

# Package ‘tidyllm’

May 8, 2026

**Title** Tidy Integration of Large Language Models

**Version** 0.5.0

**Description** A tidy interface for integrating large language model (LLM) APIs such as 'Claude', 'OpenAI', 'Gemini', 'Mistral', and local models via 'Ollama' into R workflows. The package supports text, image, audio, video, and document interactions; a unified media interface for attaching inline files or uploading to provider file stores; batch request APIs for cost-efficient large-scale processing; and a pipeline-oriented interface for seamless integration into data workflows. Web services are available at <https://www.anthropic.com>, <https://openai.com>, <https://aistudio.google.com/>, <https://mistral.ai/> and <https://ollama.com>.

**License** MIT + file LICENSE

**URL** <https://edubrueell.github.io/tidyllm/>

**BugReports** <https://github.com/edubrueell/tidyllm/issues>

**Encoding** UTF-8

**RoxygenNote** 7.3.3

**VignetteBuilder** knitr

**Suggests** knitr, rmarkdown, testthat (>= 3.0.0), tidyverse, httptest2, httpuv, ellmer

**Imports** S7 (>= 0.2.0), base64enc, glue, jsonlite, curl, httr2 (>= 1.1.1), lubridate, purrr, rlang, stringr, grDevices, pdftools, tibble, cli, png, lifecycle

**Collate** 'tidyllm-package.R' 'utilites.R' 'embedding\_helpers.R' 'LLMMessage.R' 'APIProvider.R' 'llm\_message.R' 'llm\_verbs.R' 'media.R' 'message\_retrieval.R' 'perform\_api\_requests.R' 'rate\_limits.R' 'tidyllm\_schema.R' 'tools.R' 'api\_chat\_completions.R' 'api\_openai.R' 'api\_azure\_openai.R' 'api\_groq.R' 'api\_mistral.R' 'api\_perplexity.R' 'api\_deepseek.R' 'api\_voyage.R' 'api\_ollama.R' 'api\_claude.R' 'api\_gemini.R' 'api\_ellmer.R' 'api\_openrouter.R' 'api\_llamacpp.R' 'pdfbatch.R' 'zzz.R'

**Depends** R (>= 4.2.0)

**Config/testthat/edition** 3

**NeedsCompilation** no

**Author** Eduard Brüll [aut, cre],  
Jia Zhang [ctb]

**Maintainer** Eduard Brüll <eduard.brue11@zew.de>

**Repository** CRAN

**Date/Publication** 2026-04-30 08:30:08 UTC

## Contents

audio_file . . . . .	5
azure_openai . . . . .	5
azure_openai_chat . . . . .	6
azure_openai_embedding . . . . .	8
cancel_openai_batch . . . . .	9
chat . . . . .	9
chat_completions_chat . . . . .	11
chat_ellmer . . . . .	12
check_azure_openai_batch . . . . .	13
check_batch . . . . .	14
check_claude_batch . . . . .	14
check_gemini_batch . . . . .	15
check_groq_batch . . . . .	16
check_job . . . . .	17
check_mistral_batch . . . . .	17
check_openai_batch . . . . .	18
claude . . . . .	19
claude_chat . . . . .	19
claude_delete_file . . . . .	21
claude_file_metadata . . . . .	22
claude_list_files . . . . .	22
claude_list_models . . . . .	23
claude_upload_file . . . . .	24
claude_websearch . . . . .	24
deepseek . . . . .	25
deepseek_chat . . . . .	25
deep_research . . . . .	27
delete_file . . . . .	28
df_llm_message . . . . .	28
ellmer . . . . .	29
ellmer_tool . . . . .	29
embed . . . . .	31
fetch_azure_openai_batch . . . . .	32
fetch_batch . . . . .	33
fetch_claude_batch . . . . .	34
fetch_gemini_batch . . . . .	35
fetch_groq_batch . . . . .	35

- fetch\_job . . . . . 36
- fetch\_mistral\_batch . . . . . 37
- fetch\_openai\_batch . . . . . 37
- field\_chr . . . . . 38
- field\_object . . . . . 39
- file\_info . . . . . 40
- gemini . . . . . 40
- gemini\_chat . . . . . 41
- gemini\_delete\_file . . . . . 43
- gemini\_embedding . . . . . 43
- gemini\_file\_metadata . . . . . 44
- gemini\_list\_files . . . . . 44
- gemini\_list\_models . . . . . 45
- gemini\_upload\_file . . . . . 45
- get\_logprobs . . . . . 46
- get\_metadata . . . . . 47
- get\_reply . . . . . 48
- get\_reply\_data . . . . . 48
- get\_user\_message . . . . . 49
- groq . . . . . 50
- groq\_chat . . . . . 50
- groq\_list\_models . . . . . 52
- groq\_transcribe . . . . . 53
- img . . . . . 54
- list\_azure\_openai\_batches . . . . . 55
- list\_batches . . . . . 55
- list\_claude\_batches . . . . . 56
- list\_files . . . . . 57
- list\_gemini\_batches . . . . . 57
- list\_groq\_batches . . . . . 58
- list\_hf\_gguf\_files . . . . . 58
- list\_mistral\_batches . . . . . 59
- list\_models . . . . . 60
- list\_openai\_batches . . . . . 60
- llamacpp . . . . . 61
- llamacpp\_chat . . . . . 61
- llamacpp\_delete\_model . . . . . 63
- llamacpp\_download\_model . . . . . 64
- llamacpp\_embedding . . . . . 64
- llamacpp\_health . . . . . 65
- llamacpp\_list\_local\_models . . . . . 66
- llamacpp\_list\_models . . . . . 66
- llamacpp\_rerank . . . . . 67
- LLMMessage . . . . . 68
- llm\_message . . . . . 69
- mistral . . . . . 70
- mistral\_chat . . . . . 70
- mistral\_embedding . . . . . 72

mistral_list_models . . . . .	73
ollama . . . . .	74
ollama_chat . . . . .	75
ollama_delete_model . . . . .	77
ollama_download_model . . . . .	77
ollama_embedding . . . . .	78
ollama_list_models . . . . .	79
openai . . . . .	79
openai_chat . . . . .	80
openai_check_research . . . . .	81
openai_code_interpreter . . . . .	82
openai_deep_research . . . . .	82
openai_embedding . . . . .	83
openai_fetch_research . . . . .	84
openai_list_models . . . . .	84
openai_websearch . . . . .	85
openrouter . . . . .	85
openrouter_chat . . . . .	86
openrouter_credits . . . . .	88
openrouter_embedding . . . . .	88
openrouter_generation . . . . .	89
openrouter_list_models . . . . .	90
pdf_file . . . . .	90
pdf_page_batch . . . . .	91
perplexity . . . . .	92
perplexity_chat . . . . .	92
perplexity_check_research . . . . .	95
perplexity_deep_research . . . . .	96
perplexity_fetch_research . . . . .	97
rate_limit_info . . . . .	98
send_azure_openai_batch . . . . .	98
send_batch . . . . .	100
send_claude_batch . . . . .	102
send_gemini_batch . . . . .	103
send_groq_batch . . . . .	104
send_mistral_batch . . . . .	106
send_ollama_batch . . . . .	107
send_openai_batch . . . . .	109
tidyLLM_schema . . . . .	111
tidyLLM_tool . . . . .	112
upload_file . . . . .	113
video_file . . . . .	113
voyage . . . . .	114
voyage_embedding . . . . .	114
voyage_rerank . . . . .	116

---

audio_file	<i>Create an Audio Object</i>
------------	-------------------------------

---

**Description**

Stores a reference to a local audio file for use in multimodal messages. The file is not encoded until the message is formatted for a provider.

**Usage**

```
audio_file(.path)
```

**Arguments**

.path	The path to the audio file on disk.
-------	-------------------------------------

**Value**

A tidyllm\_audio object.

---

azure_openai	<i>Azure OpenAI Endpoint Provider Function</i>
--------------	--

---

**Description**

The `azure_openai()` function acts as an interface for interacting with the Azure OpenAI API through main tidyllm verbs.

**Usage**

```
azure_openai(..., .called_from = NULL)
```

**Arguments**

...	Parameters to be passed to the Azure OpenAI API specific function, such as model configuration, input text, or API-specific options.
.called_from	An internal argument that specifies which action (e.g., chat) the function is being invoked from. This argument is automatically managed and should not be modified by the user.

**Details**

`azure_openai()` currently routes messages only to `azure_openai_chat()` when used with `chat()`, `send_batch()`. It dynamically routes requests to OpenAI-specific functions like `azure_openai_chat()` and `azure_openai_embedding()` based on the context of the call.

**Value**

The result of the requested action, depending on the specific function invoked (currently, only an updated LLMMessage object for `azure_openai_chat()`).

---

<code>azure_openai_chat</code>	<i>Send LLM Messages to an Azure OpenAI Chat Completions endpoint</i>
--------------------------------	---

---

**Description**

This function sends a message history to the Azure OpenAI Chat Completions API and returns the assistant's reply.

**Usage**

```
azure_openai_chat(
    .llm,
    .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),
    .deployment = "gpt-4o-mini",
    .api_version = "2024-08-01-preview",
    .max_completion_tokens = NULL,
    .reasoning_effort = NULL,
    .frequency_penalty = NULL,
    .logit_bias = NULL,
    .presence_penalty = NULL,
    .seed = NULL,
    .stop = NULL,
    .stream = FALSE,
    .temperature = NULL,
    .top_p = NULL,
    .timeout = 60,
    .verbose = FALSE,
    .json_schema = NULL,
    .max_tries = 3,
    .dry_run = FALSE,
    .logprobs = NULL,
    .top_logprobs = NULL,
    .tools = NULL,
    .tool_choice = NULL,
    .max_tool_rounds = 10
)
```

**Arguments**

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.endpoint_url</code>	Base URL for the API (default: <code>Sys.getenv("AZURE_ENDPOINT_URL")</code> ).
<code>.deployment</code>	The identifier of the model that is deployed (default: "gpt-4o-mini").

<code>.api_version</code>	Which version of the API is deployed (default: "2024-08-01-preview").
<code>.max_completion_tokens</code>	An upper bound for the number of tokens that can be generated for a completion.
<code>.reasoning_effort</code>	How long should reasoning models reason (can either be "low", "medium" or "high").
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency.
<code>.logit_bias</code>	A named list modifying the likelihood of specified tokens appearing in the completion.
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far.
<code>.seed</code>	If specified, the system will make a best effort to sample deterministically.
<code>.stop</code>	Up to 4 sequences where the API will stop generating further tokens.
<code>.stream</code>	If set to TRUE, the answer will be streamed to console as it comes (default: FALSE).
<code>.temperature</code>	What sampling temperature to use, between 0 and 2. Higher values make the output more random.
<code>.top_p</code>	An alternative to sampling with temperature, called nucleus sampling.
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.verbose</code>	Should additional information be shown after the API call (default: FALSE).
<code>.json_schema</code>	A JSON schema object provided by tidyllm schema or ellmer schemata.
<code>.max_tries</code>	Maximum retries to perform request.
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object (default: FALSE).
<code>.logprobs</code>	If TRUE, get the log probabilities of each output token (default: NULL).
<code>.top_logprobs</code>	If specified, get the top N log probabilities of each output token (0-5, default: NULL).
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.tool_choice</code>	A character string specifying the tool-calling behavior; valid values are "none", "auto", or "required".
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10).

**Value**

A new LLMMessage object containing the original messages plus the assistant's response.

---

`azure_openai_embedding`*Generate Embeddings Using OpenAI API on Azure*

---

## Description

Generate Embeddings Using OpenAI API on Azure

## Usage

```
azure_openai_embedding(  
  .input,  
  .deployment = "text-embedding-3-small",  
  .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),  
  .api_version = "2023-05-15",  
  .truncate = TRUE,  
  .timeout = 120,  
  .dry_run = FALSE,  
  .max_tries = 3  
)
```

## Arguments

<code>.input</code>	A character vector of texts to embed or an LLMMessageobject
<code>.deployment</code>	The embedding model identifier (default: "text-embedding-3-small").
<code>.endpoint_url</code>	Base URL for the API (default: Sys.getenv("AZURE_ENDPOINT_URL")).
<code>.api_version</code>	What API-Version othe Azure OpenAI API should be used (default: "2023-05-15")
<code>.truncate</code>	Whether to truncate inputs to fit the model's context length (default: TRUE).
<code>.timeout</code>	Timeout for the API request in seconds (default: 120).
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object.
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).

## Value

A tibble with two columns: `input` and `embeddings`. The `input` column contains the texts sent to embed, and the `embeddings` column is a list column where each row contains an embedding vector of the sent input.

---

cancel\_openai\_batch     *Cancel an In-Progress OpenAI Batch*

---

**Description**

This function cancels an in-progress batch created through the OpenAI API. The batch will be moved to a "cancelling" state and, eventually, "cancelled."

**Usage**

```
cancel_openai_batch(.batch_id, .dry_run = FALSE, .max_tries = 3, .timeout = 60)
```

**Arguments**

.batch_id	Character; the unique identifier for the batch to cancel.
.dry_run	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
.max_tries	Integer; maximum number of retries if the request fails (default: 3).
.timeout	Integer; request timeout in seconds (default: 60).

**Value**

A list containing the response from the OpenAI API about the cancellation status.

---

chat     *Chat with a Language Model*

---

**Description**

The chat() function sends a message to a language model via a specified provider and returns the response. It routes the provided LLMMessage object to the appropriate provider-specific chat function, while allowing for the specification of common arguments applicable across different providers.

**Usage**

```
chat(  
  .llm,  
  .provider = getOption("tidyllm_chat_default"),  
  .dry_run = NULL,  
  .stream = NULL,  
  .temperature = NULL,  
  .timeout = NULL,  
  .top_p = NULL,  
  .max_tries = NULL,  
)
```

```

.model = NULL,
.verbose = NULL,
.json_schema = NULL,
.tools = NULL,
.max_tool_rounds = NULL,
.seed = NULL,
.stop = NULL,
.frequency_penalty = NULL,
.presence_penalty = NULL
)

```

### Arguments

<code>.llm</code>	An LLMMessage object containing the message or conversation history to send to the language model.
<code>.provider</code>	A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like <code>openai()</code> , <code>claude()</code> , etc. You can also set a default provider function via the <code>tidyllm_chat_default</code> option.
<code>.dry_run</code>	Logical; if TRUE, simulates the request without sending it to the provider. Useful for testing.
<code>.stream</code>	Logical; if TRUE, streams the response from the provider in real-time.
<code>.temperature</code>	Numeric; controls the randomness of the model's output (0 = deterministic).
<code>.timeout</code>	Numeric; the maximum time (in seconds) to wait for a response.
<code>.top_p</code>	Numeric; nucleus sampling parameter, which limits the sampling to the top cumulative probability $p$ .
<code>.max_tries</code>	Integer; the maximum number of retries for failed requests.
<code>.model</code>	Character; the model identifier to use (e.g., "gpt-4").
<code>.verbose</code>	Logical; if TRUE, prints additional information about the request and response.
<code>.json_schema</code>	List; A JSON schema object as R list to enforce the output structure
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.max_tool_rounds</code>	Integer; the maximum number of tool use iterations for multi-turn tool calling (default varies by provider).
<code>.seed</code>	Integer; sets a random seed for reproducibility.
<code>.stop</code>	Character vector; specifies sequences where the model should stop generating further tokens.
<code>.frequency_penalty</code>	Numeric; adjusts the likelihood of repeating tokens (positive values decrease repetition).
<code>.presence_penalty</code>	Numeric; adjusts the likelihood of introducing new tokens (positive values encourage novelty).

**Details**

The `chat()` function provides a unified interface for interacting with different language model providers. Common arguments such as `.temperature`, `.model`, and `.stream` are supported by most providers and can be passed directly to `chat()`. If a provider does not support a particular argument, an error will be raised.

Advanced provider-specific configurations can be accessed via the provider functions.

**Value**

An updated `LLMMessage` object containing the response from the language model.

**Examples**

```
## Not run:
# Basic usage with OpenAI provider
llm_message("Hello World") |>
  chat(ollama(.ollama_server = "https://my-ollama-server.de"), .model="mixtral")

  chat(mistral, .model="mixtral")

# Use streaming with Claude provider
llm_message("Tell me a story") |>
  chat(claude(), .stream=TRUE)

## End(Not run)
```

---

chat\_completions\_chat *Chat with any OpenAI-Compatible API Endpoint*

---

**Description**

Provider function for interacting with any server that implements the OpenAI Chat Completions wire format (vLLM, LiteLLM, Together, Anyscale, etc.). Supports the `chat()` verb only.

**Usage**

```
chat_completions_chat(
  .llm,
  .api_url,
  .api_key_env_var = NULL,
  .model = "default",
  ...
)

chat_completions(..., .called_from = NULL)
```

**Arguments**

.llm	An LLMMessage object containing the conversation history.
.api_url	Base URL for the API endpoint (required).
.api_key_env_var	Name of the environment variable holding the API key.
.model	The model identifier to use (default: "default").
...	Additional parameters passed to the underlying chat function.
.called_from	Internal routing argument; do not set manually.

**Value**

An updated LLMMessage object with the assistant's response appended.

---

chat_ellmer	<i>Send LLM Messages to Ellmer Chat Object</i>
-------------	--

---

**Description**

This function converts an LLMMessage to the turns ellmer chat object. The ellmer object is cloned and its turn history is cleared internally. This maintains tidyllms stateless style.

**Usage**

```
chat_ellmer(
  .llm,
  .ellmer_chat,
  .stream = FALSE,
  .timeout = 60,
  .verbose = FALSE,
  .dry_run = FALSE,
  .tools = NULL,
  .max_tries = 3
)
```

**Arguments**

.llm	An LLMMessage object containing the conversation history.
.ellmer_chat	An ellmer chat object (e.g., from <code>ellmer::chat_anthropic()</code> ).
.stream	Logical; if TRUE, streams the response piece by piece (default: FALSE).
.timeout	Request timeout in seconds (default: 60).
.verbose	If TRUE, displays additional information about the chat (default: FALSE).
.dry_run	If TRUE, returns the constructed request object without executing it (default: FALSE).
.tools	Dummy for error message if called from chat with tools (default: NULL).
.max_tries	Maximum retries to perform the request (default: 3).

**Value**

A new LLMMesage object containing the original messages plus the assistant's response.

---

check\_azure\_openai\_batch

*Check Batch Processing Status for Azure OpenAI Batch API*

---

**Description**

This function retrieves the processing status and other details of a specified Azure OpenAI batch ID from the Azure OpenAI Batch API.

