# `pydantic_ai.models.gemini`

Custom interface to the `generativelanguage.googleapis.com` API using HTTPX and [Pydantic](https://docs.pydantic.dev/latest/.

The Google SDK for interacting with the `generativelanguage.googleapis.com` API `google-generativeai` reads like it was written by a Java developer who thought they knew everything about OOP, spent 30 minutes trying to learn Python, gave up and decided to build the library to prove how horrible Python is. It also doesn't use httpx for HTTP requests, and tries to implement tool calling itself, but doesn't use Pydantic or equivalent for validation.

We therefore implement support for the API directly.

Despite these shortcomings, the Gemini model is actually quite powerful and very fast.

## Setup

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

### GeminiModelName `module-attribute`

```
GeminiModelName = Literal[
    "gemini-1.5-flash",
    "gemini-1.5-flash-8b",
    "gemini-1.5-pro",
    "gemini-1.0-pro",
]
```

Named Gemini models.

See the Gemini API docs for a full list.

### GeminiModel `dataclass`

Bases: `Model`

A model that uses Gemini via `generativelanguage.googleapis.com` API.

This is implemented from scratch rather than using a dedicated SDK, good API documentation is available here.

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

> **Source code in** `pydantic_ai_slim/pydantic_ai/models/gemini.py`                                   ⌄

```
47   @dataclass(init=False)
48   class GeminiModel(Model):
49       """A model that uses Gemini via `generativelanguage.googleapis.com` API.
50
51       This is implemented from scratch rather than using a dedicated SDK, good API documentation is
52       available [here](https://ai.google.dev/api).
53
54       Apart from `__init__`, all methods are private or match those of the base class.
55       """
56
57       model_name: GeminiModelName
58       auth: AuthProtocol
59       http_client: AsyncHTTPClient
60       url: str
61
62       def __init__(
63           self,
64           model_name: GeminiModelName,
65           *,
66           api_key: str | None = None,
67           http_client: AsyncHTTPClient | None = None,
68           url_template: str = 'https://generativelanguage.googleapis.com/v1beta/models/{model}:',
69       ):
70           """Initialize a Gemini model.
71
72           Args:
73               model_name: The name of the model to use.
74               api_key: The API key to use for authentication, if not provided, the `GEMINI_API_KEY` environment variable
75                   will be used if available.
76               http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
77               url_template: The URL template to use for making requests, you shouldn't need to change this,
78                   docs [here](https://ai.google.dev/gemini-api/docs/quickstart?lang=rest#make-first-request),
79                   `model` is substituted with the model name, and `function` is added to the end of the URL.
80           """
81           self.model_name = model_name
82           if api_key is None:
83               if env_api_key := os.getenv('GEMINI_API_KEY'):
84                   api_key = env_api_key
85               else:
86                   raise exceptions.UserError('API key must be provided or set in the GEMINI_API_KEY environment variable')
87           self.auth = ApiKeyAuth(api_key)
88           self.http_client = http_client or cached_async_http_client()
89           self.url = url_template.format(model=model_name)
90
91       async def agent_model(
92           self,
93           *,
94           function_tools: list[ToolDefinition],
95           allow_text_result: bool,
96           result_tools: list[ToolDefinition],
97       ) -> GeminiAgentModel:
98           return GeminiAgentModel(
99               http_client=self.http_client,
100              model_name=self.model_name,
101              auth=self.auth,
102              url=self.url,
103              function_tools=function_tools,
104              allow_text_result=allow_text_result,
105              result_tools=result_tools,
106          )
107
108      def name(self) -> str:
109          return self.model_name
```

**__init__**

```
__init__(
    model_name: GeminiModelName,
    *,
    api_key: str | None = None,
    http_client: AsyncClient | None = None,
    url_template: str = "https://generativelanguage.googleapis.com/v1beta/models/{model}:"
)
```

Initialize a Gemini model.

