# pydantic_ai.models.openai

## Setup

For details on how to set up authentication with this model, see model configuration for OpenAI.

### OpenAIModelName `module-attribute`

```
OpenAIModelName = Union[ChatModel, str]
```

Using this more broad type for the model name instead of the ChatModel definition allows this model to be used more easily with other model types (ie, Ollama)

### OpenAIModel `dataclass`

Bases: `Model`

A model that uses the OpenAI API.

Internally, this uses the OpenAI Python client to interact with the API.

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

**Source code in** `pydantic_ai_slim/pydantic_ai/models/openai.py`

```python
53   @dataclass(init=False)
54   class OpenAIModel(Model):
55       """A model that uses the OpenAI API.
56
57       Internally, this uses the [OpenAI Python client](https://github.com/openai/openai-python) to interact with the API.
58
59       Apart from `__init__`, all methods are private or match those of the base class.
60       """
61
62       model_name: OpenAIModelName
63       client: AsyncOpenAI = field(repr=False)
64
65       def __init__(
66           self,
67           model_name: OpenAIModelName,
68           *,
69           api_key: str | None = None,
70           openai_client: AsyncOpenAI | None = None,
71           http_client: AsyncHTTPClient | None = None,
72       ):
73           """Initialize an OpenAI model.
74
75           Args:
76               model_name: The name of the OpenAI model to use. List of model names available
77                   [here](https://github.com/openai/openai-python/blob/v1.54.3/src/openai/types/chat_model.py#L7)
78                   (Unfortunately, despite being ask to do so, OpenAI do not provide `.inv` files for their API).
79               api_key: The API key to use for authentication, if not provided, the `OPENAI_API_KEY` environment variable
80                   will be used if available.
81               openai_client: An existing
82                   [`AsyncOpenAI`](https://github.com/openai/openai-python?tab=readme-ov-file#async-usage)
83                   client to use, if provided, `api_key` and `http_client` must be `None`.
84               http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
85           """
86           self.model_name: OpenAIModelName = model_name
87           if openai_client is not None:
88               assert http_client is None, 'Cannot provide both `openai_client` and `http_client`'
89               assert api_key is None, 'Cannot provide both `openai_client` and `api_key`'
90               self.client = openai_client
91           elif http_client is not None:
92               self.client = AsyncOpenAI(api_key=api_key, http_client=http_client)
93           else:
94               self.client = AsyncOpenAI(api_key=api_key, http_client=cached_async_http_client())
95
96       async def agent_model(
97           self,
98           *,
99           function_tools: list[ToolDefinition],
100          allow_text_result: bool,
101          result_tools: list[ToolDefinition],
102      ) -> AgentModel:
103          check_allow_model_requests()
104          tools = [self._map_tool_definition(r) for r in function_tools]
105          if result_tools:
106              tools += [self._map_tool_definition(r) for r in result_tools]
107          return OpenAIAgentModel(
108              self.client,
109              self.model_name,
110              allow_text_result,
111              tools,
112          )
113
114      def name(self) -> str:
115          return f'openai:{self.model_name}'
116
117      @staticmethod
118      def _map_tool_definition(f: ToolDefinition) -> chat.ChatCompletionToolParam:
119          return {
120              'type': 'function',
121              'function': {
122                  'name': f.name,
123                  'description': f.description,
124                  'parameters': f.parameters_json_schema,
125              },
126          }
```

**__init__**

```python
__init__(
    model_name: OpenAIModelName,
    *,
    api_key: str | None = None,
    openai_client: AsyncOpenAI | None = None,
    http_client: AsyncClient | None = None
)
```

Initialize an OpenAI model.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| `model_name` | `OpenAIModelName` | The name of the OpenAI model to use. List of model names available here (Unfortunately, despite being ask to do so, OpenAI do not provide `.inv` files for their API). | *required* |
| `api_key` | `str | None` | The API key to use for authentication, if not provided, the `OPENAI_API_KEY` environment variable will be used if available. | `None` |
| `openai_client` | `AsyncOpenAI | None` | An existing `AsyncOpenAI` client to use, if provided, `api_key` and `http_client` must be `None`. | `None` |
| `http_client` | `AsyncClient | None` | An existing `httpx.AsyncClient` to use for making HTTP requests. | `None` |

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

```python
65    def __init__(
66        self,
67        model_name: OpenAIModelName,
68        *,
69        api_key: str | None = None,
70        openai_client: AsyncOpenAI | None = None,
71        http_client: AsyncHTTPClient | None = None,
72    ):
73        """Initialize an OpenAI model.
74
75        Args:
76            model_name: The name of the OpenAI model to use. List of model names available
77                [here](https://github.com/openai/openai-python/blob/v1.54.3/src/openai/types/chat_model.py#L7)
78                (Unfortunately, despite being ask to do so, OpenAI do not provide `.inv` files for their API).
79            api_key: The API key to use for authentication, if not provided, the `OPENAI_API_KEY` environment variable
80                will be used if available.
81            openai_client: An existing
82                [`AsyncOpenAI`](https://github.com/openai/openai-python?tab=readme-ov-file#async-usage)
83                client to use, if provided, `api_key` and `http_client` must be `None`.
84            http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
85        """
86        self.model_name: OpenAIModelName = model_name
87        if openai_client is not None:
88            assert http_client is None, 'Cannot provide both `openai_client` and `http_client`'
89            assert api_key is None, 'Cannot provide both `openai_client` and `api_key`'
90            self.client = openai_client
91        elif http_client is not None:
92            self.client = AsyncOpenAI(api_key=api_key, http_client=http_client)
93        else:
94            self.client = AsyncOpenAI(api_key=api_key, http_client=cached_async_http_client())
```