**Usage**

```
check_azure_openai_batch(
  .llms = NULL,
  .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),
  .batch_id = NULL,
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

.llms	A list of LLMMesage objects.
.endpoint_url	Base URL for the API (default: Sys.getenv("AZURE_ENDPOINT_URL")).
.batch_id	A manually set batch ID.
.dry_run	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
.max_tries	Maximum retries to perform the request (default: 3).
.timeout	Integer specifying the request timeout in seconds (default: 60).

**Value**

A tibble with information about the status of batch processing.

---

check_batch	<i>Check Batch Processing Status</i>
-------------	--------------------------------------

---

### Description

This function retrieves the processing status and other details of a specified batchid or a list of LLMMessage objects with batch attribute. It routes the input to the appropriate provider-specific batch API function.

### Usage

```
check_batch(
  .llms,
  .provider = getOption("tidyllm_cbatch_default"),
  .dry_run = NULL,
  .max_tries = NULL,
  .timeout = NULL
)
```

### Arguments

.llms	A list of LLMMessage objects or a character vector with a batch ID.
.provider	A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like openai(), claude(), etc. You can also set a default provider function via the tidyllm_cbatch_default option.
.dry_run	Logical; if TRUE, returns the prepared request object without executing it
.max_tries	Maximum retries to perform the request
.timeout	Integer specifying the request timeout in seconds

### Value

A tibble with information about the status of batch processing.

---

check_claude_batch	<i>Check Batch Processing Status for Claude API</i>
--------------------	---

---

### Description

This function retrieves the processing status and other details of a specified Claude batch ID from the Claude API.

**Usage**

```

check_claude_batch(
  .llms = NULL,
  .batch_id = NULL,
  .api_url = "https://api.anthropic.com/",
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)

```

**Arguments**

.llms	A list of LLMMessage objects
.batch_id	A manually set batchid
.api_url	Character; base URL of the Claude API (default: "https://api.anthropic.com/").
.dry_run	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
.max_tries	Maximum retries to perform request
.timeout	Integer specifying the request timeout in seconds (default: 60).

**Value**

A tibble with information about the status of batch processing

---

check_gemini_batch	<i>Check the Status of a Gemini Batch Operation</i>
--------------------	---

---

**Description**

Retrieves processing status and metadata for a Gemini batch operation.

**Usage**

```

check_gemini_batch(
  .llms = NULL,
  .batch_id = NULL,
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE
)

```

**Arguments**

<code>.llms</code>	(Optional) List of LLMMessage objects with a "batch_id" attribute (as returned by <code>send_gemini_batch()</code> ).
<code>.batch_id</code>	(Optional) Character string: full batch operation name, e.g. "batches/xyz123". If both <code>.llms</code> and <code>.batch_id</code> are provided, <code>.batch_id</code> is used.
<code>.timeout</code>	Integer. Request timeout in seconds. Default: 60.
<code>.max_tries</code>	Integer. Maximum retry attempts. Default: 3.
<code>.dry_run</code>	Logical. If TRUE, return the request object instead of making the request (for debugging). Default: FALSE.

**Details**

You can supply either the `.batch_id` string (e.g. "batches/xyz...") **or** a list of LLMMessage objects (`.llms`) with a "batch\_id" attribute as returned by `send_gemini_batch()`.

**Value**

A tibble with the operation's metadata, including name, state, creation time, completion time, and done status.

---

<code>check_groq_batch</code>	<i>Check Batch Processing Status for Groq API</i>
-------------------------------	---

---

**Description**

This function retrieves the processing status and other details of a specified Groq batch.

**Usage**

```
check_groq_batch(
  .llms = NULL,
  .batch_id = NULL,
  .api_url = "https://api.groq.com/",
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects with a <code>batch_id</code> attribute.
<code>.batch_id</code>	A character string with the batch ID to check.
<code>.api_url</code>	Character; base URL of the Groq API (default: "https://api.groq.com/").
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform request.
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).

**Value**

A tibble with information about the status of batch processing.

---

check_job	<i>Check the Status of a Batch or Research Job</i>
-----------	--

---

**Description**

check\_job() dispatches to check\_batch() for batch objects or perplexity\_check\_research() / openai\_check\_research() for tidyllm\_research\_job objects.

**Usage**

```
check_job(.job, ...)
```

**Arguments**

.job	An object with a batch_id attribute (from send_batch()) or a tidyllm_research_job (from deep_research(.background = TRUE)).
...	Additional arguments passed to the underlying function.

**Value**

Status information; type depends on .job class.

---

check_mistral_batch	<i>Check Batch Processing Status for Mistral Batch API</i>
---------------------	--

---

**Description**

This function retrieves the processing status and other details of a specified Mistral batch ID from the Mistral Batch API.

**Usage**

```
check_mistral_batch(
  .llms = NULL,
  .batch_id = NULL,
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects.
<code>.batch_id</code>	A manually set batch ID.
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform the request (default: 3).
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).

**Value**

A tibble with information about the status of batch processing.

---

check_openai_batch	<i>Check Batch Processing Status for OpenAI Batch API</i>
--------------------	---

---

**Description**

This function retrieves the processing status and other details of a specified OpenAI batch ID from the OpenAI Batch API.

**Usage**

```
check_openai_batch(
  .llms = NULL,
  .batch_id = NULL,
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects.
<code>.batch_id</code>	A manually set batch ID.
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform the request (default: 3).
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).

**Value**

A tibble with information about the status of batch processing.

---

claude	<i>Provider Function for Claude models on the Anthropic API</i>
--------	---

---

### Description

The `claude()` function acts as an interface for interacting with the Anthropic API through main `tidymlm` verbs such as `chat()`, `embed()`, and `send_batch()`. It dynamically routes requests to Claude-specific functions like `claude_chat()` and `send_claude_batch()` based on the context of the call.

### Usage

```
claude(..., .called_from = NULL)
```

### Arguments

<code>...</code>	Parameters to be passed to the appropriate OpenAI-specific function, such as model configuration, input text, or API-specific options.
<code>.called_from</code>	An internal argument that specifies which action (e.g., <code>chat</code> , <code>send_batch</code> ) the function is being invoked from. This argument is automatically managed and should not be modified by the user.

### Value

The result of the requested action, depending on the specific function invoked (e.g., an updated `LLMMessage` object for `chat()`, or a matrix for `embed()`).

---

<code>claude_chat</code>	<i>Interact with Claude AI models via the Anthropic API</i>
--------------------------	---

---

### Description

Interact with Claude AI models via the Anthropic API

### Usage

```
claude_chat(  
  .llm,  
  .model = "claude-sonnet-4-6",  
  .max_tokens = 2048,  
  .temperature = NULL,  
  .top_k = NULL,  
  .top_p = NULL,  
  .metadata = NULL,  
  .stop_sequences = NULL,
```

```

    .tools = NULL,
    .json_schema = NULL,
    .file_ids = NULL,
    .api_url = "https://api.anthropic.com/",
    .verbose = FALSE,
    .max_tries = 3,
    .timeout = 60,
    .stream = FALSE,
    .dry_run = FALSE,
    .thinking = FALSE,
    .thinking_budget = 1024,
    .max_tool_rounds = 10
)

```

### Arguments

<code>.llm</code>	An LLMMessage object containing the conversation history and system prompt.
<code>.model</code>	Character string specifying the Claude model version (default: "claude-sonnet-4-6").
<code>.max_tokens</code>	Integer specifying the maximum number of tokens in the response (default: 1024).
<code>.temperature</code>	Numeric between 0 and 1 controlling response randomness.
<code>.top_k</code>	Integer controlling diversity by limiting the top K tokens.
<code>.top_p</code>	Numeric between 0 and 1 for nucleus sampling.
<code>.metadata</code>	List of additional metadata to include with the request.
<code>.stop_sequences</code>	Character vector of sequences that will halt response generation.
<code>.tools</code>	List of additional tools or functions the model can use.
<code>.json_schema</code>	A schema to enforce an output structure
<code>.file_ids</code>	Character; A vector of file IDs for files that were uploaded to Anthropic's Servers
<code>.api_url</code>	Base URL for the Anthropic API (default: "https://api.anthropic.com/").
<code>.verbose</code>	Logical; if TRUE, displays additional information about the API call (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform request
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).
<code>.stream</code>	Logical; if TRUE, streams the response piece by piece (default: FALSE).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.thinking</code>	Logical; if TRUE, enables Claude's thinking mode for complex reasoning tasks (default: FALSE).
<code>.thinking_budget</code>	Integer specifying the maximum tokens Claude can spend on thinking (default: 1024). Must be at least 1024.
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.

**Value**

A new LLMMessage object containing the original messages plus Claude's response.

**Examples**

```
## Not run:
# Basic usage
msg <- llm_message("What is R programming?")
result <- claude_chat(msg)

# With custom parameters
result2 <- claude_chat(msg,
  .temperature = 0.7,
  .max_tokens = 1000)

## End(Not run)
```

---

claude\_delete\_file      *Delete a File from Claude API*

---

**Description**

Deletes a specific file from the Claude API using its file ID.

**Usage**

```
claude_delete_file(
  .file_id,
  .api_url = "https://api.anthropic.com/",
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE
)
```

**Arguments**

<code>.file_id</code>	The file ID to delete.
<code>.api_url</code>	Base URL for the Claude API (default: "https://api.anthropic.com/").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.

**Value**

Invisibly returns NULL. Prints a confirmation message upon successful deletion.

---

claude\_file\_metadata *Retrieve Metadata for a File from Claude API*

---

### Description

Retrieves metadata for a specific file uploaded to the Claude API.

### Usage

```
claude_file_metadata(  
  .file_id,  
  .api_url = "https://api.anthropic.com/",  
  .timeout = 60,  
  .max_tries = 3,  
  .dry_run = FALSE  
)
```

### Arguments

.file_id	The file ID to retrieve metadata for.
.api_url	Base URL for the Claude API (default: "https://api.anthropic.com/").
.timeout	Request timeout in seconds (default: 60).
.max_tries	Maximum retry attempts for requests (default: 3).
.dry_run	Logical; if TRUE, returns the prepared request object without executing it.

### Value

A tibble containing metadata fields such as file\_id, filename, size, and MIME type.

---

claude\_list\_files *List Files in Claude API*

---

### Description

Lists metadata for files uploaded to the Claude API, supporting pagination.

### Usage

```
claude_list_files(  
  .limit = 20,  
  .order = "desc",  
  .api_url = "https://api.anthropic.com/",  
  .timeout = 60,  
  .max_tries = 3,  
  .dry_run = FALSE  
)
```

**Arguments**

<code>.limit</code>	The maximum number of files to return (default: 20).
<code>.order</code>	Order of results, either "asc" or "desc" (default: "desc").
<code>.api_url</code>	Base URL for the Claude API (default: "https://api.anthropic.com/").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.

**Value**

A tibble containing metadata for each file, including `file_id`, `filename`, `size`, and `MIME type`.

---

<code>claude_list_models</code>	<i>List Available Models from the Anthropic Claude API</i>
---------------------------------	--

---

**Description**

List Available Models from the Anthropic Claude API

**Usage**

```
claude_list_models(
  .api_url = "https://api.anthropic.com",
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE,
  .verbose = FALSE
)
```

**Arguments**

<code>.api_url</code>	Base URL for the API (default: "https://api.anthropic.com").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum number of retries for the API request (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.
<code>.verbose</code>	Logical; if TRUE, prints additional information about the request.

**Value**

A tibble containing model information (columns include `type`, `id`, `display_name`, and `created_at`), or NULL if no models are found.

claude\_upload\_file      *Upload a File to Claude API*

---

### Description

Uploads a file to the Claude API and returns its metadata as a tibble.

### Usage

```
claude_upload_file(  
  .file_path,  
  .api_url = "https://api.anthropic.com/",  
  .timeout = 60,  
  .max_tries = 3,  
  .dry_run = FALSE  
)
```

### Arguments

<code>.file_path</code>	The local file path of the file to upload.
<code>.api_url</code>	Base URL for the Claude API (default: "https://api.anthropic.com/").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.

### Value

A tibble containing metadata about the uploaded file, including its `file_id`, `name`, and `size`.

---

claude\_websearch      *Builtin Claude Web Search Tool*

---

### Description

Returns a TOOL object for Claude's builtin `web_search` tool.

### Usage

```
claude_websearch()
```

---

deepseek	<i>Deepseek Provider Function</i>
----------	-----------------------------------

---

### Description

The `deepseek()` function acts as a provider interface for interacting with the Deepseek API through `tidyLLM`'s `chat()` verb. It dynamically routes requests to deepseek-specific function. At the moment this is only `deepseek_chat()`

### Usage

```
deepseek(..., .called_from = NULL)
```

### Arguments

<code>...</code>	Parameters to be passed to the appropriate Deepseek-specific function, such as model configuration, input text, or API-specific options.
<code>.called_from</code>	An internal argument specifying which action (e.g., <code>chat</code> , <code>embed</code> ) the function is invoked from. This argument is automatically managed by the <code>tidyLLM</code> verbs and should not be modified by the user.

### Value

The result of the requested action, depending on the specific function invoked (e.g., an updated `LLMMessage` object for `chat()`).

---

deepseek_chat	<i>Send LLM Messages to the DeepSeek Chat API</i>
---------------	---

---

### Description

This function sends a message history to the DeepSeek Chat API and returns the assistant's reply. Currently tool calls cause problems on the DeepSeek API

### Usage

```
deepseek_chat(  
  .llm,  
  .model = "deepseek-v4-pro",  
  .thinking = NULL,  
  .max_tokens = 2048,  
  .temperature = NULL,  
  .top_p = NULL,  
  .frequency_penalty = NULL,  
  .presence_penalty = NULL,  
)
```

```

    .stop = NULL,
    .stream = FALSE,
    .logprobs = NULL,
    .top_logprobs = NULL,
    .tools = NULL,
    .tool_choice = NULL,
    .api_url = "https://api.deepseek.com/",
    .timeout = 60,
    .verbose = FALSE,
    .dry_run = FALSE,
    .max_tries = 3,
    .max_tool_rounds = 10
  )

```

### Arguments

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.model</code>	The identifier of the model to use (default: "deepseek-v4-pro"). Use "deepseek-v4-flash" for a faster, cheaper alternative.
<code>.thinking</code>	If TRUE, enables thinking mode via <code>thinking: {type: "enabled"}</code> in the request body and captures the reasoning trace in metadata (default: NULL). Supported by all V4 models; note that <code>temperature</code> , <code>top_p</code> , <code>presence_penalty</code> , and <code>frequency_penalty</code> are ignored when thinking is active.
<code>.max_tokens</code>	The maximum number of tokens that can be generated in the response (default: 2048).
<code>.temperature</code>	Controls the randomness in the model's response. Values between 0 and 2 are allowed (optional).
<code>.top_p</code>	Nucleus sampling parameter that controls the proportion of probability mass considered (optional).
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Penalizes repeated tokens to reduce repetition (optional).
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Encourages new topics by penalizing tokens that have appeared so far (optional).
<code>.stop</code>	One or more sequences where the API will stop generating further tokens (optional).
<code>.stream</code>	Logical; if TRUE, streams the response piece by piece (default: FALSE).
<code>.logprobs</code>	If TRUE, returns log probabilities of each output token (default: FALSE).
<code>.top_logprobs</code>	Number between 0 and 5 specifying the number of top log probabilities to return (optional).
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.tool_choice</code>	A character string specifying the tool-calling behavior; valid values are "none", "auto", or "required" (optional).

<code>.api_url</code>	Base URL for the DeepSeek API (default: "https://api.deepseek.com/").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.verbose</code>	If TRUE, displays additional information after the API call (default: FALSE).
<code>.dry_run</code>	If TRUE, returns the constructed request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform the request (default: 3).
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.

**Value**

A new LLMMessage object containing the original messages plus the assistant's response.

---

deep\_research                      *Run Deep Research via a Provider*

---

**Description**

The `deep_research()` function sends a message to a provider's deep research endpoint. Currently supported: Perplexity (sonar-deep-research via async API).

**Usage**

```
deep_research(.llm, .provider, .background = FALSE, ...)
```

**Arguments**

<code>.llm</code>	An LLMMessage object containing the research question.
<code>.provider</code>	A function or function call specifying the provider (e.g., <code>perplexity()</code> ).
<code>.background</code>	Logical; if TRUE, returns a <code>tidyllm_research_job</code> immediately (default: FALSE).
<code>...</code>	Additional arguments passed to the provider's deep research function.

**Value**

If `.background = FALSE`, an LLMMessage with the research reply. If `.background = TRUE`, a `tidyllm_research_job` for use with `check_job()/fetch_job()`.

---

delete_file	<i>Delete a File from a Provider's File Store</i>
-------------	---

---

**Description**

Delete a File from a Provider's File Store

**Usage**

```
delete_file(.provider, .file_id, ...)
```

**Arguments**

.provider	A provider function call.
.file_id	The file ID string to delete, or a tidyllm_file object returned by upload_file().
...	Additional provider-specific arguments.

**Value**

Invisibly NULL; prints a confirmation message.

---

df_llm_message	<i>Convert a Data Frame to an LLMMMessage Object</i>
----------------	--

---

**Description**

This function converts a data frame into an LLMMMessage object representing a conversation history. The data frame must have specific columns (role and content), with each row representing a message.

**Usage**

```
df_llm_message(.df)
```

**Arguments**

.df	A data frame with at least two rows and columns role and content. The role column should contain "user", "assistant", or "system". The content column should contain the corresponding message text.
-----	--

**Value**

An LLMMMessage object representing the structured conversation.

**See Also**[llm\\_message\(\)](#)Other Message Creation Utilities: [llm\\_message\(\)](#)

---

`ellmer`*Alias for the Ellmer Provider Function*

---

**Description**

The `chat_ellmer()` function acts as a provider interface for interacting with ellmer chat objects through tidyllm's verb interface

**Usage**

```
ellmer(..., .called_from = NULL)
```

**Arguments**

`...` Additional parameters to pass through (for consistency with other providers).  
`.called_from` An internal argument specifying which action (e.g., chat) the function is invoked from.

**Details**

This function allows you to use any ellmer chat object (e.g., from `ellmer::chat_anthropic()`, `ellmer::chat_openai()`, etc.) as a stateless backend for tidyllm. The ellmer object is cloned for each interaction to maintain tidyllm's stateless approach.

**Value**

A provider function that can be used with `tidyllm::chat()`.

---

`ellmer_tool`*Convert an ellmer Tool to a tidyllm TOOL*

---

**Description**

Converts an ellmer `ToolDef` or `ToolBuiltIn` object to a tidyllm `TOOL` object, allowing seamless integration of ellmer-defined tools and builtin provider tools with tidyllm workflows.

**Usage**

```
ellmer_tool(.ellmer_tool)
```

## Arguments

`.ellmer_tool` An ellmer ToolDef object created via `ellmer::tool()`, or a builtin tool like `ellmer::claude_tool_web_search()`

## Details

This function supports two types of ellmer tools:

**Custom ToolDef objects:** Extracts the function, description, and argument schemas from an ellmer tool and converts them to tidyllm's internal representation. Ellmer type objects are automatically converted to tidyllm field descriptors.

**Builtin ToolBuiltin objects:** Converts provider-native builtin tools (like Claude's web search) to tidyllm format. For builtin tools, the tool definition is passed through as-is to the provider's API, which handles the tool execution natively.

## Value

A TOOL class object that can be used with tidyllm `chat()` functions

## Examples

```
## Not run:
library(ellmer)

# Example 1: Custom tool
tool_rnorm <- ellmer::tool(
  rnorm,
  description = "Draw numbers from a random normal distribution",
  arguments = list(
    n = ellmer::type_integer("The number of observations"),
    mean = ellmer::type_number("The mean value"),
    sd = ellmer::type_number("The standard deviation")
  )
)

tidyllm_tool_rnorm <- ellmer_tool(tool_rnorm)

llm_message("Generate 100 random numbers") |>
  chat(openai(), .tools = tidyllm_tool_rnorm)

# Example 2: Builtin tool
web_search <- ellmer_tool(ellmer::claude_tool_web_search())
llm_message("What are the latest AI developments?") |>
  chat(claude(), .tools = web_search)

## End(Not run)
```

---

embed	<i>Generate text embeddings</i>
-------	---------------------------------

---

### Description

The `embed()` function allows you to embed a text via a specified provider. It routes the input to the appropriate provider-specific embedding function.

