

pydantic_ai.models.ollama

Setup

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

Example local usage

With `ollama` installed, you can run the server with the model you want to use:

terminal-run-ollama

```
ollama run llama3.2
```

(this will pull the `llama3.2` model if you don't already have it downloaded)

Then run your code, here's a minimal example:

ollama_example.py

```

from pydantic import BaseModel
from pydantic_ai import Agent

class CityLocation(BaseModel):
    city: str
    country: str

agent = Agent('ollama:llama3.2', result_type=CityLocation)

result = agent.run_sync('Where the olympics held in 2012?')
print(result.data)
#> city='London' country='United Kingdom'
print(result.cost())
#> Cost(request_tokens=56, response_tokens=8, total_tokens=64, details=None)

```

Example using a remote server

ollama_example_with_remote_server.py

```

from pydantic import BaseModel
from pydantic_ai import Agent
from pydantic_ai.models.ollama import OllamaModel

ollama_model = OllamaModel(
    model_name='qwen2.5-coder:7b', ❶
    base_url='http://192.168.1.74:11434/v1', ❷
)

class CityLocation(BaseModel):
    city: str
    country: str

agent = Agent(model=ollama_model, result_type=CityLocation)

result = agent.run_sync('Where the olympics held in 2012?')
print(result.data)
#> city='London' country='United Kingdom'
print(result.cost())
#> Cost(request_tokens=56, response_tokens=8, total_tokens=64, details=None)

```

❶ The name of the model running on the remote server

❷ The url of the remote server

See `OllamaModel` for more information

CommonOllamaModelNames module-attribute

```

CommonOllamaModelNames = Literal[
    "codellama",
    "gemma",
    "gemma2",
    "llama3",
    "llama3.1",
    "llama3.2",
    "llama3.2-vision",
    "llama3.3",
    "mistral",
    "mistral-nemo",
    "mixtral",
    "phi3",
    "qwq",
    "qwen",
    "qwen2",
    "qwen2.5",
    "starcoder2",
]

```

This contains just the most common ollama models.

For a full list see ollama.com/library.

OllamaModelName module-attribute

```

OllamaModelName = Union[CommonOllamaModelNames, str]

```

Possible ollama models.

Since Ollama supports hundreds of models, we explicitly list the most models but allow any name in the type hints.

OllamaModel `dataclass`

Bases: `Model`

A model that implements Ollama using the OpenAI API.

Internally, this uses the `OpenAI Python client` to interact with the Ollama server.

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

99 Source code in pydantic_ai Slim/pydantic_ai/models/ollama.py

```
57 @dataclass(init=False)
58 class OllamaModel(Model):
59     """A model that implements Ollama using the OpenAI API.
60
61     Internally, this uses the [OpenAI Python client](https://github.com/openai/openai-python) to interact with the Ollama server.
62
63     Apart from `__init__`, all methods are private or match those of the base class.
64     """
65
66     model_name: OllamaModelName
67     openai_model: OpenAIModel
68
69     def __init__(
70         self,
71         model_name: OllamaModelName,
72         *,
73         base_url: str | None = 'http://localhost:11434/v1/',
74         openai_client: AsyncOpenAI | None = None,
75         http_client: AsyncHTTPClient | None = None,
76     ):
77         """Initialize an Ollama model.
78
79         Ollama has built-in compatability for the OpenAI chat completions API ([source](https://ollama.com/blog/openai-compatibility)), so we reuse the
80         ['OpenAIModel'] [pydantic_ai.models.openai.OpenAIModel] here.
81
82         Args:
83             model_name: The name of the Ollama model to use. List of models available [here](https://ollama.com/library)
84             You must first download the model ('ollama pull <MODEL-NAME>') in order to use the model
85             base_url: The base url for the ollama requests. The default value is the ollama default
86             openai_client: An existing
87                 ['AsyncOpenAI'] (https://github.com/openai/openai-python?tab=readme-ov-file#async-usage)
88                 client to use, if provided, 'base_url' and 'http_client' must be 'None'.
89             http_client: An existing 'httpx.AsyncClient' to use for making HTTP requests.
90         """
91         self.model_name = model_name
92         if openai_client is not None:
93             assert base_url is None, 'Cannot provide both `openai_client` and `base_url`'
94             assert http_client is None, 'Cannot provide both `openai_client` and `http_client`'
95             self.openai_model = OpenAIModel(model_name=model_name, openai_client=openai_client)
96         else:
97             # API key is not required for ollama but a value is required to create the client
98             http_client_ = http_client or cached_async_http_client()
99             oai_client = AsyncOpenAI(base_url=base_url, api_key='ollama', http_client=http_client_)
100             self.openai_model = OpenAIModel(model_name=model_name, openai_client=oai_client)
101
102     async def agent_model(
103         self,
104         *,
105         function_tools: list[ToolDefinition],
106         allow_text_result: bool,
107         result_tools: list[ToolDefinition],
108     ) -> AgentModel:
109         return await self.openai_model.agent_model(
110             function_tools=function_tools,
111             allow_text_result=allow_text_result,
112             result_tools=result_tools,
113         )
114
115     def name(self) -> str:
116         return f'ollama:{self.model_name}'
```