**Parameters:**

| Name | Type | Description | Default |
|------|------|-------------|---------|
| model_name | GeminiModelName | The name of the model to use. | *required* |
| api_key | str | None | The API key to use for authentication, if not provided, the `GEMINI_API_KEY` environment variable will be used if available. | None |
| http_client | AsyncClient | None | An existing `httpx.AsyncClient` to use for making HTTP requests. | None |
| url_template | str | The URL template to use for making requests, you shouldn't need to change this, docs here, `model` is substituted with the model name, and `function` is added to the end of the URL. | `'https://generativelanguage.googleapis.com/v1beta/models/{model}:'` |

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

```python
62    def __init__(
63        self,
64        model_name: GeminiModelName,
65        *,
66        api_key: str | None = None,
67        http_client: AsyncHTTPClient | None = None,
68        url_template: str = 'https://generativelanguage.googleapis.com/v1beta/models/{model}:',
69    ):
70        """Initialize a Gemini model.
71
72        Args:
73            model_name: The name of the model to use.
74            api_key: The API key to use for authentication, if not provided, the `GEMINI_API_KEY` environment variable
75                will be used if available.
76            http_client: An existing `httpx.AsyncClient` to use for making HTTP requests.
77            url_template: The URL template to use for making requests, you shouldn't need to change this,
78                docs [here](https://ai.google.dev/gemini-api/docs/quickstart?lang=rest#make-first-request),
79                `model` is substituted with the model name, and `function` is added to the end of the URL.
80        """
81        self.model_name = model_name
82        if api_key is None:
83            if env_api_key := os.getenv('GEMINI_API_KEY'):
84                api_key = env_api_key
85            else:
86                raise exceptions.UserError('API key must be provided or set in the GEMINI_API_KEY environment variable')
87        self.auth = ApiKeyAuth(api_key)
88        self.http_client = http_client or cached_async_http_client()
89        self.url = url_template.format(model=model_name)
```

## AuthProtocol

Bases: `Protocol`

Abstract definition for Gemini authentication.

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

```python
112    class AuthProtocol(Protocol):
113        """Abstract definition for Gemini authentication."""
114
115        async def headers(self) -> dict[str, str]: ...
```

## ApiKeyAuth `dataclass`

Authentication using an API key for the `X-Goog-Api-Key` header.

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

```python
118    @dataclass
119    class ApiKeyAuth:
120        """Authentication using an API key for the `X-Goog-Api-Key` header."""
121
122        api_key: str
123
124        async def headers(self) -> dict[str, str]:
125            # https://cloud.google.com/docs/authentication/api-keys-use#using-with-rest
126            return {'X-Goog-Api-Key': self.api_key}
```

