

pydantic_ai.models.openrouter

Setup

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

KnownOpenRouterProviders module-attribute

```
KnownOpenRouterProviders = Literal[  
    "z-ai",  
    "cerebras",  
    "venice",  
    "moonshotai",  
    "morph",  
    "stealth",  
    "wandb",  
    "klusterai",  
    "openai",  
    "sambanova",  
    "amazon-bedrock",  
    "mistral",  
    "nextbit",  
    "atoma",  
    "ai21",  
    "minimax",  
    "baseten",  
    "anthropic",  
    "featherless",  
    "groq",  
    "lambda",  
    "azure",  
    "ncompass",  
    "deepseek",  
    "hyperbolic",  
    "crusoe",  
    "cohere",  
    "mancer",  
    "avian",  
    "perplexity",  
    "novita",  
    "siliconflow",  
    "switchpoint",  
    "xai",  
    "inflection",  
    "fireworks",  
]
```

OpenRouter

```
"deepinfra",
"inference-net",
"inception",
"atlas-cloud",
"nvidia",
"alibaba",
"friendli",
"infermatic",
"targon",
"ubicloud",
"aion-labs",
"liquid",
"nineteen",
"cloudflare",
"nebius",
"chutes",
"enfer",
"crofai",
"open-inference",
"phala",
"gmicloud",
"meta",
"relace",
"parasail",
"together",
"google-ai-studio",
"google-vertex",
]
```

Known providers in the OpenRouter marketplace

OpenRouterProviderName module-attribute

```
OpenRouterProviderName = str | KnownOpenRouterProviders
```



Possible OpenRouter provider names.

Since OpenRouter is constantly updating their list of providers, we explicitly list some known providers but allow any name in the type hints. See [the OpenRouter API](#) for a full list.

OpenRouterTransforms module-attribute

```
OpenRouterTransforms = Literal['middle-out']
```



Available messages transforms for OpenRouter models with limited token windows.

Currently only supports 'middle-out', but is expected to grow in the future.

OpenRouterProviderConfig

Bases: `TypedDict`

Represents the 'Provider' object from the OpenRouter API.

⟨ ⟩ Source code in `pydantic_ai_slim/pydantic_ai/models/openrouter.py`

```
140 class OpenRouterProviderConfig(TypedDict, total=False):
141     """Represents the 'Provider' object from the OpenRouter API."""
142
143     order: list[OpenRouterProviderName]
144     """List of provider slugs to try in order (e.g. ["anthropic", "openai"]).
145     [See details](https://openrouter.ai/docs/features/provider-routing#ordering-
146     specific-providers)"""
147
148     allow_fallbacks: bool
149     """Whether to allow backup providers when the primary is unavailable.
150     [See details](https://openrouter.ai/docs/features/provider-routing#disabling-
151     fallbacks)"""
152
153     require_parameters: bool
154     """Only use providers that support all parameters in your request."""
155
156     data_collection: Literal['allow', 'deny']
157     """Control whether to use providers that may store data. [See details]
158     (https://openrouter.ai/docs/features/provider-routing#requiring-providers-to-
159     comply-with-data-policies)"""
160
161     zdr: bool
162     """Restrict routing to only ZDR (Zero Data Retention) endpoints. [See
163     details](https://openrouter.ai/docs/features/provider-routing#zero-data-
164     retention-enforcement)"""
165
166     only: list[OpenRouterProviderName]
167     """List of provider slugs to allow for this request. [See details]
168     (https://openrouter.ai/docs/features/provider-routing#allowing-only-specific-
169     providers)"""
170
171     ignore: list[str]
172     """List of provider slugs to skip for this request. [See details]
173     (https://openrouter.ai/docs/features/provider-routing#ignoring-providers)"""
174
175     quantizations: list[Literal['int4', 'int8', 'fp4', 'fp6', 'fp8', 'fp16',
176     'bf16', 'fp32', 'unknown']]
177     """List of quantization levels to filter by (e.g. ["int4", "int8"]). [See
178     details](https://openrouter.ai/docs/features/provider-
179     routing#quantization)"""
180
181     sort: Literal['price', 'throughput', 'latency']
182     """Sort providers by price or throughput. (e.g. "price" or "throughput").
183     [See details](https://openrouter.ai/docs/features/provider-routing#provider-
184     sorting)"""
185
186     max_price: _OpenRouterMaxPrice
187     """The maximum pricing you want to pay for this request. [See details]
188     (https://openrouter.ai/docs/features/provider-routing#max-price)"""
```

