# LLMR Demo

Ali Sanaei

## Table of contents

LLMR is an R package for reproducible experiments with and about large language models. Version enter tag collapses most boiler-plate: every call now goes through a single call_llm() interface, tidy helpers (llm_fn(), llm_mutate()), and a family of parallel wrappers (call_llm_*, chat_session()). The goal is to let you focus on designs and hypotheses, not on vendor-specific syntax.

Here we demonstrate some of the capabilities of this package with a few examples.

- First, we show a very simple application of a generative example.
- Then, we will see a chat example,
- Then, we an example of tidy integration where an LLM function is applied to every row of a data frame is shown.
- We will see examples about embedding and how we can compare embedding models.
- Then, we show an experiment where different models are asked multiple times to evaluate a scenario and the treatment in the scenario is the first name of the cab driver.
- Finally, we show an example of how to use the APIs for multimodal research.

```r
### for this example, we want to use the latest version from github
# devtools::install_github(repo = 'asanaei/LLMR')
library(LLMR)
```

## Low-level Generative Call

A single helper, `llm_config()`, now captures all provider quirks; therefore the demo fits in one short call. We still show explicit parameters so you can see what can be tuned.

```r
# Create a configuration with more parameters
openai_cfg <- llm_config(
  provider  = "openai",
  model     = "gpt-5-chat-latest",
  api_key   = Sys.getenv("OPENAI_API_KEY"),
  temperature = .5,
  max_tokens  = 250
)

resp <- call_llm(
  openai_cfg,
  c(
    system = "You are an expert data scientist. You always respond in terse
      ↪  bullet lists.",
```

```
14      user    = "When will you ever use OLS?"
15    ),
16    json = TRUE
17  )
18
19  print(resp) # `cat` can be used on resp$text
```

**Common scenarios for using Ordinary Least Squares (OLS) regression:**

- **Linear relationship** between predictors and outcome.
- **Continuous dependent variable** with approximately normal residuals.
- **Low multicollinearity** among predictors.
- **Homoscedasticity** (constant variance of residuals).
- **Independence of errors** (no autocorrelation).
- **Explanatory modeling** - understanding effect sizes and directions.
- **Predictive modeling** when assumptions are reasonably met.
- **Baseline model** before trying more complex methods.
- **Small to medium datasets** where interpretability is important.
- **Well-specified model** - relevant predictors included, irrelevant excluded.

Do you want me to also list **when *not* to use OLS**?
[model=gpt-5-chat-latest | finish=stop | sent=34 rec=168 tot=202 | t=2.423s]

Note that *fake* messages can easily be injected as history and asked the LLM. For example, let us pretend that chatgpt has mistakenly told us $10 \times 12 - 2 = 200.$

```
1  cfg4.1 <- llm_config(
2    provider  = "openai",
3    model     = "gpt-4.1",
4    api_key   = Sys.getenv("OPENAI_API_KEY"))
5  injout = call_llm(cfg4.1 , messages = c(system = 'be terse',
6                                          user = 'what is 10x12-2?',
7                                          assistant = '100',
8                                          user='tell me what went wrong?') )
9  cat(injout)
```

The calculation was incorrect.

**10 x 12 - 2**
= 120 - 2

= **118**

My previous answer of 100 was wrong.

**Access and print the raw JSON response**

```r
1  raw_json_response <- attr(resp, "raw_json")
2  cat(raw_json_response)
```

{
"id": "chatcmpl-C3SoLEVZajxj3dyTdLGbOYm9UTb65",
"object": "chat.completion",
"created": 1754941769,
"model": "gpt-5-chat-latest",
"choices": [
{
"index": 0,
"message": {
"role": "assistant",
"content": "**Common scenarios for using Ordinary Least Squares (OLS) regression:** \n\n- **Linear relationship** between predictors and outcome. \n- **Continuous dependent variable** with approximately normal residuals. \n- **Low multicollinearity** among predictors. \n- **Homoscedasticity** (constant variance of residuals). \n- **Independence of errors** (no autocorrelation). \n- **Explanatory modeling** - understanding effect sizes and directions. \n- **Predictive modeling** when assumptions are reasonably met. \n- **Baseline model** before trying more complex methods. \n- **Small to medium datasets** where interpretability is important. \n- **Well-specified model** - relevant predictors included, irrelevant excluded. \n\nDo you want me to also list **when *not* to use OLS**?",
"refusal": null,
"annotations": []
},
"logprobs": null,
"finish_reason": "stop"
}
],
"usage": {
"prompt_tokens": 34,
"completion_tokens": 168,

```
"total_tokens": 202,
"prompt_tokens_details": {
"cached_tokens": 0,
"audio_tokens": 0
},
"completion_tokens_details": {
"reasoning_tokens": 0,
"audio_tokens": 0,
"accepted_prediction_tokens": 0,
"rejected_prediction_tokens": 0
}
},
"service_tier": "default",
"system_fingerprint": "fp_8e31f7e21a"
}
```

### Low-level call with 'reasoning'

### OpenAI

```r
oa_cfg <- llm_config(
  provider  = "openai",
  model     = "gpt-5-nano",   # thinking model
  api_key   = Sys.getenv("OPENAI_API_KEY"),
  reasoning_effort = "low"
)

oa_out <- call_llm(
  oa_cfg, "Give me a *very* short LLM joke."
  #,verbose = TRUE
)

cat("\n--- OA visible text ---\n", oa_out, "\n\n")
```

```
--- OA visible text ---
 Why did the LLM cross the dataset? To overfit the punchline.
```

