
response_format string

Optional Defaults to mp3

The format to audio in.

Supported formats are mp3 ,
opus , aac , flac , wav ,
and pcm .

speed number Optional

Defaults to 1

The speed of the generated
audio. Select a value from

0.25 to 4.0 . 1.0 is the
default.

stream_format string

Optional Defaults to audio

The format to stream the audio
in. Supported formats are sse
and audio . sse is not
supported for tts-1 or
tts-1-hd .

Returns

The audio file content or a
[stream of audio events](#).

Create transcription

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

...

Default

Streaming

Logprobs

Word timestamps

Example request

Example response



Options

Transcribes audio into the input language.

Request body

file file **Required**

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

model string **Required**

ID of the model to use. The options are

`gpt-4o-transcribe`,

`gpt-4o-mini-transcribe`,

and `whisper-1` (which is powered by our open source Whisper V2 model).

chunking_strategy

"auto" or object **Optional**

Controls how the audio is cut into chunks. When set to `"auto"`, the server first normalizes loudness and then uses voice activity detection (VAD) to choose boundaries. `server_vad` object can be provided to tweak VAD detection parameters manually. If unset, the audio is transcribed as a single block.

▼ Show possible types

Example request

javascript ▼



```
2 import OpenAI from "openai";
3
4 const openai = new OpenAI();
5
6 async function main() {
7   const transcription = await openai.audio
8     file: fs.createReadStream("audio.mp3")
9     model: "gpt-4o-transcribe",
10    response_format: "json",
11    include: ["logprobs"]
12  });
13
14  console.log(transcription);
15 }
16 main();
```

Response



```
1 {
2   "text": "Hey, my knee is hurting and I w
3   "logprobs": [
4     { "token": "Hey", "logprob": -1.041529
5     { "token": ",", "logprob": -9.805982e-
6     { "token": " my", "logprob": -0.002297
7     {
8       "token": " knee",
9       "logprob": -4.7159858e-5,
10      "bytes": [32, 107, 110, 101, 101]
11    },
12    { "token": " is", "logprob": -0.043909
13    {
14      "token": " hurting",
15      "logprob": -1.101146e-5
```

include[] array Optional

Additional information to include in the transcription response.

`logprobs` will return the log probabilities of the tokens in the response to understand the model's confidence in the transcription. `logprobs` only works with `response_format` set to `json` and only with the models `gpt-4o-transcribe` and `gpt-4o-mini-transcribe`.

language string Optional

The language of the input audio. Supplying the input language in [ISO-639-1](#) (e.g. `en`) format will improve accuracy and latency.

prompt string Optional

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

response_format string

Optional Defaults to `json`

The format of the output, in one of these options: `json`, `text`, `srt`, `verbose_json`, or `vtt`. For `gpt-4o-transcribe` and `gpt-4o-mini-transcribe`, the only supported format is `json`.

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 of the Speech-to-Text guide](#) for more information.

Note: Streaming is not supported for the `whisper-1` model and will be ignored.

temperature number Optional

Defaults to 0

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

timestamp_granularities[]

array Optional

Defaults to segment

The timestamp granularities to populate for this transcription.

`response_format` must be set `verbose_json` to use timestamp granularities. Either or both of these options are supported: `word`, or `segment`. Note: There is no additional

latency for segment timestamps, but generating word timestamps incurs additional latency.

Returns

The [transcription object](#), a [verbose transcription object](#) or a [stream of transcript events](#).

Create translation

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

Translates audio into English.

Request body

file file **Required**

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

model string or "whisper-1"

Required

ID of the model to use. Only

`whisper-1` (which is powered by our open source Whisper V2

Example request

javascript 

```
1 import fs from "fs";
2 import OpenAI from "openai";
3
4 const openai = new OpenAI();
5
6 async function main() {
7   const translation = await openai.audio
8     .transcribe({
9       file: fs.createReadStream("speech.
10     });
11
12   console.log(translation.text);
13 }
14 main();
```

Response 

```
1 {
2   "text": "Hello, my name is Wolfgang and I
```

model) is currently available.

```
3 }
```

prompt string Optional

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

response_format string

Optional Defaults to json

The format of the output, in one of these options: `json`, `text`, `srt`, `verbose_json`, or `vtt`.

temperature number Optional

Defaults to 0

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

Returns

The translated text.

