# `pydantic_ai.models`

Logic related to making requests to an LLM.

The aim here is to make a common interface for different LLMs, so that the rest of the code can be agnostic to the specific LLM being used.

## KnownModelName `module-attribute`

```
KnownModelName = Literal[
    "openai:gpt-4o",
    "openai:gpt-4o-mini",
    "openai:gpt-4-turbo",
    "openai:gpt-4",
    "openai:o1-preview",
    "openai:o1-mini",
    "openai:gpt-3.5-turbo",
    "groq:llama-3.1-70b-versatile",
    "groq:llama3-groq-70b-8192-tool-use-preview",
    "groq:llama3-groq-8b-8192-tool-use-preview",
    "groq:llama-3.1-70b-specdec",
    "groq:llama-3.1-8b-instant",
    "groq:llama-3.2-1b-preview",
    "groq:llama-3.2-3b-preview",
    "groq:llama-3.2-11b-vision-preview",
    "groq:llama-3.2-90b-vision-preview",
    "groq:llama3-70b-8192",
    "groq:llama3-8b-8192",
    "groq:mixtral-8x7b-32768",
    "groq:gemma2-9b-it",
    "groq:gemma-7b-it",
    "gemini-1.5-flash",
    "gemini-1.5-pro",
    "vertexai:gemini-1.5-flash",
    "vertexai:gemini-1.5-pro",
    "ollama:codellama",
    "ollama:gemma",
    "ollama:gemma2",
    "ollama:llama3",
    "ollama:llama3.1",
    "ollama:llama3.2",
    "ollama:llama3.2-vision",
    "ollama:llama3.3",
    "ollama:mistral",
    "ollama:mistral-nemo",
    "ollama:mixtral",
    "ollama:phi3",
    "ollama:qwq",
    "ollama:qwen",
    "ollama:qwen2",
    "ollama:qwen2.5",
    "ollama:starcoder2",
    "test",
]
```

Known model names that can be used with the `model` parameter of `Agent`.

`KnownModelName` is provided as a concise way to specify a model.

## Model

Bases: `ABC`

Abstract class for a model.

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

```
 77   class Model(ABC):
 78       """Abstract class for a model."""
 79
 80       @abstractmethod
 81       async def agent_model(
 82           self,
 83           *,
 84           function_tools: list[ToolDefinition],
 85           allow_text_result: bool,
 86           result_tools: list[ToolDefinition],
 87       ) -> AgentModel:
 88           """Create an agent model, this is called for each step of an agent run.
 89
 90           This is async in case slow/async config checks need to be performed that can't be done in `__init__`.
 91
 92           Args:
 93               function_tools: The tools available to the agent.
 94               allow_text_result: Whether a plain text final response/result is permitted.
 95               result_tools: Tool definitions for the final result tool(s), if any.
 96
 97           Returns:
 98               An agent model.
 99           """
100           raise NotImplementedError()
101
102       @abstractmethod
103       def name(self) -> str:
104           raise NotImplementedError()
```

### agent_model `abstractmethod` `async`

```
agent_model(
    *,
    function_tools: list[ToolDefinition],
    allow_text_result: bool,
    result_tools: list[ToolDefinition]
) -> AgentModel
```

Create an agent model, this is called for each step of an agent run.

This is async in case slow/async config checks need to be performed that can't be done in `__init__`.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| function_tools | list[ToolDefinition] | The tools available to the agent. | *required* |
| allow_text_result | bool | Whether a plain text final response/result is permitted. | *required* |
| result_tools | list[ToolDefinition] | Tool definitions for the final result tool(s), if any. | *required* |

**Returns:**

| Type | Description |
|------|-------------|
| AgentModel | An agent model. |

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

```python
 80    @abstractmethod
 81    async def agent_model(
 82        self,
 83        *,
 84        function_tools: list[ToolDefinition],
 85        allow_text_result: bool,
 86        result_tools: list[ToolDefinition],
 87    ) -> AgentModel:
 88        """Create an agent model, this is called for each step of an agent run.
 89
 90        This is async in case slow/async config checks need to be performed that can't be done in `__init__`.
 91
 92        Args:
 93            function_tools: The tools available to the agent.
 94            allow_text_result: Whether a plain text final response/result is permitted.
 95            result_tools: Tool definitions for the final result tool(s), if any.
 96
 97        Returns:
 98            An agent model.
 99        """
100        raise NotImplementedError()
```

## AgentModel

Bases: `ABC`

Model configured for each step of an Agent run.

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

```python
107    class AgentModel(ABC):
108        """Model configured for each step of an Agent run."""
109
110        @abstractmethod
111        async def request(self, messages: list[Message]) -> tuple[ModelAnyResponse, Cost]:
112            """Make a request to the model."""
113            raise NotImplementedError()
114
115        @asynccontextmanager
116        async def request_stream(self, messages: list[Message]) -> AsyncIterator[EitherStreamedResponse]:
117            """Make a request to the model and return a streaming response."""
118            raise NotImplementedError(f'Streamed requests not supported by this {self.__class__.__name__}')
119            # yield is required to make this a generator for type checking
120            # noinspection PyUnreachableCode
121            yield  # pragma: no cover
```

### request `abstractmethod` `async`

```python
request(
    messages: list[Message],
) -> tuple[ModelAnyResponse, Cost]
```

Make a request to the model.