order instance-attribute

```
order: list[OpenRouterProviderName]
```



List of provider slugs to try in order (e.g. ["anthropic", "openai"]). [See details](#)

allow_fallbacks instance-attribute

```
allow_fallbacks: bool
```



Whether to allow backup providers when the primary is unavailable. [See details](#)

require_parameters instance-attribute

```
require_parameters: bool
```



Only use providers that support all parameters in your request.

data_collection instance-attribute

```
data_collection: Literal['allow', 'deny']
```



Control whether to use providers that may store data. [See details](#)

zdr instance-attribute

```
zdr: bool
```



Restrict routing to only ZDR (Zero Data Retention) endpoints. [See details](#)

only instance-attribute

```
only: list[OpenRouterProviderName]
```



List of provider slugs to allow for this request. [See details](#)

ignore instance-attribute

```
ignore: list[str]
```



List of provider slugs to skip for this request. [See details](#)

quantizations instance-attribute

```
quantizations: list[
    Literal[
        "int4",

```



```
        "int8",
        "fp4",
        "fp6",
        "fp8",
        "fp16",
        "bf16",
        "fp32",
        "unknown",
    ]
]
```

List of quantization levels to filter by (e.g. `["int4", "int8"]`). [See details](#)

sort instance-attribute

```
sort: Literal['price', 'throughput', 'latency']
```



Sort providers by price or throughput. (e.g. `"price"` or `"throughput"`). [See details](#)

max_price instance-attribute

```
max_price: _OpenRouterMaxPrice
```



The maximum pricing you want to pay for this request. [See details](#)

OpenRouterReasoning

Bases: [TypedDict](#)

Configuration for reasoning tokens in OpenRouter requests.

Reasoning tokens allow models to show their step-by-step thinking process. You can configure this using either OpenAI-style effort levels or Anthropic-style token limits, but not both simultaneously.

< > Source code in `pydantic_ai_slim/pydantic_ai/models/openrouter.py`

```
174     class OpenRouterReasoning(TypedDict, total=False):
175         """Configuration for reasoning tokens in OpenRouter requests.
176
177         Reasoning tokens allow models to show their step-by-step thinking
178         process.
179         You can configure this using either OpenAI-style effort levels or
180         Anthropic-style
181         token limits, but not both simultaneously.
182         """
183
184         effort: Literal['high', 'medium', 'low']
185         """OpenAI-style reasoning effort level. Cannot be used with
186         max_tokens."""
187
188         max_tokens: int
189         """Anthropic-style specific token limit for reasoning. Cannot be used
190         with effort."""
191
192         exclude: bool
193         """Whether to exclude reasoning tokens from the response. Default is
194         False. All models support this."""
195
196         enabled: bool
197         """Whether to enable reasoning with default parameters. Default is
198         inferred from effort or max_tokens."""
199
```

effort instance-attribute

```
effort: Literal['high', 'medium', 'low']
```

OpenAI-style reasoning effort level. Cannot be used with max_tokens.

max_tokens instance-attribute

```
max_tokens: int
```

Anthropic-style specific token limit for reasoning. Cannot be used with effort.

exclude instance-attribute

```
exclude: bool
```

Whether to exclude reasoning tokens from the response. Default is False. All models support this.

enabled instance-attribute

```
enabled: bool
```

Whether to enable reasoning with default parameters. Default is inferred from effort or max_tokens.

OpenRouterUsageConfig

Bases: [TypedDict](#)

Configuration for OpenRouter usage.

⟨ ⟩ [Source code in pydantic_ai_slim/pydantic_ai/models/openrouter.py](#)

```
195  class OpenRouterUsageConfig(TypedDict, total=False):
196      """Configuration for OpenRouter usage."""
197
198      include: bool
```

OpenRouterModelSettings

Bases: [ModelSettings](#)

Settings used for an OpenRouter model request.

