deprecated in favor of `max_completion_tokens`, and is not compatible with [o-series models](#).

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, and querying for objects via API or the dashboard.

Keys are strings with a maximum length of 64 characters. Values are strings with a maximum length of 512 characters.

```
6  "choices": [  
7    {  
8      "index": 0,  
9      "message": {  
10       "role": "assistant",  
11       "content": "Hello! How can I assist you",  
12       "refusal": null,  
13       "annotations": []  
14     },  
15     "logprobs": null,  
16     "finish_reason": "stop"  
17   },  
18 ],  
19 "usage": {  
20   "prompt_tokens": 19,  
21   "completion_tokens": 10,  
22   "total_tokens": 29,  
23   "prompt_tokens_details": {  
24     "cached_tokens": 0,  
25     "audio_tokens": 0  
26   },  
27   "completion_tokens_details": {  
28     "reasoning_tokens": 0,  
29     "audio_tokens": 0,  
30     "accepted_prediction_tokens": 0,  
31     "rejected_prediction_tokens": 0  
32   }  
33 },  
34 "service_tier": "default"  
35 }
```


modalities array or null Optional

Output types that you would like the model to generate. Most models are capable of generating text, which is the default:

```
["text"]
```

The `gpt-4o-audio-preview` model can also be used to [generate audio](#). To request that this model generate both text and audio responses, you can use:

```
["text", "audio"]
```

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.

parallel_tool_calls boolean Optional Defaults to true

Whether to enable [parallel function calling](#) during tool use.

prediction object Optional

Configuration for a [Predicted Output](#), which can greatly improve response times when large parts of the model response are known ahead of time. This is most common when you are regenerating a file with only minor changes to most of the content.

▼ Show possible types

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.

reasoning_effort string or null Optional Defaults to medium

o-series models only

Constrains effort on reasoning for [reasoning models](#). Currently supported values are `low`, `medium`, and `high`. Reducing reasoning effort can result in faster responses and fewer tokens used on reasoning in a response.

response_format object Optional

An object specifying the format that the model must output.

Setting to `{ "type": "json_schema", "json_schema": {...} }` enables Structured Outputs which ensures the model will match your supplied JSON schema. Learn more in the [Structured Outputs guide](#).

Setting to `{ "type": "json_object" }` enables the older JSON mode, which ensures the message the model generates is valid JSON. Using `json_schema` is preferred for models that support it.

▼ Show possible types

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 auto

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', and the Project is Scale tier enabled, the system will utilize scale tier credits until they are exhausted.
- If set to 'auto', and the Project is not Scale tier enabled, the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.

- If set to 'flex', the request will be processed with the Flex Processing service tier.
[Learn more.](#)
- When not set, the default behavior is 'auto'.

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

stop string / array / null Optional Defaults to null
Not supported with latest reasoning models `o3` and `o4-mini` .

Up to 4 sequences where the API will stop generating further tokens. The returned text will not contain the stop sequence.

store boolean or null Optional Defaults to false
Whether or not to store the output of this chat completion request for use in our [model distillation](#) or [evals](#) products.

stream boolean or null Optional Defaults to false
If set to true, the model response data will be streamed to the client as it is generated using [server-sent events](#). See the [Streaming section below](#) for more information, along with the [streaming responses](#) guide for more information on how to handle the streaming events.

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.

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](#)

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](#)

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.

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 OpenAI to monitor and detect abuse. [Learn more.](#)

web_search_options object Optional

This tool searches the web for relevant results to use in a response. Learn more about the [web search tool](#).

⌵ Show properties

Returns

Returns a [chat completion](#) object, or a streamed sequence of [chat completion chunk](#) objects if the request is streamed.

Get chat completion

GET `https://api.openai.com/v1/chat/completions/{completion_id}`

Get a stored chat completion. Only Chat Completions that have been created with the `store` parameter set to `true` will be returned.

Path parameters

completion_id string Required
The ID of the chat completion to retrieve.

Returns

The [ChatCompletion](#) object matching the specified ID.