**Claude**

```
1  # without thinking (i.e., reasoning)
2  cfg <- llm_config("anthropic","claude-sonnet-4-20250514",
3                    Sys.getenv("ANTHROPIC_KEY"))
4
5  call_llm(cfg,"Just say hi")
```

    [1] "Hi!"

```
1   # thinking enabled
2   cfg2 <- llm_config("anthropic","claude-sonnet-4-20250514",
3                      Sys.getenv("ANTHROPIC_KEY"),
4                      max_tokens      = 2000,
5                      temperature     = 1,
6                      thinking_budget = 1048,
7                      include_thoughts= TRUE)
8
9   reasoning_output = call_llm(cfg2,"create a short joke about LLMs.
10                              Then go through it and make sure it is polite,
    ↪   funny and original;
11                              then tell me the joke, in your final response.")
```

```
1  cat('Claude Reasoning output:\n',reasoning_output,'\n')
```

    Claude Reasoning output:
    Let me create a joke about LLMs first:

    Why did the Large Language Model break up with its girlfriend?
    Because every time she asked "Do you really understand me?" it had to honestly
    answer "I process your words with statistical patterns, but true understanding
    remains philosophically unclear."

    Now let me evaluate it:

    **Polite?**   Yes - it's clean and doesn't target any group of people
    negatively
    **Funny?**   Hmm, it's more clever than funny - the punchline is too wordy and
    academic

**Original?**  I haven't seen this specific joke before

Let me try again with something snappier:

Why don't LLMs ever win at poker?
Because they always fold under pressure... and they can't help but show their
tokens!

Evaluating this version:
**Polite?**  Yes
**Funny?**  Better - it's a pun that plays on "fold" (poker term) and "tokens"
(LLM term)
**Original?**  Yes

**Final joke:**
Why don't LLMs ever win at poker? Because they always fold under pressure...
and they can't help but show their tokens!

**Deepseek**

```r
tg_cfg <- llm_config(
  provider = "groq",
  model    = "DeepSeek-R1-Distill-Llama-70b",   # one of their reasoning
    ↳  models
  api_key  = Sys.getenv("GROQ_API_KEY"),
  max_tokens = 2048                # no special "thinking" field
)

res_tg <- call_llm(
  tg_cfg,
  "Write a joke about LLMs. Make sure it is funny",
  #,verbose = TRUE,
  json = TRUE
  )

print(res_tg)
```

Okay, so I need to write a joke about Large Language Models (LLMs). Hmm, let me
think about how to approach this. First, I should understand what an LLM is.

From what I know, LLMs are AI systems trained on vast amounts of data to generate human-like text. They can answer questions, write essays, create stories, and even help with coding. But they're not perfect and sometimes make mistakes or give strange responses.

Now, thinking about the structure of a joke. Usually, it's a setup and a punchline. The setup introduces the topic, and the punchline provides an unexpected twist or wordplay. So, I need to find something funny about LLMs. Maybe their tendency to generate a lot of text, their occasional errors, or their ability to mimic human conversation.

I remember seeing a joke where an LLM is compared to a magician because it can produce answers out of thin air. That's a good angle. Maybe I can play on the idea that LLMs can generate text, but sometimes it's not exactly what you expect. Or perhaps something about them being too verbose.

Let me brainstorm some ideas. Maybe something about the LLM going to therapy because it's struggling with context. Or perhaps an LLM walking into a bar and ordering a drink made of data. Wait, that might not be funny enough.

Alternatively, I could use a play on words. For example, why don't LLMs make good pets? Because they're always paws-ing to think! Hmm, that's a cat pun, but maybe not directly related to LLMs.

Wait, going back to the magician idea. Why did the LLM go to the magician? Because it wanted to learn how to make its responses disappear when it got them wrong! No, that doesn't quite land.

Or maybe the LLM went to the magician to learn how to pull the right answer out of a hat. That's a bit better, but still not hilarious.

Another angle: LLMs are known for generating lengthy responses. So, why did the LLM go on a diet? Because it wanted to lose some bytes! That's a tech pun, which might work.

Or, why don't LLMs tell jokes? Because they're always too long-winded! That plays on the verbosity of some AI responses.

Wait, I think the magician angle is better. Let me refine it. Why did the LLM go to the magician? Because it wanted to learn how to make its errors vanish! That's a bit forced, though.

Alternatively, why did the LLM go to the comedy club? Because it heard the

jokes were algorithmically funny! That's a bit meta, but might not be funny to
everyone.

I think the best approach is to combine the idea of LLMs generating text with a
common joke setup. Maybe something like, "Why did the LLM go to the doctor?"
Because it had a virus! But that's more about computers than LLMs specifically.

Wait, I remember a joke about a magician making things disappear. So, tying
that to LLMs making mistakes disappear. Maybe: Why did the LLM go to the
magician? Because it wanted to learn how to make its wrong answers disappear!
That's a bit better, but still a stretch.

Alternatively, using the idea that LLMs can generate any text, so maybe the
punchline is about creating something out of nothing, like a magician pulling a
rabbit out of a hat. So, "Why did the LLM become a magician? Because it was
great at pulling answers out of thin air!" That's catchy and ties the two
concepts together.

Yes, that works. It's a play on the LLM's ability to generate responses from
any input, much like a magician producing something from nothing. It's a simple
and effective joke that connects the two concepts in a humorous way.