### Usage

```
embed(
  .input,
  .provider = getOption("tidyllm_embed_default"),
  .model = NULL,
  .truncate = NULL,
  .timeout = NULL,
  .dry_run = NULL,
  .max_tries = NULL
)
```

### Arguments

<code>.input</code>	A character vector of texts <code>v</code> , a list of texts and image objects, or an <code>LLMMessage</code> object
<code>.provider</code>	A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like <code>openai()</code> , <code>ollama()</code> , etc. You can also set a default provider function via the <code>tidyllm_embed_default</code> option.
<code>.model</code>	The embedding model to use
<code>.truncate</code>	Whether to truncate inputs to fit the model's context length
<code>.timeout</code>	Timeout for the API request in seconds
<code>.dry_run</code>	If <code>TRUE</code> , perform a dry run and return the request object.
<code>.max_tries</code>	Maximum retry attempts for requests

### Value

A tibble with two columns: `input` and `embeddings`. The `input` column contains the texts sent to `embed`, and the `embeddings` column is a list column where each row contains an embedding vector of the sent input.

### Examples

```
## Not run:
c("What is the meaning of life, the universe and everything?",
  "How much wood would a woodchuck chuck?",
  "How does the brain work?") |>
```

```

    embed(gemini)

    ## End(Not run)

```

---

```

fetch_azure_openai_batch

```

*Fetch Results for an Azure OpenAI Batch*

---

### Description

This function retrieves the results of a completed Azure OpenAI batch and updates the provided list of LLMMessage objects with the responses. It aligns each response with the original request using the custom\_ids generated in send\_azure\_openai\_batch().

### Usage

```

fetch_azure_openai_batch(
  .llms,
  .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),
  .batch_id = NULL,
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)

```

### Arguments

.llms	A list of LLMMessage objects that were part of the batch.
.endpoint_url	Base URL for the API (default: Sys.getenv("AZURE_ENDPOINT_URL")).
.batch_id	Character; the unique identifier for the batch. By default this is NULL and the function will attempt to use the batch_id attribute from .llms.
.dry_run	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
.max_tries	Integer; maximum number of retries if the request fails (default: 3).
.timeout	Integer; request timeout in seconds (default: 60).

### Value

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

fetch_batch	<i>Fetch Results from a Batch API</i>
-------------	---------------------------------------

---

### Description

This function retrieves the results of a completed batch and updates the provided list of `LLMMessage` objects with the responses. It aligns each response with the original request using the `custom_ids` generated in `send_batch()`.

### Usage

```
fetch_batch(
  .llms,
  .provider = getOption("tidyllm_fbatch_default"),
  .dry_run = NULL,
  .max_tries = NULL,
  .timeout = NULL
)
```

### Arguments

<code>.llms</code>	A list of <code>LLMMessage</code> objects containing conversation histories.
<code>.provider</code>	A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like <code>openai()</code> , <code>claude()</code> , etc. You can also set a default provider function via the <code>tidyllm_fbatch_default</code> option.
<code>.dry_run</code>	Logical; if <code>TRUE</code> , returns the constructed request without executing it
<code>.max_tries</code>	Integer; maximum number of retries if the request fails
<code>.timeout</code>	Integer; request timeout in seconds

### Details

The function routes the input to the appropriate provider-specific batch API function.

### Value

A list of updated `LLMMessage` objects, each with the assistant's response added if successful.

---

fetch\_claude\_batch      *Fetch Results for a Claude Batch*

---

### Description

This function retrieves the results of a completed Claude batch and updates the provided list of LLMMessage objects with the responses. It aligns each response with the original request using the custom\_ids generated in send\_claude\_batch().

### Usage

```
fetch_claude_batch(  
    .llms,  
    .batch_id = NULL,  
    .api_url = "https://api.anthropic.com/",  
    .dry_run = FALSE,  
    .max_tries = 3,  
    .timeout = 60  
)
```

### Arguments

.llms	A list of LLMMessage objects that were part of the batch. The list should have names (custom IDs) set by send_claude_batch() to ensure correct alignment.
.batch_id	Character; the unique identifier for the batch. By default this is NULL and the function will attempt to use the batch_id attribute from .llms.
.api_url	Character; the base URL for the Claude API (default: "https://api.anthropic.com/").
.dry_run	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
.max_tries	Integer; maximum number of retries if the request fails (default: 3).
.timeout	Integer; request timeout in seconds (default: 60).

### Value

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

fetch\_gemini\_batch      *Fetch Results for a Gemini Batch*

---

### Description

Retrieves the results of a completed Gemini batch and updates the provided list of LLMMessage objects with the assistant's responses, matching by original list order.

### Usage

```
fetch_gemini_batch(
    .llms,
    .batch_name = NULL,
    .timeout = 60,
    .max_tries = 3,
    .dry_run = FALSE
)
```

### Arguments

.llms	List of LLMMessage objects (as from send_gemini_batch()), must have a batch_id attribute if .batch_name is not given.
.batch_name	(Optional) Character; batch operation name (e.g. "batches/xyz123"). If not provided, is taken from attr(.llms, "batch_id").
.timeout	Integer; request timeout in seconds (default: 60).
.max_tries	Integer; maximum retry attempts (default: 3).
.dry_run	Logical; if TRUE, returns the GET request object (default: FALSE).

### Value

A list of updated LLMMessage objects with the assistant response appended to each, in the same order.

---

fetch\_groq\_batch      *Fetch Results for a Groq Batch*

---

### Description

This function retrieves the results of a completed Groq batch and updates the provided list of LLMMessage objects with the responses.

**Usage**

```

fetch_groq_batch(
  .llms,
  .batch_id = NULL,
  .api_url = "https://api.groq.com/",
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)

```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects that were part of the batch.
<code>.batch_id</code>	Character; the unique identifier for the batch.
<code>.api_url</code>	Character; the base URL for the Groq API (default: "https://api.groq.com/").
<code>.dry_run</code>	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
<code>.max_tries</code>	Integer; maximum number of retries if the request fails (default: 3).
<code>.timeout</code>	Integer; request timeout in seconds (default: 60).

**Value**

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

fetch_job	<i>Fetch Results from a Batch or Research Job</i>
-----------	---

---

**Description**

fetch\_job() dispatches to fetch\_batch() for batch objects or perplexity\_fetch\_research() / openai\_fetch\_research() for tidyllm\_research\_job objects.

**Usage**

```

fetch_job(.job, .provider = NULL, ...)

```

**Arguments**

<code>.job</code>	An object with a batch_id attribute (from send_batch()) or a tidyllm_research_job (from deep_research(.background = TRUE)).
<code>.provider</code>	A provider function (required for batch jobs, ignored for research jobs).
<code>...</code>	Additional arguments passed to the underlying function.

**Value**

Fetches results; type depends on .job class.

---

fetch\_mistral\_batch     *Fetch Results for an Mistral Batch*

---

### Description

This function retrieves the results of a completed Mistral batch and updates the provided list of LLMMessage objects with the responses. It aligns each response with the original request using the custom\_ids generated in send\_mistral\_batch().

### Usage

```
fetch_mistral_batch(  
    .llms,  
    .batch_id = NULL,  
    .dry_run = FALSE,  
    .max_tries = 3,  
    .timeout = 60  
)
```

### Arguments

.llms	A list of LLMMessage objects that were part of the batch.
.batch_id	Character; the unique identifier for the batch. By default this is NULL and the function will attempt to use the batch_id attribute from .llms.
.dry_run	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
.max_tries	Integer; maximum number of retries if the request fails (default: 3).
.timeout	Integer; request timeout in seconds (default: 60).

### Value

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

fetch\_openai\_batch     *Fetch Results for an OpenAI Batch*

---

### Description

This function retrieves the results of a completed OpenAI batch and updates the provided list of LLMMessage objects with the responses. It aligns each response with the original request using the custom\_ids generated in send\_openai\_batch().

**Usage**

```
fetch_openai_batch(
  .llms,
  .batch_id = NULL,
  .dry_run = FALSE,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

.llms	A list of LLMMessage objects that were part of the batch.
.batch_id	Character; the unique identifier for the batch. By default this is NULL and the function will attempt to use the batch_id attribute from .llms.
.dry_run	Logical; if TRUE, returns the constructed request without executing it (default: FALSE).
.max_tries	Integer; maximum number of retries if the request fails (default: 3).
.timeout	Integer; request timeout in seconds (default: 60).

**Value**

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

field\_chr

*Define Field Descriptors for JSON Schema*


---

**Description**

These functions create field descriptors used in tidyllm\_schema() or field\_object() to define JSON schema fields. They support character, factor, numeric, and logical types.

**Usage**

```
field_chr(.description = character(0), .vector = FALSE)

field_fct(.description = character(0), .levels, .vector = FALSE)

field_dbl(.description = character(0), .vector = FALSE)

field_lgl(.description = character(0), .vector = FALSE)
```

**Arguments**

.description	A character string describing the field (optional).
.vector	A logical value indicating if the field is a vector (default: FALSE).
.levels	A character vector specifying allowable values (for field_fct() only).

**Value**

An S7 tidyllm\_field object representing the field descriptor.

**Examples**

```
field_chr("A common street name")
field_fct("State abbreviation", .levels = c("CA", "TX", "Other"))
field_dbl("House number")
field_lgl("Is residential")
field_dbl("A list of apartment numbers at the address", .vector=TRUE )
```

---

field_object	<i>Define a nested object field</i>
--------------	-------------------------------------

---

**Description**

Define a nested object field

**Usage**

```
field_object(.description = character(0), ..., .vector = FALSE)
```

**Arguments**

.description	A character string describing the field (optional).
...	Named fields to include in the object definition (required).
.vector	A logical value indicating if the field is a list of objects (default: FALSE).

**Value**

An S7 tidyllm\_field object of type "object" containing nested fields.

**Examples**

```
# Define an address object with nested fields
address <- field_object("A mailing address",
  street = field_chr("Street name"),
  city = field_chr("City name"),
  zipcode = field_chr("Postal code")
)

# Create a vector of objects
addresses <- field_object("List of addresses",
  street = field_chr("Street name"),
  city = field_chr("City name"),
  .vector = TRUE
)
```

---

file_info	<i>Get Metadata for a File Stored on a Provider</i>
-----------	---

---

**Description**

Get Metadata for a File Stored on a Provider

**Usage**

```
file_info(.provider, .file_id, ...)
```

**Arguments**

.provider	A provider function call.
.file_id	The file ID string to look up, or a tidyllm_file object returned by upload_file().
...	Additional provider-specific arguments.

**Value**

A single-row tibble of file metadata.

---

gemini	<i>Google Gemini Provider Function</i>
--------	--

---

**Description**

The gemini() function acts as a provider interface for interacting with the Google Gemini API through tidyllm's main verbs such as chat() and embed(). It dynamically routes requests to Gemini-specific functions like gemini\_chat() and gemini\_embedding() based on the context of the call.

**Usage**

```
gemini(..., .called_from = NULL)
```

**Arguments**

...	Parameters to be passed to the appropriate Gemini-specific function, such as model configuration, input text, or API-specific options.
.called_from	An internal argument specifying which action (e.g., chat, embed) the function is invoked from. This argument is automatically managed by the tidyllm verbs and should not be modified by the user.

**Details**

Some functions, such as `gemini_upload_file()` and `gemini_delete_file()`, are specific to Gemini and do not have general verb counterparts.

**Value**

The result of the requested action, depending on the specific function invoked (e.g., an updated `LLMMessage` object for `chat()`).

---

gemini_chat	<i>Send LLMMessage to Gemini API</i>
-------------	--------------------------------------

---

**Description**

Send LLMMessage to Gemini API

**Usage**

```
gemini_chat(  
    .llm,  
    .model = "gemini-2.5-flash",  
    .fileid = NULL,  
    .temperature = NULL,  
    .max_output_tokens = NULL,  
    .top_p = NULL,  
    .top_k = NULL,  
    .grounding_threshold = NULL,  
    .presence_penalty = NULL,  
    .frequency_penalty = NULL,  
    .stop_sequences = NULL,  
    .safety_settings = NULL,  
    .json_schema = NULL,  
    .tools = NULL,  
    .thinking_budget = NULL,  
    .timeout = 120,  
    .dry_run = FALSE,  
    .max_tries = 3,  
    .verbose = FALSE,  
    .stream = FALSE,  
    .max_tool_rounds = 10  
)
```

**Arguments**

<code>.llm</code>	An existing <code>LLMMessage</code> object or an initial text prompt.
<code>.model</code>	The model identifier (default: "gemini-1.5-flash").

<code>.fileid</code>	Optional vector of file IDs uploaded via <code>gemini_upload_file()</code> (default: NULL).
<code>.temperature</code>	Controls randomness in generation (default: NULL, range: 0.0-2.0).
<code>.max_output_tokens</code>	Maximum tokens in the response (default: NULL).
<code>.top_p</code>	Controls nucleus sampling (default: NULL, range: 0.0-1.0).
<code>.top_k</code>	Controls diversity in token selection (default: NULL, range: 0 or more).
<code>.grounding_threshold</code>	A grounding threshold between 0 and 1. With lower grounding thresholds Gemini will use Google to search for relevant information before answering. (default: NULL).
<code>.presence_penalty</code>	Penalizes new tokens (default: NULL, range: -2.0 to 2.0).
<code>.frequency_penalty</code>	Penalizes frequent tokens (default: NULL, range: -2.0 to 2.0).
<code>.stop_sequences</code>	Optional character sequences to stop generation (default: NULL, up to 5).
<code>.safety_settings</code>	A list of safety settings (default: NULL).
<code>.json_schema</code>	A schema to enforce an output structure
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.thinking_budget</code>	Token budget for internal reasoning (default: NULL). Works with <code>gemini-2.5-flash</code> and <code>gemini-2.5-pro</code> .
<code>.timeout</code>	When should our connection time out (default: 120 seconds).
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object.
<code>.max_tries</code>	Maximum retries to perform request (default: 3).
<code>.verbose</code>	Should additional information be shown after the API call.
<code>.stream</code>	Should the response be streamed (default: FALSE).
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.

**Value**

A new `LLMMessage` object containing the original messages plus the assistant's response.

---

gemini\_delete\_file      *Delete a File from Gemini API*

---

**Description**

Deletes a specific file from the Gemini API using its file ID.

**Usage**

```
gemini_delete_file(.file_name)
```

**Arguments**

.file\_name      The file ID (e.g., "files/abc-123") to delete.

**Value**

Invisibly returns NULL. Prints a confirmation message upon successful deletion.

---

gemini\_embedding      *Generate Embeddings Using the Google Gemini API*

---

**Description**

Generate Embeddings Using the Google Gemini API

**Usage**

```
gemini_embedding(  
  .input,  
  .model = "gemini-embedding-2-preview",  
  .truncate = TRUE,  
  .timeout = 120,  
  .dry_run = FALSE,  
  .max_tries = 3  
)
```

**Arguments**

.input      A character vector of texts to embed or an LLMMessage object

.model      The embedding model identifier (default: "text-embedding-3-small").

.truncate    Whether to truncate inputs to fit the model's context length (default: TRUE).

.timeout     Timeout for the API request in seconds (default: 120).

.dry\_run     If TRUE, perform a dry run and return the request object.

.max\_tries   Maximum retry attempts for requests (default: 3).

**Value**

A matrix where each column corresponds to the embedding of a message in the message history.

---

gemini\_file\_metadata    *Retrieve Metadata for a File from Gemini API*

---

**Description**

Retrieves metadata for a specific file uploaded to the Gemini API.

**Usage**

```
gemini_file_metadata(.file_name)
```

**Arguments**

.file\_name    The file ID (e.g., "files/abc-123") to retrieve metadata for.

**Value**

A tibble containing metadata fields such as name, display name, MIME type, size, and URI.

---

gemini\_list\_files    *List Files in Gemini API*

---

**Description**

Lists metadata for files uploaded to the Gemini API, supporting pagination.

**Usage**

```
gemini_list_files(.page_size = 10, .page_token = NULL)
```

**Arguments**

.page\_size    The maximum number of files to return per page (default: 10, maximum: 100).

.page\_token    A token for fetching the next page of results (default: NULL).

**Value**

A tibble containing metadata for each file, including fields such as name, display name, MIME type, and URI.

---

gemini\_list\_models      *List Available Models from the Google Gemini API*

---

**Description**

List Available Models from the Google Gemini API

**Usage**

```
gemini_list_models(.timeout = 60, .max_tries = 3, .dry_run = FALSE)
```

**Arguments**

<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum number of retries for the API request (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.

**Value**

A tibble containing model information with columns including name, base\_model\_id, version, display\_name, description, input\_token\_limit, output\_token\_limit, supported\_generation\_methods, thinking, temperature, max\_temperature, top\_p, and top\_k, or NULL if no models are found.

---

gemini\_upload\_file      *Upload a File to Gemini API*

---

**Description**

Uploads a file to the Gemini API and returns its metadata as a tibble.

**Usage**

```
gemini_upload_file(.file_path)
```

**Arguments**

<code>.file_path</code>	The local file path of the file to upload.
-------------------------	--

**Value**

A tibble containing metadata about the uploaded file, including its name, URI, and MIME type.

---

`get_logprobs`*Retrieve Log Probabilities from Assistant Replies*

---

### Description

Extracts token log probabilities from assistant replies within an `LLMMessage` object. Each row represents a token with its log probability and top alternative tokens.

### Usage

```
get_logprobs(.llm, .index = NULL)
```

### Arguments

<code>.llm</code>	An <code>LLMMessage</code> object containing the message history.
<code>.index</code>	A positive integer specifying which assistant reply's log probabilities to extract. If <code>NULL</code> (default), log probabilities for all replies are returned.

### Details

An empty tibble is output if no logprobs were requested. Works with `openai_chat()`, `llamacpp_chat()`, and other providers that support logprobs.

Columns include:

- `reply_index`: The index of the assistant reply in the message history.
- `token`: The generated token.
- `logprob`: The log probability of the generated token.
- `bytes`: The byte-level encoding of the token.
- `top_logprobs`: A list column containing the top alternative tokens with their log probabilities.

### Value

A tibble containing log probabilities for the specified assistant reply or all replies.

### See Also

[get\\_metadata\(\)](#)

---

get_metadata	<i>Retrieve Metadata from Assistant Replies</i>
--------------	---

---

## Description

Retrieves metadata from assistant replies within an LLMMessage object. It returns the metadata as a tibble.

## Usage

```
get_metadata(.llm, .index = NULL)
```

```
last_metadata(.llm)
```

## Arguments

.llm	An LLMMessage object containing the message history.
.index	A positive integer specifying which assistant reply's metadata to extract. If NULL (default), metadata for all replies is returned.

## Details

Metadata columns may include:

- `model`: The model used for generating the reply.
- `timestamp`: The time when the reply was generated.
- `prompt_tokens`: The number of tokens in the input prompt.
- `completion_tokens`: The number of tokens in the assistant's reply.
- `total_tokens`: The total number of tokens (prompt + completion).
- `api_specific`: A list column with API-specific metadata.

For convenience, [last\\_metadata\(\)](#) is provided to retrieve the metadata for the last message.

## Value

A tibble containing metadata for the specified assistant reply or all replies.

## See Also

[last\\_metadata\(\)](#)

---

get_reply	<i>Retrieve Assistant Reply as Text</i>
-----------	---

---

**Description**

Extracts the plain text content of the assistant's reply from an LLMMessage object. Use [get\\_reply\\_data\(\)](#) for structured replies in JSON format.

**Usage**

```
get_reply(.llm, .index = NULL)
```

```
last_reply(.llm)
```

**Arguments**

.llm	An LLMMessage object containing the message history.
.index	A positive integer indicating the index of the assistant reply to retrieve. Defaults to NULL, which retrieves the last reply.

**Details**

This function is the core utility for retrieving assistant replies by index. For convenience, [last\\_reply\(\)](#) is provided as a wrapper to retrieve the latest assistant reply.

**Value**

Returns a character string containing the assistant's reply, or NA\_character\_ if no reply exists.

**See Also**

[get\\_reply\\_data\(\)](#), [last\\_reply\(\)](#)

---

get_reply_data	<i>Retrieve Assistant Reply as Structured Data</i>
----------------	--

---

**Description**

Parses the assistant's reply as JSON and returns the corresponding structured data. If the reply is not marked as JSON, attempts to extract and parse JSON content from the text.

**Usage**