Example request

curl ↕ 

```
1 curl https://api.openai.com/v1/chat/completions,
2 -H "Authorization: Bearer $OPENAI_API_KEY" \
3 -H "Content-Type: application/json"
```

Response



```
1 {
2   "object": "chat.completion",
3   "id": "chatcmpl-abc123",
4   "model": "gpt-4o-2024-08-06",
5   "created": 1738960610,
6   "request_id": "req_ded8ab984ec4bf840f37566c1f",
7   "tool_choice": null,
8   "usage": {
9     "total_tokens": 31,
10    "completion_tokens": 18,
11    "prompt_tokens": 13
12  },
13   "seed": 4944116822809979520,
14   "top_p": 1.0,
15   "temperature": 1.0,
16   "presence_penalty": 0.0,
17   "frequency_penalty": 0.0,
18   "system_fingerprint": "fp_50cad350e4",
19   "input_user": null,
20   "service_tier": "default",
21   "tools": null,
22   "metadata": {},
23   "choices": [
24     {
25       "index": 0,
26       "message": {
27         "content": "Mind of circuits hum, \nLe",
28         "role": "assistant",
29         "tool_calls": null,
30         "function_call": null
31       },
32       "finish_reason": "stop",
33       "logprobs": null
34     }
35   ],
36   "response_format": null
37 }
```

Get chat messages

GET https://api.openai.com/v1/chat/completions/{completion_id}/messages

Get the messages in a stored chat completion. Only Chat Completions that have been created with the `store` parameter set to `true` will be returned.

Path parameters

completion_id string Required
The ID of the chat completion to retrieve messages from.

Query parameters

after string Optional
Identifier for the last message from the previous pagination request.

limit integer Optional Defaults to 20
Number of messages to retrieve.

order string Optional Defaults to asc
Sort order for messages by timestamp. Use `asc` for ascending order or `desc` for descending order. Defaults to `asc`.

Returns

A list of `messages` for the specified chat completion.

Example requestcurl

```
1 curl https://api.openai.com/v1/chat/completions,
2   -H "Authorization: Bearer $OPENAI_API_KEY" \
3   -H "Content-Type: application/json"
```

Response

```
1  {
2    "object": "list",
3    "data": [
4      {
5        "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
6        "role": "user",
7        "content": "write a haiku about ai",
8        "name": null,
9        "content_parts": null
10     }
11  ],
12  "first_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
13  "last_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
14  "has_more": false
15  }
```

List Chat Completions

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

List stored Chat Completions. Only Chat Completions that have been stored with the `store` parameter set to `true` will be returned.

Query parameters

after string Optional
Identifier for the last chat completion from the previous pagination request.

limit integer Optional Defaults to 20
Number of Chat Completions to retrieve.

metadata map Optional
A list of metadata keys to filter the Chat Completions by. Example:

```
metadata[key1]=value1&metadata[key2]=value2
```

model string Optional
The model used to generate the Chat Completions.

order string Optional Defaults to asc

Example requestcurl

```
1 curl https://api.openai.com/v1/chat/completions
2   -H "Authorization: Bearer $OPENAI_API_KEY" \
3   -H "Content-Type: application/json"
```

Response

```
1  {
2    "object": "list",
3    "data": [
4      {
5        "object": "chat.completion",
6        "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
7        "model": "gpt-4.1-2025-04-14",
8        "created": 1738960610,
9        "request_id": "req_ded8ab984ec4bf840f375c",
10       "tool_choice": null,
11       "usage": {
12         "total_tokens": 31,
13         "completion_tokens": 18,
14         "prompt_tokens": 13
15       },
16       "seed": 4944116822809979520,
17     }
18   ]
19 }
```

Sort order for Chat Completions by timestamp. Use `asc` for ascending order or `desc` for descending order. Defaults to `asc`.

Returns

A list of `Chat Completions` matching the specified filters.

```
18     "top_p": 1.0,
19     "temperature": 1.0,
20     "presence_penalty": 0.0,
21     "frequency_penalty": 0.0,
22     "system_fingerprint": "fp_50cad350e4",
23     "input_user": null,
24     "service_tier": "default",
25     "tools": null,
26     "metadata": {},
27     "choices": [
28       {
29         "index": 0,
30         "message": {
31           "content": "Mind of circuits hum,
32           "role": "assistant",
33           "tool_calls": null,
34           "function_call": null
35         },
36         "finish_reason": "stop",
37         "logprobs": null
38       }
39     ],
40     "response_format": null
41   },
42   "first_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobMf1mj",
43   "last_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobMf1mj",
44   "has_more": false
45 }
```

Update chat completion

POST `https://api.openai.com/v1/chat/completions/{completion_id}`

Modify a stored chat completion. Only Chat Completions that have been created with the `store` parameter set to `true` can be modified. Currently, the only supported modification is to update the `metadata` field.

Path parameters

completion_id string Required
The ID of the chat completion to update.

Request body

metadata map Required
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, and querying for objects via API or the dashboard.

Keys are strings with a maximum length of 64 characters. Values are strings with a maximum length of 512 characters.