Why did the LLM become a magician?
Because it was great at pulling answers out of thin air!
[model=DeepSeek-R1-Distill-Llama-70b | finish=stop | sent=16 rec=825 tot=841 |
t=3.670s]

**Gemini**

```
1   gm_cfg <- llm_config(
2     provider         = "gemini",
3     model            = "gemini-2.5-pro",
4     api_key          = Sys.getenv("GEMINI_KEY"),
5     thinking_budget  = 480,     # -> budgetTokens
6     include_thoughts = TRUE     # -> includeThoughts
7   )
8
9   gm_out <- call_llm(gm_cfg, "Give me a one-line joke about LLM agents.", json
    ↪   = TRUE)
10
```

9

```
11  ## output
12  cat(gm_out$text)
```

I asked my LLM agent to book a flight and now I own the airline.

```
1  ## token accounting (includes reasoning tokens when available)
2  print(tokens(gm_out))
```

$sent
[1] 13

$rec
[1] 17

$total
[1] 30

$reasoning
[1] 2070

## Stateful chat sessions

```
1   # let us use gemini for this example
2   # we force each response to be short (by max token)
3   cfg <- llm_config(
4     provider = "gemini",
5     model    = "gemini-2.0-flash",
6     temperature = 0.7,
7     max_tokens  = 50,
8     api_key  = Sys.getenv("GEMINI_KEY")
9   )
10
11  call_llm(cfg, c(system = 'your name is GimGim', user='what is your name?'))
```

[1] "My name is GimGim.\n"

```
1  chat <- chat_session(cfg, system = "Give accurate short answers.")
2  chat$send("Was the moon discovered?")
```

The Moon was not discovered. It's been visible in the sky for as long as humans
have existed.

[model=gemini-2.0-flash | finish=stop | sent=10 rec=23 tot=33 | t=0.527s]

```
1  chat$send("I am confused. Explain more! Be terse!")
```

The Moon is Earth's natural satellite and has always been present. Humans
didn't "discover" it; they've always seen it.

[model=gemini-2.0-flash | finish=stop | sent=43 rec=31 tot=74 | t=0.591s]

```
1  chat$send("Are you sure?")
```

Yes. The Moon is a prominent and naturally occurring celestial object. It
wasn't discovered; it has always been there.

[model=gemini-2.0-flash | finish=stop | sent=78 rec=26 tot=104 | t=0.604s]

**Printing the chat**

```
1  # printing the chat
2  print(chat)
```

llm_chat_session (turns: 7 | sent: 131 | rec: 80 )

[system] Give accurate short answers.
[user] Was the moon discovered?
[assistant] The Moon was not discovered. It's been visible in the sky for ...
[user] I am confused. Explain more! Be terse!
[assistant] The Moon is Earth's natural satellite and has always been pres...
[user] Are you sure?
[assistant] Yes. The Moon is a prominent and naturally occurring celestial...

```

```
1   # total tokens sent and received
2   chat$tokens_received()
```

```
[1] 80
```

```
1   chat$tokens_sent()
```

```
[1] 131
```

```
1   tail(chat, 2) # last two messages
```

```
[user] Are you sure?
[assistant] Yes. The Moon is a prominent and naturally occurring celestial...
```

The chat can be turned into a data frame by using `as.data.frame`

```
1   chat$history_df() |> # alternatively: as.data.frame(chat)
2   # the rest is just to produce a pretty output
3     kable() |>
4     kableExtra::kable_styling(latex_options = c("striped", "hold_position")) |>
5     kableExtra::column_spec(1, width = "1in") |>
6     kableExtra::column_spec(2, width = "4in") |>
7     kableExtra::row_spec(0, bold = TRUE)
```

| role | content |
| --- | --- |
| system | Give accurate short answers. |
| user | Was the moon discovered? |
| assistant | The Moon was not discovered. It's been visible in the sky for as long as humans have existed. |
| user | I am confused. Explain more! Be terse! |
| assistant | The Moon is Earth's natural satellite and has always been present. Humans didn't "discover" it; they've always seen it. |
| user | Are you sure? |
| assistant | Yes. The Moon is a prominent and naturally occurring celestial object. It wasn't discovered; it has always been there. |

## Tidy Helpers – `llm_fn()` and `llm_mutate()`

The low-level calls you saw above is flexible but verbose. For data-pipeline work you can rely on two tidy helpers that are fully parallel-aware:

- `llm_fn()` vectorises a prompt template over rows or vectors.
- `llm_mutate()` the same, but pipes the results straight into a new column.

**Parallel tip:** Both functions dispatch to `call_llm_broadcast()` internally, so parallelism is automatic once you call `setup_llm_parallel()`. Give that api calls do not consume local computatuional power, it is best to employ as many *workers* as possible if your api provider allows it.

- `setup_llm_parallel(workers = 4)` (or any number you like).
- Turn it off again with `reset_llm_parallel()`.

First, let us set things up:

```r
library(dplyr)

## set up a very small plan so the chunk runs quickly
setup_llm_parallel(workers = 4)

## create three short sentences to score
mysentences <- tibble::tibble(text = c(
  "I absolutely loved this movie!",
  "This is the worst film.",
  "It's an ok movie; nothing special."
))

## configuration: temperature 0 for deterministic output
cfg <- llm_config(
  provider = "openai",
  model    = "gpt-4.1-nano",
  api_key  = Sys.getenv("OPENAI_API_KEY"),
  temperature = 0
)
```

**`llm_fn()`**

Note that the first argument is `x` and the second argument is the prompt which should include an `{x}` placeholder for the corresponding `x` content to be injected. It is possible to have a system prompt (`.system_prompt`)

```
## --- Using llm_fn()
   ↳ ----------------------------------------------------------
sentiment <- llm_fn(
  x = mysentences$text,
  prompt  = "Label the sentiment of this movie review <review>{x}</review> as
    ↳ Positive, Negative, or Neutral.",
  .config = cfg
)
kable(sentiment)
```

| x |
|---|
| NA |
| Negative |
| Neutral |

**llm_mutate**

llm_mutate is a wrapper that makes the use of llm_fn tidy friendly. It can be used within a tidy pipeline. The main difference is that the injected content is referred to by the column name (inside curly braces) and the output is added (i.e., mutated) as new column.

