

OpenAI Platform

 Copy page

Structured model outputs

Ensure text responses from the model adhere to a JSON schema you define.

[Overview](#)

[Use case examples](#)

[API Usage](#)

[Streaming](#)

[Supported schemas](#)

JSON is one of the most widely used formats in the world for applications to exchange data.

Structured Outputs is a feature that ensures the model will always generate responses that adhere to your supplied [JSON Schema](#), so you don't need to worry about the model omitting a required key, or hallucinating an invalid enum value.

Some benefits of Structured Outputs include:

- 1 **Reliable type-safety:** No need to validate or retry incorrectly formatted responses
- 2 **Explicit refusals:** Safety-based model refusals are now programmatically detectable
- 3 **Simpler prompting:** No need for strongly worded prompts to achieve consistent formatting

In addition to supporting JSON Schema in the REST API, the OpenAI SDKs for [Python](#) and [JavaScript](#) also make it easy to define object schemas using [Pydantic](#) and [Zod](#) respectively. Below, you can see how to extract information from unstructured text that conforms to a schema defined in code.

Getting a structured response

python 

```
1 from pydantic import BaseModel
2 from openai import OpenAI
3
4 client = OpenAI()
5
6 class CalendarEvent(BaseModel):
```

```
7     name: str
8     date: str
9     participants: list[str]
10
11 completion = client.chat.completions.parse(
12     model="gpt-4o-2024-08-06",
13     messages=[
14         {"role": "system", "content": "Extract the names of Alice and Bob from the following text."},
15         {"role": "user", "content": "Alice and Bob are going to meet at the park."}
16     ],
17     response_format=CalendarEvent,
18 )
19
20 event = completion.choices[0].message.parsed
```

Supported models

Structured Outputs is available in our [latest large language models](#), starting with GPT-4o. Older models like `gpt-4-turbo` and earlier may use [JSON mode](#) instead.

When to use Structured Outputs via function calling vs via response_format

Structured Outputs is available in two forms in the OpenAI API:

- 1 When using [function calling](#)
- 2 When using a `json_schema` response format

Function calling is useful when you are building an application that bridges the models and functionality of your application.

For example, you can give the model access to functions that query a database in order to build an AI assistant that can help users with their orders, or functions that can interact with the UI.

Conversely, Structured Outputs via `response_format` are

more suitable when you want to indicate a structured schema for use when the model responds to the user, rather than when the model calls a tool.

For example, if you are building a math tutoring application, you might want the assistant to respond to your user using a specific JSON Schema so that you can generate a UI that displays different parts of the model's output in distinct ways.

Put simply:

If you are connecting the model to tools, functions, data, etc. in your system, then you should use function calling - If you want to structure the model's output when it responds to the user, then you should use a structured `response_format`

The remainder of this guide will focus on non-function calling use cases in the Chat Completions API. To learn

- ⓘ more about how to use Structured Outputs with function calling, check out the Function Calling guide.

Structured Outputs vs JSON mode

Structured Outputs is the evolution of JSON mode. While both ensure valid JSON is produced, only Structured Outputs ensure schema adherence. Both Structured Outputs and JSON mode are supported in the Responses API, Chat Completions API, Assistants API, Fine-tuning API and Batch API.

We recommend always using Structured Outputs instead of JSON mode when possible.

However, Structured Outputs with

`response_format: {type: "json_schema", ...}` is only supported with the `gpt-4o-mini`, `gpt-4o-mini-2024-07-18`, and `gpt-4o-2024-08-06` model snapshots and later.

	STRUCTURED OUTPUTS	JSON MODE
Outputs valid JSON	Yes	Yes
Adheres to schema	Yes (see supported schemas)	No
Compatible models	gpt-4o-mini, gpt-4o-2024-08-06, and later	gpt-3.5-turbo, gpt-4-* and gpt-4o-* models
Enabling	response_format: { type: "json_schema", json_schema: {"strict": true, "schema": ...} }	response_format: { type: "json_object" }

Examples

[Chain of thought](#) [Structured data extraction](#) [UI generation](#)

Chain of thought

You can ask the model to output an answer in a structured, step-by-step way, to guide the user through the solution.

```
Structured Outputs for chain-of-thought... python ⚙️ 🐦
1 from pydantic import BaseModel
2 from openai import OpenAI
3
4 client = OpenAI()
5
6 class Step(BaseModel):
7     explanation: str
8     output: str
9
10 class MathReasoning(BaseModel):
11     steps: list[Step]
12     final_answer: str
13
14 completion = client.chat.completions.parse(
15     model="gpt-4o-2024-08-06",
16     messages=[
17         {"role": "system", "content": "You are a helpful assistant."},
18         {"role": "user", "content": "How can I solve this equation?"}
```

```
19     ],
20     response_format=MathReasoning,
21 )
22
23 math_reasoning = completion.choices[0].message.par
```

