

# Predicted Outputs

[Copy page](#)

Reduce latency for model responses where much of the response is known ahead of time.

**Predicted Outputs** enable you to speed up API responses from [Chat Completions](#) when many of the output tokens are known ahead of time. This is most common when you are regenerating a text or code file with minor modifications. You can provide your prediction using the `prediction` request parameter in [Chat Completions](#).

Predicted Outputs are available today using the latest `gpt-4o` and `gpt-4o-mini` models. Read on to learn how to use Predicted Outputs to reduce latency in your applications.

## Code refactoring example

Predicted Outputs are particularly useful for regenerating text documents and code files with small modifications. Let's say you want the [GPT-4o model](#) to refactor a piece of TypeScript code, and convert the `username` property of the `User` class to be `email` instead:

```
1 class User {  
2   firstName: string = "";  
3   lastName: string = "";  
4   username: string = "";  
5 }  
6  
7 export default User;
```

Most of the file will be unchanged, except for line 4 above. If you use the current text of the code file as your prediction, you can regenerate the entire file with lower latency. These time savings add up quickly for larger files.

Below is an example of using the `prediction` parameter in our SDKs to predict that the final output of the model will be very similar to our original code file, which we use as the prediction text.

Refactor a TypeScript class with a Predicted Output

javascript ↕ [Copy](#)

```
1 import OpenAI from "openai";  
2  
3 const code = `  
4 class User {  
5   firstName: string = "";  
6   lastName: string = "";
```

```

7   username: string = "";
8   }
9
10  export default User;
11  `.trim();
12
13  const openai = new OpenAI();
14
15  const refactorPrompt = `
16  Replace the "username" property with an "email" property. Respond only
17  with code, and with no markdown formatting.
18  `;
19
20  const completion = await openai.chat.completions.create({
21    model: "gpt-4o",
22    messages: [
23      {
24        role: "user",
25        content: refactorPrompt
26      },
27      {
28        role: "user",
29        content: code
30      }
31    ],
32    store: true,
33    prediction: {
34      type: "content",
35      content: code
36    }
37  });
38
39  // Inspect returned data
40  console.log(completion);
41  console.log(completion.choices[0].message.content);

```

In addition to the refactored code, the model response will contain data that looks something like this:

```

1  {
2    id: 'chatcmpl-xxx',
3    object: 'chat.completion',
4    created: 1730918466,
5    model: 'gpt-4o-2024-08-06',
6    choices: [ /* ...actual text response here... */ ],
7    usage: {
8      prompt_tokens: 81,
9      completion_tokens: 39,
10     total_tokens: 120,
11     prompt_tokens_details: { cached_tokens: 0, audio_tokens: 0 },
12     completion_tokens_details: {
13       reasoning_tokens: 0,
14       audio_tokens: 0,
15       accepted_prediction_tokens: 18,

```

```
      rejected_prediction_tokens: 10
    },
    system_fingerprint: 'fp_159d8341cc'
  }
}
```

Note both the `accepted_prediction_tokens` and `rejected_prediction_tokens` in the `usage` object. In this example, 18 tokens from the prediction were used to speed up the response, while 10 were rejected.

**i** Note that any rejected tokens are still billed like other completion tokens generated by the API, so Predicted Outputs can introduce higher costs for your requests.

## Streaming example

The latency gains of Predicted Outputs are even greater when you use streaming for API responses. Here is an example of the same code refactoring use case, but using streaming in the OpenAI SDKs instead.

