
✨ AI SDK 5 is available now.

[View Announcement](#)

> Menu

AI SDK Providers > OpenAI

 Copy markdown

OpenAI Provider

The [OpenAI](#) provider contains language model support for the OpenAI responses, chat, and completion APIs, as well as embedding model support for the OpenAI embeddings API.

Setup

The OpenAI provider is available in the `@ai-sdk/openai` module. You can install it with

pnpm npm yarn bun

```
$ pnpm add @ai-sdk/openai
```



Provider Instance

You can import the default provider instance `openai` from `@ai-sdk/openai`:

```
1 import { openai } from '@ai-sdk/openai';
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

Deny

Accept all

Consent Settings

`openAI` from `@ai-sdk/openai` and

```

1  import { createOpenAI } from '@ai-sdk/openai';
2
3  const openai = createOpenAI({
4    // custom settings, e.g.
5    headers: {
6      'header-name': 'header-value',
7    },
8  });

```

You can use the following optional settings to customize the OpenAI provider instance:

- **baseUrl** *string*

Use a different URL prefix for API calls, e.g. to use proxy servers. The default prefix is

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

- **apiKey** *string*

API key that is being sent using the `Authorization` header. It defaults to the

`OPENAI_API_KEY` environment variable.

- **name** *string*

The provider name. You can set this when using OpenAI compatible providers to change the model provider property. Defaults to `openai`.

- **organization** *string*

OpenAI Organization.

- **project** *string*

OpenAI project.

- **headers** *Record<string,string>*

Custom headers to include in the requests.

- **fetch** *(input: RequestInfo, init?: RequestInit) => Promise<Response>*

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

al `fetch` function. You can use it as custom fetch implementation for e.g.

Language Models

The OpenAI provider instance is a function that you can invoke to create a language model:

```
1  const model = openai('gpt-5');
```

It automatically selects the correct API based on the model id. You can also pass additional settings in the second argument:

```
1  const model = openai('gpt-5', {  
2    // additional settings  
3  });
```

The available options depend on the API that's automatically chosen for the model (see below). If you want to explicitly select a specific model API, you can use `.chat` or `.completion`.

Example

You can use OpenAI language models to generate text with the `generateText` function:

```
1  import { openai } from '@ai-sdk/openai';  
2  import { generateText } from 'ai';  
3  
4  const { text } = await generateText({  
5    model: openai('gpt-5'),  
6    prompt: 'Write a vegetarian lasagna recipe for 4 people.',  
7  });
```

OpenAI language models can also be used in the `streamText`, `generateObject`, and `streamObject` functions (see [AI SDK Core](#)).

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

Using the `.chat()` factory method.
AI chat models support tool calls

and some have multi-modal capabilities.

```
1  const model = openai.chat('gpt-5');
```

OpenAI chat models support also some model specific provider options that are not part of the [standard call settings](#). You can pass them in the `providerOptions` argument:

```
1  const model = openai.chat('gpt-5');
2
3  await generateText({
4    model,
5    providerOptions: {
6      openai: {
7        logitBias: {
8          // optional likelihood for specific tokens
9          '50256': -100,
10       },
11       user: 'test-user', // optional unique user identifier
12     },
13   },
14 });
```

The following optional provider options are available for OpenAI chat models:

- **logitBias** *Record<number, number>*

Modifies the likelihood of specified tokens appearing in the completion.

Accepts a JSON object that maps tokens (specified by their token ID in the GPT tokenizer) to an associated bias value from -100 to 100. You can use this tokenizer tool to convert text to token IDs. Mathematically, the bias is added to the logits generated by

result in a ban or exclusive selection of the relevant token.

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

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

logprobs will increase the response can be useful to better understand

how the model is behaving.

Setting to true will return the log probabilities of the tokens that were generated.

Setting to a number will return the log probabilities of the top n tokens that were generated.

- **parallelToolCalls** *boolean*

Whether to enable parallel function calling during tool use. Defaults to `true`.

- **user** *string*

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

- **reasoningEffort** *'minimal' / 'low' / 'medium' / 'high'*

Reasoning effort for reasoning models. Defaults to `medium`. If you use `providerOptions` to set the `reasoningEffort` option, this model setting will be ignored.

- **structuredOutputs** *boolean*

Whether to use [structured outputs](#). Defaults to `true`.

When enabled, tool calls and object generation will be strict and follow the provided schema.

- **maxCompletionTokens** *number*

Maximum number of completion tokens to generate. Useful for reasoning models.

- **store** *boolean*

Whether to enable persistence in Responses API.

- **metadata** *Record<string, string>*

Metadata to associate with the request.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

Service tier for the request. Set to 'flex' for 50% cheaper processing at the cost of increased latency. Only available for o3, o4-mini, and gpt-5 models. Defaults to 'auto'.

- **strictJsonSchema** *boolean*

Whether to use strict JSON schema validation. Defaults to `false`.

- **textVerbosity** *'low' | 'medium' | 'high'*

Controls the verbosity of the model's responses. Lower values will result in more concise responses, while higher values will result in more verbose responses.

- **promptCacheKey** *string*

A cache key for manual prompt caching control. Used by OpenAI to cache responses for similar requests to optimize your cache hit rates.

- **safetyIdentifier** *string*

A stable identifier used to help detect users of your application that may be violating OpenAI's usage policies. The IDs should be a string that uniquely identifies each user.

Reasoning

OpenAI has introduced the `o1`, `o3`, and `o4` series of [reasoning models](#)⁷. Currently, `o4-mini`, `o3`, `o3-mini`, and `o1` are available via both the chat and responses APIs. The models `codex-mini-latest` and `computer-use-preview` are available only via the [responses API](#).

Reasoning models currently only generate text, have several limitations, and are only supported using `generateText` and `streamText`.

They support additional settings and response metadata:

- You can use `providerOptions` to set
 - the `reasoningEffort` option (or alternatively the `reasoningEffort` model setting), which determines the amount of reasoning the model performs.
- You can use response `providerMetadata` to access the number of reasoning tokens

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```

3
4  const { text, usage, providerMetadata } = await generateText({
5    model: openai('gpt-5'),
6    prompt: 'Invent a new holiday and describe its traditions.',
7    providerOptions: {
8      openai: {
9        reasoningEffort: 'low',
10     },
11   },
12 });
13
14 console.log(text);
15 console.log('Usage:', {
16   ...usage,
17   reasoningTokens: providerMetadata?.openai?.reasoningTokens,
18 });

```



System messages are automatically converted to OpenAI developer messages for reasoning models when supported.



Reasoning models require additional runtime inference to complete their reasoning phase before generating a response. This introduces longer latency compared to other models.



`maxOutputTokens` is automatically mapped to `max_completion_tokens` for reasoning models.

Structured Outputs

Structured outputs are enabled by default. You can disable them by setting the `structuredOutputs` option to `false`.

```

1  import { openai } from '@ai-sdk/openai';
2  import { generateObject } from 'ai';
3  import { z } from 'zod';
4
5  const result = await generateObject({

```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```
12     schemaName: 'recipe',
13     schemaDescription: 'A recipe for lasagna.',
14     schema: z.object({
15       name: z.string(),
16       ingredients: z.array(
17         z.object({
18           name: z.string(),
19           amount: z.string(),
20         }),
21       ),
22       steps: z.array(z.string()),
23     }),
24     prompt: 'Generate a lasagna recipe.',
25   });
26
27   console.log(JSON.stringify(result.object, null, 2));
```

OpenAI structured outputs have several [limitations](#)[↗], in particular around the [supported schemas](#)[↗], and are therefore opt-in.



For example, optional schema properties are not supported. You need to change Zod `.nullish()` and `.optional()` to `.nullable()`.

Logprobs

OpenAI provides logprobs information for completion/chat models. You can access it in the `providerMetadata` object.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.


```

1  import { openai } from '@ai-sdk/openai';
2  import { generateText } from 'ai';
3
4  const result = await generateText({
5    model: openai('gpt-5'),
6    prompt: 'Write a vegetarian lasagna recipe for 4 people.',
7    providerOptions: {
8      openai: {
9        // this can also be a number,
10        // refer to logprobs provider options section for more
11        logprobs: true,
12      },
13    },
14  });
15
16  const openaiMetadata = (await result.providerMetadata)?.openai;
17
18  const logprobs = openaiMetadata?.logprobs;

```

Image Support

The OpenAI Chat API supports Image inputs for appropriate models. You can pass Image files as part of the message content using the 'image' type:

```

1  const result = await generateText({
2    model: openai('gpt-5'),
3    messages: [
4      {
5        role: 'user',
6        content: [
7          {
8            type: 'text',
9            text: 'Please describe the image.',
10          },
11          {
12            type: 'image',
13            image: fs.readFileSync('./data/image.png'),
14          },
15        ],
16      },
17    ],
18  });

```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

and to questions about it. The image

You can also pass the URL of an image.