Example response

```
1  {
2   "steps": [
3     {
4       "explanation": "Start with the equation  $8x - 7 = -23$ ",
5       "output": " $8x + 7 = -23$ "
6     },
7     {
8       "explanation": "Subtract 7 from both sides to get  $8x = -23 - 7$ ",
9       "output": " $8x = -30$ "
10    },
11    {
12      "explanation": "Simplify the right side of the equation to get  $8x = -30$ ",
13      "output": " $8x = -30$ "
14    },
15    {
16      "explanation": "Divide both sides by 8 to solve for x to get  $x = -30 / 8$ ",
17      "output": " $x = -30 / 8$ "
18    },
19    {
20      "explanation": "Simplify the fraction.",
21      "output": " $x = -15 / 4$ "
22    }
23  ],
24  "final_answer": " $x = -15 / 4$ "
25 }
```

How to use Structured Outputs with `response_format`

You can use Structured Outputs with the new SDK helper to parse the model's output into your desired format, or you can specify the JSON schema directly.

Note: for fine tuned models, the first request you make with any schema will have additional latency as our API processes the schema, but subsequent requests with the same schema will not have additional latency. Other

- ⓘ

models do not have this limitation.

SDK objects Manual schema

- › Step 1: Define your object
 - › Step 2: Supply your object in the API call
 - › Step 3: Handle edge cases

Refusals with Structured Outputs

When using Structured Outputs with user-generated input, OpenAI models may occasionally refuse to fulfill the request for safety reasons. Since a refusal does not necessarily follow the schema you have supplied in `response_format`, the API response will include a new field called `refusal` to indicate that the model refused to fulfill the request.

When the `refusal` property appears in your output object, you might present the refusal in your UI, or include conditional logic in code that consumes the response to handle the case of a refused request.

python ◊

```
15 response_format=MathReasoning,
16 )
17
18 math_reasoning = completion.choices[0].message
19
20 # If the model refuses to respond, you will get a
21
22 if (math_reasoning.refusal):
23     print(math_reasoning.refusal)
24 else:
25     print(math_reasoning.parsed)
```

The API response from a refusal will look something like this:

```
json ◊
1 {
2     "id": "chatcmpl-9nYAG9LPNonX8DAyrkwYfemr3C8HC",
3     "object": "chat.completion",
4     "created": 1721596428,
5     "model": "gpt-4o-2024-08-06",
6     "choices": [
7         {
8             "index": 0,
9             "message": {
10                 "role": "assistant",
11                 "refusal": "I'm sorry, I cannot assist with that request.",
12             },
13             "logprobs": null,
14             "finish_reason": "stop"
15         }
16     ],
17     "usage": {
18         "prompt_tokens": 81,
19         "completion_tokens": 11,
20         "total_tokens": 92,
21         "completion_tokens_details": {
22             "reasoning_tokens": 0,
23             "accepted_prediction_tokens": 0,
24             "rejected_prediction_tokens": 0
25         }
26     },
27     "system_fingerprint": "fp_3407719c7f"
28 }
```

Tips and best practices

Handling user-generated input

If your application is using user-generated input, make sure your prompt includes instructions on how to handle situations where the input cannot result in a valid response.

The model will always try to adhere to the provided schema, which can result in hallucinations if the input is completely unrelated to the schema.

You could include language in your prompt to specify that you want to return empty parameters, or a specific sentence, if the model detects that the input is incompatible with the task.

Handling mistakes

Structured Outputs can still contain mistakes. If you see mistakes, try adjusting your instructions, providing examples in the system instructions, or splitting tasks into simpler subtasks. Refer to the [prompt engineering guide](#) for more guidance on how to tweak your inputs.

Avoid JSON schema divergence

To prevent your JSON Schema and corresponding types in your programming language from diverging, we strongly recommend using the native Pydantic/zod sdk support.

If you prefer to specify the JSON schema directly, you could add CI rules that flag when either the JSON schema or underlying data objects are edited, or add a CI step that auto-generates the JSON Schema from type definitions (or vice-versa).

Streaming

You can use streaming to process model responses or

function call arguments as they are being generated, and parse them as structured data.