## OpenAIAgentModel `dataclass`

Bases: `AgentModel`

Implementation of `AgentModel` for OpenAI models.

```python
129     @dataclass
130     class OpenAIAgentModel(AgentModel):
131         """Implementation of `AgentModel` for OpenAI models."""
132
133         client: AsyncOpenAI
134         model_name: OpenAIModelName
135         allow_text_result: bool
136         tools: list[chat.ChatCompletionToolParam]
137
138         async def request(self, messages: list[Message]) -> tuple[ModelAnyResponse, result.Cost]:
139             response = await self._completions_create(messages, False)
140             return self._process_response(response), _map_cost(response)
141
142         @asynccontextmanager
143         async def request_stream(self, messages: list[Message]) -> AsyncIterator[EitherStreamedResponse]:
144             response = await self._completions_create(messages, True)
145             async with response:
146                 yield await self._process_streamed_response(response)
147
148         @overload
149         async def _completions_create(
150             self, messages: list[Message], stream: Literal[True]
151         ) -> AsyncStream[ChatCompletionChunk]:
152             pass
153
154         @overload
155         async def _completions_create(self, messages: list[Message], stream: Literal[False]) -> chat.ChatCompletion:
156             pass
157
158         async def _completions_create(
159             self, messages: list[Message], stream: bool
160         ) -> chat.ChatCompletion | AsyncStream[ChatCompletionChunk]:
161             # standalone function to make it easier to override
162             if not self.tools:
163                 tool_choice: Literal['none', 'required', 'auto'] | None = None
164             elif not self.allow_text_result:
165                 tool_choice = 'required'
166             else:
167                 tool_choice = 'auto'
168
169             openai_messages = [self._map_message(m) for m in messages]
170             return await self.client.chat.completions.create(
171                 model=self.model_name,
172                 messages=openai_messages,
173                 n=1,
174                 parallel_tool_calls=True if self.tools else NOT_GIVEN,
175                 tools=self.tools or NOT_GIVEN,
176                 tool_choice=tool_choice or NOT_GIVEN,
177                 stream=stream,
178                 stream_options={'include_usage': True} if stream else NOT_GIVEN,
179             )
180
181         @staticmethod
182         def _process_response(response: chat.ChatCompletion) -> ModelAnyResponse:
183             """Process a non-streamed response, and prepare a message to return."""
184             timestamp = datetime.fromtimestamp(response.created, tz=timezone.utc)
185             choice = response.choices[0]
186             if choice.message.tool_calls is not None:
187                 return ModelStructuredResponse(
188                     [ToolCall.from_json(c.function.name, c.function.arguments, c.id) for c in choice.message.tool_calls],
189                     timestamp=timestamp,
190                 )
191             else:
192                 assert choice.message.content is not None, choice
193                 return ModelTextResponse(choice.message.content, timestamp=timestamp)
194
195         @staticmethod
196         async def _process_streamed_response(response: AsyncStream[ChatCompletionChunk]) -> EitherStreamedResponse:
197             """Process a streamed response, and prepare a streaming response to return."""
198             timestamp: datetime | None = None
199             start_cost = Cost()
200             # the first chunk may contain enough information so we iterate until we get either `tool_calls` or `content`
201             while True:
202                 try:
203                     chunk = await response.__anext__()
204                 except StopAsyncIteration as e:
205                     raise UnexpectedModelBehavior('Streamed response ended without content or tool calls') from e
206
207                 timestamp = timestamp or datetime.fromtimestamp(chunk.created, tz=timezone.utc)
208                 start_cost += _map_cost(chunk)
209
210                 if chunk.choices:
211                     delta = chunk.choices[0].delta
212
213                     if delta.content is not None:
214                         return OpenAIStreamTextResponse(delta.content, response, timestamp, start_cost)
215                     elif delta.tool_calls is not None:
216                         return OpenAIStreamStructuredResponse(
217                             response,
218                             {c.index: c for c in delta.tool_calls},
219                             timestamp,
220                             start_cost,
221                         )
222                     # else continue until we get either delta.content or delta.tool_calls
223
224         @staticmethod
225         def _map_message(message: Message) -> chat.ChatCompletionMessageParam:
226             """Just maps a `pydantic_ai.Message` to a `openai.types.ChatCompletionMessageParam`."""
227             if message.role == 'system':
228                 # SystemPrompt ->
229                 return chat.ChatCompletionSystemMessageParam(role='system', content=message.content)
230             elif message.role == 'user':
231                 # UserPrompt ->
232                 return chat.ChatCompletionUserMessageParam(role='user', content=message.content)
233             elif message.role == 'tool-return':
234                 # ToolReturn ->
235                 return chat.ChatCompletionToolMessageParam(
236                     role='tool',
237                     tool_call_id=_guard_tool_call_id(message),
238                     content=message.model_response_str(),
239                 )
240             elif message.role == 'retry-prompt':
241                 # RetryPrompt ->
242                 if message.tool_name is None:
243                     return chat.ChatCompletionUserMessageParam(role='user', content=message.model_response())
244                 else:
245                     return chat.ChatCompletionToolMessageParam(
246                         role='tool',
247                         tool_call_id=_guard_tool_call_id(message),
248                         content=message.model_response(),
249                     )
250             elif message.role == 'model-text-response':
251                 # ModelTextResponse ->
252                 return chat.ChatCompletionAssistantMessageParam(role='assistant', content=message.content)
253             elif message.role == 'model-structured-response':
254                 assert (
255                     message.role == 'model-structured-response'
```

```python
256                    ), f'Expected role to be "llm-tool-calls", got {message.role}'
257                    # ModelStructuredResponse ->
258                    return chat.ChatCompletionAssistantMessageParam(
259                        role='assistant',
260                        tool_calls=[_map_tool_call(t) for t in message.calls],
261                    )
262                else:
263                    assert_never(message)
```

