groq    Personal ⇕    /    Default Project ⇕          Playground    API Keys    Dashboard    **Docs**    ⚙

# Text Generation

[ Copy page | ⌄ ]

Generating text with Groq's Chat Completions API enables you to have natural, conversational interactions with Groq's large language models. It processes a series of messages and generates human-like responses that can be used for various applications including conversational agents, content generation, task automation, and generating structured data outputs like JSON for your applications.

## Chat Completions

Chat completions allow your applications to have dynamic interactions with Groq's models. You can send messages that include user inputs and system instructions, and receive responses that match the conversational context.

Chat models can handle both multi-turn discussions (conversations with multiple back-and-forth exchanges) and single-turn tasks where you need just one response.

For details about all available parameters, visit the API reference page.

### Getting Started with Groq SDK

To start using Groq's Chat Completions API, you'll need to install the Groq SDK and set up your API key.

Python    **JavaScript**

```shell
npm install --save groq-sdk
```

## Performing a Basic Chat Completion

The example below shows how to send messages to the model and receive a complete response. The messages array should include user queries and can optionally include system messages to guide the model's behavior.

```JavaScript
import Groq from "groq-sdk";

const groq = new Groq();

export async function main() {
  const completion = await getGroqChatCompletion();
  console.log(completion.choices[0]?.message?.content || "");
}

export const getGroqChatCompletion = async () => {
  return groq.chat.completions.create({
    messages: [
      // Set an optional system message. This sets the behavior of the
```

```
14        // assistant and can be used to provide specific instructions for
15        // how it should behave throughout the conversation.
16      {
17        role: "system",
18        content: "You are a helpful assistant.",
19      },
20      // Set a user message for the assistant to respond to.
21      {
22        role: "user",
23        content: "Explain the importance of fast language models",
24      },
25    ],
26    model: "openai/gpt-oss-20b",
27  });
28 };
29
30 main();
```

Streaming allows you to receive and process the model's response token by token as it's being generated. This creates a more interactive experience and can reduce perceived latency.

To stream a completion, set the parameter `stream=true` and process the response chunks as they arrive.

JavaScript

```
1  import Groq from "groq-sdk";
2
3  const groq = new Groq();
4
5  export async function main() {
6    const stream = await getGroqChatStream();
7    for await (const chunk of stream) {
8      // Print the completion returned by the LLM.
9      process.stdout.write(chunk.choices[0]?.delta?.content || "");
10   }
11 }
12
13 export async function getGroqChatStream() {
14   return groq.chat.completions.create({
15     //
16     // Required parameters
17     //
18     messages: [
19       // Set an optional system message. This sets the behavior of the
20       // assistant and can be used to provide specific instructions for
21       // how it should behave throughout the conversation.
22       {
23         role: "system",
24         content: "You are a helpful assistant.",
25       },
26       // Set a user message for the assistant to respond to.
27       {
28         role: "user",
29         content: "Explain the importance of fast language models",
30       },
31     ],
32
33     // The language model which will generate the completion.
34     model: "openai/gpt-oss-20b",
35
```

JavaScript

```
1  import Groq from "groq-sdk";
2                                                          m completions.
3  const groq = new Groq();                               become deterministic
```

```javascript
 4
 5  export async function main() {
 6    const stream = await getGroqChatStream();
 7    for await (const chunk of stream) {
 8      // Print the completion returned by the LLM.
 9      process.stdout.write(chunk.choices[0]?.delta?.content || "");
10    }
11  }
12
13  export async function getGroqChatStream() {
14    return groq.chat.completions.create({
15      //
16      // Required parameters
17      //
18      messages: [
19        // Set an optional system message. This sets the behavior of the
20        // assistant and can be used to provide specific instructions for
21        // how it should behave throughout the conversation.
22        {
23          role: "system",
24          content: "You are a helpful assistant.",
25        },
26        // Set a user message for the assistant to respond to.
27        {
28          role: "user",
29          content:
30            "Start at 1 and count to 10.  Separate each number with a comma and a space",
31        },
32      ],
33
34      // The language model which will generate the completion.
35      model: "llama-3.3-70b-versatile",
36
37      //
38      // Optional parameters
39      //
40
41      // Controls randomness: lowering results in less random completions.
42      // As the temperature approaches zero, the model will become deterministic
43      // and repetitive.
44      temperature: 0.5,
45
46      // The maximum number of tokens to generate. Requests can use up to
47      // 2048 tokens shared between prompt and completion.
48      max_completion_tokens: 1024,
49
50      // Controls diversity via nucleus sampling: 0.5 means half of all
51      // likelihood-weighted options are considered.
52      top_p: 1,
53
54      // A stop sequence is a predefined or user-specified text string that
55      // signals an AI to stop generating content, ensuring its responses
56      // remain focused and concise. Examples include punctuation marks and
57      // markers like "[end]".
58      //
59      // For this example, we will use ", 6" so that the llm stops counting at 5.
60      // If multiple stop values are needed, an array of string may be passed,
61      // stop: [", 6", ", six", ", Six"]
62      stop: ", 6",
63
64      // If set, partial message deltas will be sent.
65      stream: true,
66    });
67  }
68
69  main();
```