That way, you don't have to wait for the entire response to complete before handling it. This is particularly useful if you would like to display JSON fields one by one, or handle function call arguments as soon as they are available.

We recommend relying on the SDKs to handle streaming with Structured Outputs.

You can find an example of how to stream function call arguments without the SDK `stream` helper in the [function calling guide](#).

Here is how you can stream a model response with the `stream` helper:

```
python ◊
```

```
1 from typing import List
2 from pydantic import BaseModel
3 from openai import OpenAI
4
5 class EntitiesModel(BaseModel):
6     attributes: List[str]
7     colors: List[str]
8     animals: List[str]
9
10 client = OpenAI()
11
12 with client.beta.chat.completions.stream(
13     model="gpt-4.1",
14     messages=[
15         {"role": "system", "content": "Extract entities from this sentence."},
16         {
17             "role": "user",
18             "content": "The quick brown fox jumps over the fence."
19         },
20     ],
21     response_format=EntitiesModel,
22 ) as stream:
23     for event in stream:
24         if event.type == "content.delta":
25             if event.parsed is not None:
26                 # Print the parsed data as JSON
27                 print("content.delta parsed:", event.parsed)
28         elif event.type == "content.done":
```

```
29         print("content.done")
30     elif event.type == "error":
31         print("Error in stream:", event.error)
32
33 final_completion = stream.get_final_completion()
34 print("Final completion:", final_completion)
```

You can also use the `stream` helper to parse function call arguments:

```
python ◊
```

```
1 from pydantic import BaseModel
2 import openai
3 from openai import OpenAI
4
5 class GetWeather(BaseModel):
6     city: str
7     country: str
8
9     client = OpenAI()
10
11 with client.beta.chat.completions.stream(
12     model="gpt-4.1",
13     messages=[
14         {
15             "role": "user",
16             "content": "What's the weather like in "
17         },
18     ],
19     tools=[
20         openai.pydantic_function_tool(GetWeather,
21     ],
22     parallel_tool_calls=True,
23 ) as stream:
24     for event in stream:
25         if event.type == "tool_calls.function.argv":
26             print(event)
27
28 print(stream.get_final_completion())
```

Supported schemas

Structured Outputs supports a subset of the [JSON Schema](#) language.

Supported types

The following types are supported for Structured Outputs:

String

Number

Boolean

Integer

Object

Array

Enum

anyOf

Supported properties

In addition to specifying the type of a property, you can specify a selection of additional constraints:

Supported `string` properties:

`pattern` — A regular expression that the string must match.

`format` — Predefined formats for strings. Currently supported:

`date-time`

`time`

`date`

`duration`

`email`

`hostname`

`ipv4`

`ipv6`

`uuid`

Supported `number` properties:

`multipleOf` — The number must be a multiple of this value.

`maximum` — The number must be less than or equal to this value.

`exclusiveMaximum` — The number must be less than this value.

`minimum` — The number must be greater than or equal to this value.

`exclusiveMinimum` — The number must be greater than this value.

Supported `array` properties:

`minItems` — The array must have at least this many items.

`maxItems` — The array must have at most this many items.

Here are some examples on how you can use these type restrictions:

[String Restrictions](#)

[Number Restrictions](#)

```
1  {
2      "name": "user_data",
3      "strict": true,
4      "schema": {
5          "type": "object",
6          "properties": {
7              "name": {
8                  "type": "string",
9                  "description": "The name of the user"
10             },
11             "username": {
12                 "type": "string",
13                 "description": "The username of the user"
14                 "pattern": "^[a-zA-Z0-9_]+$"
15             },
16             "email": {
17                 "type": "string",
18                 "description": "The email of the user"
19                 "format": "email"
20             }
21         },
22         "additionalProperties": false,
```

```
23     "required": [
24         "name", "username", "email"
25     ]
26 }
27 }
```

ⓘ Note these constraints are not yet supported for fine-tuned models.