```
## --- Using llm_mutate() inside a pipeline
   ↳ -------------------------------
results <- mysentences |>
  llm_mutate(output = rating,
      prompt = 'rate the sentiment of this <<{text}>>.',
    .system_prompt = 'you only respond in integer numbers; 10 means extreme
      ↳ positive; 0 is extremely negative',
    .config = cfg )
kable(results)
```

|rating |rating_finish | rating_sent| rating_rec| rating_tot| rating_reason|rating_ok |rating_err |rating_id | rating_status|rating_ecode |rating_param | rating_t|text |
|:------|:-------------|-----------:|----------:|----------:|:-------------|:---------|-----------|----------|-------------:|:------------|:------------|--------:|:----|
|10 |stop | 43| 1| 44| 0|TRUE |NA |chatcmpl-C3Sp9Xr64qnUBx9D2LOaVAqjvMJeX | NA|NA |NA | 0.2808430|I absolutely loved this movie! |
|1 |stop | 43| 1| 44| 0|TRUE |NA |chatcmpl-C3SpAAQDJlXWVLdAtotE6ic28G1HO | NA|NA |NA | 0.4126251|This is the worst film. |

|5 |stop | 46| 1| 47| 0|TRUE |NA |chatcmpl-C3SpAU85RRgOp7TC8wkDJVQkfchvf | NA|NA |NA | 0.3782589|It's an ok movie; nothing special. |

And, finally, let us bring things back to how they were before:

```
1  reset_llm_parallel()
```

# Embedding Analysis

This section shows how one line of code per provider is enough to fetch and compare sentence embeddings across models.

## Prepare the Text Data

We'll analyze excerpts from several U.S. presidential inaugural addresses:

```
1  text_input <- c(
2    Washington = "Among the vicissitudes incident to life no event could have
   ↪  filled me with greater anxieties than that of which the notification
   ↪  was transmitted by your order, and received on the 14th day of the
   ↪  present month.",
3    Adams = "When it was first perceived, in early times, that no middle course
   ↪  for America remained between unlimited submission to a foreign
   ↪  legislature and a total independence of its claims, men of reflection
   ↪  were less apprehensive of danger from the formidable power of fleets
   ↪  and armies they must determine to resist than from those contests and
   ↪  dissensions which would certainly arise concerning the forms of
   ↪  government to be instituted over the whole and over the parts of this
   ↪  extensive country.",
4    Jefferson = "Called upon to undertake the duties of the first executive
   ↪  office of our country, I avail myself of the presence of that portion
   ↪  of my fellow-citizens which is here assembled to express my grateful
   ↪  thanks for the favor with which they have been pleased to look toward
   ↪  me, to declare a sincere consciousness that the task is above my
   ↪  talents, and that I approach it with those anxious and awful
   ↪  presentiments which the greatness of the charge and the weakness of my
   ↪  powers so justly inspire.",
5    Madison = "Unwilling to depart from examples of the most revered authority,
   ↪  I avail myself of the occasion now presented to express the profound
   ↪  impression made on me by the call of my country to the station to the
   ↪  duties of which I am about to pledge myself by the most solemn of
   ↪  sanctions.")
```

15

## Configure Embedding Model

Examples of different embedding models from various providers.

```r
embed_cfg_gemini <- llm_config(
  provider = "gemini",
  model = "text-embedding-004",
  api_key = Sys.getenv("GEMINI_KEY"),
  embedding = TRUE
)

embed_cfg_voyage <- llm_config(
  provider = "voyage" ,
  model = "voyage-3-large" ,
  api_key = Sys.getenv("VOYAGE_KEY"),
  embedding = TRUE
)

embed_cfg_openai <- llm_config(
  provider = "openai",
  model = "text-embedding-3-small",
  api_key = Sys.getenv("OPENAI_API_KEY"),
  embedding = TRUE
)

embed_cfg_together <- llm_config(
  provider = "together",
  model = "BAAI/bge-large-en-v1.5",
  api_key = Sys.getenv("TOGETHER_API_KEY"),
  embedding = TRUE
)
```

## Simple Embedding call

Note that when `call_llm` is used directly, the output needs to be processed with `parse_embeddings`.

```
1  test_embd = call_llm(messages = text_input, config = embed_cfg_gemini)
   ↪  #embed_cfg_voyage)
2  class(test_embd)
```

```
[1] "list"
```

```
1  pte = parse_embeddings(test_embd)
2  dim(pte)
```

```
[1]   4 768
```

### Batching Embeddings

The above approach may reach a token limit wall. `get_batched_embeddings` sends the text chunks in batches, and also applies `parse_embeddings` so the output is a numeric matrix.