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

```python
110    @abstractmethod
111    async def request(self, messages: list[Message]) -> tuple[ModelAnyResponse, Cost]:
112        """Make a request to the model."""
113        raise NotImplementedError()
```

### request_stream `async`

```python
request_stream(
    messages: list[Message],
) -> AsyncIterator[EitherStreamedResponse]
```

Make a request to the model and return a streaming response.

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

```python
115    @asynccontextmanager
116    async def request_stream(self, messages: list[Message]) -> AsyncIterator[EitherStreamedResponse]:
117        """Make a request to the model and return a streaming response."""
118        raise NotImplementedError(f'Streamed requests not supported by this {self.__class__.__name__}')
119        # yield is required to make this a generator for type checking
120        # noinspection PyUnreachableCode
121        yield  # pragma: no cover
```

## StreamTextResponse

Bases: `ABC`

Streamed response from an LLM when returning text.

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

```python
124    class StreamTextResponse(ABC):
125        """Streamed response from an LLM when returning text."""
126
127        def __aiter__(self) -> AsyncIterator[None]:
128            """Stream the response as an async iterable, building up the text as it goes.
129
130            This is an async iterator that yields `None` to avoid doing the work of validating the input and
131            extracting the text field when it will often be thrown away.
132            """
133            return self
134
135        @abstractmethod
136        async def __anext__(self) -> None:
137            """Process the next chunk of the response, see above for why this returns `None`."""
138            raise NotImplementedError()
139
140        @abstractmethod
141        def get(self, *, final: bool = False) -> Iterable[str]:
142            """Returns an iterable of text since the last call to `get()` — e.g. the text delta.
143
144            Args:
145                final: If True, this is the final call, after iteration is complete, the response should be fully validated
146                    and all text extracted.
147            """
148            raise NotImplementedError()
149
150        @abstractmethod
151        def cost(self) -> Cost:
152            """Return the cost of the request.
153
154            NOTE: this won't return the ful cost until the stream is finished.
155            """
156            raise NotImplementedError()
157
158        @abstractmethod
159        def timestamp(self) -> datetime:
160            """Get the timestamp of the response."""
161            raise NotImplementedError()
```

### __aiter__

```
__aiter__() -> AsyncIterator[None]
```

Stream the response as an async iterable, building up the text as it goes.

This is an async iterator that yields `None` to avoid doing the work of validating the input and extracting the text field when it will often be thrown away.

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

```python
127    def __aiter__(self) -> AsyncIterator[None]:
128        """Stream the response as an async iterable, building up the text as it goes.
129
130        This is an async iterator that yields `None` to avoid doing the work of validating the input and
131        extracting the text field when it will often be thrown away.
132        """
133        return self
```

### __anext__ `abstractmethod` `async`

```
__anext__() -> None
```

Process the next chunk of the response, see above for why this returns `None`.

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

```python
135    @abstractmethod
136    async def __anext__(self) -> None:
137        """Process the next chunk of the response, see above for why this returns `None`."""
138        raise NotImplementedError()
```

### get `abstractmethod`

```
get(*, final: bool = False) -> Iterable[str]
```

Returns an iterable of text since the last call to `get()` — e.g. the text delta.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| `final` | `bool` | If True, this is the final call, after iteration is complete, the response should be fully validated and all text extracted. | `False` |

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

```python
140    @abstractmethod
141    def get(self, *, final: bool = False) -> Iterable[str]:
142        """Returns an iterable of text since the last call to `get()` — e.g. the text delta.
143
144        Args:
145            final: If True, this is the final call, after iteration is complete, the response should be fully validated
146                and all text extracted.
147        """
148        raise NotImplementedError()
```

### cost `abstractmethod`

```
cost() -> Cost
```

Return the cost of the request.

NOTE: this won't return the ful cost until the stream is finished.

```python
150    @abstractmethod
151    def cost(self) -> Cost:
152        """Return the cost of the request.
153
154        NOTE: this won't return the ful cost until the stream is finished.
155        """
156        raise NotImplementedError()
```

### timestamp `abstractmethod`

```
timestamp() -> datetime
```

Get the timestamp of the response.

```python
158    @abstractmethod
159    def timestamp(self) -> datetime:
160        """Get the timestamp of the response."""
161        raise NotImplementedError()
```

