

pydantic_ai.models.vertexai

Custom interface to the `*-aiplatform.googleapis.com` API for Gemini models.

This model uses `GeminiAgentModel` with just the URL and auth method changed from `GeminiModel`, it relies on the VertexAI `generateContent` and `streamGenerateContent` function endpoints having the same schemas as the equivalent `Gemini endpoints`.

Setup

For details on how to set up authentication with this model as well as a comparison with the `generativelanguage.googleapis.com` API used by `GeminiModel`, see [model configuration for Gemini via VertexAI](#).

Example Usage

With the default google project already configured in your environment using "application default credentials":

vertex_example_env.py

```

from pydantic_ai import Agent
from pydantic_ai.models.vertexai import VertexAIModel

model = VertexAIModel('gemini-1.5-flash')
agent = Agent(model)
result = agent.run_sync('Tell me a joke.')
print(result.data)
#> Did you hear about the toothpaste scandal? They called it Colgate.

```

Or using a service account JSON file:

vertex_example_service_account.py

```

from pydantic_ai import Agent
from pydantic_ai.models.vertexai import VertexAIModel

model = VertexAIModel(
    'gemini-1.5-flash',
    service_account_file='path/to/service-account.json',
)
agent = Agent(model)
result = agent.run_sync('Tell me a joke.')
print(result.data)
#> Did you hear about the toothpaste scandal? They called it Colgate.

```

VERTEX_AI_URL_TEMPLATE module-attribute

```

VERTEX_AI_URL_TEMPLATE = "https://{region}-aiplatform.googleapis.com/v1/projects/{project_id}/locations/{region}/publishers/{model_publisher}/models/{model}:"

```

URL template for Vertex AI.

See [generateContent docs](#) and [streamGenerateContent docs](#) for more information.

The template is used thus:

- `region` is substituted with the `region` argument, see [available regions](#)
- `model_publisher` is substituted with the `model_publisher` argument
- `model` is substituted with the `model_name` argument
- `project_id` is substituted with the `project_id` from auth/credentials
- `function` (`generateContent` or `streamGenerateContent`) is added to the end of the URL