### OpenAIStreamTextResponse `dataclass`

Bases: `StreamTextResponse`

Implementation of `StreamTextResponse` for OpenAI models.

```python
266  @dataclass
267  class OpenAIStreamTextResponse(StreamTextResponse):
268      """Implementation of `StreamTextResponse` for OpenAI models."""
269
270      _first: str | None
271      _response: AsyncStream[ChatCompletionChunk]
272      _timestamp: datetime
273      _cost: result.Cost
274      _buffer: list[str] = field(default_factory=list, init=False)
275
276      async def __anext__(self) -> None:
277          if self._first is not None:
278              self._buffer.append(self._first)
279              self._first = None
280              return None
281
282          chunk = await self._response.__anext__()
283          self._cost += _map_cost(chunk)
284          try:
285              choice = chunk.choices[0]
286          except IndexError:
287              raise StopAsyncIteration()
288
289          # we don't raise StopAsyncIteration on the last chunk because usage comes after this
290          if choice.finish_reason is None:
291              assert choice.delta.content is not None, f'Expected delta with content, invalid chunk: {chunk!r}'
292          if choice.delta.content is not None:
293              self._buffer.append(choice.delta.content)
294
295      def get(self, *, final: bool = False) -> Iterable[str]:
296          yield from self._buffer
297          self._buffer.clear()
298
299      def cost(self) -> Cost:
300          return self._cost
301
302      def timestamp(self) -> datetime:
303          return self._timestamp
```

### OpenAIStreamStructuredResponse `dataclass`

Bases: `StreamStructuredResponse`

Implementation of `StreamStructuredResponse` for OpenAI models.

```python
306  @dataclass
307  class OpenAIStreamStructuredResponse(StreamStructuredResponse):
308      """Implementation of `StreamStructuredResponse` for OpenAI models."""
309
310      _response: AsyncStream[ChatCompletionChunk]
311      _delta_tool_calls: dict[int, ChoiceDeltaToolCall]
312      _timestamp: datetime
313      _cost: result.Cost
314
315      async def __anext__(self) -> None:
316          chunk = await self._response.__anext__()
317          self._cost += _map_cost(chunk)
318          try:
319              choice = chunk.choices[0]
320          except IndexError:
321              raise StopAsyncIteration()
322
323          if choice.finish_reason is not None:
324              raise StopAsyncIteration()
325
326          assert choice.delta.content is None, f'Expected tool calls, got content instead, invalid chunk: {chunk!r}'
327
328          for new in choice.delta.tool_calls or []:
329              if current := self._delta_tool_calls.get(new.index):
330                  if current.function is None:
331                      current.function = new.function
332                  elif new.function is not None:
333                      current.function.name = _utils.add_optional(current.function.name, new.function.name)
334                      current.function.arguments = _utils.add_optional(current.function.arguments, new.function.arguments)
335              else:
336                  self._delta_tool_calls[new.index] = new
337
338      def get(self, *, final: bool = False) -> ModelStructuredResponse:
339          calls: list[ToolCall] = []
340          for c in self._delta_tool_calls.values():
341              if f := c.function:
342                  if f.name is not None and f.arguments is not None:
343                      calls.append(ToolCall.from_json(f.name, f.arguments, c.id))
344
345          return ModelStructuredResponse(calls, timestamp=self._timestamp)
346
347      def cost(self) -> Cost:
348          return self._cost
349
350      def timestamp(self) -> datetime:
351          return self._timestamp
```