## StreamStructuredResponse

Bases: `ABC`

Streamed response from an LLM when calling a tool.

```python
164    class StreamStructuredResponse(ABC):
165        """Streamed response from an LLM when calling a tool."""
166
167        def __aiter__(self) -> AsyncIterator[None]:
168            """Stream the response as an async iterable, building up the tool call as it goes.
169
170            This is an async iterator that yields `None` to avoid doing the work of building the final tool call when
171            it will often be thrown away.
172            """
173            return self
174
175        @abstractmethod
176        async def __anext__(self) -> None:
177            """Process the next chunk of the response, see above for why this returns `None`."""
178            raise NotImplementedError()
179
180        @abstractmethod
181        def get(self, *, final: bool = False) -> ModelStructuredResponse:
182            """Get the `ModelStructuredResponse` at this point.
183
184            The `ModelStructuredResponse` may or may not be complete, depending on whether the stream is finished.
185
186            Args:
187                final: If True, this is the final call, after iteration is complete, the response should be fully validated.
188            """
189            raise NotImplementedError()
190
191        @abstractmethod
192        def cost(self) -> Cost:
193            """Get the cost of the request.
194
195            NOTE: this won't return the full cost until the stream is finished.
196            """
197            raise NotImplementedError()
198
199        @abstractmethod
200        def timestamp(self) -> datetime:
201            """Get the timestamp of the response."""
202            raise NotImplementedError()
```

### __aiter__

```
__aiter__() -> AsyncIterator[None]
```

Stream the response as an async iterable, building up the tool call as it goes.

This is an async iterator that yields `None` to avoid doing the work of building the final tool call when it will often be thrown away.

```python
167    def __aiter__(self) -> AsyncIterator[None]:
168        """Stream the response as an async iterable, building up the tool call as it goes.
169
170        This is an async iterator that yields `None` to avoid doing the work of building the final tool call when
171        it will often be thrown away.
172        """
173        return self
```

### __anext__ `abstractmethod` `async`

```
__anext__() -> None
```

Process the next chunk of the response, see above for why this returns `None`.

```
175   @abstractmethod
176   async def __anext__(self) -> None:
177       """Process the next chunk of the response, see above for why this returns `None`."""
178       raise NotImplementedError()
```

**get** `abstractmethod`

```
get(*, final: bool = False) -> ModelStructuredResponse
```

Get the `ModelStructuredResponse` at this point.

The `ModelStructuredResponse` may or may not be complete, depending on whether the stream is finished.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| final | bool | If True, this is the final call, after iteration is complete, the response should be fully validated. | False |

```
180   @abstractmethod
181   def get(self, *, final: bool = False) -> ModelStructuredResponse:
182       """Get the `ModelStructuredResponse` at this point.
183
184       The `ModelStructuredResponse` may or may not be complete, depending on whether the stream is finished.
185
186       Args:
187           final: If True, this is the final call, after iteration is complete, the response should be fully validated.
188       """
189       raise NotImplementedError()
```

**cost** `abstractmethod`

```
cost() -> Cost
```

Get the cost of the request.

NOTE: this won't return the full cost until the stream is finished.

```
191   @abstractmethod
192   def cost(self) -> Cost:
193       """Get the cost of the request.
194
195       NOTE: this won't return the full cost until the stream is finished.
196       """
197       raise NotImplementedError()
```

**timestamp** `abstractmethod`

```
timestamp() -> datetime
```

Get the timestamp of the response.

```
199   @abstractmethod
200   def timestamp(self) -> datetime:
201       """Get the timestamp of the response."""
202       raise NotImplementedError()
```

## ALLOW_MODEL_REQUESTS `module-attribute`

```
ALLOW_MODEL_REQUESTS = True
```

Whether to allow requests to models.

This global setting allows you to disable request to most models, e.g. to make sure you don't accidentally make costly requests to a model during tests.

The testing models `TestModel` and `FunctionModel` are no affected by this setting.

## check_allow_model_requests

```
check_allow_model_requests() -> None
```

Check if model requests are allowed.

If you're defining your own models that have cost or latency associated with their use, you should call this in `Model.agent_model`.

**Raises:**

| Type | Description |
|------|-------------|
| RuntimeError | If model requests are not allowed. |

```
219   def check_allow_model_requests() -> None:
220       """Check if model requests are allowed.
221
222       If you're defining your own models that have cost or latency associated with their use, you should call this in
223       [`Model.agent_model`][pydantic_ai.models.Model.agent_model].
224
225       Raises:
226           RuntimeError: If model requests are not allowed.
227       """
228       if not ALLOW_MODEL_REQUESTS:
229           raise RuntimeError('Model requests are not allowed, since ALLOW_MODEL_REQUESTS is False')
```

## override_allow_model_requests

```
override_allow_model_requests(
    allow_model_requests: bool,
) -> Iterator[None]
```

Context manager to temporarily override `ALLOW_MODEL_REQUESTS`.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| allow_model_requests | bool | Whether to allow model requests within the context. | *required* |

```
232   @contextmanager
233   def override_allow_model_requests(allow_model_requests: bool) -> Iterator[None]:
234       """Context manager to temporarily override [`ALLOW_MODEL_REQUESTS`][pydantic_ai.models.ALLOW_MODEL_REQUESTS].
235
236       Args:
237           allow_model_requests: Whether to allow model requests within the context.
238       """
239       global ALLOW_MODEL_REQUESTS
240       old_value = ALLOW_MODEL_REQUESTS
241       ALLOW_MODEL_REQUESTS = allow_model_requests  # pyright: ignore[reportConstantRedefinition]
242       try:
243           yield
244       finally:
245           ALLOW_MODEL_REQUESTS = old_value  # pyright: ignore[reportConstantRedefinition]
```