```
1  {  
2    type: 'image',  
3    image: 'https://sample.edu/image.png',  
4  }
```

PDF support

The OpenAI Chat API supports reading PDF files. You can pass PDF files as part of the message content using the `file` type:

```
1  const result = await generateText({  
2    model: openai('gpt-5'),  
3    messages: [  
4      {  
5        role: 'user',  
6        content: [  
7          {  
8            type: 'text',  
9            text: 'What is an embedding model?',  
10           },  
11          {  
12            type: 'file',  
13            data: fs.readFileSync('./data/ai.pdf'),  
14            mediaType: 'application/pdf',  
15            filename: 'ai.pdf', // optional  
16          },  
17        ],  
18      },  
19    ],  
20  });
```

The model will have access to the contents of the PDF file and respond to questions about it. The PDF file should be passed using the `data` field, and the `mediaType` should be set to `'application/pdf'`.

You can also pass a file-id from the OpenAI Files API.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```
1  {
2    type: 'file',
3    data: 'file-8EFBcWHsQxZV7YGezBC1fq',
4    mediaType: 'application/pdf',
5  }
```

You can also pass the URL of a PDF.

```
1  {
2    type: 'file',
3    data: 'https://sample.edu/example.pdf',
4    mediaType: 'application/pdf',
5    filename: 'ai.pdf', // optional
6  }
```

Predicted Outputs

OpenAI supports [predicted outputs](#) for `gpt-4o` and `gpt-4o-mini`. Predicted outputs help you reduce latency by allowing you to specify a base text that the model should modify. You can enable predicted outputs by adding the `prediction` option to the `providerOptions.openai` object:

```
1  const result = streamText({
2    model: openai('gpt-5'),
3    messages: [
4      {
5        role: 'user',
6        content: 'Replace the Username property with an Email property.',
7      },
8      {
9        role: 'user',
10       content: existingCode,
11     },
12   ],
13   providerOptions: {
14     openai: {
15       prediction: {
16         type: 'content',
17         content: existingCode
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

OpenAI provides usage information for predicted outputs (`acceptedPredictionTokens` and `rejectedPredictionTokens`). You can access it in the `providerMetadata` object.

```
1  const openaiMetadata = (await result.providerMetadata)?.openai;
2
3  const acceptedPredictionTokens = openaiMetadata?.acceptedPredictionTokens;
4  const rejectedPredictionTokens = openaiMetadata?.rejectedPredictionTokens;
```



OpenAI Predicted Outputs have several [limitations](#)⁷, e.g. unsupported API parameters and no tool calling support.

Image Detail

You can use the `openai` provider option to set the [image input detail](#)⁷ to `high`, `low`, or `auto`:

```
1  const result = await generateText({
2    model: openai('gpt-5'),
3    messages: [
4      {
5        role: 'user',
6        content: [
7          { type: 'text', text: 'Describe the image in detail.' },
8          {
9            type: 'image',
10           image:
11             'https://github.com/vercel/ai/blob/main/examples/ai-core/data/comic-c
12
13           // OpenAI specific options - image detail:
14           providerOptions: {
15             openai: { imageDetail: 'low' },
16           },
17         },
18       ],
19     },
20   ],
21 });
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.



`useChat`) does not support the `messages` first before passing the

messages to functions like `generateText` or `streamText`. For more details on `providerOptions` usage, see [here](#).

Distillation

OpenAI supports model distillation for some models. If you want to store a generation for use in the distillation process, you can add the `store` option to the `providerOptions.openai` object. This will save the generation to the OpenAI platform for later use in distillation.

```
1 import { openai } from '@ai-sdk/openai';
2 import { generateText } from 'ai';
3 import 'dotenv/config';
4
5 async function main() {
6   const { text, usage } = await generateText({
7     model: openai('gpt-4o-mini'),
8     prompt: 'Who worked on the original macintosh?',
9     providerOptions: {
10       openai: {
11         store: true,
12         metadata: {
13           custom: 'value',
14         },
15       },
16     },
17   });
18
19   console.log(text);
20   console.log();
21   console.log('Usage:', usage);
22 }
23
24 main().catch(console.error);
```

Prompt Caching

OpenAI has introduced [Prompt Caching](#) for supported models including `gpt-4o` and `gpt-4o-mini`.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

models, when the prompt is 1024 tokens or less. Prompt caching is only available for models with a maximum context length of 128,000 tokens, when the prompt is 1024 tokens or less.