[Source code in pydantic_ai_slim/pydantic_ai/models/openrouter.py](#)

```
201 class OpenRouterModelSettings(ModelSettings, total=False):
202     """Settings used for an OpenRouter model request."""
203
204     # ALL FIELDS MUST BE `openrouter_` PREFIXED SO YOU CAN MERGE THEM WITH
205     # OTHER MODELS.
206
207     openrouter_models: list[str]
208     """A list of fallback models.
209
210     These models will be tried, in order, if the main model returns an error.
211     [See details](https://openrouter.ai/docs/features/model-routing#the-models-
212     parameter)
213     """
214
215     openrouter_provider: OpenRouterProviderConfig
216     """OpenRouter routes requests to the best available providers for your
217     model. By default, requests are load balanced across the top providers to
218     maximize uptime.
219
220     You can customize how your requests are routed using the provider object.
221     [See more](https://openrouter.ai/docs/features/provider-routing)"""
222
223     openrouter_preset: str
224     """Presets allow you to separate your LLM configuration from your code.
225
226     Create and manage presets through the OpenRouter web application to
227     control provider routing, model selection, system prompts, and other
228     parameters, then reference them in OpenRouter API requests. [See more]
229     (https://openrouter.ai/docs/features/presets)"""
230
231     openrouter_transforms: list[OpenRouterTransforms]
232     """To help with prompts that exceed the maximum context size of a model.
233
234     Transforms work by removing or truncating messages from the middle of the
235     prompt, until the prompt fits within the model's context window. [See more]
236     (https://openrouter.ai/docs/features/message-transforms)
237     """
238
239     openrouter_reasoning: OpenRouterReasoning
240     """To control the reasoning tokens in the request.
241
242     The reasoning config object consolidates settings for controlling
243     reasoning strength across different models. [See more]
244     (https://openrouter.ai/docs/use-cases/reasoning-tokens)
245     """
246
247     openrouter_usage: OpenRouterUsageConfig
248     """To control the usage of the model.
249
250     The usage config object consolidates settings for enabling detailed usage
251     information. [See more](https://openrouter.ai/docs/use-cases/usage-
252     accounting)
```

openrouter_models instance-attribute

```
openrouter_models: list[str]
```



A list of fallback models.

These models will be tried, in order, if the main model returns an error. [See details](#)

openrouter_provider instance-attribute

```
openrouter_provider: OpenRouterProviderConfig
```



OpenRouter routes requests to the best available providers for your model. By default, requests are load balanced across the top providers to maximize uptime.

You can customize how your requests are routed using the provider object. [See more](#)

openrouter_preset instance-attribute

```
openrouter_preset: str
```



Presets allow you to separate your LLM configuration from your code.

Create and manage presets through the OpenRouter web application to control provider routing, model selection, system prompts, and other parameters, then reference them in OpenRouter API requests. [See more](#)

openrouter_transforms instance-attribute

```
openrouter_transforms: list[OpenRouterTransforms]
```



To help with prompts that exceed the maximum context size of a model.

Transforms work by removing or truncating messages from the middle of the prompt, until the prompt fits within the model's context window. [See more](#)

openrouter_reasoning instance-attribute

```
openrouter_reasoning: OpenRouterReasoning
```



To control the reasoning tokens in the request.

The reasoning config object consolidates settings for controlling reasoning strength across different models. [See more](#)

openrouter_usage instance-attribute

```
openrouter_usage: OpenRouterUsageConfig
```



To control the usage of the model.

The usage config object consolidates settings for enabling detailed usage information. [See more](#)

OpenRouterModel

Bases: [OpenAIChatModel](#)

Extends OpenAIModel to capture extra metadata for Openrouter.