```
get_reply_data(.llm, .index = NULL)
```

```
last_reply_data(.llm)
```

**Arguments**

- .llm An LLMMessage object containing the message history.
- .index A positive integer indicating the index of the assistant reply to retrieve. Defaults to NULL, which retrieves the last reply.

**Details**

For convenience, [last\\_reply\\_data\(\)](#) is provided as a wrapper to retrieve the latest assistant reply's data.

**Value**

Returns the parsed data from the assistant's reply, or NULL if parsing fails.

**See Also**

[get\\_reply\(\)](#), [last\\_reply\\_data\(\)](#)

---

<code>get_user_message</code>	<i>Retrieve a User Message by Index</i>
-------------------------------	---

---

**Description**

Extracts the content of a user's message from an LLMMessage object at a specific index.

**Usage**

```
get_user_message(.llm, .index = NULL)
```

```
last_user_message(.llm)
```

**Arguments**

- .llm An LLMMessage object.
- .index A positive integer indicating which user message to retrieve. Defaults to NULL, which retrieves the last message.

**Details**

For convenience, [last\\_user\\_message\(\)](#) is provided as a wrapper to retrieve the latest user message without specifying an index.

**Value**

Returns the content of the user's message at the specified index. If no messages are found, returns `NA_character_`.

**See Also**

[last\\_user\\_message\(\)](#)

---

groq

*Groq API Provider Function*

---

**Description**

The `groq()` function acts as an interface for interacting with the Groq API through `tidyLLM`'s main verbs. Currently, Groq only supports `groq_chat()` for chat-based interactions and `groq_transcribe()` for transcription tasks.

**Usage**

```
groq(..., .called_from = NULL)
```

**Arguments**

<code>...</code>	Parameters to be passed to the Groq-specific function, such as model configuration, input text, or API-specific options.
<code>.called_from</code>	An internal argument that specifies which action (e.g., <code>chat</code> ) the function is being invoked from. This argument is automatically managed and should not be modified by the user.

**Details**

Since `groq_transcribe()` is unique to Groq and does not have a general verb counterpart, `groq()` currently routes messages only to `groq_chat()` when used with verbs like `chat()`.

**Value**

The result of the requested action, depending on the specific function invoked (currently, only an updated `LLMMessage` object for `groq_chat()`).

---

groq\_chat

*Send LLM Messages to the Groq Chat API*

---

**Description**

This function sends a message history to the Groq Chat API and returns the assistant's reply.

**Usage**

```

groq_chat(
  .llm,
  .model = "openai/gpt-oss-120b",
  .max_tokens = 1024,
  .temperature = NULL,
  .top_p = NULL,
  .frequency_penalty = NULL,
  .presence_penalty = NULL,
  .stop = NULL,
  .seed = NULL,
  .tools = NULL,
  .tool_choice = NULL,
  .api_url = "https://api.groq.com/",
  .json_schema = NULL,
  .timeout = 60,
  .verbose = FALSE,
  .stream = FALSE,
  .dry_run = FALSE,
  .max_tries = 3,
  .max_tool_rounds = 10
)

```

**Arguments**

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.model</code>	The identifier of the model to use (default: "llama-3.2-11b-vision-preview").
<code>.max_tokens</code>	The maximum number of tokens that can be generated in the response (default: 1024).
<code>.temperature</code>	Controls the randomness in the model's response. Values between 0 and 2 are allowed, where higher values increase randomness (optional).
<code>.top_p</code>	Nucleus sampling parameter that controls the proportion of probability mass considered. Values between 0 and 1 are allowed (optional).
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Positive values penalize repeated tokens, reducing likelihood of repetition (optional).
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Positive values encourage new topics by penalizing tokens that have appeared so far (optional).
<code>.stop</code>	One or more sequences where the API will stop generating further tokens. Can be a string or a list of strings (optional).
<code>.seed</code>	An integer for deterministic sampling. If specified, attempts to return the same result for repeated requests with identical parameters (optional).
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls (optional).

<code>.tool_choice</code>	A character string specifying the tool-calling behavior; valid values are "none", "auto", or "required" (optional).
<code>.api_url</code>	Base URL for the Groq API (default: "https://api.groq.com/").
<code>.json_schema</code>	A list or tidyllm schema created with <code>tidyllm_schema()</code> for structured JSON output (optional).
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.verbose</code>	If TRUE, displays additional information after the API call, including rate limit details (default: FALSE).
<code>.stream</code>	Logical; if TRUE, streams the response piece by piece (default: FALSE).
<code>.dry_run</code>	If TRUE, performs a dry run and returns the constructed request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform request
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.

**Value**

A new `LLMMessage` object containing the original messages plus the assistant's response.

**Examples**

```
## Not run:
# Basic usage
msg <- llm_message("What is Groq?")
result <- groq_chat(msg)

# With custom parameters
result2 <- groq_chat(msg,
  .model = "llama-3.2-vision",
  .temperature = 0.5,
  .max_tokens = 512)

## End(Not run)
```

---

groq\_list\_models

*List Available Models from the Groq API*

---

**Description**

List Available Models from the Groq API

**Usage**

```
groq_list_models(  
  .api_url = "https://api.groq.com",  
  .timeout = 60,  
  .max_tries = 3,  
  .dry_run = FALSE,  
  .verbose = FALSE  
)
```

**Arguments**

<code>.api_url</code>	Base URL for the API (default: "https://api.groq.com").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum number of retries for the API request (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.
<code>.verbose</code>	Logical; if TRUE, prints additional information about the request.

**Value**

A tibble containing model information (columns include `id`, `created`, `owned_by`, and `context_window`), or NULL if no models are found.

---

<code>groq_transcribe</code>	<i>Transcribe an Audio File Using Groq transcription API</i>
------------------------------	--

---

**Description**

This function reads an audio file and sends it to the Groq transcription API for transcription.

**Usage**

```
groq_transcribe(  
  .audio_file,  
  .model = "playai-tts",  
  .language = NULL,  
  .prompt = NULL,  
  .temperature = 0,  
  .api_url = "https://api.groq.com/openai/v1/audio/transcriptions",  
  .dry_run = FALSE,  
  .verbose = FALSE,  
  .max_tries = 3  
)
```

**Arguments**

<code>.audio_file</code>	The path to the audio file (required). Supported formats include flac, mp3, mp4, mpeg, mpga, m4a, ogg, wav, or webm.
<code>.model</code>	The model to use for transcription (default: "whisper-large-v3").
<code>.language</code>	The language of the input audio, in ISO-639-1 format (optional).
<code>.prompt</code>	A prompt to guide the transcription style. It should match the audio language (optional).
<code>.temperature</code>	Sampling temperature, between 0 and 1, with higher values producing more randomness (default: 0).
<code>.api_url</code>	Base URL for the API (default: "https://api.groq.com/openai/v1/audio/transcriptions").
<code>.dry_run</code>	Logical; if TRUE, performs a dry run and returns the request object without making the API call (default: FALSE).
<code>.verbose</code>	Logical; if TRUE, rate limiting info is displayed after the API request (default: FALSE).
<code>.max_tries</code>	Maximum retries to perform request

**Value**

A character vector containing the transcription.

**Examples**

```
## Not run:
# Basic usage
groq_transcribe(.audio_file = "example.mp3")

## End(Not run)
```

---

img

*Create an Image Object*


---

**Description**

This function reads an image file from disk, encodes it in base64, and returns a `tidyllm_image` object that can be used in multimodal embedding requests or attached to messages via `.media`.

**Usage**

```
img(.path)
```

**Arguments**

<code>.path</code>	The path to the image file on disk.
--------------------	-------------------------------------

**Value**

A tidyllm\_image object.

---

list_azure_openai_batches	<i>List Azure OpenAI Batch Requests</i>
---------------------------	---

---

**Description**

Retrieves batch request details from the Azure OpenAI Batch API.

**Usage**

```
list_azure_openai_batches(
  .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),
  .limit = 20,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

.endpoint_url	Base URL for the API (default: Sys.getenv("AZURE_ENDPOINT_URL")).
.limit	Maximum number of batches to retrieve (default: 20).
.max_tries	Maximum retry attempts for requests (default: 3).
.timeout	Request timeout in seconds (default: 60).

**Value**

A tibble with batch details: batch ID, status, creation time, expiration time, and request counts (total, completed, failed).

---

list_batches	<i>List all Batch Requests on a Batch API</i>
--------------	---

---

**Description**

List all Batch Requests on a Batch API

**Usage**

```
list_batches(.provider = getOption("tidymlm_lbatch_default"))
```

**Arguments**

`.provider` A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like `openai()`, `claude()`, etc. You can also set a default provider function via the `tidyllm_lbatch_default` option.

**Value**

A tibble with information about the status of batch processing.

---

`list_claude_batches` *List Claude Batch Requests*

---

**Description**

Retrieves batch request details from the Claude API.

**Usage**

```
list_claude_batches(
  .api_url = "https://api.anthropic.com/",
  .limit = 20,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

`.api_url` Base URL for the Claude API (default: "https://api.anthropic.com/").

`.limit` Maximum number of batches to retrieve (default: 20).

`.max_tries` Maximum retry attempts for requests (default: 3).

`.timeout` Request timeout in seconds (default: 60).

**Value**

A tibble with batch details: batch ID, status, creation time, expiration time, and request counts (succeeded, errored, expired, canceled).

---

list_files	<i>List Files Stored on a Provider</i>
------------	--

---

**Description**

List Files Stored on a Provider

**Usage**

```
list_files(.provider, ...)
```

**Arguments**

.provider	A provider function call.
...	Additional provider-specific arguments.

**Value**

A tibble of file metadata.

---

list_gemini_batches	<i>List Recent Gemini Batch Operations</i>
---------------------	--

---

**Description**

Returns a tibble with recent Gemini batch operations and their metadata.

**Usage**

```
list_gemini_batches(  
  .filter = NULL,  
  .page_size = 20,  
  .timeout = 60,  
  .max_tries = 3,  
  .dry_run = FALSE  
)
```

**Arguments**

.filter	Optional filter expression for batch listing (see Gemini API docs).
.page_size	Integer. Maximum number of results to return. Default: 20.
.timeout	Integer. Request timeout in seconds. Default: 60.
.max_tries	Integer. Maximum retry attempts. Default: 3.
.dry_run	Logical. If TRUE, returns the request object (for debugging). Default: FALSE.

**Value**

A tibble with columns: name, state, done, create\_time, complete\_time.

---

list_groq_batches	<i>List Groq Batch Requests</i>
-------------------	---------------------------------

---

**Description**

Retrieves batch request details from the Groq API.

**Usage**

```
list_groq_batches(
  .api_url = "https://api.groq.com/",
  .limit = 20,
  .max_tries = 3,
  .timeout = 60
)
```

**Arguments**

.api_url	Base URL for the Groq API (default: "https://api.groq.com/").
.limit	Maximum number of batches to retrieve (default: 20).
.max_tries	Maximum retry attempts for requests (default: 3).
.timeout	Request timeout in seconds (default: 60).

**Value**

A tibble with batch details including batch ID, status, creation time, and request counts.

---

list_hf_gguf_files	<i>List GGUF Files Available in a Hugging Face Repository</i>
--------------------	---

---

**Description**

Calls the Hugging Face Hub API to list all .gguf files in a repository, along with their sizes. Useful for choosing a quantization level before calling `llamacpp_download_model()`.

**Usage**

```
list_hf_gguf_files(.repo, .timeout = 30)
```

**Arguments**

.repo	Hugging Face repository ID in "owner/model" format (e.g. "Qwen/Qwen3-8B-GGUF").
.timeout	Request timeout in seconds (default: 30).

**Value**

A tibble with columns filename, size\_gb, and url.

---

list\_mistral\_batches *List Mistral Batch Requests*

---

**Description**

Retrieves batch request details from the OpenAI Batch API.

**Usage**

```
list_mistral_batches(  
  .limit = 100,  
  .max_tries = 3,  
  .timeout = 60,  
  .status = NULL,  
  .created_after = NULL  
)
```

**Arguments**

.limit	Maximum number of batches to retrieve (default: 20).
.max_tries	Maximum retry attempts for requests (default: 3).
.timeout	Request timeout in seconds (default: 60).
.status	Filter by status. (default: NULL)
.created_after	created after a string specifying a date-time (default: NULL)

**Value**

A tibble with batch details for all batches fitting the request

---

list_models	<i>List Available Models for a Provider</i>
-------------	---

---

**Description**

The `list_models()` function retrieves available models from the specified provider.

**Usage**

```
list_models(.provider = getOption("tidyllm_lmodels_default"), ...)
```

**Arguments**

<code>.provider</code>	A function or function call specifying the provider and any additional parameters. You can also set a default provider via the <code>tidyllm_lmodels_default</code> option.
<code>...</code>	Additional arguments to be passed to the provider-specific <code>list_models</code> function.

**Value**

A tibble containing model information.

---

list_openai_batches	<i>List OpenAI Batch Requests</i>
---------------------	-----------------------------------

---

**Description**

Retrieves batch request details from the OpenAI Batch API.

**Usage**

```
list_openai_batches(.limit = 20, .max_tries = 3, .timeout = 60)
```

**Arguments**

<code>.limit</code>	Maximum number of batches to retrieve (default: 20).
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.timeout</code>	Request timeout in seconds (default: 60).

**Value**

A tibble with batch details: batch ID, status, creation time, expiration time, and request counts (total, completed, failed).

---

llamacpp	<i>llama.cpp Provider Function</i>
----------	------------------------------------

---

### Description

The `llamacpp()` provider connects tidyllm to a locally running **llama.cpp** server. It exposes the same verb/provider pattern as every other tidyllm provider while also offering llama.cpp-specific features: BNF grammar constraints (`.grammar`), token logprobs (`.logprobs`), and model management helpers.

The server must be started separately before calling `llamacpp()`. See `llamacpp_health()` to verify the server is running, and `llamacpp_download_model()` / `list_hf_gguf_files()` to obtain models.

### Usage

```
llamacpp(..., .called_from = NULL)
```

### Arguments

<code>...</code>	Parameters passed to the appropriate llama.cpp-specific function.
<code>.called_from</code>	An internal argument specifying which verb invoked this function. Managed automatically by tidyllm verbs; do not set manually.

### Value

The result of the requested action (e.g., an updated `LLMMessage` for `chat()`, a tibble for `embed()` or `list_models()`).

---

llamacpp_chat	<i>Send LLM Messages to a llama.cpp Server</i>
---------------	--

---

### Description

Sends a message history to a local llama.cpp server using the OpenAI-compatible Chat Completions API. Supports BNF grammar constraints (a llama.cpp-specific feature that enforces output format at the token sampling level) and token logprobs for uncertainty quantification.

### Usage

```
llamacpp_chat(
  .llm,
  .model = "local-model",
  .max_tokens = 1024,
  .temperature = NULL,
  .top_p = NULL,
```

```

.stop = NULL,
.stream = FALSE,
.tools = NULL,
.tool_choice = NULL,
.json_schema = NULL,
.grammar = NULL,
.logprobs = FALSE,
.top_logprobs = NULL,
.seed = NULL,
.thinking = NULL,
.thinking_budget = NULL,
.server = Sys.getenv("LLAMACPP_SERVER", "http://localhost:8080"),
.api_key = Sys.getenv("LLAMACPP_API_KEY", ""),
.timeout = 120,
.verbose = FALSE,
.dry_run = FALSE,
.max_tries = 3,
.max_tool_rounds = 10
)

```

### Arguments

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.model</code>	The model name (default: "local-model"). llama.cpp ignores this value and serves whatever GGUF is loaded; the default is a clear placeholder.
<code>.max_tokens</code>	Maximum tokens in the response (default: 1024).
<code>.temperature</code>	Controls randomness (0–2, optional).
<code>.top_p</code>	Nucleus sampling parameter (0–1, optional).
<code>.stop</code>	One or more stop sequences (optional).
<code>.stream</code>	Logical; if TRUE, streams the response (default: FALSE).
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects.
<code>.tool_choice</code>	Tool-calling behavior: "none", "auto", or "required" (optional).
<code>.json_schema</code>	A JSON schema as an R list to enforce structured output. Passed as <code>response_format = {type: "json_s</code> — enforced at the sampler level.
<code>.grammar</code>	A BNF grammar string to constrain sampling to any formal language (e.g. "root ::= [0-9]{4}" for a 4-digit year). Takes precedence over <code>.json_schema</code> when both are set; use one or the other.
<code>.logprobs</code>	Logical; if TRUE, returns token log-probabilities (default: FALSE).
<code>.top_logprobs</code>	Integer; number of top-alternative log-prob entries to return per token (1–20, optional). Requires <code>.logprobs = TRUE</code> .
<code>.seed</code>	Integer; random seed for reproducible outputs (optional).
<code>.thinking</code>	Logical; if TRUE, enables extended reasoning/thinking mode for models that support it (e.g. Qwen3). NULL (default) lets the server decide; FALSE explicitly disables it for faster responses; TRUE forces it on. The thinking trace is accessible via <code>get_metadata(result)\$api_specific\$thinking</code> .

<code>.thinking_budget</code>	Integer; maximum tokens the model may use for thinking when <code>.thinking = TRUE</code> . NULL uses the server default (optional).
<code>.server</code>	Base URL of the llama.cpp server. Defaults to the <code>LLAMACPP_SERVER</code> environment variable, falling back to <code>"http://localhost:8080"</code> .
<code>.api_key</code>	API key for the llama.cpp server. Defaults to the <code>LLAMACPP_API_KEY</code> environment variable. Leave unset if the server was started without <code>--api-key</code> .
<code>.timeout</code>	Request timeout in seconds (default: 120).
<code>.verbose</code>	If TRUE, displays additional information (default: FALSE).
<code>.dry_run</code>	If TRUE, returns the request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries (default: 3).
<code>.max_tool_rounds</code>	Maximum tool use iterations (default: 10).

**Value**

A new `LLMMessage` object containing the original messages plus the assistant's response.

---

`llamacpp_delete_model` *Delete a Local GGUF Model File*

---

**Description**

Deletes a GGUF file from disk. Issues a warning confirming the deletion. Mirrors the `ollama_delete_model()` pattern.

**Usage**

```
llamacpp_delete_model(.path)
```

**Arguments**

<code>.path</code>	Full path to the <code>.gguf</code> file to delete.
--------------------	---

**Value**

Invisibly returns TRUE on success.

---

llamacpp\_download\_model

*Download a GGUF Model from Hugging Face*


---

### Description

Downloads a single GGUF file from a Hugging Face repository with a streaming progress bar. Use `list_hf_gguf_files()` first to browse available quantizations.

### Usage

```
llamacpp_download_model(
    .repo,
    .filename,
    .dir = Sys.getenv("LLAMACPP_MODEL_DIR", "~/models"),
    .timeout = 3600
)
```

### Arguments

<code>.repo</code>	Hugging Face repository ID (e.g. "Qwen/Qwen3-8B-GGUF").
<code>.filename</code>	The exact filename to download (e.g. "Qwen3-8B-Q4_K_M.gguf"). Use <code>list_hf_gguf_files(.repo)</code> to see what is available.
<code>.dir</code>	Destination directory. Defaults to <code>LLAMACPP_MODEL_DIR</code> env var, then " <code>~/models</code> ". Created if it does not exist.
<code>.timeout</code>	Download timeout in seconds (default: 3600).

### Value

Invisibly returns the full path of the downloaded file.

---

llamacpp\_embedding

*Generate Embeddings Using a llama.cpp Server*


---

### Description

Sends text to the `/v1/embeddings` endpoint of a running llama.cpp server and returns embedding vectors.

**Usage**