## GeminiAgentModel `dataclass`

Bases: `AgentModel`

Implementation of `AgentModel` for Gemini models.

```python
@dataclass(init=False)
class GeminiAgentModel(AgentModel):
    """Implementation of `AgentModel` for Gemini models."""

    http_client: AsyncHTTPClient
    model_name: GeminiModelName
    auth: AuthProtocol
    tools: _GeminiTools | None
    tool_config: _GeminiToolConfig | None
    url: str

    def __init__(
        self,
        http_client: AsyncHTTPClient,
        model_name: GeminiModelName,
        auth: AuthProtocol,
        url: str,
        function_tools: list[ToolDefinition],
        allow_text_result: bool,
        result_tools: list[ToolDefinition],
    ):
        check_allow_model_requests()
        tools = [_function_from_abstract_tool(t) for t in function_tools]
        if result_tools:
            tools += [_function_from_abstract_tool(t) for t in result_tools]

        if allow_text_result:
            tool_config = None
        else:
            tool_config = _tool_config([t['name'] for t in tools])

        self.http_client = http_client
        self.model_name = model_name
        self.auth = auth
        self.tools = _GeminiTools(function_declarations=tools) if tools else None
        self.tool_config = tool_config
        self.url = url

    async def request(self, messages: list[Message]) -> tuple[ModelAnyResponse, result.Cost]:
        async with self._make_request(messages, False) as http_response:
            response = _gemini_response_ta.validate_json(await http_response.aread())
        return self._process_response(response), _metadata_as_cost(response)

    @asynccontextmanager
    async def request_stream(self, messages: list[Message]) -> AsyncIterator[EitherStreamedResponse]:
        async with self._make_request(messages, True) as http_response:
            yield await self._process_streamed_response(http_response)

    @asynccontextmanager
    async def _make_request(self, messages: list[Message], streamed: bool) -> AsyncIterator[HTTPResponse]:
        contents: list[_GeminiContent] = []
        sys_prompt_parts: list[_GeminiTextPart] = []
        for m in messages:
            either_content = self._message_to_gemini(m)
            if left := either_content.left:
                sys_prompt_parts.append(left.value)
            else:
                contents.append(either_content.right)

        request_data = _GeminiRequest(contents=contents)
        if sys_prompt_parts:
            request_data['system_instruction'] = _GeminiTextContent(role='user', parts=sys_prompt_parts)
        if self.tools is not None:
            request_data['tools'] = self.tools
        if self.tool_config is not None:
            request_data['tool_config'] = self.tool_config

        url = self.url + ('streamGenerateContent' if streamed else 'generateContent')

        headers = {
            'Content-Type': 'application/json',
            'User-Agent': get_user_agent(),
            **await self.auth.headers(),
        }

        request_json = _gemini_request_ta.dump_json(request_data, by_alias=True)

        async with self.http_client.stream('POST', url, content=request_json, headers=headers) as r:
            if r.status_code != 200:
                await r.aread()
                raise exceptions.UnexpectedModelBehavior(f'Unexpected response from gemini {r.status_code}', r.text)
            yield r

    @staticmethod
    def _process_response(response: _GeminiResponse) -> ModelAnyResponse:
        either = _extract_response_parts(response)
        if left := either.left:
            return _structured_response_from_parts(left.value)
        else:
            return ModelTextResponse(content=''.join(part['text'] for part in either.right))

    @staticmethod
    async def _process_streamed_response(http_response: HTTPResponse) -> EitherStreamedResponse:
        """Process a streamed response, and prepare a streaming response to return."""
        aiter_bytes = http_response.aiter_bytes()
        start_response: _GeminiResponse | None = None
        content = bytearray()

        async for chunk in aiter_bytes:
            content.extend(chunk)
            responses = _gemini_streamed_response_ta.validate_json(
                content,
                experimental_allow_partial='trailing-strings',
            )
            if responses:
                last = responses[-1]
                if last['candidates'] and last['candidates'][0]['content']['parts']:
                    start_response = last
                    break

        if start_response is None:
            raise UnexpectedModelBehavior('Streamed response ended without content or tool calls')

        if _extract_response_parts(start_response).is_left():
            return GeminiStreamStructuredResponse(_content=content, _stream=aiter_bytes)
        else:
            return GeminiStreamTextResponse(_json_content=content, _stream=aiter_bytes)

    @staticmethod
    def _message_to_gemini(m: Message) -> _utils.Either[_GeminiTextPart, _GeminiContent]:
        """Convert a message to a _GeminiTextPart for "system_instructions" or _GeminiContent for "contents"."""
        if m.role == 'system':
            # SystemPrompt ->
            return _utils.Either(left=_GeminiTextPart(text=m.content))
        elif m.role == 'user':
            # UserPrompt ->
            return _utils.Either(right=_content_user_text(m.content))
```

```
256          elif m.role == 'tool-return':
257              # ToolReturn ->
258              return _utils.Either(right=_content_function_return(m))
259          elif m.role == 'retry-prompt':
260              # RetryPrompt ->
261              return _utils.Either(right=_content_function_retry(m))
262          elif m.role == 'model-text-response':
263              # ModelTextResponse ->
264              return _utils.Either(right=_content_model_text(m.content))
265          elif m.role == 'model-structured-response':
266              # ModelStructuredResponse ->
267              return _utils.Either(right=_content_function_call(m))
268          else:
269              assert_never(m)
```

## GeminiStreamTextResponse `dataclass`

Bases: `StreamTextResponse`

Implementation of `StreamTextResponse` for the Gemini model.

> Source code in `pydantic_ai_slim/pydantic_ai/models/gemini.py`

```
272  @dataclass
273  class GeminiStreamTextResponse(StreamTextResponse):
274      """Implementation of `StreamTextResponse` for the Gemini model."""
275
276      _json_content: bytearray
277      _stream: AsyncIterator[bytes]
278      _position: int = 0
279      _timestamp: datetime = field(default_factory=_utils.now_utc, init=False)
280      _cost: result.Cost = field(default_factory=result.Cost, init=False)
281
282      async def __anext__(self) -> None:
283          chunk = await self._stream.__anext__()
284          self._json_content.extend(chunk)
285
286      def get(self, *, final: bool = False) -> Iterable[str]:
287          if final:
288              all_items = pydantic_core.from_json(self._json_content)
289              new_items = all_items[self._position :]
290              self._position = len(all_items)
291              new_responses = _gemini_streamed_response_ta.validate_python(new_items)
292          else:
293              all_items = pydantic_core.from_json(self._json_content, allow_partial=True)
294              new_items = all_items[self._position : -1]
295              self._position = len(all_items) - 1
296              new_responses = _gemini_streamed_response_ta.validate_python(
297                  new_items, experimental_allow_partial='trailing-strings'
298              )
299          for r in new_responses:
300              self._cost += _metadata_as_cost(r)
301              parts = r['candidates'][0]['content']['parts']
302              if _all_text_parts(parts):
303                  for part in parts:
304                      yield part['text']
305              else:
306                  raise UnexpectedModelBehavior(
307                      'Streamed response with unexpected content, expected all parts to be text'
308                  )
309
310      def cost(self) -> result.Cost:
311          return self._cost
312
313      def timestamp(self) -> datetime:
314          return self._timestamp
```