the number of prompt tokens that

- Note that caching behavior is dependent on load on OpenAI's infrastructure. Prompt prefixes generally remain in the cache following 5-10 minutes of inactivity before they are evicted, but during off-peak periods they may persist for up to an hour.

```
1 import { openai } from '@ai-sdk/openai';
2 import { generateText } from 'ai';
3
4 const { text, usage, providerMetadata } = await generateText({
5   model: openai('gpt-4o-mini'),
6   prompt: `A 1024-token or longer prompt...`,
7 });
8
9 console.log(`usage:`, {
10   ...usage,
11   cachedPromptTokens: providerMetadata?.openai?.cachedPromptTokens,
12 });
```

To improve cache hit rates, you can manually control caching using the `promptCacheKey` option:

```
1 import { openai } from '@ai-sdk/openai';
2 import { generateText } from 'ai';
3
4 const { text, usage, providerMetadata } = await generateText({
5   model: openai('gpt-5'),
6   prompt: `A 1024-token or longer prompt...`,
7   providerOptions: {
8     openai: {
9       promptCacheKey: 'my-custom-cache-key-123',
10     },
11   },
12 });
13
14 console.log(`usage:`, {
15   ...usage,
16   cachedPromptTokens: providerMetadata?.openai?.cachedPromptTokens,
17 });
```

Audio Input

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

audio files to the model.



The `gpt-4o-audio-preview` model is currently in preview and requires at least some audio inputs. It will not work with non-audio data.

```
1 import { openai } from '@ai-sdk/openai';
2 import { generateText } from 'ai';
3
4 const result = await generateText({
5   model: openai('gpt-4o-audio-preview'),
6   messages: [
7     {
8       role: 'user',
9       content: [
10        { type: 'text', text: 'What is the audio saying?' },
11        {
12          type: 'file',
13          mediaType: 'audio/mpeg',
14          data: fs.readFileSync('./data/galileo.mp3'),
15        },
16      ],
17    },
18  ],
19 });
```

Responses Models

You can use the OpenAI responses API with the `openai.responses(modelId)` factory method.

```
1 const model = openai.responses('gpt-5');
```

Further configuration can be done using OpenAI provider options. You can validate the provider options using the `OpenAIResponsesProviderOptions` type.

```
1 import { openai, OpenAIResponsesProviderOptions } from '@ai-sdk/openai';
2 import { generateText } from 'ai';
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```
3   parallelToolCalls: false,
```

```

9      store: false,
10     user: 'user_123',
11     // ...
12   } satisfies OpenAIResponsesProviderOptions,
13 },
14 // ...
15 });

```

The following provider options are available:

- **parallelToolCalls** *boolean* Whether to use parallel tool calls. Defaults to `true`.

- **store** *boolean*

Whether to store the generation. Defaults to `true`.

When using reasoning models (o1, o3, o4-mini) with multi-step tool calls and `store: false`, include `['reasoning.encrypted_content']` in the `include` option to ensure reasoning content is available across conversation steps.

- **metadata** *Record<string, string>* Additional metadata to store with the generation.
- **previousResponseId** *string* The ID of the previous response. You can use it to continue a conversation. Defaults to `undefined`.
- **instructions** *string* Instructions for the model. They can be used to change the system or developer message when continuing a conversation using the `previousResponseId` option. Defaults to `undefined`.
- **user** *string* A unique identifier representing your end-user, which can help OpenAI to monitor and detect abuse. Defaults to `undefined`.
- **reasoningEffort** *'minimal' | 'low' | 'medium' | 'high'* Reasoning effort for reasoning models. Defaults to `medium`. If you use `providerOptions` to set the `reasoningEffort` option, this model setting will be ignored.
- **reasoningSummary** *'auto' | 'detailed'* Controls whether the model returns its reasoning process. Set to `'auto'` for a condensed summary, `'detailed'` for more

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

to reasoning summaries). When as events with type `'reasoning'` ng field.