Root objects must not be anyOf and must be an object

Note that the root level object of a schema must be an object, and not use `anyOf`. A pattern that appears in Zod (as one example) is using a discriminated union, which produces an `anyOf` at the top level. So code such as the following won't work:

javascript ◊

```
1 import { z } from 'zod';
2 import { zodResponseFormat } from 'openai/helpers',
3
4 const BaseResponseSchema = z.object({/* ... */});
5 const UnsuccessfulResponseSchema = z.object({/* .. .
6
7 const finalSchema = z.discriminatedUnion('status',
8 BaseResponseSchema,
9 UnsuccessfulResponseSchema,
10 ]);
11
12 // Invalid JSON Schema for Structured Outputs
13 const json = zodResponseFormat(finalSchema, 'final')
```

All fields must be required

To use Structured Outputs, all fields or function parameters must be specified as `required`.

json ◊

```
1 {
2     "name": "get_weather",
3     "description": "Fetches the weather in the gi
```

```
4     "strict": true,
5     "parameters": {
6         "type": "object",
7         "properties": {
8             "location": {
9                 "type": "string",
10                "description": "The location to get the weather for"
11            },
12            "unit": {
13                "type": "string",
14                "description": "The unit to return temperature in",
15                "enum": ["F", "C"]
16            }
17        },
18        "additionalProperties": false,
19        "required": ["location", "unit"]
20    }
21 }
```

Although all fields must be required (and the model will return a value for each parameter), it is possible to emulate an optional parameter by using a union type with `null`.



```
1  {
2      "name": "get_weather",
3      "description": "Fetches the weather in the given location",
4      "strict": true,
5      "parameters": {
6          "type": "object",
7          "properties": {
8              "location": {
9                  "type": "string",
10                 "description": "The location to get the weather for"
11             },
12             "unit": {
13                 "type": ["string", "null"],
14                 "description": "The unit to return the temperature in",
15                 "enum": ["F", "C"]
16             }
17         },
18         "additionalProperties": false,
19         "required": [
20             "location", "unit"
21         ]
22     }
23 }
```

Objects have limitations on nesting depth and size

A schema may have up to 5000 object properties total, with up to 10 levels of nesting.

Limitations on total string size

In a schema, total string length of all property names, definition names, enum values, and const values cannot exceed 120,000 characters.

Limitations on enum size

A schema may have up to 1000 enum values across all enum properties.

For a single enum property with string values, the total string length of all enum values cannot exceed 15,000 characters when there are more than 250 enum values.

additionalProperties: false must always be set in objects

`additionalProperties` controls whether it is allowable for an object to contain additional keys / values that were not defined in the JSON Schema.

Structured Outputs only supports generating specified keys / values, so we require developers to set `additionalProperties: false` to opt into Structured Outputs.

```
json ⚙️  
1  {  
2      "name": "get_weather",  
3      "description": "Fetches the weather in the given location",  
4      "strict": true,  
5      "schema": {  
6          "type": "object",  
7          "properties": {  
8              "location": {  
9                  "type": "string",  
10                 "description": "The location to get the weather for"},  
11             },  
12             "unit": {  
13                 "type": "string",  
14                 "description": "The unit to return the temperature in",  
15                 "enum": ["F", "C"]  
16             }  
17         },  
18         "additionalProperties": false,  
19         "required": [  
20             "location", "unit"  
21         ]  
22     }  
23 }
```

Key ordering

When using Structured Outputs, outputs will be produced in the same order as the ordering of keys in the schema.

Some type-specific keywords are not yet supported

Composition: `allOf`, `not`, `dependentRequired`,
`dependentSchemas`, `if`, `then`, `else`

For fine-tuned models, we additionally do not support the following:

For strings: `minLength`, `maxLength`, `pattern`,
`format`

For numbers: `minimum`, `maximum`, `multipleOf`

For objects: `patternProperties`

For arrays: `minItems`, `maxItems`

If you turn on Structured Outputs by supplying `strict: true` and call the API with an unsupported JSON Schema, you will receive an error.

For anyOf, the nested schemas must each be a valid JSON Schema per this subset

Here's an example supported anyOf schema:

```
json ⚙️

1  {
2      "type": "object",
3      "properties": {
4          "item": {
5              "anyOf": [
6                  {
7                      "type": "object",
8                      "description": "The user object",
9                      "properties": {
10                          "name": {
11                              "type": "string",
12                              "description": "The name of the user"
13                          },
14                          "age": {
15                              "type": "number",
16                              "description": "The age of the user"
17                          }
18                      },
19                      "additionalProperties": false,
20                      "required": [
21                          "name",
22                          "age"
23                      ]
24                  }
25              ]
26          }
27      }
28  }
```

```
24     },
25     {
26         "type": "object",
27         "description": "The address of the location",
28         "properties": {
29             "number": {
30                 "type": "string",
31                 "description": "The number of the building"
32             },
33             "street": {
34                 "type": "string",
35                 "description": "The street name"
36             },
37             "city": {
38                 "type": "string",
39                 "description": "The city name"
40             }
41         },
42         "additionalProperties": false,
43         "required": [
44             "number",
45             "street",
46             "city"
47         ]
48     }
49 ],
50 }
51 },
52 "additionalProperties": false,
53 "required": [
54     "item"
55 ]
56 }
```

