Documentation

| Text | ⇕ |

## Chat Completion Models

The Groq Chat Completions API processes a series of messages and generates output responses. These models can perform multi-turn discussions or tasks that require only one interaction.

For details about the parameters, visit the reference page.

### JSON mode *(beta)*

JSON mode is a beta feature that guarantees all chat completions are valid JSON.
Usage:

1. Set `"response_format": {"type": "json_object"}` in your chat completion request

2. Add a description of the desired JSON structure within the system prompt (see below for example system prompts)

Recommendations for best beta results:

- Mixtral performs best at generating JSON, followed by Gemma, then Llama

- Use pretty-printed JSON instead of compact JSON

- Keep prompts concise

Beta Limitations:

- Does not support streaming

- Does not support stop sequences

Error Code:

- Groq will return a 400 error with an error code of `json_validate_failed` if JSON generation fails.

Example system prompts:

```
You are a legal advisor who summarizes documents in JSON
```

```
You are a data analyst API capable of sentiment analysis that responds in JSON.  The JSON schema should in
{
  "sentiment_analysis": {
    "sentiment": "string (positive, negative, neutral)",
    "confidence_score": "number (0-1)"
    # Include additional fields as required
  }
}
```

### Generating Chat Completions with groq SDK

Code Overview

**Python**   JavaScript

```
pip install groq
```

## Performing a basic Chat Completion

```
1   from groq import Groq
2
3   client = Groq()
4
5   chat_completion = client.chat.completions.create(
6       #
7       # Required parameters
8       #
9       messages=[
10          # Set an optional system message. This sets the behavior of the
11          # assistant and can be used to provide specific instructions for
12          # how it should behave throughout the conversation.
13          {
14              "role": "system",
15              "content": "you are a helpful assistant."
16          },
17          # Set a user message for the assistant to respond to.
18          {
19              "role": "user",
20              "content": "Explain the importance of fast language models",
21          }
22      ],
23
24      # The language model which will generate the completion.
25      model="llama-3.3-70b-versatile",
26
27      #
28      # Optional parameters
29      #
30
31      # Controls randomness: lowering results in less random completions.
32      # As the temperature approaches zero, the model will become deterministic
33      # and repetitive.
34      temperature=0.5,
35
36      # The maximum number of tokens to generate. Requests can use up to
37      # 32,768 tokens shared between prompt and completion.
38      max_tokens=1024,
39
40      # Controls diversity via nucleus sampling: 0.5 means half of all
41      # likelihood-weighted options are considered.
42      top_p=1,
43
44      # A stop sequence is a predefined or user-specified text string that
45      # signals an AI to stop generating content, ensuring its responses
46      # remain focused and concise. Examples include punctuation and newline
47      # markers like "[end]".
48      stop=None,
49
50      # If set, partial message deltas will be sent.
51      stream=False,
52  )
```

## Streaming a Chat Completion

To stream a completion, simply set the parameter `stream=True`. Then the completion function will return an iterator of completion deltas rather than a single, full completion.

```
from groq import Groq

client = Groq()

stream = client.chat.completions.create(
    #
    # Required parameters
    #
    messages=[
        # Set an optional system message. This sets the behavior of the
        # assistant and can be used to provide specific instructions for
        # how it should behave throughout the conversation.
        {
            "role": "system",
            "content": "you are a helpful assistant."
        },
        # Set a user message for the assistant to respond to.
        {
            "role": "user",
            "content": "Explain the importance of fast language models",
        }
    ],

    # The language model which will generate the completion.
    model="llama-3.3-70b-versatile",
```

## Performing a Chat Completion with a stop sequence

```
from groq import Groq

client = Groq()

chat_completion = client.chat.completions.create(
    #
    # Optional parameters
    #
```

```
29  # Print the completion returned by the LLM.
30  print(chat_completion.choices[0].message.content)
31      # Controls randomness: lowering results in less random completions.
32      # As the temperature approaches zero, the model will become deterministic
33      # and repetitive.
34      temperature=0.5,
35
36      # The maximum number of tokens to generate. Requests can use up to
37      # 2048 tokens shared between prompt and completion.
38      max_completion_tokens=1024,
39
40      # Controls diversity via nucleus sampling: 0.5 means half of all
41      # likelihood-weighted options are considered.
42      top_p=1,
43
44      # A stop sequence is a predefined or user-specified text string that
45      # signals an AI to stop generating content, ensuring its responses
46      # remain focused and concise. Examples include punctuation marks and
47      # markers like "[end]".
48      stop=None,
49
50      # If set, partial message deltas will be sent.
51      stream=False,
52  )
```

## Performing an Async Chat Completion

Simply use the Async client to enable asyncio

```
import asyncio

from groq import AsyncGroq

async def main():
    client = AsyncGroq()

    # Print the completion returned by the LLM.
    for chunk in chat_completion.choices[0].delta.content, end="")

    # Required parameters

    # The maximum number of tokens to generate. Requests can use up to
    # 2048 tokens shared between prompt and completion.
    max_completion_tokens=1024,

    # Controls diversity via nucleus sampling: 0.5 means half of all
    # likelihood-weighted options are considered.
    top_p=1,
```

## Streaming an Async Chat Completion

```
import asyncio

from groq import AsyncGroq

async def main():
    client = AsyncGroq()

    stream = await client.chat.completions.create(
        # Required parameters

# Print the completion returned by the LLM.
print(chunk.choices[0]...)
```

## JSON Mode

```python
from typing import List, Optional
import json
from pydantic import BaseModel
from groq import Groq

groq = Groq()

# Data model for LLM to generate
class Ingredient(BaseModel):
    name: str
    quantity: str
    quantity_unit: Optional[str]

class Recipe(BaseModel):
    recipe_name: str
    ingredients: List[Ingredient]
    directions: List[str]

def get_recipe(recipe_name: str) -> Recipe:
    chat_completion = groq.chat.completions.create(
        messages=[
            {
                "role": "system",
                "content": "You are a recipe database that outputs recipes in JSON.\n"
                # Pass the json schema to the model. Pretty printing improves results.
                f"The JSON object must use the schema: {json.dumps(Recipe.model_json_schema(), indent
            },
            {
                "role": "user",
                "content": f"Fetch a recipe for {recipe_name}",
            },
        ],
        model="llama3-70b-8192",
        temperature=0,
        # Streaming is not supported in JSON mode
        stream=False,
        # Enable JSON mode by setting the response format
        response_format={"type": "json_object"},
    )
    return Recipe.model_validate_json(chat_completion.choices[0].message.content)


def print_recipe(recipe: Recipe):
    print("Recipe:", recipe.recipe_name)

    print("\nIngredients:")
    for ingredient in recipe.ingredients:
        print(
            f"- {ingredient.name}: {ingredient.quantity} {ingredient.quantity_unit or ''}"
        )
    print("\nDirections:")
    for step, direction in enumerate(recipe.directions, start=1):
        print(f"{step}. {direction}")


recipe = get_recipe("apple pie")
print_recipe(recipe)
```