

pydantic_ai.models.groq

Setup

For details on how to set up authentication with this model, see [model configuration for Groq](#).

GroqModelName module-attribute

```
GroqModelName = Literal[
    "llama-3.1-70b-versatile",
    "llama3-groq-70b-8192-tool-use-preview",
    "llama3-groq-8b-8192-tool-use-preview",
    "llama-3.1-70b-specdec",
    "llama-3.1-8b-instant",
    "llama-3.2-1b-preview",
    "llama-3.2-3b-preview",
    "llama-3.2-11b-vision-preview",
    "llama-3.2-90b-vision-preview",
    "llama3-70b-8192",
    "llama3-8b-8192",
    "mixtral-8x7b-32768",
    "gemma2-9b-it",
    "gemma-7b-it",
]
```

Named Groq models.

See [the Groq docs](#) for a full list.

GroqModel dataclass

Bases: `Model`

A model that uses the Groq API.

Internally, this uses the [Groq Python client](#) to interact with the API.

Apart from `__init__`, all methods are private or match those of the base class.

```
68 @dataclass(init=False)
69 class GroqModel(Model):
70     """A model that uses the Groq API.
71
72     Internally, this uses the [Groq Python client](https://github.com/groq/groq-python) to interact with the API.
73
74     Apart from `__init__`, all methods are private or match those of the base class.
75     """
76
77     model_name: GroqModelName
78     client: AsyncGroq = field(repr=False)
79
80     def __init__(
81         self,
82         model_name: GroqModelName,
83         *,
84         api_key: str | None = None,
85         groq_client: AsyncGroq | None = None,
86         http_client: AsyncHTTPClient | None = None,
87     ):
88         """Initialize a Groq model.
89
90         Args:
91             model_name: The name of the Groq model to use. List of model names available
92                 [here](https://console.groq.com/docs/models).
93             api_key: The API key to use for authentication, if not provided, the `GROQ_API_KEY` environment variable
94                 will be used if available.
95             groq_client: An existing
96                 [AsyncGroq](https://github.com/groq/groq-python?tab=readme-ov-file#async-usage)
97                 client to use, if provided, `api_key` and `http_client` must be `None`.
98             http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
99
100         """
101         self.model_name = model_name
102         if groq_client is not None:
103             assert http_client is None, 'Cannot provide both `groq_client` and `http_client`'
104             assert api_key is None, 'Cannot provide both `groq_client` and `api_key`'
105             self.client = groq_client
106         elif http_client is not None:
107             self.client = AsyncGroq(api_key=api_key, http_client=http_client)
108         else:
109             self.client = AsyncGroq(api_key=api_key, http_client=cached_async_http_client())
110
111     async def agent_model(
112         self,
113         *,
114         function_tools: list[ToolDefinition],
115         allow_text_result: bool,
116         result_tools: list[ToolDefinition],
117     ) -> AgentModel:
118         check_allow_model_requests()
119         tools = [self._map_tool_definition(r) for r in function_tools]
120         if result_tools:
121             tools += [self._map_tool_definition(r) for r in result_tools]
122         return GroqAgentModel(
123             self.client,
124             self.model_name,
125             allow_text_result,
126             tools,
127         )
128
129     def name(self) -> str:
130         return f'groq:{self.model_name}'
131
132     @staticmethod
133     def _map_tool_definition(f: ToolDefinition) -> chat.ChatCompletionToolParam:
134         return {
135             'type': 'function',
136             'function': {
137                 'name': f.name,
138                 'description': f.description,
139                 'parameters': f.parameters_json_schema,
140             },
141         }
```

__init__

```
__init__(
    model_name: GroqModelName,
    *,
    api_key: str | None = None,
    groq_client: AsyncGroq | None = None,
    http_client: AsyncClient | None = None
)
```

Initialize a Groq model.

Parameters:

Name	Type	Description	Default
model_name	GroqModelName	The name of the Groq model to use. List of model names available here .	required
api_key	str None	The API key to use for authentication, if not provided, the <code>GROQ_API_KEY</code> environment variable will be used if available.	None
groq_client	AsyncGroq None	An existing <code>AsyncGroq</code> client to use, if provided, <code>api_key</code> and <code>http_client</code> must be <code>None</code> .	None
http_client	AsyncClient None	An existing <code>httpx.AsyncClient</code> to use for making HTTP requests.	None

```

80     def __init__(
81         self,
82         model_name: GroqModelName,
83         *,
84         api_key: str | None = None,
85         groq_client: AsyncGroq | None = None,
86         http_client: AsyncHTTPClient | None = None,
87     ):
88         """Initialize a Groq model.
89
90     Args:
91         model_name: The name of the Groq model to use. List of model names available
92             [here](https://console.groq.com/docs/models).
93         api_key: The API key to use for authentication, if not provided, the `GROQ_API_KEY` environment variable
94             will be used if available.
95         groq_client: An existing
96             [`AsyncGroq`](https://github.com/groq/groq-python?tab=readme-ov-file#async-usage)
97             client to use, if provided, `api_key` and `http_client` must be `None`.
98         http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
99     """
100     self.model_name = model_name
101     if groq_client is not None:
102         assert http_client is None, 'Cannot provide both `groq_client` and `http_client`'
103         assert api_key is None, 'Cannot provide both `groq_client` and `api_key`'
104         self.client = groq_client
105     elif http_client is not None:
106         self.client = AsyncGroq(api_key=api_key, http_client=http_client)
107     else:
108         self.client = AsyncGroq(api_key=api_key, http_client=cached_async_http_client())