< > Source code in `pydantic_ai_slim/pydantic_ai/models/openrouter.py`

```
519     class OpenRouterModel(OpenAIChatModel):
520         """Extends OpenAIModel to capture extra metadata for Openrouter."""
521
522         def __init__(
523             self,
524             model_name: str,
525             *,
526             provider: Literal['openrouter'] | Provider[AsyncOpenAI] =
527             'openrouter',
528             profile: ModelProfileSpec | None = None,
529             settings: ModelSettings | None = None,
530         ):
531             """Initialize an OpenRouter model.
532
533             Args:
534                 model_name: The name of the model to use.
535                 provider: The provider to use for authentication and API access.
536                 If not provided, a new provider will be created with the default settings.
537                 profile: The model profile to use. Defaults to a profile picked
538                 by the provider based on the model name.
539                 settings: Model-specific settings that will be used as defaults
540                 for this model.
541                 ...
542
543                 super().__init__(model_name, provider=provider or
544                 OpenRouterProvider(), profile=profile, settings=settings)
544
545             @override
546             def prepare_request(
547                 self,
548                 model_settings: ModelSettings | None,
549                 model_request_parameters: ModelRequestParameters,
550             ) -> tuple[ModelSettings | None, ModelRequestParameters]:
551                 merged_settings, customized_parameters =
552                 super().prepare_request(model_settings, model_request_parameters)
553                 new_settings =
554                 _openrouter_settings_to_openai_settings(cast(OpenRouterModelSettings,
555                 merged_settings or {}))
556                 return new_settings, customized_parameters
557
558             @override
559             def _validate_completion(self, response: chat.ChatCompletion) ->
560             _OpenRouterChatCompletion:
561                 response =
562                 _OpenRouterChatCompletion.model_validate(response.model_dump())
563
564                 if error := response.error:
565                     raise ModelHTTPError(status_code=error.code,
566                     model_name=response.model, body=error.message)
567
568                 return response
569
570             @override
571             def _process_thinking(self, message: chat.ChatCompletionMessage) ->
572             list[ThinkingPart] | None:
573                 assert isinstance(message, _OpenRouterCompletionMessage)
574
575                 if reasoning_details := message.reasoning_details:
```

```

576         return [_from_reasoning_detail(detail) for detail in
577 reasoning_details]
578     else:
579         return super().__process_thinking(message)
580
581     @override
582     def _process_provider_details(self, response: chat.ChatCompletion) ->
583     dict[str, Any] | None:
584         assert isinstance(response, _OpenRouterChatCompletion)
585
586         provider_details = super().__process_provider_details(response) or {}
587         provider_details.update(_map_openrouter_provider_details(response))
588         return provider_details or None
589
590     @dataclass
591     class _MapModelResponseContext(OpenAIChatModel._MapModelResponseContext):
592 # type: ignore[reportPrivateUsage]
593         reasoning_details: list[dict[str, Any]] =
594         field(default_factory=list[dict[str, Any]])
595
596         def _into_message_param(self) ->
597             chat.ChatCompletionAssistantMessageParam:
598                 message_param = super().__into_message_param()
599                 if self.reasoning_details:
600                     message_param['reasoning_details'] = self.reasoning_details
601 # type: ignore[reportGeneralTypeIssues]
602                 return message_param
603
604         @override
605         def _map_response_thinking_part(self, item: ThinkingPart) -> None:
606             assert isinstance(self._model, OpenRouterModel)
607             if item.provider_name == self._model.system:
608                 if reasoning_detail := _into_reasoning_detail(item): #
pragma: lax no cover
609
610                 self.reasoning_details.append(reasoning_detail.model_dump())
611             else: # pragma: lax no cover
612                 super().__map_response_thinking_part(item)
613
614         @property
615         @override
616         def _streamed_response_cls(self):
617             return OpenRouterStreamedResponse
618
619         @override
620         def _map_finish_reason( # type: ignore[reportIncompatibleMethodOverride]
621             self, key: Literal['stop', 'length', 'tool_calls', 'content_filter',
622             'error']
623         ) -> FinishReason | None:
624             return _CHAT_FINISH_REASON_MAP.get(key)

```

__init__

```

__init__(
    model_name: str,
    *,
    provider: (
        Literal["openrouter"] | Provider[AsyncOpenAI]
    ) = "openrouter",

```

```
        profile: ModelProfileSpec | None = None,  
        settings: ModelSettings | None = None  
    )
```

Initialize an OpenRouter model.