- **strictJsonSchema** *boolean* Whether to use strict JSON schema validation. Defaults to `false`.
- **serviceTier** *'auto' | 'flex' | 'priority'* Service tier for the request. Set to 'flex' for 50% cheaper processing at the cost of increased latency (available for o3, o4-mini, and gpt-5 models). Set to 'priority' for faster processing with Enterprise access (available for gpt-4, gpt-5, gpt-5-mini, o3, o4-mini; gpt-5-nano is not supported). Defaults to 'auto'.
- **textVerbosity** *'low' | 'medium' | 'high'* Controls the verbosity of the model's response. Lower values result in more concise responses, while higher values result in more verbose responses. Defaults to `'medium'`.
- **include** *Array<string>* Specifies additional content to include in the response. Supported values: `['reasoning.encrypted_content']` for accessing reasoning content across conversation steps, and `['file_search_call.results']` for including file search results in responses. Defaults to `undefined`.
- **promptCacheKey** *string* A cache key for manual prompt caching control. Used by OpenAI to cache responses for similar requests to optimize your cache hit rates.
- **safetyIdentifier** *_string_0* A stable identifier used to help detect users of your application that may be violating OpenAI's usage policies. The IDs should be a string that uniquely identifies each user.

The OpenAI responses provider also returns provider-specific metadata:

```
1  const { providerMetadata } = await generateText({
2    model: openai.responses('gpt-5'),
3  });
4
5  const openaiMetadata = providerMetadata?.openai;
```

The following OpenAI-specific metadata is returned:

- **responseId** *string* The ID of the response. Can be used to continue a conversation.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

not tokens that were a cache hit.

tokens that the model generated.

Web Search

The OpenAI responses API supports web search through the

`openai.tools.webSearchPreview` tool.

You can force the use of the web search tool by setting the `toolChoice` parameter to `{ type: 'tool', toolName: 'web_search_preview' }`.

```
1  const result = await generateText({
2    model: openai.responses('gpt-5'),
3    prompt: 'What happened in San Francisco last week?',
4    tools: {
5      web_search_preview: openai.tools.webSearchPreview({
6        // optional configuration:
7        searchContextSize: 'high',
8        userLocation: {
9          type: 'approximate',
10         city: 'San Francisco',
11         region: 'California',
12       },
13     },
14   },
15   // Force web search tool:
16   toolChoice: { type: 'tool', toolName: 'web_search_preview' },
17 });
18
19 // URL sources
20 const sources = result.sources;
```

Reasoning Output

For reasoning models like `gpt-5`, you can enable reasoning summaries to see the model's thought process. Different models support different summarizers—for example, `o4-mini` supports detailed summaries. Set `reasoningSummary: "auto"` to automatically receive the richest level available.

```
1  import { openai } from '@ai-sdk/openai';
2  import { streamText } from 'ai';
3
4  const result = streamText({
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

ebate in San Francisco.',

for condensed or 'detailed' for con

```

11     },
12   });
13
14   for await (const part of result.fullStream) {
15     if (part.type === 'reasoning') {
16       console.log(`Reasoning: ${part.textDelta}`);
17     } else if (part.type === 'text-delta') {
18       process.stdout.write(part.textDelta);
19     }
20   }

```

For non-streaming calls with `generateText`, the reasoning summaries are available in the `reasoning` field of the response:

```

1  import { openai } from '@ai-sdk/openai';
2  import { generateText } from 'ai';
3
4  const result = await generateText({
5    model: openai.responses('gpt-5'),
6    prompt: 'Tell me about the Mission burrito debate in San Francisco.',
7    providerOptions: {
8      openai: {
9        reasoningSummary: 'auto',
10      },
11    },
12  });
13  console.log('Reasoning:', result.reasoning);

```

Learn more about reasoning summaries in the [OpenAI documentation](#).

Verbosity Control

You can control the length and detail of model responses using the `textVerbosity` parameter:

```

1  import { openai } from '@ai-sdk/openai';
2  import { generateText } from 'ai';
3
4  const result = await generateText({

```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```

    prompt: 'What is the best pet dog.',

```

```

    use, 'medium' (default), or 'high' if

```

```

10    },

```

```
11     },
12   });
```

The `textVerbosity` parameter scales output length without changing the underlying prompt:

- `'low'`: Produces terse, minimal responses
- `'medium'`: Balanced detail (default)
- `'high'`: Verbose responses with comprehensive detail


File Search

The OpenAI responses API supports file search through the `openai.tools.fileSearch` tool.