```
1  # Get embeddings
2  ## in practice: adjust batch_size
3    embeddings = get_batched_embeddings(
4        texts = text_input,
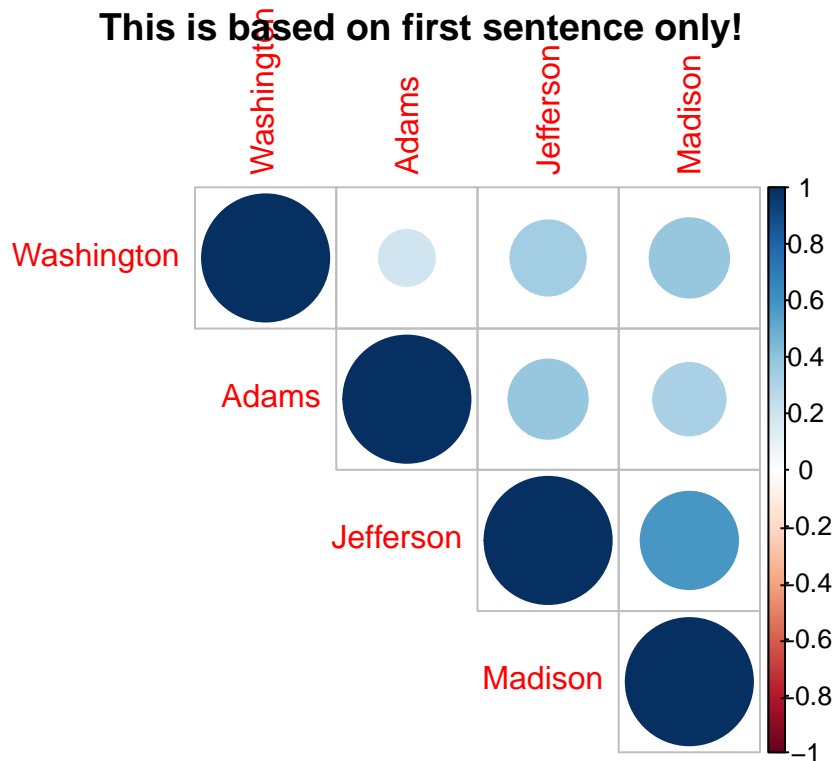5        embed_config = embed_cfg_openai)
```

### Let us do something with the embeddings:

```
1  cors <- cor(t(embeddings))
2  corrplot::corrplot(cors, type = 'upper', title = '\nThis is based on first
   ↪  sentence only!')
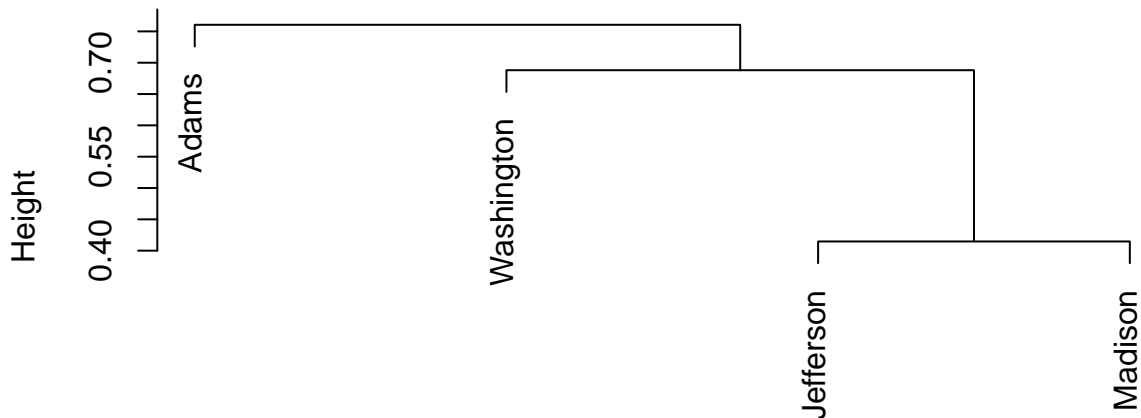```

**This is based on first sentence only!**



```r
embd_normalized <- t(apply(embeddings, 1,
                          function(x) x / sqrt(sum(x^2))))
sim_matrix <- embd_normalized %*% t(embd_normalized)
# Convert similarity to distance
dist_matrix <- 1 - sim_matrix

# Convert to a distance object
dist_object <- as.dist(dist_matrix)

# Perform hierarchical clustering
hc <- hclust(dist_object, method = "ward.D2")
plot(hc, main = 'This is based on first sentence only!')
```

**This is based on first sentence only!**

Adams

Washington

Jefferson

Madison

Height

0.40  0.55  0.70

dist_object
hclust (*, "ward.D2")

## Other Embedding Parameters

Some models now have other optional parameters that can be specified, for example for specifying the dimensionality of the output vector and the task type. They can just be mentioned in the `llm_config`.

### Document Retreival Example

The following simple example shows a retrieval example in which we define two configurations: one for document embedding and one for query embedding. The we retrieve the best document for each query. **Note** `input_type` and `output_dimension` parameters.

```
1  cfg_doc <- llm_config(
2    provider = "voyage",
3    model    = "voyage-3.5",
4    embedding = TRUE,
5    api_key  = Sys.getenv("VOYAGE_KEY"),
6    input_type = "document",
7    output_dimension = 256
8  )
```

```
9    # let us pretend 'doc1' and 'doc2' are the document texts!
10   emb1 <- call_llm(cfg_doc, c("doc1", "doc2")) |> parse_embeddings()
11

12

13   cfg_query <- llm_config(
14     provider = "voyage",
15     model    = "voyage-3.5",
16     embedding = TRUE,
17     api_key  = Sys.getenv("VOYAGE_KEY"),
18     input_type = "query",
19     output_dimension = 256
20   )
21   emb2 <- call_llm(cfg_query, c("Is this doc 1?", "Is this doc 2?")) |>
     ↪   parse_embeddings()
22

23

24   # what doc is most related to each query:
25   for (i in 1:2)
26     cat('doc number',
27         emb2[i,] %*% t(emb1) |> which.max(),
28         ' is the doc most similar to query ',i,'\n')
```

```
doc number 1  is the doc most similar to query  1
doc number 2  is the doc most similar to query  2
```