`__init__`

```
__init__(
    model_name: OllamaModelName,
    *,
    base_url: str | None = "http://localhost:11434/v1/",
    openai_client: AsyncOpenAI | None = None,
    http_client: AsyncClient | None = None
)
```

Initialize an Ollama model.

Ollama has built-in compatability for the OpenAI chat completions API ([source](#)), so we reuse the `OpenAIModel` here.

Parameters:

Name	Type	Description	Default
<code>model_name</code>	<code>OllamaModelName</code>	The name of the Ollama model to use. List of models available here You must first download the model ('ollama pull <MODEL-NAME>') in order to use the model	<i>required</i>
<code>base_url</code>	<code>str</code> <code>None</code>	The base url for the ollama requests. The default value is the ollama default	<code>'http://localhost:11434/v1/'</code>
<code>openai_client</code>	<code>AsyncOpenAI</code> <code>None</code>	An existing <code>AsyncOpenAI</code> client to use, if provided, <code>base_url</code> and <code>http_client</code> must be <code>None</code> .	<code>None</code>
<code>http_client</code>	<code>AsyncClient</code> <code>None</code>	An existing <code>httpx.AsyncClient</code> to use for making HTTP requests.	<code>None</code>

```

69     def __init__(
70         self,
71         model_name: OllamaModelName,
72         *,
73         base_url: str | None = 'http://localhost:11434/v1/',
74         openai_client: AsyncOpenAI | None = None,
75         http_client: AsyncHTTPClient | None = None,
76     ):
77         """Initialize an Ollama model.
78
79         Ollama has built-in compatability for the OpenAI chat completions API ([source](https://ollama.com/blog/openai-compatibility)), so we reuse the
80         ['OpenAIModel']([pydantic_ai.models.openai.OpenAIModel]) here.
81
82         Args:
83             model_name: The name of the Ollama model to use. List of models available [here](https://ollama.com/library)
84                 You must first download the model ('ollama pull <MODEL-NAME>') in order to use the model
85             base_url: The base url for the ollama requests. The default value is the ollama default
86             openai_client: An existing
87                 ['AsyncOpenAI']([https://github.com/openai/openai-python?tab=readme-ov-file#async-usage])
88                 client to use, if provided, 'base_url' and 'http_client' must be 'None'.
89             http_client: An existing 'httpx.AsyncClient' to use for making HTTP requests.
90         """
91         self.model_name = model_name
92         if openai_client is not None:
93             assert base_url is None, 'Cannot provide both 'openai_client' and 'base_url'
94             assert http_client is None, 'Cannot provide both 'openai_client' and 'http_client'
95             self.openai_model = OpenAIModel(model_name=model_name, openai_client=openai_client)
96         else:
97             # API key is not required for ollama but a value is required to create the client
98             http_client_ = http_client or cached_async_http_client()
99             oai_client = AsyncOpenAI(base_url=base_url, api_key='ollama', http_client=http_client_)
100             self.openai_model = OpenAIModel(model_name=model_name, openai_client=oai_client)

```