You can force the use of the file search tool by setting the `toolChoice` parameter to `{ type: 'tool', toolName: 'file_search' }`.

```
1  const result = await generateText({
2    model: openai.responses('gpt-5'),
3    prompt: 'What does the document say about user authentication?',
4    tools: {
5      file_search: openai.tools.fileSearch({
6        // optional configuration:
7        vectorStoreIds: ['vs_123', 'vs_456'],
8        maxNumResults: 10,
9        ranking: {
10         ranker: 'auto',
11       },
12       filters: {
13         type: 'and',
14         filters: [
15           { key: 'author', type: 'eq', value: 'John Doe' },
16           { key: 'date', type: 'gte', value: '2023-01-01' },
17         ],
18       },
19     }),
20   },
21   // Force file search tool:
22   toolChoice: { type: 'tool', toolName: 'file_search' },
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

 The tool must be named `file_search` when using OpenAI's file search functionality. This name is required by OpenAI's API specification and cannot be customized.

Code Interpreter

The OpenAI responses API supports the code interpreter tool through the `openai.tools.codeInterpreter` tool. This allows models to write and execute Python code.

```
1  import { openai } from '@ai-sdk/openai';
2  import { generateText } from 'ai';
3
4  const result = await generateText({
5    model: openai.responses('gpt-5'),
6    prompt: 'Write and run Python code to calculate the factorial of 10',
7    tools: {
8      code_interpreter: openai.tools.codeInterpreter({
9        // optional configuration:
10       container: {
11         fileIds: ['file-123', 'file-456'], // optional file IDs to make available
12       },
13     },
14   },
15 });
```

The code interpreter tool can be configured with:

- **container:** Either a container ID string or an object with `fileIds` to specify uploaded files that should be available to the code interpreter


 The tool must be named `code_interpreter` when using OpenAI's code interpreter functionality. This name is required by OpenAI's API specification and cannot be customized.

Image Support

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

appropriate models. You can pass
'type':

```

2    model: openai.responses('gpt-5'),
3    messages: [
4      {
5        role: 'user',
6        content: [
7          {
8            type: 'text',
9            text: 'Please describe the image.',
10           },
11           {
12             type: 'image',
13             image: fs.readFileSync('./data/image.png'),
14           },
15         ],
16       },
17     ],
18   });

```

The model will have access to the image and will respond to questions about it. The image should be passed using the `image` field.

You can also pass a file-id from the OpenAI Files API.

```

1  {
2    type: 'image',
3    image: 'file-8EFBcWHsQxZV7YGezBC1fq'
4  }

```

You can also pass the URL of an image.

```

1  {
2    type: 'image',
3    image: 'https://sample.edu/image.png',
4  }

```

PDF support

The OpenAI Responses API supports reading PDF files. You can pass PDF files as part of the message content using the `file` type:

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```

4      {
5        role: 'user',
6        content: [
7          {
8            type: 'text',
9            text: 'What is an embedding model?',
10         },
11         {
12           type: 'file',
13           data: fs.readFileSync('./data/ai.pdf'),
14           mediaType: 'application/pdf',
15           filename: 'ai.pdf', // optional
16         },
17       ],
18     },
19   ],
20 });

```

You can also pass a file-id from the OpenAI Files API.

```

1  {
2    type: 'file',
3    data: 'file-8EFBcWHsQxZV7YGezBC1fq',
4    mediaType: 'application/pdf',
5  }

```

You can also pass the URL of a pdf.

```

1  {
2    type: 'file',
3    data: 'https://sample.edu/example.pdf',
4    mediaType: 'application/pdf',
5    filename: 'ai.pdf', // optional
6  }

```

The model will have access to the contents of the PDF file and respond to questions about it. The PDF file should be passed using the `data` field, and the `mediaType` should be set to `'application/pdf'`.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

You can enforce structured outputs a `schema` option. Additionally, you

can pass a Zod or JSON Schema object to the `experimental_output` option when using `generateText` or `streamText`.

```
1 // Using generateObject
2 const result = await generateObject({
3   model: openai.responses('gpt-4.1'),
4   schema: z.object({
5     recipe: z.object({
6       name: z.string(),
7       ingredients: z.array(
8         z.object({
9           name: z.string(),
10          amount: z.string(),
11        }),
12      ),
13      steps: z.array(z.string()),
14    }),
15  }),
16   prompt: 'Generate a lasagna recipe.',
17 });
18
19 // Using generateText
20 const result = await generateText({
21   model: openai.responses('gpt-4.1'),
22   prompt: 'How do I make a pizza?',
23   experimental_output: Output.object({
24     schema: z.object({
25       ingredients: z.array(z.string()),
26       steps: z.array(z.string()),
27     }),
28   }),
29 });
```

Completion Models

You can create models that call the [OpenAI completions API](#) using the `.completion()` factory method. The first argument is the model id. Currently only `gpt-3.5-turbo-instruct` is supported.

```
1 const model = openai.completion('gpt-3.5-turbo-instruct');
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

specific settings that are not part of
ns argument:


```

1  const model = openai.completion('gpt-3.5-turbo-instruct');
2
3  await model.doGenerate({
4    providerOptions: {
5      openai: {
6        echo: true, // optional, echo the prompt in addition to the completion
7        logitBias: {
8          // optional likelihood for specific tokens
9          '50256': -100,
10       },
11       suffix: 'some text', // optional suffix that comes after a completion of in
12       user: 'test-user', // optional unique user identifier
13     },
14   },
15 });

```

The following optional provider options are available for OpenAI completion models:

- **echo:** *boolean*

Echo back the prompt in addition to the completion.

- **logitBias** *Record<number, number>*

Modifies the likelihood of specified tokens appearing in the completion.

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

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

- **logprobs** *boolean / number*

Return the log probabilities of the tokens. Including logprobs will increase the response

This site uses tracking technologies. You may opt in or opt out of the use of these technologies. can be useful to better understand

tokens that were generated.

Setting to a number will return the log probabilities of the top n tokens that were generated.

- **suffix string**

The suffix that comes after a completion of inserted text.

- **user string**

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

Model Capabilities

Model	Image Input	Audio Input	Object Generation	Tool Usage
<code>gpt-4.1</code>	✓	✗	✓	✓
<code>gpt-4.1-mini</code>	✓	✗	✓	✓
<code>gpt-4.1-nano</code>	✓	✗	✓	✓
<code>gpt-4o</code>	✓	✗	✓	✓
<code>gpt-4o-mini</code>	✓	✗	✓	✓
<code>gpt-4o-audio-preview</code>	✗	✓	✓	✓
<code>gpt-4-turbo</code>	✓	✗	✓	✓
<code>gpt-4</code>	✗	✗	✓	✓
<code>gpt-3.5-turbo</code>	✗	✗	✓	✓
<code>o1</code>	✓	✗	✓	✓
<code>o3-mini</code>	✗	✗	✓	✓
<code>o3</code>	✓	✗	✓	✓
<code>o4-mini</code>	✓	✗	✓	✓
This site uses tracking technologies. You may opt in or opt out of the use of these technologies.			✗	✗
			✓	✓

Model	Image Input	Audio Input	Object Generation	Tool Usage
<code>gpt-5-mini</code>	✓	✗	✓	✓
<code>gpt-5-nano</code>	✓	✗	✓	✓
<code>gpt-5-chat-latest</code>	✓	✗	✗	✗



The table above lists popular models. Please see the [OpenAI docs](#) [↗] for a full list of available models. The table above lists popular models. You can also pass any available provider model ID as a string if needed.

Embedding Models

You can create models that call the [OpenAI embeddings API](#) [↗] using the `.textEmbedding()` factory method.

```
1  const model = openai.textEmbedding('text-embedding-3-large');
```

OpenAI embedding models support several additional provider options. You can pass them as an options argument:

```
1  import { openai } from '@ai-sdk/openai';
2  import { embed } from 'ai';
3
4  const { embedding } = await embed({
5    model: openai.textEmbedding('text-embedding-3-large'),
6    value: 'sunny day at the beach',
7    providerOptions: {
8      openai: {
9        dimensions: 512, // optional, number of dimensions for the embedding
10       user: 'test-user', // optional unique user identifier
11     }
12   });
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

The following optional provider options are available for OpenAI embedding models:

- **dimensions:** *number*

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

- **user** *string*

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


Model Capabilities

Model	Default Dimensions	Custom Dimensions
text-embedding-3-large	3072	✓
text-embedding-3-small	1536	✓
text-embedding-ada-002	1536	×

Image Models

You can create models that call the [OpenAI image generation API ↗](#) using the `.image()` factory method.

```
1  const model = openai.image('dall-e-3');
```

 Dall-E models do not support the `aspectRatio` parameter. Use the `size` parameter instead.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

Model	Sizes
<code>gpt-image-1</code>	1024x1024, 1536x1024, 1024x1536
<code>dall-e-3</code>	1024x1024, 1792x1024, 1024x1792
<code>dall-e-2</code>	256x256, 512x512, 1024x1024

You can pass optional `providerOptions` to the image model. These are prone to change by OpenAI and are model dependent. For example, the `gpt-image-1` model supports the `quality` option:

```
1  const { image, providerMetadata } = await generateImage({
2    model: openai.image('gpt-image-1'),
3    prompt: 'A salamander at sunrise in a forest pond in the Seychelles.',
4    providerOptions: {
5      openai: { quality: 'high' },
6    },
7  });
```

For more on `generateImage()` see [Image Generation](#).