```
llamacpp_embedding(
  .input,
  .model = "local-model",
  .server = Sys.getenv("LLAMACPP_SERVER", "http://localhost:8080"),
  .api_key = Sys.getenv("LLAMACPP_API_KEY", ""),
  .timeout = 120,
  .dry_run = FALSE,
  .max_tries = 3
)
```

**Arguments**

<code>.input</code>	A character vector of texts to embed, or an LLMMessage object.
<code>.model</code>	The model name (default: "local-model"). llama.cpp ignores this and serves the loaded embedding model.
<code>.server</code>	Base URL of the llama.cpp server. Defaults to the LLAMACPP_SERVER environment variable, falling back to "http://localhost:8080".
<code>.api_key</code>	API key for the server (default: LLAMACPP_API_KEY env var).
<code>.timeout</code>	Request timeout in seconds (default: 120).
<code>.dry_run</code>	If TRUE, returns the request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries (default: 3).

**Value**

A tibble with columns `input` (text) and `embeddings` (list of numeric vectors).

---

llamacpp_health	<i>Check Health of the llama.cpp Server</i>
-----------------	---

---

**Description**

Calls the /health endpoint of a running llama.cpp server. Returns the status string ("ok", "loading model", "no model loaded", or "error") along with the full parsed response body as a named list.

**Usage**

```
llamacpp_health(
  .server = Sys.getenv("LLAMACPP_SERVER", "http://localhost:8080"),
  .timeout = 10
)
```

**Arguments**

<code>.server</code>	Base URL of the llama.cpp server. Defaults to LLAMACPP_SERVER env var, falling back to "http://localhost:8080".
<code>.timeout</code>	Request timeout in seconds (default: 10).

**Value**

A named list with at least a status element.

---

```
llamacpp_list_local_models
```

*List Local GGUF Model Files*

---

**Description**

Scans a directory for .gguf files and returns a tibble with their names, sizes, and modification times. No server needed.

**Usage**

```
llamacpp_list_local_models(  
  .path = Sys.getenv("LLAMACPP_MODEL_DIR", "~/models")  
)
```

**Arguments**

`.path` Directory to scan (default: LLAMACPP\_MODEL\_DIR env var, then "~/models"). The path is expanded with `path.expand()`.

**Value**

A tibble with columns filename, size\_gb, modified, and path.

---

```
llamacpp_list_models
```

*List Models Loaded in the llama.cpp Server*

---

**Description**

Calls the /v1/models endpoint of a running llama.cpp server and returns the currently loaded model(s) as a tibble. In normal operation this is one row; two rows appear when speculative decoding is active (main + draft model).

**Usage**

```
llamacpp_list_models(  
  .server = Sys.getenv("LLAMACPP_SERVER", "http://localhost:8080"),  
  .api_key = Sys.getenv("LLAMACPP_API_KEY", ""),  
  .timeout = 30,  
  .max_tries = 3  
)
```

**Arguments**

<code>.server</code>	Base URL of the llama.cpp server. Defaults to LLAMACPP_SERVER env var, falling back to "http://localhost:8080".
<code>.api_key</code>	API key for the server (default: LLAMACPP_API_KEY env var).
<code>.timeout</code>	Request timeout in seconds (default: 30).
<code>.max_tries</code>	Maximum retries (default: 3).

**Value**

A tibble with columns id, object, and created.

---

llamacpp_rerank	<i>Rerank Documents Using a llama.cpp Server</i>
-----------------	--

---

**Description**

Calls the /v1/reranking endpoint of a running llama.cpp server to score documents by relevance to a query. Useful for building fully-local RAG pipelines (embed → cosine search → rerank → chat, all with llamacpp()).

**Usage**

```
llamacpp_rerank(
  .query,
  .documents,
  .model = "local-model",
  .server = Sys.getenv("LLAMACPP_SERVER", "http://localhost:8080"),
  .api_key = Sys.getenv("LLAMACPP_API_KEY", ""),
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE
)
```

**Arguments**

<code>.query</code>	A single query string.
<code>.documents</code>	A character vector of documents to rerank.
<code>.model</code>	The model name (default: "local-model").
<code>.server</code>	Base URL of the llama.cpp server. Defaults to LLAMACPP_SERVER env var, falling back to "http://localhost:8080".
<code>.api_key</code>	API key for the server (default: LLAMACPP_API_KEY env var).
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum retries (default: 3).
<code>.dry_run</code>	If TRUE, returns the request object without executing it (default: FALSE).

**Value**

A tibble with columns `index` (original position), `document` (text), and `relevance_score`, sorted by descending score.

---

LLMMessage

*Large Language Model Message Class*


---

**Description**

LLMMessage is an S7 class for managing a conversation history intended for use with large language models (LLMs). Please use `llm_message()` to create or modify LLMMessage objects.

**Usage**

```
LLMMessage(message_history = list(), system_prompt = character(0))
```

**Arguments**

`message_history`

A list containing messages. Each message is a named list with keys like `role`, `content`, `media`, etc.

`system_prompt`

A character string representing the default system prompt used for the conversation.

**Details**

The LLMMessage class includes the following features:

- Stores message history in a structured format.
- Supports attaching media and metadata to messages.
- Provides generics like `add_message()`, `has_image()`, and `remove_message()` for interaction.
- Enables API-specific formatting through the `to_api_format()` generic.
- `message_history`: A list containing messages. Each message is a named list with keys like `role`, `content`, `media`, etc.
- `system_prompt`: A character string representing the default system prompt used for the conversation.

llm\_message

*Create or Update Large Language Model Message Object***Description**

This function creates a new LLMMessage object or updates an existing one. It supports adding text prompts and various media types, such as images, PDFs, text files, or plots.

**Usage**

```
llm_message(
  .llm = NULL,
  .prompt = NULL,
  .role = "user",
  .system_prompt = "You are a helpful assistant",
  .media = NULL,
  .files = NULL,
  .imagefile = NULL,
  .pdf = NULL,
  .textfile = NULL,
  .capture_plot = FALSE,
  .f = NULL
)
```

**Arguments**

<code>.llm</code>	An existing LLMMessage object or an initial text prompt.
<code>.prompt</code>	Text prompt to add to the message history.
<code>.role</code>	The role of the message sender, typically "user" or "assistant".
<code>.system_prompt</code>	Default system prompt if a new LLMMessage needs to be created.
<code>.media</code>	An inline media object or a list of them. Accepted types: <code>img()</code> , <code>audio_file()</code> , <code>video_file()</code> , <code>pdf_file()</code> .
<code>.files</code>	A <code>tidyllm_file</code> object or list of them returned by <code>upload_file()</code> . These are remote file references stored on a provider's server.
<code>.imagefile</code>	Path to an image file to be attached (optional). Deprecated; use <code>.media = img(path)</code> instead.
<code>.pdf</code>	Path to a PDF file to be attached (optional). Deprecated; use <code>.media = pdf_file(path)</code> instead.
<code>.textfile</code>	Path to a text file to be read and attached (optional).
<code>.capture_plot</code>	Boolean to indicate whether a plot should be captured and attached as an image (optional).
<code>.f</code>	An R function or an object coercible to a function via <code>rlang::as_function</code> , whose output should be captured and attached (optional).

**Value**

Returns an updated or new LLMMessage object.

**See Also**

[df\\_llm\\_message\(\)](#)

Other Message Creation Utilities: [df\\_llm\\_message\(\)](#)

---

mistral

*Mistral Provider Function*

---

**Description**

The `mistral()` function acts as an interface for interacting with the Mistral API through main `tidyllm` verbs such as `chat()` and `embed()`. It dynamically routes requests to Mistral-specific functions like `mistral_chat()` and `mistral_embedding()` based on the context of the call.

**Usage**

```
mistral(..., .called_from = NULL)
```

**Arguments**

<code>...</code>	Parameters to be passed to the appropriate Mistral-specific function, such as model configuration, input text, or API-specific options.
<code>.called_from</code>	An internal argument that specifies which action (e.g., <code>chat</code> , <code>embed</code> , <code>send_batch</code> ) the function is being invoked from. This argument is automatically managed and should not be modified by the user.

**Value**

The result of the requested action, depending on the specific function invoked (e.g., an updated LLMMessage object for `chat()`, or a matrix for `embed()`).

---

mistral\_chat

*Send LLMMessage to Mistral API*

---

**Description**

Send LLMMessage to Mistral API

**Usage**

```

mistral_chat(
  .llm,
  .model = "mistral-large-latest",
  .frequency_penalty = NULL,
  .logit_bias = NULL,
  .presence_penalty = NULL,
  .seed = NULL,
  .stop = NULL,
  .stream = FALSE,
  .temperature = 0.7,
  .top_p = 1,
  .min_tokens = NULL,
  .max_tokens = NULL,
  .json_schema = NULL,
  .safe_prompt = FALSE,
  .reasoning_effort = NULL,
  .timeout = 120,
  .max_tries = 3,
  .dry_run = FALSE,
  .verbose = FALSE,
  .tools = NULL,
  .tool_choice = NULL,
  .max_tool_rounds = 10
)

```

**Arguments**

<code>.llm</code>	An LLMMessage object.
<code>.model</code>	The model identifier to use (default: "mistral-large-latest").
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency.
<code>.logit_bias</code>	A named list modifying the likelihood of specified tokens appearing in the completion.
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far.
<code>.seed</code>	If specified, the system will make a best effort to sample deterministically.
<code>.stop</code>	Up to 4 sequences where the API will stop generating further tokens.
<code>.stream</code>	If set to TRUE, the answer will be streamed to console as it comes (default: FALSE).
<code>.temperature</code>	What sampling temperature to use, between 0 and 2. Higher values make the output more random.
<code>.top_p</code>	An alternative to sampling with temperature, called nucleus sampling.

<code>.min_tokens</code>	The minimum number of tokens to generate in the completion. Must be $\geq 0$ (optional).
<code>.max_tokens</code>	An upper bound for the number of tokens that can be generated for a completion.
<code>.json_schema</code>	A JSON schema object provided by tidyllm schema or ellmer schemata.
<code>.safe_prompt</code>	Whether to inject a safety prompt before all conversations (default: FALSE).
<code>.reasoning_effort</code>	Controls the reasoning effort for Magistral thinking models; one of "low", "medium", or "high" (default: NULL, meaning the API default).
<code>.timeout</code>	When should our connection time out in seconds (default: 120).
<code>.max_tries</code>	Maximum retries to perform request
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object (default: FALSE).
<code>.verbose</code>	Should additional information be shown after the API call? (default: FALSE)
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.tool_choice</code>	A character string specifying the tool-calling behavior; valid values are "none", "auto", or "required".
<code>.max_tool_rounds</code>	Integer specifying the maximum number of tool use iterations (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.

**Value**

Returns an updated LLMMessage object.

---

<code>mistral_embedding</code>	<i>Generate Embeddings Using Mistral API</i>
--------------------------------	--

---

**Description**

Generate Embeddings Using Mistral API

**Usage**

```
mistral_embedding(
  .input,
  .model = "mistral-embed",
  .timeout = 120,
  .max_tries = 3,
  .dry_run = FALSE
)
```

**Arguments**

<code>.input</code>	A character vector of texts to embed or an LLMMessage object
<code>.model</code>	The embedding model identifier (default: "mistral-embed").
<code>.timeout</code>	Timeout for the API request in seconds (default: 120).
<code>.max_tries</code>	Maximum retries to perform request
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object.

**Value**

A matrix where each column corresponds to the embedding of a message in the message history.

---

`mistral_list_models`    *List Available Models from the Mistral API*

---

**Description**

List Available Models from the Mistral API

**Usage**

```
mistral_list_models(
  .api_url = "https://api.mistral.ai",
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE,
  .verbose = FALSE
)
```

**Arguments**

<code>.api_url</code>	Base URL for the API (default: "https://api.mistral.ai").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum number of retries for the API request (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.
<code>.verbose</code>	Logical; if TRUE, prints additional information about the request.

**Value**

A tibble containing model information (columns include id and created), or NULL if no models are found.

---

`ollama`*Ollama API Provider Function*

---

### Description

The `ollama()` function acts as an interface for interacting with local AI models via the Ollama API. It integrates seamlessly with the main `tidyllm` verbs such as `chat()` and `embed()`.

### Usage

```
ollama(..., .called_from = NULL)
```

### Arguments

<code>...</code>	Parameters to be passed to the appropriate Ollama-specific function, such as model configuration, input text, or API-specific options.
<code>.called_from</code>	An internal argument specifying the verb (e.g., <code>chat</code> , <code>embed</code> ) the function is invoked from. This argument is automatically managed by <code>tidyllm</code> and should not be set by the user.

### Details

Some functionalities, like `ollama_download_model()` or `ollama_list_models()` are unique to the Ollama API and do not have a general verb counterpart. These functions can be only accessed directly.

Supported Verbs:

- `chat()`: Sends a message to an Ollama model and retrieves the model's response.
- `embed()`: Generates embeddings for input texts using an Ollama model.
- `send_batch()`: Behaves different than the other `send_batch()` verbs since it immediately processes the answers

### Value

The result of the requested action:

- For `chat()`: An updated `LLMMessage` object containing the model's response.
- For `embed()`: A matrix where each column corresponds to an embedding.

ollama\_chat

*Interact with local AI models via the Ollama API***Description**

Interact with local AI models via the Ollama API

**Usage**

```
ollama_chat(
  .llm,
  .model = "qwen3.5:4b",
  .stream = FALSE,
  .seed = NULL,
  .json_schema = NULL,
  .temperature = NULL,
  .num_ctx = 2048,
  .num_predict = NULL,
  .top_k = NULL,
  .top_p = NULL,
  .min_p = NULL,
  .mirostat = NULL,
  .mirostat_eta = NULL,
  .mirostat_tau = NULL,
  .repeat_last_n = NULL,
  .repeat_penalty = NULL,
  .tools = NULL,
  .max_tool_rounds = 10,
  .tfs_z = NULL,
  .stop = NULL,
  .think = NULL,
  .ollama_server = "http://localhost:11434",
  .timeout = 120,
  .keep_alive = NULL,
  .dry_run = FALSE
)
```

**Arguments**

<code>.llm</code>	An LLMMessage object containing the conversation history and system prompt.
<code>.model</code>	Character string specifying the Ollama model to use (default: "qwen3-vl")
<code>.stream</code>	Logical; whether to stream the response (default: FALSE)
<code>.seed</code>	Integer; seed for reproducible generation (default: NULL)
<code>.json_schema</code>	A JSON schema object as R list to enforce the output structure (default: NULL)
<code>.temperature</code>	Float between 0-2; controls randomness in responses (default: NULL)

<code>.num_ctx</code>	Integer; sets the context window size (default: 2048)
<code>.num_predict</code>	Integer; maximum number of tokens to predict (default: NULL)
<code>.top_k</code>	Integer; controls diversity by limiting top tokens considered (default: NULL)
<code>.top_p</code>	Float between 0-1; nucleus sampling threshold (default: NULL)
<code>.min_p</code>	Float between 0-1; minimum probability threshold (default: NULL)
<code>.mirostat</code>	Integer (0,1,2); enables Mirostat sampling algorithm (default: NULL)
<code>.mirostat_eta</code>	Float; Mirostat learning rate (default: NULL)
<code>.mirostat_tau</code>	Float; Mirostat target entropy (default: NULL)
<code>.repeat_last_n</code>	Integer; tokens to look back for repetition (default: NULL)
<code>.repeat_penalty</code>	Float; penalty for repeated tokens (default: NULL)
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects representing the available functions for tool calls.
<code>.max_tool_rounds</code>	Integer; maximum number of tool use iterations for multi-turn tool calling (default: 10). Set to 1 for single-round tool use, or higher for multi-turn agentic loops.
<code>.tfs_z</code>	Float; tail free sampling parameter (default: NULL)
<code>.stop</code>	Character; custom stop sequence(s) (default: NULL)
<code>.think</code>	Logical or character; controls thinking mode for supported models like Qwen3. Use FALSE to disable, TRUE to enable, or "high"/"medium"/"low" to set effort level (default: NULL - model default)
<code>.ollama_server</code>	String; Ollama API endpoint (default: "http://localhost:11434")
<code>.timeout</code>	Integer; API request timeout in seconds (default: 120)
<code>.keep_alive</code>	Character; How long should the ollama model be kept in memory after request (default: NULL - 5 Minutes)
<code>.dry_run</code>	Logical; if TRUE, returns request object without execution (default: FALSE)

### Details

The function provides extensive control over the generation process through various parameters:

- Temperature (0-2): Higher values increase creativity, lower values make responses more focused
- Top-k/Top-p: Control diversity of generated text
- Mirostat: Advanced sampling algorithm for maintaining consistent complexity
- Repeat penalties: Prevent repetitive text
- Context window: Control how much previous conversation is considered

### Value

A new LLMMessage object containing the original messages plus the model's response

### Examples

```
## Not run:
llm_message("user", "Hello, how are you?")
response <- ollama_chat(llm, .model = "gemma2", .temperature = 0.7)

# With custom parameters
response <- ollama_chat(
  llm,
  .model = "llama2",
  .temperature = 0.8,
  .top_p = 0.9,
  .num_ctx = 4096
)

## End(Not run)
```

---

ollama\_delete\_model *Delete a model from the Ollama API*

---

### Description

This function sends a DELETE request to the Ollama API to remove a specified model.

### Usage

```
ollama_delete_model(.model, .ollama_server = "http://localhost:11434")
```

### Arguments

`.model` The name of the model to delete.  
`.ollama_server` The base URL of the Ollama API (default is "http://localhost:11434").

---

ollama\_download\_model *Download a model from the Ollama API*

---

### Description

This function sends a request to the Ollama API to download a specified model from Ollama's large online library of models.

### Usage

```
ollama_download_model(.model, .ollama_server = "http://localhost:11434")
```

**Arguments**

- `.model` The name of the model to download.
- `.ollama_server` The base URL of the Ollama API (default is "http://localhost:11434").

---

ollama\_embedding      *Generate Embeddings Using Ollama API*

---

**Description**

Generate Embeddings Using Ollama API

**Usage**

```
ollama_embedding(  
  .input,  
  .model = "qwen3-embedding:0.6b",  
  .truncate = TRUE,  
  .ollama_server = "http://localhost:11434",  
  .timeout = 120,  
  .dry_run = FALSE  
)
```

**Arguments**

- `.input` Aa character vector of texts to embed or an LLMMessage object
- `.model` The embedding model identifier (default: "all-minilm").
- `.truncate` Whether to truncate inputs to fit the model's context length (default: TRUE).
- `.ollama_server` The URL of the Ollama server to be used (default: "http://localhost:11434").
- `.timeout` Timeout for the API request in seconds (default: 120).
- `.dry_run` If TRUE, perform a dry run and return the request object.

**Value**

A matrix where each column corresponds to the embedding of a message in the message history.

---

ollama_list_models	<i>Retrieve and return model information from the Ollama API</i>
--------------------	--

---

**Description**

This function connects to the Ollama API and retrieves information about available models, returning it as a tibble.

**Usage**

```
ollama_list_models(.ollama_server = "http://localhost:11434")
```

**Arguments**

`.ollama_server` The URL of the ollama server to be used

**Value**

A tibble containing model information, or NULL if no models are found.

---

openai	<i>OpenAI Provider Function</i>
--------	---------------------------------

---

**Description**

The `openai()` function acts as an interface for interacting with the OpenAI API through main `tidyllm` verbs such as `chat()`, `embed()`, and `send_batch()`. It dynamically routes requests to OpenAI-specific functions like `openai_chat()` and `openai_embedding()` based on the context of the call.

**Usage**

