

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 ir
{
  "sentiment_analysis": {
    "sentiment": "string (positive, negative, neutral)",
    "confidence_score": "number (0-1)"
    # Include additional fields as required
  }
}
```



Generating Chat Completions with groq SDK

```
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 )
```

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.

```
24 # The language model which will generate the completion.
25 model="llama3-8b-8192",
26
27 from groq import Groq
28 # Optional parameters
29 client = Groq()
30
31 stream=True # Controls randomness in generated results in less random completions.
32 # As the temperature approaches zero, the model will become deterministic
33 # and predictable.
34 # Required parameters
35 temperature=0.5,
36 messages=[
37     # Set an optional system message. This sets the behavior of the
38     # assistant and can be used to provide specific instructions for
39     # how it should behave throughout the conversation.
40     {
41         "role": "system",
42         "content": "you are a helpful assistant."
43     },
44     # Set a user message for the assistant to respond to.
45     {
46         "role": "user",
47         "content": "Explain the importance of fast language models",
48     }
49 ]
```

```

18 # A stop sequence is a predefined or user-specified text string that
19 # signals an AI to stop generating content, ensuring its responses
20 # remain focused and concise. Examples include punctuation marks and

```

Performing a Chat Completion with a stop sequence

```

21 # markers like "[end]".
22 stop=None,
23 from groq import Groq
24 # If a page token is returned by the LLM,
25 client = Groq(api_key="gsk-8b-8192",
26 )
27 chat_completion = client.chat.completions.create(
28 # Parameters returned by the LLM.
29 print(required_parameters[0].message.content)
30 #
31 # Controls randomness: lowering results in less random completions.
32 # As the temperature approaches zero, the model will become deterministic
33 # and predictable and can be used to provide specific instructions for
34 # how the model should behave throughout the conversation.
35 {
36 # The maximum number of tokens to generate. Requests can use up to
37 # 2048 tokens shared between prompt and completion.
38 max_tokens=1024,
39 # Set a user message for the assistant to respond to.
40 # Controls diversity via nucleus sampling: 0.5 means half of all
41 # likelihood-weighted options are considered.

```

Performing an Async Chat Completion

Simply use the Async client to enable asyncio

```

22 # A stop sequence is a predefined or user-specified text string that
23 # signals an AI to stop generating content, ensuring its responses
24 # remain focused and concise. Examples include punctuation marks and
25 # markers like "[end]".
26 stop=None,
27 from groq import AsyncGroq
28 # Optional parameters: message deltas will be sent.
29 stream=True,
30 async def main():
31 # Controls randomness: lowering results in less random completions.
32 # As the temperature approaches zero, the model will become deterministic
33 # and predictable and can be used to provide specific instructions for
34 # how the model should behave throughout the conversation.
35 {
36 # The maximum number of tokens to generate. Requests can use up to
37 # 2048 tokens shared between prompt and completion.
38 max_tokens=1024,
39 # Set a user message for the assistant to respond to.
40 # Controls diversity via nucleus sampling: 0.5 means half of all
41 # likelihood-weighted options are considered.
42 top_p=1, "role": "system",
43 "content": "you are a helpful assistant."

```

Streaming an Async Chat Completion

```

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 # For this example, we will use the first stop sequence found.
49 from groq import AsyncGroq
50 # If multiple stop values are needed, an array of string may be passed,
51 stop=[" ", "6", " ", "six", " ", "Six"]
52 # Set a user message for the assistant to respond to.
53 async def main():
54 # Create a language model which will generate the completion.
55 client = AsyncGroq()

```

```

28 # If set, partial message deltas will be sent.
29 stream=False client.chat.completions.create(
30 )
31 # Optional parameters
32 # Print the completion returned by the LLM.
33 print(chat_completion.choices[0].message.content)
34 # Controls randomness: lowering results in less random completions. As the temperature approaches zero, the model will generate deterministic outputs for
35 # deterministic behavior throughout the conversation.
36 temperature=0.5,
37 "role": "system",
38 # The maximum number of tokens to generate. Requests can use up to
39 # 2048 tokens shared between prompt and completion.
40 max_tokens=1024,
41 # Controls diversity via nucleus sampling: 0.5 means half of all
42 # likelihood-weighted options are considered.
43 top_p=1,
44 from pydantic import BaseModel
45 # A stop sequence is a predefined or user-specified text string that
46 # signals the model to stop generating content. Examples include punctuation marks and
47 # markers like "[end]".
48 stop=None,
49 # Data model for LLM to generate
50 class Ingredient(BaseModel):
51     name:str,
52     quantity:str,
53     unit:Optional[str],
54     # As the temperature approaches zero, the model will become
55     # more deterministic and repeatable.
56     print(chat_completion.choices[0].message.content)
57 class Recipe(BaseModel):
58     recipe_name:str,
59     ingredients:list[Ingredient],
60     directions:str,
61     # Controls diversity via nucleus sampling: 0.5 means half of all
62     # likelihood-weighted options are considered.
63     chat_completion = groq.chat.completions.create(
64         messages=[
65             # A stop sequence is a predefined or user-specified text string that
66             # signals the model to stop generating content, ensuring its responses
67             # remain focused and concise. Examples include punctuation marks and
68             # markers like "[end]".
69             stop=None,
70             # The JSON object must use the schema: {json.dumps(Recipe.model_json_schema(), indent
71             },
72             # If set, partial message deltas will be sent.
73             stream=True,
74             "role": "user",
75             "content": f"Fetch a recipe for {recipe_name}",
76         ],
77     ),
78     # Print the incremental deltas returned by the LLM.
79     async for chunk in groq.chat.completions.create(
80         model="llama3-8b-8192",
81         temperature=0,
82         print(chunk.choices[0].delta.content, end="")
83         # Streaming is not supported in JSON mode
84     ):
85         asyncio.StreamWriter.write,
86         # Enable JSON mode by setting the response format
87         response_format={"type": "json_object"},
88     )

```

```
44     return Recipe.model_validate_json(chat_completion.choices[0].message.content)
45
46
47 def print_recipe(recipe: Recipe):
48     print("Recipe:", recipe.recipe_name)
49
50     print("\nIngredients:")
51     for ingredient in recipe.ingredients:
52         print(
53             f"- {ingredient.name}: {ingredient.quantity} {ingredient.quantity_unit or ''}"
54         )
55     print("\nDirections:")
56     for step, direction in enumerate(recipe.directions, start=1):
57         print(f"{step}. {direction}")
58
59
60 recipe = get_recipe("apple pie")
61 print_recipe(recipe)
```