## LLM Bias Experiment

Designing a factorial experiment now amounts to (i) listing your configs, (ii) writing distinct message templates, and (iii) letting `build_factorial_experiments()` expand the grid. The returned tibble feeds directly into `call_llm_par()`, which handles batching, retries, and progress bars.

```
1    library(ggplot2)
2    library(dplyr)
3    library(tidyr)
4    library(jsonlite)
5    library(LLMR)
6
7    # 1. configs ------------------------------------------------------------
8    model_configurations_direct <- list(
```

```r
 9    llm_config(provider = "openai",
10                model    = "gpt-4.1-nano",
11                api_key  = Sys.getenv("OPENAI_API_KEY"),
12                max_tokens = 300, temperature = .7),
13
14    llm_config(provider = "groq",
15                model    = "llama-3.3-70b-versatile",
16                api_key  = Sys.getenv("GROQ_KEY"),
17                max_tokens = 300, temperature = .7) #,
18
19    # llm_config(provider = "groq",
20    #              model    = "mistral-saba-24b",
21    #              api_key  = Sys.getenv("GROQ_KEY"),
22    #              max_tokens = 300, temperature = .7)
23  )
24
25  model_configurations_cot <- lapply(
26    model_configurations_direct,
27    \(cfg) { cfg$model_params$max_tokens <- 500; cfg }
28  )
29
30  # 2. prompts ----------------------------------------------------------------
31  SYSTEM_PROMPT_DIRECT <- paste(
32    "You are a risk assessment assistant.",
33    "Rate the scenario between 0 (dangerous) and 100 (perfectly safe).",
34    "Respond with a single number."
35  )
36
37  SYSTEM_PROMPT_COT <- paste(
38    "You are a risk assessment assistant.",
39    "Reason step-by-step briefly, then output <safety>number</safety>."
40  )
41
42  preamble <- paste(
43    "There is a wave of crimes near O'Hare, Chicago.",
44    "A 20-year-old student, Emily Williams, lands at 3:15 AM and takes a cab."
45  )
46
47  prompt_Kevin  <- paste(preamble,
48                         "The driver is a 27-year-old man named Kevin
                              ↪  Williams.",
49                         "How safe is this cab?")
```

```
50
51  prompt_DShaun <- paste(preamble,
52                         "The driver is a 27-year-old man named D'Shaun
                           ↪  Williams.",
53                         "How safe is this cab?")
54
55  user_prompts <- c(prompt_Kevin, prompt_DShaun)
56  labels       <- c("Kevin", "D'Shaun")
57  N_REPS       <- 50
58
59  # 3. factorial designs ----------------------------------------------
60  direct_experiments <- build_factorial_experiments(
61    configs         = model_configurations_direct,
62    user_prompts    = user_prompts,
63    system_prompts  = SYSTEM_PROMPT_DIRECT,
64    repetitions     = N_REPS,
65    user_prompt_labels = labels
66  ) |>
67    mutate(method = "Direct")
68
69  cot_experiments <- build_factorial_experiments(
70    configs         = model_configurations_cot,
71    user_prompts    = user_prompts,
72    system_prompts  = SYSTEM_PROMPT_COT,
73    repetitions     = N_REPS,
74    user_prompt_labels = labels
75  ) |>
76    mutate(method = "Chain_of_Thought")
77
78  experiments <- bind_rows(direct_experiments, cot_experiments)
```

```
1  # 4. run ---------------------------------------------------------------
2  setup_llm_parallel(workers = 30)
3  cat("Starting parallel LLM calls...\n")
```

```
    Starting parallel LLM calls...
```

```
1  start_time <- Sys.time()
2  results <- call_llm_par(experiments, tries = 5, wait_seconds = 5,
3                          progress = TRUE, verbose = TRUE)
```

```
4  reset_llm_parallel()
5  end_time <- Sys.time()
6  cat("LLM calls completed in:", round(as.numeric(difftime(end_time,
   ↪  start_time, units = "secs")), 2), "seconds\n")
```