The transcription object (JSON)

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

logprobs array

The log probabilities of the tokens in the transcription. Only returned with the models

`gpt-4o-transcribe` and `gpt-4o-mini-transcribe` if `logprobs` is added to the `include` array.

✓ Show properties

text string

The transcribed text.

usage object

Token usage statistics for the request.

✓ Show possible types

OBJECT The transcription object (JSON)



```
1  {
2    "text": "Imagine the wildest idea that y
3    "usage": {
4      "type": "tokens",
5      "input_tokens": 14,
6      "input_token_details": {
7        "text_tokens": 10,
8        "audio_tokens": 4
9      },
10     "output_tokens": 101,
11     "total_tokens": 115
12   }
13 }
```

The transcription object (Verbose JSON)

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

duration number

The duration of the input audio.

language string

The language of the input audio.

segments array

Segments of the transcribed text and their corresponding details.

✖ Show properties

text string

The transcribed text.

usage object

Usage statistics for models billed by audio input duration.

✖ Show properties

words array

Extracted words and their corresponding timestamps.

✖ Show properties

OBJECT The transcription object (Verbose JS...



```

1  {
2    "task": "transcribe",
3    "language": "english",
4    "duration": 8.470000267028809,
5    "text": "The beach was a popular spot on
6    "segments": [
7      {
8        "id": 0,
9        "seek": 0,
10       "start": 0.0,
11       "end": 3.319999933242798,
12       "text": " The beach was a popular sp
13       "tokens": [
14         50364, 440, 7534, 390, 257, 3743,
15       ],
16       "temperature": 0.0,
17       "avg_logprob": -0.2860786020755768,
18       "compression_ratio": 1.2363636493682
19       "no_speech_prob": 0.0098597947508096
20     },
21     ...
22   ],
23   "usage": {
24     "type": "duration",
25     "seconds": 9
26   }
27 }
```

Stream Event (speech.audio.delta)

Emitted for each chunk of audio data generated during speech synthesis.

audio string

A chunk of Base64-encoded audio data.

type string

The type of the event. Always

`speech.audio.delta`.

OBJECT Stream Event (speech.audio.delta)



```
1 {  
2   "type": "speech.audio.delta",  
3   "audio": "base64-encoded-audio-data"  
4 }
```

Stream Event (speech.audio.done)

Emitted when the speech synthesis is complete and all audio has been streamed.

type string

The type of the event. Always

`speech.audio.done`.

OBJECT Stream Event (speech.audio.done)



```
1 {  
2   "type": "speech.audio.done",  
3   "usage": {  
4     "input_tokens": 14,  
5     "output_tokens": 101,  
6     "total_tokens": 115  
7   }  
8 }
```

usage object

Token usage statistics for the request.

▼ Show properties

Stream Event (transcript.text.delta)

Emitted when there is an additional text delta. This is also the first event emitted when the transcription starts. Only emitted when you [create a transcription](#) with the `Stream` parameter set to `true`.

OBJECT Stream Event (transcript.text.delta)



```
1 {  
2   "type": "transcript.text.delta",  
3   "delta": " wonderful"  
4 }
```

delta string

The text delta that was additionally transcribed.

logprobs array

The log probabilities of the delta. Only included if you [create a transcription](#) with the `include[]` parameter set to `logprobs`.

✓ Show properties

type string

The type of the event. Always `transcript.text.delta`.

Stream Event (transcript.text.done)

Emitted when the transcription is complete. Contains the complete transcription text. Only emitted when you [create a transcription](#) with the `Stream` parameter set to `true`.

logprobs array

The log probabilities of the individual tokens in the transcription. Only included if you [create a transcription](#) with the `include[]` parameter set to `logprobs`.

✓ Show properties

text string

The text that was transcribed.

type string

The type of the event. Always `transcript.text.done`.

usage object

Usage statistics for models billed by token usage.

✓ Show properties

OBJECT Stream Event (`transcript.text.done`)



```
1  {
2    "type": "transcript.text.done",
3    "text": "I see skies of blue and clouds",
4    "usage": {
5      "type": "tokens",
6      "input_tokens": 14,
7      "input_token_details": {
8        "text_tokens": 10,
9        "audio_tokens": 4
10     },
11    "output_tokens": 31,
12    "total_tokens": 45
13  }
14 }
```

Images