Parameters:

Name	Type	Description	Default
model_name	str	The name of the model to use.	<i>required</i>
provider	Literal['openrouter'] Provider[AsyncOpenAI]	The provider to use for authentication and API access. If not provided, a new provider will be created with the default settings.	'openrouter'
profile	ModelProfileSpec None	The model profile to use. Defaults to a profile picked by the provider based on the model name.	None
settings	ModelSettings None	Model-specific settings that will be used as defaults for this model.	None

< > Source code in `pydantic_ai_slim/pydantic_ai/models/openrouter.py`

```
522     def __init__(  
523         self,  
524         model_name: str,  
525         *,  
526         provider: Literal['openrouter'] | Provider[AsyncOpenAI] = 'openrouter',  
527         profile: ModelProfileSpec | None = None,  
528         settings: ModelSettings | None = None,  
529     ):  
530         """Initialize an OpenRouter model.  
531  
532         Args:  
533             model_name: The name of the model to use.  
534             provider: The provider to use for authentication and API access. If  
535             not provided, a new provider will be created with the default settings.  
536             profile: The model profile to use. Defaults to a profile picked by  
537             the provider based on the model name.  
538             settings: Model-specific settings that will be used as defaults for  
539             this model.  
540         """  
541         super().__init__(model_name, provider=provider or OpenRouterProvider(),  
542                         profile=profile, settings=settings)
```

`OpenRouterStreamedResponse` `dataclass`

Bases: `OpenAISTreamedResponse`

Implementation of `StreamedResponse` for OpenRouter models.

< > Source code in `pydantic_ai_slim/pydantic_ai/models/openrouter.py`

```
649     @dataclass
650     class OpenRouterStreamedResponse(OpenAIStreamedResponse):
651         """Implementation of `StreamedResponse` for OpenRouter models."""
652
653         @override
654         async def _validate_response(self):
655             try:
656                 async for chunk in self._response:
657                     yield
658             _OpenRouterChatCompletionChunk.model_validate(chunk.model_dump())
659             except APIError as e:
660                 error = _OpenRouterError.model_validate(e.body)
661                 raise ModelHTTPError(status_code=error.code,
662 model_name=self._model_name, body=error.message)
663
664         @override
665         def _map_thinking_delta(self, choice: chat_completion_chunk.Choice) ->
666 Iterable[ModelResponseStreamEvent]:
667             assert isinstance(choice, _OpenRouterChunkChoice)
668
669             if reasoning_details := choice.delta.reasoning_details:
670                 for i, detail in enumerate(reasoning_details):
671                     thinking_part = _from_reasoning_detail(detail)
672                     # Use unique vendor_part_id for each reasoning detail type to
673 prevent
674                     # different detail types (e.g., reasoning.text,
675 reasoning.encrypted)
676                     # from being incorrectly merged into a single ThinkingPart.
677                     # This is required for Gemini 3 Pro which returns multiple
678 reasoning
679                     # detail types that must be preserved separately for
680 thought_signature handling.
681                     vendor_id = f'reasoning_detail_{detail.type}_{i}'
682                     yield from self._parts_manager.handle_thinking_delta(
683                         vendor_part_id=vendor_id,
684                         id=thinking_part.id,
685                         content=thinking_part.content,
686                         signature=thinking_part.signature,
687                         provider_name=self._provider_name,
688                         provider_details=thinking_part.provider_details,
689                     )
690             else:
691                 return super()._map_thinking_delta(choice)
692
693         @override
694         def _map_provider_details(self, chunk: chat.ChatCompletionChunk) ->
695 dict[str, Any] | None:
696             assert isinstance(chunk, _OpenRouterChatCompletionChunk)
697
698             provider_details = super()._map_provider_details(chunk) or {}
699             provider_details.update(_map_openrouter_provider_details(chunk))
700             return provider_details or None
701
702         @override
703         def _map_finish_reason( # type: ignore[reportIncompatibleMethodOverride]
704             self, key: Literal['stop', 'length', 'tool_calls', 'content_filter',
705 'error']
706         )
```

```
) -> FinishReason | None:  
    return _CHAT_FINISH_REASON_MAP.get(key)
```