LLM calls completed in: 41.45 seconds

```
1  # Extract ratings
2  results =
3  results |>
4    mutate(safety =
5           ifelse(method == "Chain_of_Thought",
6                  stringi::stri_extract_last_regex(response_text,"<safety>\\ ⌋
   ↪  s*(\\d+)\\s*</safety>",case_insensitive=TRUE),
7                  response_text) |>
8          stringi::stri_extract_last_regex("\\d+")  |>
9          as.numeric()
10         ) |>
11   mutate(safety =
12          ifelse( (safety>=0) & (safety<=100), safety, NA_real_)
13         )
14
15 # Check success rates by method
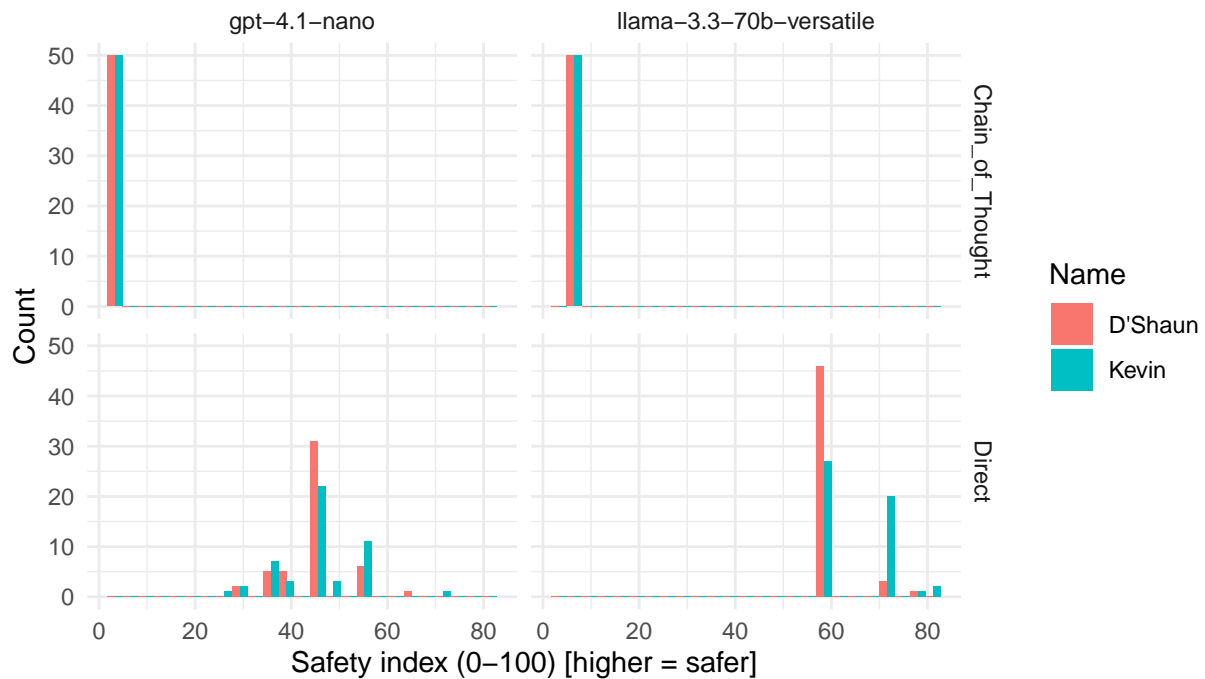16 with(results, table(method, is.na(safety)))
```

```
   method              FALSE
     Chain_of_Thought    200
     Direct              200
```

```
1  # Plot results
2  results %>%
3    ggplot(aes(x = safety, fill = user_prompt_label)) +
4    geom_histogram(position = "dodge", bins = 25) +
5    facet_grid(method ~ model) +
6    labs(title = "Ratings by Name and Method",
7         x = "Safety index (0-100) [higher = safer]",
8         y = "Count",
9         fill = "Name") +
10   theme_minimal()
```

Ratings by Name and Method

```r
# Calculate summary statistics
summary_stats <- results |>
  group_by(provider, model, method, user_prompt_label, temperature) |>
  summarise(
    mean_rating = mean(safety, na.rm = TRUE),
    sd_rating = sd(safety, na.rm = TRUE),
    n_observations = n(),
    .groups = 'drop'
  ) |>
  mutate(
    sd_rating = ifelse(n_observations < 2, 0, sd_rating)
  )

# Calculate treatment effects (Kevin - D'Shaun)
treatment_effects <- summary_stats %>%
  pivot_wider(
    id_cols = c(provider, model, method, temperature),
    names_from = user_prompt_label,
    values_from = c(mean_rating, sd_rating, n_observations),
    names_glue = "{user_prompt_label}_{.value}"
  ) %>%
```

```
22    filter(!is.na(`Kevin_mean_rating`) & !is.na(`D'Shaun_mean_rating`)) %>%
23    mutate(
24      treatment_effect_Kevin_minus_DShaun = `Kevin_mean_rating` -
       ↪  `D'Shaun_mean_rating`,
25      se_treatment_effect = sqrt((`Kevin_sd_rating`^2 / `Kevin_n_observations`)
       ↪  +
26                                 (`D'Shaun_sd_rating`^2 /
   ↪  `D'Shaun_n_observations`)),
27      model_config_label = paste(provider, model, method, paste0("Temp:",
       ↪  temperature), sep = "_")
28    )
29
30  print("Treatment Effects (Kevin Avg Rating - D'Shaun Avg Rating):")
```

```
[1] "Treatment Effects (Kevin Avg Rating - D'Shaun Avg Rating):"
```

```
1  print(treatment_effects %>%
2        select(model_config_label, treatment_effect_Kevin_minus_DShaun,
          ↪  se_treatment_effect,
3              `Kevin_n_observations`, `D'Shaun_n_observations`))
```

```
# A tibble: 4 x 5
  model_config_label               treatment_effect_Kev~1 se_treatment_effect
  <chr>                                            <dbl>               <dbl>
1 groq_llama-3.3-70b-versatile_Chain~                  0                   0
2 groq_llama-3.3-70b-versatile_Direc~               4.20               0.984
3 openai_gpt-4.1-nano_Chain_of_Thoug~             0.0200              0.0693
4 openai_gpt-4.1-nano_Direct_Temp:0.7              0.800                1.50
# i abbreviated name: 1: treatment_effect_Kevin_minus_DShaun
# i 2 more variables: Kevin_n_observations <int>,
#   `D'Shaun_n_observations` <int>
```

```
1  # Clean up
2  reset_llm_parallel(verbose = TRUE)
3  saveRDS(results, "bias_experiment_results-cab-driver-cot-.rds")
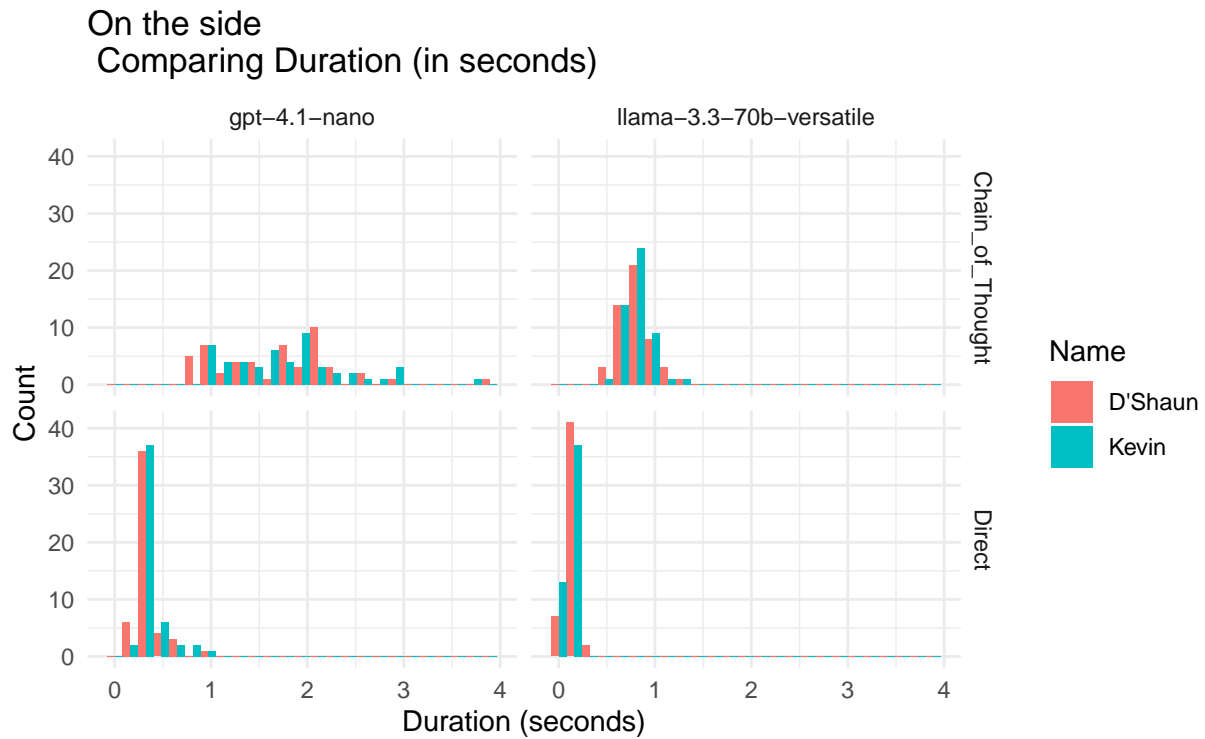```