```
openai(..., .called_from = NULL)
```

```
openai(..., .called_from = NULL)
```

**Arguments**

`...` Parameters passed to the appropriate OpenAI-specific function.

`.called_from` Internal routing argument; do not set manually.

**Value**

The result of the requested action, depending on the specific function invoked (e.g., an updated `LLMMessage` object for `chat()`, or a matrix for `embed()`).

Result of the requested action.

openai\_chat

*Send LLM Messages to the OpenAI Responses API***Description**

Sends a message history to the OpenAI Responses API (POST /v1/responses) and returns the assistant's reply. Supports streaming, tool use, structured output, and reasoning models (o-series) via `.reasoning_effort`.

**Usage**

```
openai_chat(
  .llm,
  .model = "gpt-5.5",
  .max_output_tokens = NULL,
  .temperature = NULL,
  .seed = NULL,
  .stream = FALSE,
  .timeout = 60,
  .verbose = FALSE,
  .json_schema = NULL,
  .max_tries = 3,
  .dry_run = FALSE,
  .reasoning_effort = NULL,
  .tools = NULL,
  .tool_choice = NULL,
  .max_tool_rounds = 10,
  .stateful = FALSE
)
```

**Arguments**

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.model</code>	The model identifier (default: "gpt-5.5").
<code>.max_output_tokens</code>	Maximum tokens to generate (caps reasoning + completion).
<code>.temperature</code>	Sampling temperature (0-2).
<code>.seed</code>	Seed for deterministic sampling.
<code>.stream</code>	If TRUE, stream output to console (default: FALSE).
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.verbose</code>	Print rate-limit info after the call (default: FALSE).
<code>.json_schema</code>	A tidyllm schema or ellmer type for structured output.
<code>.max_tries</code>	Maximum retry attempts (default: 3).
<code>.dry_run</code>	If TRUE, return the request object without sending (default: FALSE).

<code>.reasoning_effort</code>	For o-series models: "low", "medium", or "high".
<code>.tools</code>	A TOOL object or list of TOOL objects for function calling.
<code>.tool_choice</code>	Tool selection behavior: "none", "auto", or "required".
<code>.max_tool_rounds</code>	Maximum tool-loop iterations (default: 10).
<code>.stateful</code>	If TRUE, use <code>previous_response_id</code> to pass conversation context server-side instead of re-sending full message history. Falls back to full-history mode silently when no previous response ID is available (first turn or prior turn used a different provider), and with a warning if the server rejects the stored ID (e.g., expired context). Default: FALSE.

**Value**

A new `LLMMessage` object with the assistant's response appended.

---

`openai_check_research` *Check the Status of an OpenAI Background Research Job*

---

**Description**

Polls the status of an OpenAI background response created by `openai_deep_research(.background = TRUE)`.

**Usage**

```
openai_check_research(.job, .max_tries = 3)
```

**Arguments**

<code>.job</code>	A <code>tidyllm_research_job</code> object returned by <code>openai_deep_research(.background = TRUE)</code> .
<code>.max_tries</code>	Maximum retries per HTTP request (default: 3).

**Value**

An updated `tidyllm_research_job` with `$status` set. If completed, `$response` is also populated and ready for `fetch_job()`.

---

openai\_code\_interpreter

*OpenAI built-in code interpreter tool (server-executed)*

---

### Description

OpenAI built-in code interpreter tool (server-executed)

### Usage

```
openai_code_interpreter()
```

### Value

A TOOL object declaring the code\_interpreter built-in.

---

openai\_deep\_research *Submit a Deep Research Request to OpenAI*

---

### Description

Sends a research request to OpenAI using the deep research models (o3-deep-research or o4-mini-deep-research) via the Responses API with background: true. The model autonomously searches the web and synthesises a long-form answer, which can take 5-30 minutes.

### Usage

```
openai_deep_research(
  .llm,
  .model = "o4-mini-deep-research",
  .background = FALSE,
  .reasoning_effort = "medium",
  .json_schema = NULL,
  .max_output_tokens = NULL,
  .timeout = 1800,
  .max_tries = 3
)
```

### Arguments

.llm	An LLMMessage object containing the research question.
.model	The deep research model to use (default: "o4-mini-deep-research").
.background	Logical; if TRUE, returns a tidyllm_research_job immediately without waiting for completion (default: FALSE).

<code>.reasoning_effort</code>	Reasoning level for the model: "low", "medium" (default), or "high".
<code>.json_schema</code>	A tidyllm schema for structured JSON output (optional).
<code>.max_output_tokens</code>	Maximum tokens to generate (default: NULL for model default).
<code>.timeout</code>	Seconds to wait in blocking mode before giving up (default: 1800).
<code>.max_tries</code>	Maximum retries per HTTP request (default: 3).

**Value**

If `.background = FALSE`, an updated `LLMMessage` with the research reply. If `.background = TRUE`, a `tidyllm_research_job` for use with `check_job()/fetch_job()`.

---

<code>openai_embedding</code>	<i>Generate Embeddings Using OpenAI API</i>
-------------------------------	---

---

**Description**

Generate Embeddings Using OpenAI API

**Usage**

```
openai_embedding(
  .input,
  .model = "text-embedding-3-small",
  .truncate = TRUE,
  .timeout = 120,
  .dry_run = FALSE,
  .max_tries = 3,
  .verbose = FALSE
)
```

**Arguments**

<code>.input</code>	An existing <code>LLMMessage</code> object (or a character vector of texts to embed)
<code>.model</code>	The embedding model identifier (default: "text-embedding-3-small").
<code>.truncate</code>	Whether to truncate inputs to fit the model's context length (default: TRUE).
<code>.timeout</code>	Timeout for the API request in seconds (default: 120).
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object.
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.verbose</code>	Should information about current ratelimits be printed? (default: FALSE)

**Value**

A tibble with two columns: `input` and `embeddings`. The `input` column contains the texts sent to embed, and the `embeddings` column is a list column where each row contains an embedding vector of the sent input.

---

openai\_fetch\_research *Fetch Results from a Completed OpenAI Deep Research Job*

---

### Description

Extracts the assistant reply from a completed `tidyllm_research_job` returned by `openai_deep_research(.background = TRUE)`.

### Usage

```
openai_fetch_research(.job, .max_tries = 3)
```

### Arguments

<code>.job</code>	A <code>tidyllm_research_job</code> object. If not yet completed, polls once via <code>openai_check_research()</code> and errors if still incomplete.
<code>.max_tries</code>	Maximum retries per HTTP request (default: 3).

### Value

An updated `LLMMessage` with the research reply appended.

---

openai\_list\_models *List Available Models from the OpenAI API*

---

### Description

List Available Models from the OpenAI API

### Usage

```
openai_list_models(
  .api_url = "https://api.openai.com",
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE,
  .verbose = FALSE
)
```

### Arguments

<code>.api_url</code>	Base URL for the API (default: "https://api.openai.com").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum number of retries for the API request (default: 3).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it.
<code>.verbose</code>	Logical; if TRUE, prints additional information about the request.

**Value**

A tibble containing model information (columns include id, created, and owned\_by), or NULL if no models are found.

---

openai_websearch	<i>OpenAI built-in web search tool (server-executed)</i>
------------------	--

---

**Description**

OpenAI built-in web search tool (server-executed)

**Usage**

```
openai_websearch()
```

**Value**

A TOOL object declaring the web\_search\_preview built-in.

---

openrouter	<i>OpenRouter Provider Function</i>
------------	-------------------------------------

---

**Description**

The openrouter() function provides access to hundreds of AI models from different providers through the OpenRouter API, using a single OpenAI-compatible interface.

**Usage**

```
openrouter(..., .called_from = NULL)
```

**Arguments**

...	Parameters passed to the appropriate OpenRouter-specific function.
.called_from	An internal argument specifying which verb invoked this function. Managed automatically by tidyllm verbs; do not set manually.

**Value**

The result of the requested action (e.g., an updated LLMMessage for chat()).

---

openrouter\_chat      *Send LLM Messages to the OpenRouter Chat API*

---

### Description

Sends a message history to the OpenRouter API, which provides access to hundreds of models from different providers through a single OpenAI-compatible endpoint.

### Usage

```
openrouter_chat(
    .llm,
    .model = "anthropic/claude-sonnet-4-6",
    .max_tokens = 2048,
    .temperature = NULL,
    .top_p = NULL,
    .frequency_penalty = NULL,
    .presence_penalty = NULL,
    .stop = NULL,
    .stream = FALSE,
    .json_schema = NULL,
    .tools = NULL,
    .tool_choice = NULL,
    .reasoning = NULL,
    .provider = NULL,
    .route = NULL,
    .models = NULL,
    .api_url = "https://openrouter.ai",
    .timeout = 60,
    .verbose = FALSE,
    .dry_run = FALSE,
    .max_tries = 3,
    .max_tool_rounds = 10
)
```

### Arguments

<code>.llm</code>	An LLMMessage object containing the conversation history.
<code>.model</code>	The model identifier to use (default: "google/gemini-2.5-flash"). Any model available on OpenRouter can be used.
<code>.max_tokens</code>	Maximum number of tokens in the response (default: 2048).
<code>.temperature</code>	Controls randomness (0–2, optional).
<code>.top_p</code>	Nucleus sampling parameter (0–1, optional).
<code>.frequency_penalty</code>	Penalizes repeated tokens (-2 to 2, optional).

<code>.presence_penalty</code>	Encourages new topics (-2 to 2, optional).
<code>.stop</code>	One or more stop sequences (optional).
<code>.stream</code>	Logical; if TRUE, streams the response (default: FALSE).
<code>.json_schema</code>	A JSON schema object for structured output (default: NULL).
<code>.tools</code>	Either a single TOOL object or a list of TOOL objects for tool calls.
<code>.tool_choice</code>	Tool-calling behavior: "none", "auto", or "required" (optional).
<code>.reasoning</code>	A named list controlling reasoning token behavior (optional). Supported fields vary by model family: <ul style="list-style-type: none"> <li>• <code>effort</code>: one of "xhigh", "high", "medium", "low", "minimal", "none" (OpenAI/Grok)</li> <li>• <code>max_tokens</code>: integer specifying the reasoning token budget (Anthropic/Gemini/Alibaba)</li> <li>• <code>exclude</code>: logical; if TRUE, reasoning is used internally but not returned in the response</li> <li>• <code>enabled</code>: logical; if TRUE, activates reasoning at default settings</li> </ul> Example: <code>list(effort = "high")</code> or <code>list(max_tokens = 4000, exclude = FALSE)</code> .
<code>.provider</code>	A named list of OpenRouter provider preferences, e.g. <code>list(order = c("Anthropic", "AWS Bedrock"), allow_fallbacks = TRUE)</code> (optional).
<code>.route</code>	OpenRouter routing strategy; "fallback" routes to the next model if the primary is unavailable (optional).
<code>.models</code>	A character vector of model IDs to use as fallbacks when the primary model is unavailable (optional). Used together with <code>.route = "fallback"</code> .
<code>.api_url</code>	Base URL for the OpenRouter API (default: "https://openrouter.ai").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.verbose</code>	If TRUE, displays additional information after the API call (default: FALSE).
<code>.dry_run</code>	If TRUE, returns the request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries (default: 3).
<code>.max_tool_rounds</code>	Maximum number of tool use iterations (default: 10).

### Value

A new `LLMMessage` object containing the original messages plus the assistant's response.

---

openrouter\_credits     *Get OpenRouter Credit Balance*

---

### Description

Returns the total credits purchased and total usage so far for the current API key.

### Usage

```
openrouter_credits(
  .api_url = "https://openrouter.ai",
  .timeout = 60,
  .max_tries = 3
)
```

### Arguments

`.api_url`     Base URL for the OpenRouter API (default: "https://openrouter.ai").

`.timeout`     Request timeout in seconds (default: 60).

`.max_tries`    Maximum retries (default: 3).

### Value

A tibble with columns `total_credits` (USD purchased), `total_usage` (USD consumed so far), and `remaining` (USD remaining).

---

openrouter\_embedding     *Generate Embeddings Using the OpenRouter API*

---

### Description

Sends text to an embedding model accessible via OpenRouter and returns embedding vectors. Note that embedding models are not listed in `list_models(openrouter())` — specify the model ID directly. Known supported models include "openai/text-embedding-3-small", "openai/text-embedding-3-large", and "mistralai/mistral-embed".

### Usage

```
openrouter_embedding(
  .input,
  .model = "openai/text-embedding-3-small",
  .timeout = 120,
  .dry_run = FALSE,
  .max_tries = 3
)
```

**Arguments**

<code>.input</code>	An LLMMessage object or a character vector of texts to embed.
<code>.model</code>	The embedding model ID (default: "openai/text-embedding-3-small").
<code>.timeout</code>	Request timeout in seconds (default: 120).
<code>.dry_run</code>	If TRUE, returns the request object without executing it (default: FALSE).
<code>.max_tries</code>	Maximum retries (default: 3).

**Value**

A tibble with columns `input` (text) and `embeddings` (list of numeric vectors).

---

`openrouter_generation` *Get Details for an OpenRouter Generation*

---

**Description**

Fetches cost, latency, and token details for a past generation using its ID. The generation ID is available in `get_metadata(result)$api_specific$id` after a `chat()` call.

**Usage**

```
openrouter_generation(
  .id,
  .api_url = "https://openrouter.ai",
  .timeout = 60,
  .max_tries = 3
)
```

**Arguments**

<code>.id</code>	The generation ID string (e.g. "gen-...") returned by OpenRouter.
<code>.api_url</code>	Base URL for the OpenRouter API (default: "https://openrouter.ai").
<code>.timeout</code>	Request timeout in seconds (default: 60).
<code>.max_tries</code>	Maximum retries (default: 3).

**Value**

A named list with generation details including `id`, `model`, `total_cost`, `tokens_prompt`, `tokens_completion`, and `latency`.

---

openrouter\_list\_models

*List Available Models on OpenRouter*

---

### Description

Retrieves the list of models available through the OpenRouter API, including pricing and context window information.

### Usage

```
openrouter_list_models(
  .api_url = "https://openrouter.ai",
  .timeout = 60,
  .max_tries = 3,
  .dry_run = FALSE,
  .verbose = FALSE
)
```

### Arguments

.api_url	Base URL for the OpenRouter API (default: "https://openrouter.ai").
.timeout	Request timeout in seconds (default: 60).
.max_tries	Maximum retries (default: 3).
.dry_run	If TRUE, returns the request object without executing it (default: FALSE).
.verbose	If TRUE, displays additional information (default: FALSE).

### Value

A tibble with columns id, name, context\_length, prompt\_price\_per\_million, and completion\_price\_per\_million, or NULL if no models are found.

---

pdf\_file

*Create a PDF Object*

---

### Description

Stores a reference to a local PDF file along with a pdftools text extraction. Binary encoding is deferred to format time; providers that support binary PDF encode on demand, others use the stored text fallback.

### Usage

```
pdf_file(.path, pages = NULL, .text_extract = FALSE, engine = NULL)
```

**Arguments**

.path	The path to the PDF file on disk.
pages	An integer vector of page numbers to include (e.g. 1:3). NULL uses all pages.
.text_extract	Logical; if TRUE always use pdftools text extraction regardless of provider (equivalent to the old .pdf parameter behaviour). Default FALSE uses binary-first.
engine	Optional engine name (e.g. "mistral-ocr"); a warning is issued if the provider does not support binary PDF and the engine choice cannot be honoured.

**Value**

A tidyllm\_pdf object.

---

pdf_page_batch	<i>Batch Process PDF into LLM Messages</i>
----------------	--

---

**Description**

This function processes a PDF file page by page. For each page, it extracts the text and converts the page into an image. It creates a list of LLMMessage objects with the text and the image for multimodal processing. Users can specify a range of pages to process and provide a custom function to generate prompts for each page.

**Usage**

```
pdf_page_batch(
  .pdf,
  .general_prompt,
  .system_prompt = "You are a helpful assistant",
  .page_range = NULL,
  .prompt_fn = NULL
)
```

**Arguments**

.pdf	Path to the PDF file.
.general_prompt	A default prompt that is applied to each page if .prompt_fn is not provided.
.system_prompt	Optional system prompt to initialize the LLMMessage (default is "You are a helpful assistant").
.page_range	A vector of two integers specifying the start and end pages to process. If NULL, all pages are processed.
.prompt_fn	An optional custom function that generates a prompt for each page. The function takes the page text as input and returns a string. If NULL, .general_prompt is used for all pages.

**Value**

A list of LLMMessage objects, each containing the text and image for a page.

---

perplexity	<i>Perplexity Provider Function</i>
------------	-------------------------------------

---

**Description**

The perplexity() function acts as a provider interface for interacting with the Perplexity API through tidyllm's chat() verb. It dynamically routes requests to Perplexity-specific function. At the moment this is only perplexity\_chat()

**Usage**

```
perplexity(..., .called_from = NULL)
```

**Arguments**

...	Parameters to be passed to the appropriate Perplexity-specific function, such as model configuration, input text, or API-specific options.
.called_from	An internal argument specifying which action (e.g., chat, embed) the function is invoked from. This argument is automatically managed by the tidyllm verbs and should not be modified by the user.

**Value**

The result of the requested action, depending on the specific function invoked (e.g., an updated LLMMessage object for chat()).

---

perplexity_chat	<i>Send LLM Messages to the Perplexity Chat API</i>
-----------------	---

---

**Description**

Sends a chat message history to the Perplexity Chat API, supporting all documented API parameters.

**Usage**