Returns

The `ChatCompletion` object matching the specified ID.

Example request curl 📄

```
1 curl -X POST https://api.openai.com/v1/chat/completions/{completion_id}
2 -H "Authorization: Bearer $OPENAI_API_KEY" \
3 -H "Content-Type: application/json" \
4 -d '{"metadata": {"foo": "bar"}}'
```

Response 📄

```
1 {
2   "object": "chat.completion",
3   "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobMf1mj",
4   "model": "gpt-4o-2024-08-06",
5   "created": 1738960610,
6   "request_id": "req_ded8ab984ec4bf840f37566c1",
7   "tool_choice": null,
8   "usage": {
9     "total_tokens": 31,
10    "completion_tokens": 18,
11    "prompt_tokens": 13
12  },
13   "seed": 4944116822809979520,
14   "top_p": 1.0,
15   "temperature": 1.0,
16   "presence_penalty": 0.0,
17   "frequency_penalty": 0.0,
18   "system_fingerprint": "fp_50cad350e4",
19   "input_user": null,
20   "service_tier": "default",
21   "tools": null,
22   "metadata": {
23     "foo": "bar"
24   },
25   "choices": [
26     {
27       "index": 0,
28       "message": {
29         "content": "Mind of circuits hum,
30         "role": "assistant",
31         "tool_calls": null,
32         "function_call": null
33       },
34       "finish_reason": "stop",
35       "logprobs": null
36     }
37   ],
38   "response_format": null
39 }
```

```
29     "index": 0,
30     "message": {
31       "content": "Mind of circuits hum, \nLe
32       "role": "assistant",
33       "tool_calls": null,
34       "function_call": null
35     },
36     "finish_reason": "stop",
37     "logprobs": null
38   }
39 ],
  "response_format": null
}
```

Delete chat completion

DELETE `https://api.openai.com/v1/chat/completions/{completion_id}`

Delete a stored chat completion. Only Chat Completions that have been created with the `store` parameter set to `true` can be deleted.

Path parameters

completion_id string Required
The ID of the chat completion to delete.

Returns

A deletion confirmation object.

Example request

curl ↕

```
1 curl -X DELETE https://api.openai.com/v1/chat/c
2 -H "Authorization: Bearer $OPENAI_API_KEY" \
3 -H "Content-Type: application/json"
```

Response

```
1 {
2   "object": "chat.completion.deleted",
3   "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobMflmj2'
4   "deleted": true
5 }
```

The chat completion object

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

choices array
A list of chat completion choices. Can be more than one if `n` is greater than 1.
Show properties

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

id string
A unique identifier for the chat completion.

model string
The model used for the chat completion.

object string
The object type, which is always `chat.completion`.

service_tier string or 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', and the Project is Scale tier enabled, the system will utilize scale tier credits until they are exhausted.

OBJECT The chat completion object

```
1 {
2   "id": "chatcmpl-B9MHDbslfkBeAs8l4bebGdF0J6Pe
3   "object": "chat.completion",
4   "created": 1741570283,
5   "model": "gpt-4o-2024-08-06",
6   "choices": [
7     {
8       "index": 0,
9       "message": {
10         "role": "assistant",
11         "content": "The image shows a wooden bo
12         "refusal": null,
13         "annotations": []
14       },
15       "logprobs": null,
16       "finish_reason": "stop"
17     }
18   ],
19   "usage": {
20     "prompt_tokens": 1117,
21     "completion_tokens": 46,
22     "total_tokens": 1163,
23     "prompt_tokens_details": {
24       "cached_tokens": 0,
25       "audio_tokens": 0
26     }
27   },
```


- If set to 'auto', and the Project is not Scale tier enabled, the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'flex', the request will be processed with the Flex Processing service tier. [Learn more](#).
- When not set, the default behavior is 'auto'.

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

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.

usage object

Usage statistics for the completion request.

⌵ Show properties

The chat completion list object

An object representing a list of Chat Completions.

data array

An array of chat completion objects.