```

GroqAgentModel `dataclass`

Bases: `AgentModel`

Implementation of `AgentModel` for Groq models.

```

143 @dataclass
144 class GroqAgentModel(AgentModel):
145     """Implementation of 'AgentModel' for Groq models."""
146
147     client: AsyncGroq
148     model_name: str
149     allow_text_result: bool
150     tools: list[chat.ChatCompletionToolParam]
151
152     async def request(self, messages: list[Message]) -> tuple[ModelAnyResponse, result.Cost]:
153         response = await self._completions_create(messages, False)
154         return self._process_response(response), _map_cost(response)
155
156     @asynctxcontextmanager
157     async def request_stream(self, messages: list[Message]) -> AsyncIterator[EitherStreamedResponse]:
158         response = await self._completions_create(messages, True)
159         async with response:
160             yield await self._process_streamed_response(response)
161
162     @overload
163     async def _completions_create(
164         self, messages: list[Message], stream: Literal[True]
165     ) -> AsyncStream[ChatCompletionChunk]:
166         pass
167
168     @overload
169     async def _completions_create(self, messages: list[Message], stream: Literal[False]) -> chat.ChatCompletion:
170         pass
171
172     async def _completions_create(
173         self, messages: list[Message], stream: bool
174     ) -> chat.ChatCompletion | AsyncStream[ChatCompletionChunk]:
175         # standalone function to make it easier to override
176         if not self.tools:
177             tool_choice: Literal['none', 'required', 'auto'] | None = None
178         elif not self.allow_text_result:
179             tool_choice = 'required'
180         else:
181             tool_choice = 'auto'
182
183         groq_messages = [self._map_message(m) for m in messages]
184         return await self.client.chat.completions.create(
185             model=self.model_name,
186             messages=groq_messages,
187             temperature=0.0,
188             n=1,
189             parallel_tool_calls=True if self.tools else NOT_GIVEN,
190             tools=self.tools or NOT_GIVEN,
191             tool_choice=tool_choice or NOT_GIVEN,
192             stream=stream,
193         )
194
195     @staticmethod
196     def _process_response(response: chat.ChatCompletion) -> ModelAnyResponse:
197         """Process a non-streamed response, and prepare a message to return."""
198         timestamp = datetime.fromtimestamp(response.created, tz=tztimezon.utc)
199         choice = response.choices[0]
200         if choice.message.tool_calls is not None:
201             return ModelStructuredResponse(
202                 [ToolCall.from_json(c.function.name, c.function.arguments, c.id) for c in choice.message.tool_calls],
203                 timestamp=timestamp,
204             )
205         else:
206             assert choice.message.content is not None, choice
207             return ModelTextResponse(choice.message.content, timestamp=timestamp)
208
209     @staticmethod
210     async def _process_streamed_response(response: AsyncStream[ChatCompletionChunk]) -> EitherStreamedResponse:
211         """Process a streamed response, and prepare a streaming response to return."""
212         timestamp: datetime | None = None
213         start_cost = Cost()
214         # the first chunk may contain enough information so we iterate until we get either 'tool_calls' or 'content'
215         while True:
216             try:
217                 chunk = await response.__anext__()
218             except StopAsyncIteration as e:
219                 raise UnexpectedModelBehavior('Streamed response ended without content or tool calls') from e
220             timestamp = timestamp or datetime.fromtimestamp(chunk.created, tz=tztimezon.utc)
221             start_cost += _map_cost(chunk)
222
223             if chunk.choices:
224                 delta = chunk.choices[0].delta
225
226                 if delta.content is not None:
227                     return GroqStreamTextResponse(delta.content, response, timestamp, start_cost)
228                 elif delta.tool_calls is not None:
229                     return GroqStreamStructuredResponse(
230                         response,
231                         {c.index: c for c in delta.tool_calls},
232                         timestamp,
233                         start_cost,
234                     )
235
236     @staticmethod
237     def _map_message(message: Message) -> chat.ChatCompletionMessageParam:
238         """Just maps a 'pydantic_ai.Message' to a 'groq.types.ChatCompletionMessageParam'."""
239         if message.role == 'system':
240             # SystemPrompt ->
241             return chat.ChatCompletionSystemMessageParam(role='system', content=message.content)
242         elif message.role == 'user':
243             # UserPrompt ->
244             return chat.ChatCompletionUserMessageParam(role='user', content=message.content)
245         elif message.role == 'tool-return':
246             # ToolReturn ->
247             return chat.ChatCompletionToolMessageParam(
248                 role='tool',
249                 tool_call_id=_guard_tool_call_id(message),
250                 content=message.model_response_str(),
251             )
252         elif message.role == 'retry-prompt':
253             # RetryPrompt ->
254             if message.tool_name is None:
255                 return chat.ChatCompletionUserMessageParam(role='user', content=message.model_response())
256             else:
257                 return chat.ChatCompletionToolMessageParam(
258                     role='tool',
259                     tool_call_id=_guard_tool_call_id(message),
260                     content=message.model_response(),
261                 )
262         elif message.role == 'model-text-response':
263             # ModelTextResponse ->
264             return chat.ChatCompletionAssistantMessageParam(role='assistant', content=message.content)
265         elif message.role == 'model-structured-response':
266             assert (
267                 message.role == 'model-structured-response'
268             ), f'Expected role to be "llm-tool-calls", got {message.role}'
269             # ModelStructuredResponse ->