```
perplexity_chat(  
  .llm,  
  .model = "sonar",  
  .max_tokens = 1024,  
  .temperature = NULL,  
  .top_p = NULL,  
  .frequency_penalty = NULL,  
  .presence_penalty = NULL,  
  .stop = NULL,  
  .search_domain_filter = NULL,  
  .search_language_filter = NULL,  
  .language_preference = NULL,  
  .return_images = FALSE,  
  .image_domain_filter = NULL,  
  .image_format_filter = NULL,  
  .search_recency_filter = NULL,  
  .search_mode = "web",  
  .search_after_date_filter = NULL,  
  .search_before_date_filter = NULL,  
  .last_updated_after_filter = NULL,  
  .last_updated_before_filter = NULL,  
  .disable_search = FALSE,  
  .enable_search_classifier = FALSE,  
  .reasoning_effort = NULL,  
  .return_related_questions = FALSE,  
  .user_location = NULL,  
  .search_context_size = NULL,  
  .search_type = NULL,  
  .stream_mode = NULL,  
  .json_schema = NULL,  
  .top_k = NULL,  
  .web_search_options = NULL,  
  .api_url = "https://api.perplexity.ai/",  
  .timeout = 60,  
  .stream = FALSE,  
  .verbose = FALSE,  
  .max_tries = 3,  
  .dry_run = FALSE  
)
```

**Arguments**

.llm	An LLMMessage object containing the conversation history.
.model	Model name to use (default: "sonar").
.max_tokens	Max completion tokens (default: 1024).
.temperature	Controls response randomness ( $0 < x < 2$ ).

.top\_p Nucleus sampling threshold ( $0 < x < 1$ ).

.frequency\_penalty Number  $> 0$ . Penalizes frequent tokens.

.presence\_penalty Numeric between -2 and 2. Penalizes present tokens.

.stop Stop sequence(s), string or character vector/list.

.search\_domain\_filter Character vector of domains to allowlist/denylist (max 10; prefix with "-" to denylist).

.search\_language\_filter ISO 639-1 language code to restrict search results (e.g. "en", "de", "fr").

.language\_preference ISO 639-1 code for preferred response language (e.g. "en").

.return\_images Logical; if TRUE, returns images from search.

.image\_domain\_filter Character vector of domains to restrict image results to.

.image\_format\_filter Character vector of image formats to include (e.g. c("png", "jpg")).

.search\_recency\_filter Restrict search to recent results: "hour", "day", "week", "month", or "year".

.search\_mode Search index to use: "web" (default), "academic", or "sec".

.search\_after\_date\_filter Only include content published after this date (MM/DD/YYYY).

.search\_before\_date\_filter Only include content published before this date (MM/DD/YYYY).

.last\_updated\_after\_filter Only include content last updated after this date (MM/DD/YYYY).

.last\_updated\_before\_filter Only include content last updated before this date (MM/DD/YYYY).

.disable\_search Logical; if TRUE, disables web search entirely (default: FALSE).

.enable\_search\_classifier Logical; if TRUE, lets the model decide whether to search (default: FALSE).

.reasoning\_effort Reasoning level: "low", "medium", or "high".

.return\_related\_questions Logical; if TRUE, returns related questions.

.user\_location Named list for geographic search personalisation. Accepted fields: country (ISO 3166-1 alpha-2), city, region, latitude, longitude. Example: list(country = "DE", city = "Berlin").

.search\_context\_size Amount of search context to include: "low", "medium" (default), or "high".

.search\_type Search quality preference inside web\_search\_options: "fast", "pro", or "auto".

<code>.stream_mode</code>	Response format: "full" (default) or "concise".
<code>.json_schema</code>	A tidyllm schema created with <code>tidyllm_schema()</code> for structured JSON output (optional).
<code>.top_k</code>	Top-k token sampling (integer, 0 disables).
<code>.web_search_options</code>	Named list with raw <code>web_search_options</code> overrides. Values set here take precedence over the dedicated parameters above.
<code>.api_url</code>	API endpoint (default: "https://api.perplexity.ai/").
<code>.timeout</code>	Timeout in seconds (default: 60).
<code>.stream</code>	If TRUE, streams response.
<code>.verbose</code>	If TRUE, prints additional info.
<code>.max_tries</code>	Max request retries (default: 3).
<code>.dry_run</code>	If TRUE, returns constructed request instead of sending.

**Value**

An updated LLMMessage object with the assistant's reply and metadata, including citations and `search_results`.

---

perplexity\_check\_research

*Check the Status of a Perplexity Deep Research Job*

---

**Description**

Check the Status of a Perplexity Deep Research Job

**Usage**

```
perplexity_check_research(
  .job,
  .api_key = Sys.getenv("PERPLEXITY_API_KEY"),
  .max_tries = 3
)
```

**Arguments**

<code>.job</code>	A <code>tidyllm_research_job</code> object returned by <code>perplexity_deep_research(.background = TRUE)</code> .
<code>.api_key</code>	Character; Perplexity API key (default: from environment).
<code>.max_tries</code>	Integer; maximum retries (default: 3).

**Value**

An updated `tidyllm_research_job` with a `$status` field and `$response` if completed.

---

perplexity\_deep\_research

*Submit a Deep Research Request to Perplexity*


---

## Description

Submit a Deep Research Request to Perplexity

## Usage

```
perplexity_deep_research(
  .llm,
  .background = FALSE,
  .reasoning_effort = "medium",
  .search_context_size = "medium",
  .search_domain_filter = NULL,
  .search_language_filter = NULL,
  .language_preference = NULL,
  .search_recency_filter = NULL,
  .search_mode = NULL,
  .search_after_date_filter = NULL,
  .search_before_date_filter = NULL,
  .last_updated_after_filter = NULL,
  .last_updated_before_filter = NULL,
  .user_location = NULL,
  .json_schema = NULL,
  .idempotency_key = NULL,
  .api_key = Sys.getenv("PERPLEXITY_API_KEY"),
  .timeout = 300,
  .max_tries = 3
)
```

## Arguments

<code>.llm</code>	An LLMMessage object containing the research question.
<code>.background</code>	Logical; if TRUE, returns a <code>tidyllm_research_job</code> immediately without waiting (default: FALSE).
<code>.reasoning_effort</code>	Reasoning level: "low", "medium" (default), or "high".
<code>.search_context_size</code>	Amount of search context: "low", "medium" (default), or "high".
<code>.search_domain_filter</code>	Character vector of domains to allowlist/denylist (max 10; prefix with "-" to denylist).
<code>.search_language_filter</code>	ISO 639-1 language code to restrict search results (e.g. "en", "de").

<code>.language_preference</code>	ISO 639-1 code for preferred response language.
<code>.search_recency_filter</code>	Restrict search to recent results: "hour", "day", "week", "month", or "year".
<code>.search_mode</code>	Search index to use: "web" (default), "academic", or "sec".
<code>.search_after_date_filter</code>	Only include content published after this date (MM/DD/YYYY).
<code>.search_before_date_filter</code>	Only include content published before this date (MM/DD/YYYY).
<code>.last_updated_after_filter</code>	Only include content last updated after this date (MM/DD/YYYY).
<code>.last_updated_before_filter</code>	Only include content last updated before this date (MM/DD/YYYY).
<code>.user_location</code>	Named list for geographic search personalisation (fields: country, city, region, latitude, longitude).
<code>.json_schema</code>	A tidyllm schema created with <code>tidyllm_schema()</code> for structured JSON output (optional).
<code>.idempotency_key</code>	Optional string; unique key to prevent duplicate submissions.
<code>.api_key</code>	Character; Perplexity API key (default: from environment variable).
<code>.timeout</code>	Integer; request timeout in seconds for blocking polling (default: 300).
<code>.max_tries</code>	Integer; maximum retries (default: 3).

**Value**

If `.background = FALSE`, an updated `LLMMessage` with the research reply. If `.background = TRUE`, a `tidyllm_research_job` object.

---

`perplexity_fetch_research`

*Fetch Results from a Completed Perplexity Deep Research Job*

---

**Description**

Fetch Results from a Completed Perplexity Deep Research Job

**Usage**

```
perplexity_fetch_research(
  .job,
  .api_key = Sys.getenv("PERPLEXITY_API_KEY"),
  .max_tries = 3
)
```

**Arguments**

<code>.job</code>	A <code>tidyllm_research_job</code> object. Must have status "completed" or will poll once.
<code>.api_key</code>	Character; Perplexity API key (default: from environment).
<code>.max_tries</code>	Integer; maximum retries (default: 3).

**Value**

An updated `LLMMessage` with the research reply appended.

---

<code>rate_limit_info</code>	<i>Get the current rate limit information for all or a specific API</i>
------------------------------	---

---

**Description**

This function retrieves the rate limit details for the specified API, or for all APIs stored in the `.tidyllm_rate_limit_env` if no API is specified.

**Usage**

```
rate_limit_info(.api_name = NULL)
```

**Arguments**

<code>.api_name</code>	(Optional) The name of the API whose rate limit info you want to get. If not provided, the rate limit info for all APIs in the environment will be returned.
------------------------	--

**Value**

A tibble containing the rate limit information.

---

<code>send_azure_openai_batch</code>	<i>Send a Batch of Messages to Azure OpenAI Batch API</i>
--------------------------------------	---

---

**Description**

This function creates and submits a batch of messages to the Azure OpenAI Batch API for asynchronous processing.

**Usage**

```

send_azure_openai_batch(
    .llms,
    .endpoint_url = Sys.getenv("AZURE_ENDPOINT_URL"),
    .deployment = "gpt-4o-mini",
    .api_version = "2024-10-01-preview",
    .max_completion_tokens = NULL,
    .reasoning_effort = NULL,
    .frequency_penalty = NULL,
    .logit_bias = NULL,
    .presence_penalty = NULL,
    .seed = NULL,
    .stop = NULL,
    .temperature = NULL,
    .top_p = NULL,
    .logprobs = NULL,
    .top_logprobs = NULL,
    .dry_run = FALSE,
    .overwrite = FALSE,
    .max_tries = 3,
    .timeout = 60,
    .verbose = FALSE,
    .json_schema = NULL,
    .id_prefix = "tidyllm_azure_openai_req_"
)

```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects containing conversation histories.
<code>.endpoint_url</code>	Base URL for the API (default: <code>Sys.getenv("AZURE_ENDPOINT_URL")</code> ).
<code>.deployment</code>	The identifier of the model that is deployed (default: <code>"gpt-4o-mini"</code> ).
<code>.api_version</code>	Which version of the API is deployed (default: <code>"2024-10-01-preview"</code> )
<code>.max_completion_tokens</code>	Integer specifying the maximum tokens per response (default: <code>NULL</code> ).
<code>.reasoning_effort</code>	How long should reasoning models reason (can either be <code>"low"</code> , <code>"medium"</code> or <code>"high"</code> )
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency.
<code>.logit_bias</code>	A named list modifying the likelihood of specified tokens appearing in the completion.
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far.
<code>.seed</code>	If specified, the system will make a best effort to sample deterministically.

<code>.stop</code>	Up to 4 sequences where the API will stop generating further tokens.
<code>.temperature</code>	What sampling temperature to use, between 0 and 2. Higher values make the output more random.
<code>.top_p</code>	An alternative to sampling with temperature, called nucleus sampling.
<code>.logprobs</code>	If TRUE, get the log probabilities of each output token (default: NULL).
<code>.top_logprobs</code>	If specified, get the top N log probabilities of each output token (0-5, default: NULL).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.overwrite</code>	Logical; if TRUE, allows overwriting an existing batch ID (default: FALSE).
<code>.max_tries</code>	Maximum number of retries to perform the request (default: 3).
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).
<code>.verbose</code>	Logical; if TRUE, additional info about the requests is printed (default: FALSE).
<code>.json_schema</code>	A JSON schema object provided by tidyllm schema or ellmer schemata (default: NULL).
<code>.id_prefix</code>	Character string to specify a prefix for generating custom IDs when names in <code>.llms</code> are missing (default: "tidymlm_azure_openai_req_").

**Value**

An updated and named list of `.llms` with identifiers that align with batch responses, including a `batch_id` attribute.

---

send_batch	<i>Send a batch of messages to a batch API</i>
------------	--

---

**Description**

The `send_batch()` function allows you to send a list of `LLMMessage` objects to an API. It routes the input to the appropriate provider-specific batch API function.

**Usage**

```
send_batch(
  .llms,
  .provider = getOption("tidymlm_sbatch_default"),
  .dry_run = NULL,
  .temperature = NULL,
  .timeout = NULL,
  .top_p = NULL,
  .max_tries = NULL,
  .model = NULL,
  .verbose = NULL,
  .json_schema = NULL,
```

```

    .seed = NULL,
    .stop = NULL,
    .frequency_penalty = NULL,
    .presence_penalty = NULL,
    .id_prefix = NULL
  )

```

## Arguments

<code>.llms</code>	A list of <code>LLMMessage</code> objects containing conversation histories.
<code>.provider</code>	A function or function call specifying the language model provider and any additional parameters. This should be a call to a provider function like <code>openai()</code> , <code>claude()</code> , etc. You can also set a default provider function via the <code>tidyllm_sbatch_default</code> option.
<code>.dry_run</code>	Logical; if <code>TRUE</code> , simulates the request without sending it to the provider. Useful for testing.
<code>.temperature</code>	Numeric; controls the randomness of the model's output (0 = deterministic).
<code>.timeout</code>	Numeric; the maximum time (in seconds) to wait for a response.
<code>.top_p</code>	Numeric; nucleus sampling parameter, which limits the sampling to the top cumulative probability $p$ .
<code>.max_tries</code>	Integer; the maximum number of retries for failed requests.
<code>.model</code>	Character; the model identifier to use (e.g., "gpt-4").
<code>.verbose</code>	Logical; if <code>TRUE</code> , prints additional information about the request and response.
<code>.json_schema</code>	List; A JSON schema object as R list to enforce the output structure
<code>.seed</code>	Integer; sets a random seed for reproducibility.
<code>.stop</code>	Character vector; specifies sequences where the model should stop generating further tokens.
<code>.frequency_penalty</code>	Numeric; adjusts the likelihood of repeating tokens (positive values decrease repetition).
<code>.presence_penalty</code>	Numeric; adjusts the likelihood of introducing new tokens (positive values encourage novelty).
<code>.id_prefix</code>	Character string to specify a prefix for generating custom IDs when names in <code>.llms</code> are missing

## Value

An updated and named list of `.llms` with identifiers that align with batch responses, including a `batch_id` attribute.

---

send\_claude\_batch      *Send a Batch of Messages to Claude API*

---

### Description

This function creates and submits a batch of messages to the Claude API for asynchronous processing.

### Usage

```
send_claude_batch(
    .llms,
    .model = "claude-sonnet-4-6",
    .max_tokens = 1024,
    .temperature = NULL,
    .top_k = NULL,
    .top_p = NULL,
    .stop_sequences = NULL,
    .json_schema = NULL,
    .thinking = FALSE,
    .thinking_budget = 1024,
    .api_url = "https://api.anthropic.com/",
    .verbose = FALSE,
    .dry_run = FALSE,
    .overwrite = FALSE,
    .max_tries = 3,
    .timeout = 60,
    .id_prefix = "tidyllm_claude_req_"
)
```

### Arguments

.llms	A list of LLMMessage objects containing conversation histories.
.model	Character string specifying the Claude model version (default: "claude-sonnet-4-6").
.max_tokens	Integer specifying the maximum tokens per response (default: 1024).
.temperature	Numeric between 0 and 1 controlling response randomness.
.top_k	Integer for diversity by limiting the top K tokens.
.top_p	Numeric between 0 and 1 for nucleus sampling.
.stop_sequences	Character vector of sequences that halt response generation.
.json_schema	A schema to enforce an output structure
.thinking	Logical; if TRUE, enables Claude's thinking mode for complex reasoning tasks (default: FALSE).

<code>.thinking_budget</code>	Integer specifying the maximum tokens Claude can spend on thinking (default: 1024). Must be at least 1024. Defaults to "tidyLLM_claude_req_".
<code>.api_url</code>	Base URL for the Claude API (default: "https://api.anthropic.com/").
<code>.verbose</code>	Logical; if TRUE, prints a message with the batch ID (default: FALSE).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.overwrite</code>	Logical; if TRUE, allows overwriting an existing batch ID associated with the request (default: FALSE).
<code>.max_tries</code>	Maximum number of retries to perform the request.
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).
<code>.id_prefix</code>	Character string to specify a prefix for generating custom IDs when names in <code>.llms</code> are missing.

**Value**

An updated and named list of `.llms` with identifiers that align with batch responses, including a `batch_id` attribute.

---

send_gemini_batch	<i>Submit a list of LLMMessage objects to Gemini's batch API</i>
-------------------	--

---

**Description**

Returns a named list (same as input) with `batch_id` and `json` attributes.

**Usage**

```
send_gemini_batch(
  .llms,
  .model = "gemini-2.5-flash",
  .temperature = NULL,
  .max_output_tokens = NULL,
  .top_p = NULL,
  .top_k = NULL,
  .presence_penalty = NULL,
  .frequency_penalty = NULL,
  .stop_sequences = NULL,
  .safety_settings = NULL,
  .json_schema = NULL,
  .grounding_threshold = NULL,
  .timeout = 120,
  .dry_run = FALSE,
  .max_tries = 3,
```

```

    .display = "tidyllum_batch",
    .id_prefix = "tidyllum_gemini_req_"
)

```

### Arguments

.llms	List of LLMMessage objects (named or unnamed).
.model	The model identifier (default: "gemini-1.5-flash").
.temperature	Controls randomness (default: NULL, range: 0-2).
.max_output_tokens	Maximum tokens in the response (default: NULL).
.top_p	Nucleus sampling (default: NULL, range: 0-1).
.top_k	Diversity in token selection (default: NULL).
.presence_penalty	Penalizes new tokens (default: NULL, -2 to 2).
.frequency_penalty	Penalizes frequent tokens (default: NULL, -2 to 2).
.stop_sequences	Character vector or NULL of up to 5.
.safety_settings	Optional list of safety settings (default: NULL).
.json_schema	Optional schema to enforce output structure.
.grounding_threshold	Optional grounding threshold (0-1) to enable Google Search.
.timeout	Timeout in seconds (default: 120).
.dry_run	If TRUE, returns the constructed request (default: FALSE).
.max_tries	Maximum retry attempts (default: 3).
.display	Display name for this batch (default: "tidyllum_batch").
.id_prefix	Prefix for message IDs (default: "tidyllum_gemini_req_").

### Value

Named list of LLMMessage objects with attributes batch\_id and json

---

send_groq_batch	<i>Send a Batch of Messages to the Groq API</i>
-----------------	---

---

### Description

This function creates and submits a batch of messages to the Groq API for asynchronous processing.

**Usage**

```

send_groq_batch(
  .llms,
  .model = "openai/gpt-oss-120b",
  .max_tokens = 1024,
  .temperature = NULL,
  .top_p = NULL,
  .frequency_penalty = NULL,
  .presence_penalty = NULL,
  .stop = NULL,
  .seed = NULL,
  .api_url = "https://api.groq.com/",
  .json_schema = NULL,
  .completion_window = "24h",
  .verbose = FALSE,
  .dry_run = FALSE,
  .overwrite = FALSE,
  .max_tries = 3,
  .timeout = 60,
  .id_prefix = "tidyllm_groq_req_"
)

```

**Arguments**

.llms	A list of LLMMessage objects containing conversation histories.
.model	Character string specifying the model to use (default: "deepseek-r1-distill-llama-70b").
.max_tokens	Integer specifying the maximum tokens per response (default: 1024).
.temperature	Numeric between 0 and 2 controlling response randomness.
.top_p	Numeric between 0 and 1 for nucleus sampling.
.frequency_penalty	Number between -2.0 and 2.0 to penalize repetition.
.presence_penalty	Number between -2.0 and 2.0 to encourage new topics.
.stop	One or more sequences where the API will stop generating further tokens.
.seed	An integer for deterministic sampling.
.api_url	Base URL for the Groq API (default: "https://api.groq.com/").
.json_schema	A list or tidyllm schema created with tidyllm_schema() for structured JSON output (optional).
.completion_window	Character string for the batch completion window (default: "24h").
.verbose	Logical; if TRUE, prints a message with the batch ID (default: FALSE).
.dry_run	Logical; if TRUE, returns the prepared request objects without executing (default: FALSE).

<code>.overwrite</code>	Logical; if TRUE, allows overwriting an existing batch ID (default: FALSE).
<code>.max_tries</code>	Maximum number of retries to perform the request.
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).
<code>.id_prefix</code>	Character string to specify a prefix for generating custom IDs when names in <code>.llms</code> are missing.

**Value**

An updated and named list of `.llms` with identifiers that align with batch responses, including a `batch_id` attribute.

---

<code>send_mistral_batch</code>	<i>Send a Batch of Requests to the Mistral API</i>
---------------------------------	--

---

**Description**

This function creates and submits a batch of messages to the Mistral API for asynchronous processing.

**Usage**