## GeminiStreamStructuredResponse `dataclass`

Bases: `StreamStructuredResponse`

Implementation of `StreamStructuredResponse` for the Gemini model.

```python
317  @dataclass
318  class GeminiStreamStructuredResponse(StreamStructuredResponse):
319      """Implementation of `StreamStructuredResponse` for the Gemini model."""
320
321      _content: bytearray
322      _stream: AsyncIterator[bytes]
323      _timestamp: datetime = field(default_factory=_utils.now_utc, init=False)
324      _cost: result.Cost = field(default_factory=result.Cost, init=False)
325
326      async def __anext__(self) -> None:
327          chunk = await self._stream.__anext__()
328          self._content.extend(chunk)
329
330      def get(self, *, final: bool = False) -> ModelStructuredResponse:
331          """Get the `ModelStructuredResponse` at this point.
332
333          NOTE: It's not clear how the stream of responses should be combined because Gemini seems to always
334          reply with a single response, when returning a structured data.
335
336          I'm therefore assuming that each part contains a complete tool call, and not trying to combine data from
337          separate parts.
338          """
339          responses = _gemini_streamed_response_ta.validate_json(
340              self._content,
341              experimental_allow_partial='off' if final else 'trailing-strings',
342          )
343          combined_parts: list[_GeminiFunctionCallPart] = []
344          self._cost = result.Cost()
345          for r in responses:
346              self._cost += _metadata_as_cost(r)
347              candidate = r['candidates'][0]
348              parts = candidate['content']['parts']
349              if _all_function_call_parts(parts):
350                  combined_parts.extend(parts)
351              elif not candidate.get('finish_reason'):
352                  # you can get an empty text part along with the finish_reason, so we ignore that case
353                  raise UnexpectedModelBehavior(
354                      'Streamed response with unexpected content, expected all parts to be function calls'
355                  )
356          return _structured_response_from_parts(combined_parts, timestamp=self._timestamp)
357
358      def cost(self) -> result.Cost:
359          return self._cost
360
361      def timestamp(self) -> datetime:
362          return self._timestamp
```

### get

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

Get the `ModelStructuredResponse` at this point.

NOTE: It's not clear how the stream of responses should be combined because Gemini seems to always reply with a single response, when returning a structured data.

I'm therefore assuming that each part contains a complete tool call, and not trying to combine data from separate parts.

```python
330  def get(self, *, final: bool = False) -> ModelStructuredResponse:
331      """Get the `ModelStructuredResponse` at this point.
332
333      NOTE: It's not clear how the stream of responses should be combined because Gemini seems to always
334      reply with a single response, when returning a structured data.
335
336      I'm therefore assuming that each part contains a complete tool call, and not trying to combine data from
337      separate parts.
338      """
339      responses = _gemini_streamed_response_ta.validate_json(
340          self._content,
341          experimental_allow_partial='off' if final else 'trailing-strings',
342      )
343      combined_parts: list[_GeminiFunctionCallPart] = []
344      self._cost = result.Cost()
345      for r in responses:
346          self._cost += _metadata_as_cost(r)
347          candidate = r['candidates'][0]
348          parts = candidate['content']['parts']
349          if _all_function_call_parts(parts):
350              combined_parts.extend(parts)
351          elif not candidate.get('finish_reason'):
352              # you can get an empty text part along with the finish_reason, so we ignore that case
353              raise UnexpectedModelBehavior(
354                  'Streamed response with unexpected content, expected all parts to be function calls'
355              )
356      return _structured_response_from_parts(combined_parts, timestamp=self._timestamp)
```