```

```

270     return chat.ChatCompletionAssistantMessageParam(
271         role='assistant',
272         tool_calls=[_map_tool_call(t) for t in message.calls],
273     )
274     else:
275         assert_never(message)

```

GroqStreamTextResponse `dataclass`

Bases: `StreamTextResponse`

Implementation of `StreamTextResponse` for Groq models.

Source code in `pydantic_ai_slim/pydantic_ai/models/groq.py`

```

278 @dataclass
279 class GroqStreamTextResponse(StreamTextResponse):
280     """Implementation of 'StreamTextResponse' for Groq models."""
281
282     _first: str | None
283     _response: AsyncStream[ChatCompletionChunk]
284     _timestamp: datetime
285     _cost: result.Cost
286     _buffer: list[str] = field(default_factory=list, init=False)
287
288     async def __anext__(self) -> None:
289         if self._first is not None:
290             self._buffer.append(self._first)
291             self._first = None
292             return None
293
294         chunk = await self._response.__anext__()
295         self._cost = _map_cost(chunk)
296
297         try:
298             choice = chunk.choices[0]
299         except IndexError:
300             raise StopAsyncIteration()
301
302         # we don't raise StopAsyncIteration on the last chunk because usage comes after this
303         if choice.finish_reason is None:
304             assert choice.delta.content is not None, f'Expected delta with content, invalid chunk: {chunk!r}'
305         if choice.delta.content is not None:
306             self._buffer.append(choice.delta.content)
307
308     def get(self, *, final: bool = False) -> Iterable[str]:
309         yield from self._buffer
310         self._buffer.clear()
311
312     def cost(self) -> Cost:
313         return self._cost
314
315     def timestamp(self) -> datetime:
316         return self._timestamp

```

GroqStreamStructuredResponse `dataclass`

Bases: `StreamStructuredResponse`

Implementation of `StreamStructuredResponse` for Groq models.

Source code in `pydantic_ai_slim/pydantic_ai/models/groq.py`

```

319 @dataclass
320 class GroqStreamStructuredResponse(StreamStructuredResponse):
321     """Implementation of 'StreamStructuredResponse' for Groq models."""
322
323     _response: AsyncStream[ChatCompletionChunk]
324     _delta_tool_calls: dict[int, ChoiceDeltaToolCall]
325     _timestamp: datetime
326     _cost: result.Cost
327
328     async def __anext__(self) -> None:
329         chunk = await self._response.__anext__()
330         self._cost = _map_cost(chunk)
331
332         try:
333             choice = chunk.choices[0]
334         except IndexError:
335             raise StopAsyncIteration()
336
337         if choice.finish_reason is not None:
338             raise StopAsyncIteration()
339
340         assert choice.delta.content is None, f'Expected tool calls, got content instead, invalid chunk: {chunk!r}'
341
342         for new in choice.delta.tool_calls or []:
343             if current := self._delta_tool_calls.get(new.index):
344                 if current.function is None:
345                     current.function = new.function
346                 elif new.function is not None:
347                     current.function.name = _utils.add_optional(current.function.name, new.function.name)
348                     current.function.arguments = _utils.add_optional(current.function.arguments, new.function.arguments)
349             else:
350                 self._delta_tool_calls[new.index] = new
351
352     def get(self, *, final: bool = False) -> ModelStructuredResponse:
353         calls: list[ToolCall] = []
354         for c in self._delta_tool_calls.values():
355             if f := c.function:
356                 if f.name is not None and f.arguments is not None:
357                     calls.append(ToolCall.from_json(f.name, f.arguments, c.id))
358
359         return ModelStructuredResponse(calls, timestamp=self._timestamp)
360
361     def cost(self) -> Cost:
362         return self._cost
363
364     def timestamp(self) -> datetime:
365         return self._timestamp

```