Predicted Outputs with streaming

javascript  

```
1  import OpenAI from "openai";
2
3  const code = `
4  class User {
5    firstName: string = "";
6    lastName: string = "";
7    username: string = "";
8  }
9
10 export default User;
11 `.trim();
12
13 const openai = new OpenAI();
14
15 const refactorPrompt = `
16 Replace the "username" property with an "email" property. Respond only
17 with code, and with no markdown formatting.
18 `;
19
20 const completion = await openai.chat.completions.create({
21   model: "gpt-4o",
22   messages: [
23     {
24       role: "user",
25       content: refactorPrompt
26     },
27     {
28       role: "user",
29       content: code
30     }
31   ]
32 })
```

```

    ],
    store: true,
    prediction: {
      type: "content",
      content: code
    },
    stream: true
  });

  // Inspect returned data
  for await (const chunk of stream) {
    process.stdout.write(chunk.choices[0]?.delta?.content || "");
  }

```

## Position of predicted text in response

When providing prediction text, your prediction can appear anywhere within the generated response, and still provide latency reduction for the response. Let's say your predicted text is the simple [Hono](#) server shown below:

```

1  import { serveStatic } from "@hono/node-server/serve-static";
2  import { serve } from "@hono/node-server";
3  import { Hono } from "hono";
4
5  const app = new Hono();
6
7  app.get("/api", (c) => {
8    return c.text("Hello Hono!");
9  });
10
11 // You will need to build the client code first `pnpm run ui:build`
12 app.use(
13   "/*",
14   serveStatic({
15     rewriteRequestPath: (path) => `./dist${path}`,
16   })
17 );
18
19 const port = 3000;
20 console.log(`Server is running on port ${port}`);
21
22 serve({
23   fetch: app.fetch,
24   port,
25 });

```

You could prompt the model to regenerate the file with a prompt like:

```

1 Add a get route to this application that responds with
2 the text "hello world". Generate the entire application

```

file again with this route added, and with no other markdown formatting.

The response to the prompt might look something like this:

```
1  import { serveStatic } from "@hono/node-server/serve-static";
2  import { serve } from "@hono/node-server";
3  import { Hono } from "hono";
4
5  const app = new Hono();
6
7  app.get("/api", (c) => {
8    return c.text("Hello Hono!");
9  });
10
11 app.get("/hello", (c) => {
12   return c.text("hello world");
13 });
14
15 // You will need to build the client code first `pnpm run ui:build`
16 app.use(
17   "/*",
18   serveStatic({
19     rewriteRequestPath: (path) => `./dist${path}`,
20   })
21 );
22
23 const port = 3000;
24 console.log(`Server is running on port ${port}`);
25
26 serve({
27   fetch: app.fetch,
28   port,
29 });
```

You would still see accepted prediction tokens in the response, even though the prediction text appeared both before and after the new content added to the response:

```
1  {
2    id: 'chatcmpl-xxx',
3    object: 'chat.completion',
4    created: 1731014771,
5    model: 'gpt-4o-2024-08-06',
6    choices: [ /* completion here... */ ],
7    usage: {
8      prompt_tokens: 203,
9      completion_tokens: 159,
10     total_tokens: 362,
11     prompt_tokens_details: { cached_tokens: 0, audio_tokens: 0 },
12     completion_tokens_details: {
13       reasoning_tokens: 0,
14       audio_tokens: 0,
15       accepted_prediction_tokens: 60,
16       rejected_prediction_tokens: 0
```

```
}  
},  
system_fingerprint: 'fp_9ee9e968ea'  
}
```

This time, there were no rejected prediction tokens, because the entire content of the file we predicted was used in the final response. Nice! 🔥

## Limitations

When using Predicted Outputs, you should consider the following factors and limitations.

Predicted Outputs are only supported with the GPT-4o and GPT-4o-mini series of models.

When providing a prediction, any tokens provided that are not part of the final completion are still charged at completion token rates. See the `rejected_prediction_tokens` property of the `usage` object to see how many tokens are not used in the final response.

The following [API parameters](#) are not supported when using Predicted Outputs:

`n` : values higher than 1 are not supported

`logprobs` : not supported

`presence_penalty` : values greater than 0 are not supported

`frequency_penalty` : values greater than 0 are not supported

`audio` : Predicted Outputs are not compatible with [audio inputs and outputs](#)

`modalities` : Only `text` modalities are supported

`max_completion_tokens` : not supported

`tools` : Function calling is not currently supported with Predicted Outputs