⌵ Show properties

first_id string

The identifier of the first chat completion in the data array.

has_more boolean

Indicates whether there are more Chat Completions available.

last_id string

The identifier of the last chat completion in the data array.

object string

The type of this object. It is always set to "list".

```
28     "completion_tokens_details": {
29       "reasoning_tokens": 0,
30       "audio_tokens": 0,
31       "accepted_prediction_tokens": 0,
32       "rejected_prediction_tokens": 0
33     }
34   },
35   "service_tier": "default",
36   "system_fingerprint": "fp_fc9f1d7035"
}
```

OBJECT The chat completion list object



```
1  {
2    "object": "list",
3    "data": [
4      {
5        "object": "chat.completion",
6        "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
7        "model": "gpt-4o-2024-08-06",
8        "created": 1738960610,
9        "request_id": "req_ded8ab984ec4bf840f3750",
10       "tool_choice": null,
11       "usage": {
12         "total_tokens": 31,
13         "completion_tokens": 18,
14         "prompt_tokens": 13
15       },
16       "seed": 4944116822809979520,
17       "top_p": 1.0,
18       "temperature": 1.0,
19       "presence_penalty": 0.0,
20       "frequency_penalty": 0.0,
21       "system_fingerprint": "fp_50cad350e4",
22       "input_user": null,
23       "service_tier": "default",
24       "tools": null,
25       "metadata": {},
26       "choices": [
27         {
28           "index": 0,
29           "message": {
30             "content": "Mind of circuits hum,",
31             "role": "assistant",
32             "tool_calls": null,
33             "function_call": null
34           },
35           "finish_reason": "stop",
36           "logprobs": null
37         }
38       ],
39       "response_format": null
40     ],
41     "response_format": null
}
```

```
42     }
43   ],
44   "first_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
45   "last_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
46   "has_more": false
47 }
```

The chat completion message list object

An object representing a list of chat completion messages.

data	array
An array of chat completion message objects.	
⌵ Show properties	
first_id	string
The identifier of the first chat message in the data array.	
has_more	boolean
Indicates whether there are more chat messages available.	
last_id	string
The identifier of the last chat message in the data array.	
object	string
The type of this object. It is always set to "list".	

OBJECT The chat completion message list object



```
1  {
2    "object": "list",
3    "data": [
4      {
5        "id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
6        "role": "user",
7        "content": "write a haiku about ai",
8        "name": null,
9        "content_parts": null
10     }
11   ],
12   "first_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
13   "last_id": "chatcmpl-AyPNinnUqUDYo9SAdA52NobM",
14   "has_more": false
15 }
```

Streaming

Stream Chat Completions in real time. Receive chunks of completions returned from the model using server-sent events. [Learn more](#).

The chat completion chunk object

Represents a streamed chunk of a chat completion response returned by the model, based on the provided input. [Learn more](#).

choices	array
A list of chat completion choices. Can contain more than one elements if <code>n</code> is greater than 1. Can also be empty for the last chunk if you set <code>stream_options: {"include_usage": true}</code> .	
⌵ Show properties	
created	integer
The Unix timestamp (in seconds) of when the chat completion was created. Each chunk has the same timestamp.	
id	string
A unique identifier for the chat completion. Each chunk has the same ID.	
model	string
The model to generate the completion.	

OBJECT The chat completion chunk object



```
1  {"id":"chatcmpl-123","object":"chat.completion.chunk",
2
3  {"id":"chatcmpl-123","object":"chat.completion.chunk",
4
5  ....
6
7  {"id":"chatcmpl-123","object":"chat.completion.chunk",
```

object string

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

service_tier string or 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', and the Project is Scale tier enabled, the system will utilize scale tier credits until they are exhausted.
- If set to 'auto', and the Project is not Scale tier enabled, the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'default', the request will be processed using the default service tier with a lower uptime SLA and no latency guarentee.
- If set to 'flex', the request will be processed with the Flex Processing service tier. [Learn more](#).
- When not set, the default behavior is 'auto'.

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

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.

usage object or null

Usage statistics for the completion request.

⌵ Show properties

Realtime Beta

Communicate with a GPT-4o class model in real time using WebRTC or WebSockets. Supports text and audio inputs and ouputs, along with audio transcriptions. [Learn more about the Realtime API](#).

Session tokens

REST API endpoint to generate ephemeral session tokens for use in client-side applications.

Create session

`POST https://api.openai.com/v1/realtime/sessions`

Create an ephemeral API token for use in client-side applications with the Realtime API. Can be configured with the same session parameters as the `session.update` client event.

It responds with a session object, plus a `client_secret` key which contains a usable ephemeral API token that can be used to authenticate browser clients for the Realtime API.

Request body

input_audio_format string Optional Defaults to pcm16

Example request curl ⌵ 📄

```
1 curl -X POST https://api.openai.com/v1/realtime,
2   -H "Authorization: Bearer $OPENAI_API_KEY" \
3   -H "Content-Type: application/json" \
4   -d '{
5     "model": "gpt-4o-realtime-preview",
6     "modalities": ["audio", "text"],
7     "instructions": "You are a friendly assistan
8   }'
```

Response 📄