Definitions are supported

You can use definitions to define subschemas which are referenced throughout your schema. The following is a simple example.

```
1  {
2      "type": "object",
3      "properties": {
4          "steps": {
5              "type": "array",
6              "items": {
```

```
7             "$ref": "#/$defs/step"
8         }
9     },
10    "final_answer": {
11        "type": "string"
12    }
13},
14 "$defs": {
15     "step": {
16         "type": "object",
17         "properties": {
18             "explanation": {
19                 "type": "string"
20             },
21             "output": {
22                 "type": "string"
23             }
24         },
25         "required": [
26             "explanation",
27             "output"
28         ],
29         "additionalProperties": false
30     }
31 },
32 "required": [
33     "steps",
34     "final_answer"
35 ],
36 "additionalProperties": false
37 }
```

Recursive schemas are supported

Sample recursive schema using `#` to indicate root recursion.

```
1  {
2      "name": "ui",
3      "description": "Dynamically generated UI",
4      "strict": true,
5      "schema": {
6          "type": "object",
7          "properties": {
8              "type": {
9                  "type": "string",
10                 "description": "The type of the U"
```

```
1          "enum": ["div", "button", "header"
2      },
3      "label": {
4          "type": "string",
5          "description": "The label of the component"
6      },
7      "children": {
8          "type": "array",
9          "description": "Nested UI components"
10         "items": {
11             "$ref": "#"
12         }
13     },
14     "attributes": {
15         "type": "array",
16         "description": "Arbitrary attributes"
17         "items": {
18             "type": "object",
19             "properties": {
20                 "name": {
21                     "type": "string",
22                     "description": "The name of the attribute"
23                 },
24                 "value": {
25                     "type": "string",
26                     "description": "The value of the attribute"
27                 }
28             },
29             "additionalProperties": false,
30             "required": ["name", "value"]
31         }
32     }
33 },
34 "required": ["type", "label", "children"],
35 "additionalProperties": false
36 }
37 }
```

Sample recursive schema using explicit recursion:

```
1  {
2      "type": "object",
3      "properties": {
4          "linked_list": {
5              "$ref": "#/$defs/linked_list_node"
6          }
7      },
8      "definitions": {
9          "linked_list_node": {
10              "type": "object",
11              "properties": {
12                  "value": "string",
13                  "next": {
14                      "$ref": "#/$defs/linked_list_node"
15                  }
16              }
17          }
18      }
19  }
```

```
8     "$defs": {
9         "linked_list_node": {
10             "type": "object",
11             "properties": {
12                 "value": {
13                     "type": "number"
14                 },
15                 "next": {
16                     "anyOf": [
17                         {
18                             "$ref": "#/$defs/linked_list"
19                         },
20                         {
21                             "type": "null"
22                         }
23                     ]
24                 }
25             },
26             "additionalProperties": false,
27             "required": [
28                 "next",
29                 "value"
30             ]
31         },
32         "additionalProperties": false,
33         "required": [
34             "linked_list"
35         ]
36     }
37 }
```

JSON mode

JSON mode is a more basic version of the Structured Outputs feature. While JSON mode ensures that model output is valid JSON, Structured Outputs reliably matches the model's output to the schema you specify. We recommend you use Structured Outputs if it is supported for your use case.

When JSON mode is turned on, the model's output is ensured to be valid JSON, except for in some edge cases that you should detect and handle appropriately.

To turn on JSON mode with the Chat Completions, set the `response_format` to `{ "type": "json_object" }`. If you

are using function calling, JSON mode is always turned on.

Important notes:

When using JSON mode, you must always instruct the model to produce JSON via some message in the conversation, for example via your system message. If you don't include an explicit instruction to generate JSON, the model may generate an unending stream of whitespace and the request may run continually until it reaches the token limit. To help ensure you don't forget, the API will throw an error if the string "JSON" does not appear somewhere in the context.

JSON mode will not guarantee the output matches any specific schema, only that it is valid and parses without errors. You should use Structured Outputs to ensure it matches your schema, or if that is not possible, you should use a validation library and potentially retries to ensure that the output matches your desired schema.

Your application must detect and handle the edge cases that can result in the model output not being a complete JSON object (see below)

› Handling edge cases

Resources

To learn more about Structured Outputs, we recommend browsing the following resources:

Check out our [introductory cookbook](#) on Structured Outputs

Learn [how to build multi-agent systems](#) with Structured Outputs