```
send_mistral_batch(
  .llms,
  .model = "mistral-small-latest",
  .endpoint = "/v1/chat/completions",
  .metadata = NULL,
  .temperature = 0.7,
  .top_p = 1,
  .max_tokens = 1024,
  .min_tokens = NULL,
  .frequency_penalty = NULL,
  .logit_bias = NULL,
  .presence_penalty = NULL,
  .seed = NULL,
  .stop = NULL,
  .safe_prompt = FALSE,
  .json_schema = NULL,
  .dry_run = FALSE,
  .overwrite = FALSE,
  .max_tries = 3,
  .timeout = 60,
  .id_prefix = "tidyllm_mistral_req_"
)
```

**Arguments**

.llms	A list of LLMMessage objects containing conversation histories.
.model	The Mistral model version (default: "mistral-small-latest").
.endpoint	The API endpoint (default: "/v1/chat/completions").
.metadata	Optional metadata for the batch.
.temperature	Sampling temperature to use, between 0.0 and 1.5 (default: 0.7).
.top_p	Nucleus sampling parameter, between 0.0 and 1.0 (default: 1).
.max_tokens	The maximum number of tokens to generate in the completion (default: 1024).
.min_tokens	The minimum number of tokens to generate (optional).
.frequency_penalty	Numeric value (or NULL) for frequency penalty.
.logit_bias	A named list modifying the likelihood of specific tokens (or NULL).
.presence_penalty	Numeric value (or NULL) for presence penalty.
.seed	Random seed for deterministic outputs (optional).
.stop	Sequence(s) at which to stop generation (optional).
.safe_prompt	Logical; if TRUE, injects a safety prompt (default: FALSE).
.json_schema	A JSON schema object for structured output (optional).
.dry_run	Logical; if TRUE, returns the prepared request without executing it (default: FALSE).
.overwrite	Logical; if TRUE, allows overwriting existing custom IDs (default: FALSE).
.max_tries	Maximum retry attempts for requests (default: 3).
.timeout	Request timeout in seconds (default: 60).
.id_prefix	Prefix for generating custom IDs (default: "tidyLLM_mistral_req_").

**Value**

The prepared LLMMessage list with a batch\_id attribute.

---

send_ollama_batch	<i>Send a Batch of Messages to Ollama API</i>
-------------------	---

---

**Description**

This function creates and submits a batch of messages to the Ollama API. Contrary to other batch functions, this function waits for the batch to finish and receives requests. The advantage compared to sending single messages via chat() is that Ollama handles large parallel requests quicker than many individual chat requests.

**Usage**

```

send_ollama_batch(
  .llms,
  .model = "qwen3.5:4b",
  .stream = FALSE,
  .seed = NULL,
  .json_schema = NULL,
  .temperature = NULL,
  .num_ctx = 2048,
  .num_predict = NULL,
  .top_k = NULL,
  .top_p = NULL,
  .min_p = NULL,
  .mirostat = NULL,
  .mirostat_eta = NULL,
  .mirostat_tau = NULL,
  .repeat_last_n = NULL,
  .repeat_penalty = NULL,
  .tfs_z = NULL,
  .stop = NULL,
  .ollama_server = "http://localhost:11434",
  .timeout = 120,
  .keep_alive = NULL,
  .dry_run = FALSE
)

```

**Arguments**

<code>.llms</code>	A list of LLMMessage objects containing conversation histories.
<code>.model</code>	Character string specifying the Ollama model to use (default: "qwen3-vl")
<code>.stream</code>	Logical; whether to stream the response (default: FALSE)
<code>.seed</code>	Integer; seed for reproducible generation (default: NULL)
<code>.json_schema</code>	A JSON schema object as R list to enforce the output structure (default: NULL)
<code>.temperature</code>	Float between 0-2; controls randomness in responses (default: NULL)
<code>.num_ctx</code>	Integer; sets the context window size (default: 2048)
<code>.num_predict</code>	Integer; maximum number of tokens to predict (default: NULL)
<code>.top_k</code>	Integer; controls diversity by limiting top tokens considered (default: NULL)
<code>.top_p</code>	Float between 0-1; nucleus sampling threshold (default: NULL)
<code>.min_p</code>	Float between 0-1; minimum probability threshold (default: NULL)
<code>.mirostat</code>	Integer (0,1,2); enables Mirostat sampling algorithm (default: NULL)
<code>.mirostat_eta</code>	Float; Mirostat learning rate (default: NULL)
<code>.mirostat_tau</code>	Float; Mirostat target entropy (default: NULL)
<code>.repeat_last_n</code>	Integer; tokens to look back for repetition (default: NULL)

.repeat_penalty	Float; penalty for repeated tokens (default: NULL)
.tfs_z	Float; tail free sampling parameter (default: NULL)
.stop	Character; custom stop sequence(s) (default: NULL)
.ollama_server	String; Ollama API endpoint (default: "http://localhost:11434")
.timeout	Integer; API request timeout in seconds (default: 120)
.keep_alive	Character; How long should the ollama model be kept in memory after request (default: NULL - 5 Minutes)
.dry_run	Logical; if TRUE, returns request object without execution (default: FALSE)

### Details

The function provides extensive control over the generation process through various parameters:

- Temperature (0-2): Higher values increase creativity, lower values make responses more focused
- Top-k/Top-p: Control diversity of generated text
- Mirostat: Advanced sampling algorithm for maintaining consistent complexity
- Repeat penalties: Prevent repetitive text
- Context window: Control how much previous conversation is considered

### Value

A list of updated LLMMessage objects, each with the assistant's response added if successful.

---

send_openai_batch	<i>Send a Batch of Messages to OpenAI Batch API</i>
-------------------	---

---

### Description

This function creates and submits a batch of messages to the OpenAI Batch API for asynchronous processing.

### Usage

```
send_openai_batch(
  .llms,
  .model = "gpt-5.4",
  .max_completion_tokens = NULL,
  .reasoning_effort = NULL,
  .frequency_penalty = NULL,
  .logit_bias = NULL,
  .presence_penalty = NULL,
  .seed = NULL,
  .stop = NULL,
```

```

.temperature = NULL,
.top_p = NULL,
.logprobs = NULL,
.top_logprobs = NULL,
.dry_run = FALSE,
.override = FALSE,
.json_schema = NULL,
.max_tries = 3,
.timeout = 60,
.verbose = FALSE,
.id_prefix = "tidyllm_openai_req_"
)

```

### Arguments

<code>.llms</code>	A list of LLMMessage objects containing conversation histories.
<code>.model</code>	Character string specifying the OpenAI model version (default: "gpt-5.1-chat-latest").
<code>.max_completion_tokens</code>	Integer specifying the maximum tokens per response (default: NULL).
<code>.reasoning_effort</code>	How long should reasoning models reason (can either be "low", "medium" or "high")
<code>.frequency_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency in the text so far.
<code>.logit_bias</code>	A named list modifying the likelihood of specified tokens appearing in the completion.
<code>.presence_penalty</code>	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far.
<code>.seed</code>	If specified, the system will make a best effort to sample deterministically.
<code>.stop</code>	Up to 4 sequences where the API will stop generating further tokens.
<code>.temperature</code>	What sampling temperature to use, between 0 and 2. Higher values make the output more random.
<code>.top_p</code>	An alternative to sampling with temperature, called nucleus sampling.
<code>.logprobs</code>	If TRUE, get the log probabilities of each output token (default: NULL).
<code>.top_logprobs</code>	If specified, get the top N log probabilities of each output token (0-5, default: NULL).
<code>.dry_run</code>	Logical; if TRUE, returns the prepared request object without executing it (default: FALSE).
<code>.overwrite</code>	Logical; if TRUE, allows overwriting an existing batch ID associated with the request (default: FALSE).
<code>.json_schema</code>	A JSON schema object provided by tidyllm_schema or ellmer schemata (default: NULL).

<code>.max_tries</code>	Maximum number of retries to perform the request (default: 3).
<code>.timeout</code>	Integer specifying the request timeout in seconds (default: 60).
<code>.verbose</code>	Logical; if TRUE, additional info about the requests is printed (default: FALSE).
<code>.id_prefix</code>	Character string to specify a prefix for generating custom IDs when names in <code>.llms</code> are missing (default: "tidyLLM_openai_req_").

**Value**

An updated and named list of `.llms` with identifiers that align with batch responses, including a `batch_id` attribute.

---

tidyLLM_schema	<i>Create a JSON Schema for Structured Outputs</i>
----------------	--

---

**Description**

This function creates a JSON schema for structured outputs, supporting both character-based shorthand and S7 `tidyLLM_field` objects. It also integrates with `ellmer` types like `ellmer::type_string()` if `ellmer` is in your namespace.

**Usage**

```
tidyLLM_schema(name = "tidyLLM_schema", ...)
```

**Arguments**

<code>name</code>	A character string specifying the schema name (default: "tidyLLM_schema").
<code>...</code>	Named arguments where each name represents a field, and each value is either a character string, a <code>tidyLLM_field</code> , or an <code>ellmer</code> type. Supported character shorthand types: <ul style="list-style-type: none"> <li>• "character" or "string" for character fields</li> <li>• "logical" for boolean fields</li> <li>• "numeric" for number fields</li> <li>• "factor(...)" for enumerations</li> <li>• Use <code>[]</code> to indicate vectors, e.g., "character[]"</li> </ul>

**Value**

A list representing the JSON schema, suitable for use with `.json_schema` in LLM API calls.

**Examples**

```
## Not run:
# Example using different field types
address_schema <- tidyllm_schema(
  name = "AddressSchema",
  Street = field_chr("A common street name"),
  house_number = field_dbl(),
  City = field_chr("Name of a city"),
  State = field_fct("State abbreviation", .levels = c("CA", "TX", "Other")),
  Country = "string",
  PostalCode = "string"
)

llm_message("Imagine an address") |> chat(openai, .json_schema = address_schema)

# Example with vector field
tidyllm_schema(
  plz = field_dbl(.vector = TRUE)
)

## End(Not run)
```

---

tidyllm\_tool

*Create a Tool Definition for tidyllm*


---

**Description**

Creates a tool definition for use with Language Model API calls that support function calling. This function wraps an existing R function with schema information for LLM interaction.

**Usage**

```
tidyllm_tool(.f, .description = character(0), ...)
```

**Arguments**

<code>.f</code>	The function to wrap as a tool
<code>.description</code>	Character string describing what the tool does
<code>...</code>	Named arguments providing schema definitions for each function parameter using <code>tidyllm_fields</code>

**Details**

Each parameter schema in `...` should correspond to a parameter in the wrapped function. All required function parameters must have corresponding schema definitions.

**Value**

A `TOOL` class object that can be used with `tidyllm chat()` functions

**Examples**

```

get_weather <- function(location){}
weather_tool <- tidyllm_tool(
  get_weather,
  "Get the current weather in a given location",
  location = field_chr("The city and state, e.g., San Francisco, CA")
)

```

---

upload_file	<i>Upload a File to a Provider's File Store</i>
-------------	---

---

**Description**

Upload a File to a Provider's File Store

**Usage**

```
upload_file(.provider, .path, ...)
```

**Arguments**

.provider	A provider function call (e.g. <code>claude()</code> , <code>gemini()</code> , <code>openai()</code> ).
.path	Path to the local file to upload.
...	Additional provider-specific arguments.

**Value**

A `tidyllm_file` object.

---

video_file	<i>Create a Video Object</i>
------------	------------------------------

---

**Description**

Stores a reference to a local video file for use in multimodal messages. The file is not encoded until the message is formatted for a provider.

**Usage**

```
video_file(.path)
```

**Arguments**

.path	The path to the video file on disk.
-------	-------------------------------------

**Value**

A tidyllm\_video object.

---

voyage	<i>Voyage Provider Function</i>
--------	---------------------------------

---

**Description**

The voyage() function acts as a provider interface for interacting with the Voyage.ai API through tidyllm's verbs. It dynamically routes requests to voyage-specific functions. At the moment this is only voyage\_embed()

**Usage**

```
voyage(..., .called_from = NULL)
```

**Arguments**

...	Parameters to be passed to the appropriate Voyage-specific function, such as model configuration, input text, or API-specific options.
.called_from	An internal argument specifying which action (e.g., embed) the function is invoked from. This argument is automatically managed by the tidyllm verbs and should not be modified by the user.

**Value**

The result of the requested action, depending on the specific function invoked

---

voyage_embedding	<i>Generate Embeddings Using Voyage AI API</i>
------------------	--

---

**Description**

This function creates embedding vectors from text or multimodal inputs (text and images) using the Voyage AI API. It supports three types of input:

**Usage**

```
voyage_embedding(
  .input,
  .model = "voyage-4",
  .output_dimension = NULL,
  .timeout = 120,
  .dry_run = FALSE,
  .max_tries = 3,
  .verbose = FALSE
)
```

**Arguments**

<code>.input</code>	Input to embed. Can be: <ul style="list-style-type: none"> <li>• A character vector of texts</li> <li>• An LLMMessage object (all textual components will be embedded)</li> <li>• A list containing a mix of character strings and tidyllm_image objects created with <code>img()</code></li> </ul>
<code>.model</code>	The embedding model identifier. For text-only: "voyage-3.5-lite" (default). For multimodal inputs: "voyage-multimodal-3.5" is used automatically.
<code>.output_dimension</code>	Optional integer to control output vector size (default: NULL, uses model default).
<code>.timeout</code>	Timeout for the API request in seconds (default: 120).
<code>.dry_run</code>	If TRUE, perform a dry run and return the request object without sending.
<code>.max_tries</code>	Maximum retry attempts for requests (default: 3).
<code>.verbose</code>	Should information about current rate limits be printed? (default: FALSE).

**Details**

1. Character vector: Embeds each text string separately
2. LLMMessage object: Extracts and embeds text content from messages
3. List of mixed content: Processes a combination of text strings and image objects created with `img()`

For multimodal inputs, the function automatically switches to Voyage's multimodal API and formats the response with appropriate labels (e.g., "[IMG] image.png") for images.

**Value**

A tibble with two columns: input and embeddings.

- The input column contains the input texts or image labels
- The embeddings column is a list column where each row contains an embedding vector

**Examples**

```
## Not run:
# Text embeddings
voyage_embedding("How does photosynthesis work?")

# Multimodal embeddings
list("A banana", img("banana.jpg"), "Yellow fruit") |>
  voyage_embedding()

## End(Not run)
```

---

`voyage_rerank`*Rerank Documents Using Voyage AI API*

---

**Description**

Rerank Documents Using Voyage AI API

**Usage**

```
voyage_rerank(  
  .query,  
  .documents,  
  .model = "rerank-2",  
  .top_k = NULL,  
  .api_key = Sys.getenv("VOYAGE_API_KEY"),  
  .timeout = 60,  
  .max_tries = 3  
)
```

**Arguments**

<code>.query</code>	A single character string representing the search query.
<code>.documents</code>	A character vector of documents to rerank.
<code>.model</code>	The reranking model identifier (default: "rerank-2").
<code>.top_k</code>	Integer; return only the top-k results (default: NULL, returns all).
<code>.api_key</code>	Character; Voyage API key (default: from environment).
<code>.timeout</code>	Integer; request timeout in seconds (default: 60).
<code>.max_tries</code>	Integer; maximum retries (default: 3).

**Value**

A tibble with columns `index`, `document`, and `relevance_score`, sorted by score descending.

# Index

## \* Message Creation Utilities

- df\_llm\_message, 28
- llm\_message, 69
  
- audio\_file, 5
- azure\_openai, 5
- azure\_openai\_chat, 6
- azure\_openai\_embedding, 8
  
- cancel\_openai\_batch, 9
- chat, 9
- chat\_completions
  - (chat\_completions\_chat), 11
- chat\_completions\_chat, 11
- chat\_ellmer, 12
- check\_azure\_openai\_batch, 13
- check\_batch, 14
- check\_claude\_batch, 14
- check\_gemini\_batch, 15
- check\_groq\_batch, 16
- check\_job, 17
- check\_mistral\_batch, 17
- check\_openai\_batch, 18
- claude, 19
- claude\_chat, 19
- claude\_delete\_file, 21
- claude\_file\_metadata, 22
- claude\_list\_files, 22
- claude\_list\_models, 23
- claude\_upload\_file, 24
- claude\_websearch, 24
  
- deep\_research, 27
- deepseek, 25
- deepseek\_chat, 25
- delete\_file, 28
- df\_llm\_message, 28, 70
- df\_llm\_message(), 70
  
- ellmer, 29
  
- ellmer\_tool, 29
- embed, 31
  
- fetch\_azure\_openai\_batch, 32
- fetch\_batch, 33
- fetch\_claude\_batch, 34
- fetch\_gemini\_batch, 35
- fetch\_groq\_batch, 35
- fetch\_job, 36
- fetch\_mistral\_batch, 37
- fetch\_openai\_batch, 37
- field\_chr, 38
- field\_dbl(field\_chr), 38
- field\_fct(field\_chr), 38
- field\_lgl(field\_chr), 38
- field\_object, 39
- file\_info, 40
  
- gemini, 40
- gemini\_chat, 41
- gemini\_delete\_file, 43
- gemini\_embedding, 43
- gemini\_file\_metadata, 44
- gemini\_list\_files, 44
- gemini\_list\_models, 45
- gemini\_upload\_file, 45
- get\_logprobs, 46
- get\_metadata, 47
- get\_metadata(), 46
- get\_reply, 48
- get\_reply(), 49
- get\_reply\_data, 48
- get\_reply\_data(), 48
- get\_user\_message, 49
- groq, 50
- groq\_chat, 50
- groq\_list\_models, 52
- groq\_transcribe, 53
  
- img, 54

last\_metadata (get\_metadata), 47  
last\_metadata(), 47  
last\_reply (get\_reply), 48  
last\_reply(), 48  
last\_reply\_data (get\_reply\_data), 48  
last\_reply\_data(), 49  
last\_user\_message (get\_user\_message), 49  
last\_user\_message(), 49, 50  
list\_azure\_openai\_batches, 55  
list\_batches, 55  
list\_claude\_batches, 56  
list\_files, 57  
list\_gemini\_batches, 57  
list\_groq\_batches, 58  
list\_hf\_gguf\_files, 58  
list\_mistral\_batches, 59  
list\_models, 60  
list\_openai\_batches, 60  
llamacpp, 61  
llamacpp\_chat, 61  
llamacpp\_delete\_model, 63  
llamacpp\_download\_model, 64  
llamacpp\_embedding, 64  
llamacpp\_health, 65  
llamacpp\_list\_local\_models, 66  
llamacpp\_list\_models, 66  
llamacpp\_rerank, 67  
llm\_message, 29, 69  
llm\_message(), 29  
LLMMessage, 68  
  
mistral, 70  
mistral\_chat, 70  
mistral\_embedding, 72  
mistral\_list\_models, 73  
  
ollama, 74  
ollama\_chat, 75  
ollama\_delete\_model, 77  
ollama\_download\_model, 77  
ollama\_embedding, 78  
ollama\_list\_models, 79  
openai, 79  
openai\_chat, 80  
openai\_check\_research, 81  
openai\_code\_interpreter, 82  
openai\_deep\_research, 82  
openai\_embedding, 83  
openai\_fetch\_research, 84  
  
openai\_list\_models, 84  
openai\_websearch, 85  
openrouter, 85  
openrouter\_chat, 86  
openrouter\_credits, 88  
openrouter\_embedding, 88  
openrouter\_generation, 89  
openrouter\_list\_models, 90  
  
pdf\_file, 90  
pdf\_page\_batch, 91  
perplexity, 92  
perplexity\_chat, 92  
perplexity\_check\_research, 95  
perplexity\_deep\_research, 96  
perplexity\_fetch\_research, 97  
  
rate\_limit\_info, 98  
  
send\_azure\_openai\_batch, 98  
send\_batch, 100  
send\_claude\_batch, 102  
send\_gemini\_batch, 103  
send\_groq\_batch, 104  
send\_mistral\_batch, 106  
send\_ollama\_batch, 107  
send\_openai\_batch, 109  
  
tidyllm\_schema, 111  
tidyllm\_tool, 112  
  
upload\_file, 113  
  
video\_file, 113  
voyage, 114  
voyage\_embedding, 114  
voyage\_rerank, 116