VertexAIModel dataclass

Bases: `Model`

A model that uses Gemini via the `*-aiplatform.googleapis.com` VertexAI API.

```

54 @dataclass(init=False)
55 class VertexAIModel(Model):
56     """A model that uses Gemini via the `*-aiplatform.googleapis.com` VertexAI API."""
57
58     model_name: GeminiModelName
59     service_account_file: Path | str | None
60     project_id: str | None
61     region: VertexAiRegion
62     model_publisher: Literal['google']
63     http_client: AsyncHTTPClient
64     url_template: str
65
66     auth: BearerTokenAuth | None
67     url: str | None
68
69     # TODO __init__ can be removed once we drop 3.9 and we can set kw_only correctly on the dataclass
70     def __init__(
71         self,
72         model_name: GeminiModelName,
73         *,
74         service_account_file: Path | str | None = None,
75         project_id: str | None = None,
76         region: VertexAiRegion = 'us-central1',
77         model_publisher: Literal['google'] = 'google',
78         http_client: AsyncHTTPClient | None = None,
79         url_template: str = VERTEX_AI_URL_TEMPLATE,
80     ):
81         """Initialize a Vertex AI Gemini model.
82
83         Args:
84             model_name: The name of the model to use. I couldn't find a list of supported Google models, in VertexAI
85                 so for now this uses the same models as the [Gemini model][pydantic_ai.models.gemini.GeminiModel].
86             service_account_file: Path to a service account file.
87                 If not provided, the default environment credentials will be used.
88             project_id: The project ID to use, if not provided it will be taken from the credentials.
89             region: The region to make requests to.
90             model_publisher: The model publisher to use, I couldn't find a good list of available publishers,
91                 and from trial and error it seems non-google models don't work with the 'generateContent' and
92                 'streamGenerateContent' functions, hence only 'google' is currently supported.
93                 Please create an issue or PR if you know how to use other publishers.
94             http_client: An existing 'httpx.AsyncClient' to use for making HTTP requests.
95             url_template: URL template for Vertex AI, see
96                 ['VERTEX_AI_URL_TEMPLATE' docs][pydantic_ai.models.vertexai.VERTEX_AI_URL_TEMPLATE]
97                 for more information.
98
99             self.model_name = model_name
100             self.service_account_file = service_account_file
101             self.project_id = project_id
102             self.region = region
103             self.model_publisher = model_publisher
104             self.http_client = http_client or cached_async_http_client()
105             self.url_template = url_template
106
107             self.auth = None
108             self.url = None
109
110         async def agent_model(
111             self,
112             *,
113             function_tools: list[ToolDefinition],
114             allow_text_result: bool,
115             result_tools: list[ToolDefinition],
116         ) -> GeminiAgentModel:
117             url, auth = await self._ainit()
118             return GeminiAgentModel(
119                 http_client=self.http_client,
120                 model_name=self.model_name,
121                 auth=auth,
122                 url=url,
123                 function_tools=function_tools,
124                 allow_text_result=allow_text_result,
125                 result_tools=result_tools,
126             )
127
128         async def _ainit(self) -> tuple[str, BearerTokenAuth]:
129             if self.url is not None and self.auth is not None:
130                 return self.url, self.auth
131
132             if self.service_account_file is not None:
133                 creds: BaseCredentials | ServiceAccountCredentials = _creds_from_file(self.service_account_file)
134                 assert creds.project_id is None or isinstance(creds.project_id, str)
135                 creds_project_id: str | None = creds.project_id
136                 creds_source = 'service account file'
137             else:
138                 creds, creds_project_id = await _async_google_auth()
139                 creds_source = 'google.auth.default()'
140
141             if self.project_id is None:
142                 if creds_project_id is None:
143                     raise UserError(f'No project_id provided and none found in {creds_source}')
144                 project_id = creds_project_id
145             else:
146                 if creds_project_id is not None and self.project_id != creds_project_id:
147                     raise UserError(
148                         f'The project_id you provided does not match the one from {creds_source}: '
149                         f'{self.project_id!r} != {creds_project_id!r}'
150                     )
151                 project_id = self.project_id
152
153             self.url = url = self.url_template.format(
154                 region=self.region,
155                 project_id=project_id,
156                 model_publisher=self.model_publisher,
157                 model=self.model_name,
158             )
159             self.auth = auth = BearerTokenAuth(creds)
160             return url, auth
161
162         def name(self) -> str:
163             return f'vertexai:{self.model_name}'

```

__init__

```

__init__(
    model_name: GeminiModelName,
    *,
    service_account_file: Path | str | None = None,
    project_id: str | None = None,
    region: VertexAiRegion = "us-central1",
    model_publisher: Literal["google"] = "google",
    http_client: AsyncClient | None = None,
    url_template: str = VERTEX_AI_URL_TEMPLATE
)

```

Initialize a Vertex AI Gemini model.

Parameters:

Name	Type	Description	Default
<code>model_name</code>	<code>GeminiModelName</code>	The name of the model to use. I couldn't find a list of supported Google models, in VertexAI so for now this uses the same models as the <code>Gemini model</code> .	<i>required</i>
<code>service_account_file</code>	<code>Path</code> <code>str</code> <code>None</code>	Path to a service account file. If not provided, the default environment credentials will be used.	<code>None</code>
<code>project_id</code>	<code>str</code> <code>None</code>	The project ID to use, if not provided it will be taken from the credentials.	<code>None</code>
<code>region</code>	<code>VertexAiRegion</code>	The region to make requests to.	<code>'us-central1'</code>
<code>model_publisher</code>	<code>Literal['google']</code>	The model publisher to use, I couldn't find a good list of available publishers, and from trial and error it seems non-google models don't work with the <code>generateContent</code> and <code>streamGenerateContent</code> functions, hence only <code>google</code> is currently supported. Please create an issue or PR if you know how to use other publishers.	<code>'google'</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>
<code>url_template</code>	<code>str</code>	URL template for Vertex AI, see <code>VERTEX_AI_URL_TEMPLATE docs</code> for more information.	<code>VERTEX_AI_URL_TEMPLATE</code>

```
99 Source code in pydantic_ai_slim/pydantic_ai/models/vertexai.py
70 def __init__(
71     self,
72     model_name: GeminiModelName,
73     *,
74     service_account_file: Path | str | None = None,
75     project_id: str | None = None,
76     region: VertexAiRegion = 'us-central1',
77     model_publisher: Literal['google'] = 'google',
78     http_client: AsyncHTTPClient | None = None,
79     url_template: str = VERTEX_AI_URL_TEMPLATE,
80 ):
81     """Initialize a Vertex AI Gemini model.
82
83     Args:
84         model_name: The name of the model to use. I couldn't find a list of supported Google models, in VertexAI
85             so for now this uses the same models as the [Gemini model][pydantic_ai.models.gemini.GeminiModel].
86         service_account_file: Path to a service account file.
87             If not provided, the default environment credentials will be used.
88         project_id: The project ID to use, if not provided it will be taken from the credentials.
89         region: The region to make requests to.
90         model_publisher: The model publisher to use, I couldn't find a good list of available publishers,
91             and from trial and error it seems non-google models don't work with the 'generateContent' and
92             'streamGenerateContent' functions, hence only 'google' is currently supported.
93             Please create an issue or PR if you know how to use other publishers.
94         http_client: An existing 'httpx.AsyncClient' to use for making HTTP requests.
95         url_template: URL template for Vertex AI, see
96             ['VERTEX_AI_URL_TEMPLATE' docs][pydantic_ai.models.vertexai.VERTEX_AI_URL_TEMPLATE]
97             for more information.
98     """
99     self.model_name = model_name
100     self.service_account_file = service_account_file
101     self.project_id = project_id
102     self.region = region
103     self.model_publisher = model_publisher
104     self.http_client = http_client or cached_async_http_client()
105     self.url_template = url_template
106
107     self.auth = None
108     self.url = None
```

BearerTokenAuth dataclass

Authentication using a bearer token generated by google-auth.

```
99 Source code in pydantic_ai_slim/pydantic_ai/models/vertexai.py
184 @dataclass
185 class BearerTokenAuth:
186     """Authentication using a bearer token generated by google-auth."""
187
188     credentials: BaseCredentials | ServiceAccountCredentials
189     token_created: datetime | None = field(default=None, init=False)
190
191     async def headers(self) -> dict[str, str]:
192         if self.credentials.token is None or self._token_expired():
193             await run_in_executor(self._refresh_token)
194             self.token_created = datetime.now()
195             return {'Authorization': f'Bearer {self.credentials.token}'}
196
197     def _token_expired(self) -> bool:
198         if self.token_created is None:
199             return True
200         else:
201             return (datetime.now() - self.token_created) > MAX_TOKEN_AGE
202
203     def _refresh_token(self) -> str:
204         self.credentials.refresh(Request())
205         assert isinstance(self.credentials.token, str), f'Expected token to be a string, got {self.credentials.token}'
206         return self.credentials.token
```

VertexAiRegion module-attribute

```
VertexAiRegion = Literal[
    "us-central1",
    "us-east1",
    "us-east4",
    "us-south1",
    "us-west1",
    "us-west2",
    "us-west3",
    "us-west4",
]
```

```
"us-east5",
"europa-central2",
"europa-north1",
"europa-southwest1",
"europa-west1",
"europa-west2",
"europa-west3",
"europa-west4",
"europa-west6",
"europa-west8",
"europa-west9",
"europa-west12",
"africa-south1",
"asia-east1",
"asia-east2",
"asia-northeast1",
"asia-northeast2",
"asia-northeast3",
"asia-south1",
"asia-southeast1",
"asia-southeast2",
"australia-southeast1",
"australia-southeast2",
"me-central1",
"me-central2",
"me-west1",
"northamerica-northeast1",
"northamerica-northeast2",
"southamerica-east1",
"southamerica-west1",
]
```

Regions available for Vertex AI.

More details [here](#).