OpenAI's image models may return a revised prompt for each image. It can be access at `providerMetadata.openai.images[0]?.revisedPrompt`.

For more information on the available OpenAI image model options, see the [OpenAI API reference](#).

Transcription Models

You can create models that call the [OpenAI transcription API](#) using the `.transcription()` factory method.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

```
1  const model = openai.transcription('whisper-1');
```

You can also pass additional provider-specific options using the `providerOptions` argument. For example, supplying the input language in ISO-639-1 (e.g. `en`) format will improve accuracy and latency.

```
1  import { experimental_transcribe as transcribe } from 'ai';
2  import { openai } from '@ai-sdk/openai';
3
4  const result = await transcribe({
5    model: openai.transcription('whisper-1'),
6    audio: new Uint8Array([1, 2, 3, 4]),
7    providerOptions: { openai: { language: 'en' } },
8  });
```

To get word-level timestamps, specify the granularity:

```
1  import { experimental_transcribe as transcribe } from 'ai';
2  import { openai } from '@ai-sdk/openai';
3
4  const result = await transcribe({
5    model: openai.transcription('whisper-1'),
6    audio: new Uint8Array([1, 2, 3, 4]),
7    providerOptions: {
8      openai: {
9        //timestampGranularities: ['word'],
10       timestampGranularities: ['segment'],
11     },
12   },
13 });
14
15 // Access word-level timestamps
16 console.log(result.segments); // Array of segments with startSecond/endSecond
```

The following provider options are available:

- **timestampGranularities** *string[]* The granularity of the timestamps in the transcription.

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

`d']`, `['segment']`, and `['word']`,
segment timestamps, but
y.

- **language** *string* The language of the input audio. Supplying the input language in ISO-639-1 format (e.g. 'en') will improve accuracy and latency. Optional.
- **prompt** *string* An optional text to guide the model's style or continue a previous audio segment. The prompt should match the audio language. Optional.
- **temperature** *number* The sampling temperature, between 0 and 1. Higher values like 0.8 will make the output more random, while lower values like 0.2 will make it more focused and deterministic. If set to 0, the model will use log probability to automatically increase the temperature until certain thresholds are hit. Defaults to 0. Optional.
- **include** *string[]* Additional information to include in the transcription response.

Model Capabilities

Model	Transcription	Duration	Segments	Language
<code>whisper-1</code>	✓	✓	✓	✓
<code>gpt-4o-mini-transcribe</code>	✓	✗	✗	✗
<code>gpt-4o-transcribe</code>	✓	✗	✗	✗

Speech Models

You can create models that call the [OpenAI speech API](#) using the `.speech()` factory method.

The first argument is the model id e.g. `tts-1`.

```
1  const model = openai.speech('tts-1');
```

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.

ing the `providerOptions` generated audio.

```

1  import { experimental_generateSpeech as generateSpeech } from 'ai';
2  import { openai } from '@ai-sdk/openai';
3
4  const result = await generateSpeech({
5    model: openai.speech('tts-1'),
6    text: 'Hello, world!',
7    providerOptions: { openai: {} },
8  });

```

- **instructions** *string* Control the voice of your generated audio with additional instructions e.g. "Speak in a slow and steady tone". Does not work with `tts-1` or `tts-1-hd`. Optional.
- **response_format** *string* The format to audio in. Supported formats are `mp3`, `opus`, `aac`, `flac`, `wav`, and `pcm`. Defaults to `mp3`. Optional.
- **speed** *number* The speed of the generated audio. Select a value from 0.25 to 4.0. Defaults to 1.0. Optional.

Model Capabilities

Model	Instructions
<code>tts-1</code>	✓
<code>tts-1-hd</code>	✓
<code>gpt-4o-mini-tts</code>	✓

This site uses tracking technologies. You may opt in or opt out of the use of these technologies.



Resources

[Docs](#)

[Cookbook](#)

[Providers](#)

[Showcase](#)

[GitHub](#)

[Discussions](#)

More

[Playground](#)

[VX](#)

[Contact Sales](#)

About Vercel

[Next.js + Vercel](#)

[Open Source Software](#)

[GitHub](#)

[X](#)

Legal

[Privacy Policy](#)

© 2025 Vercel, Inc.



This site uses tracking technologies. You may opt in or opt out of the use of these technologies.