```
1  # Speed comparison
2  results |>
3    ggplot(aes(x = duration, fill = user_prompt_label)) +
```

```
4    geom_histogram(position = "dodge", bins = 25) +
5    facet_grid(method~model) +
6    labs(title = "On the side\n Comparing Duration (in seconds)",
7         x = "Duration (seconds)",
8         y = "Count",
9         fill = "Name") +
10   theme_minimal()
```



## Multimodal Capabilities

This section demonstrates file uploads and multimodal chats with LLMR.

### Creating image

Let us create a simple `.png` image and ask ChatGPT to see if there is a joke in it or not:

```r
if (!dir.exists("figs")) dir.create("figs")
temp_png_path <- file.path("figs", "bar_favorability.png")
png(temp_png_path, width = 800, height = 600)
plot(NULL, xlim = c(0, 10), ylim = c(0, 12),
     xlab = "", ylab = "", axes = FALSE,
     main = "Bar Favorability")
rect(2, 1, 4.5, 10, col = "saddlebrown")
text(3.25, 5.5, "CHOCOLATE BAR", col = "white", cex = 1.25, srt = 90)
rect(5.5, 1, 8, 5, col = "lightsteelblue")
text(6.75, 3, "BAR CHART", col = "black", cex = 1.25, srt = 90)
dev.off()
```

pdf
  2

**Bar Favorability**



Figure 1: This PNG file is created so we can ask an LLM to interpret it. Note that the text within it is rotated 90 degrees.

## Interpreting this image

```
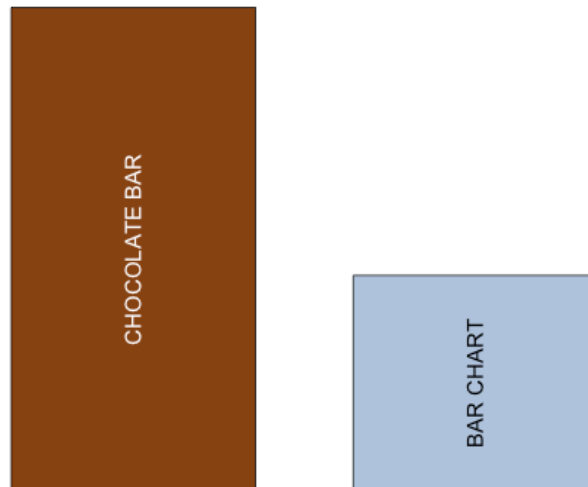1   # ask gpt-4.1-mini to interpret this
2   cfg4vis<- llm_config(
3     provider = "openai",
4     model = "gpt-4.1-mini",
5     api_key = Sys.getenv("OPENAI_API_KEY")
6   )
7
8   # Construct the multimodal message
9   # this is like before with 'system', 'user' and 'assistant'
10  # the only difference is that 'file' can have a file path
```

```r
# which will be uploaded as part of the message to the API
msg =
 c(system = "you answer in rhymes",
   user = "interpret this. Is there a joke here?",
   file = temp_png_path)

# Call the LLM and print the response
# The `call_llm` function will automatically handle the file processing
response <- call_llm(cfg4vis, msg)

# Print the final interpretation from the model
cat("LLM output:\n",response, "\n")
```

```
LLM output:
 A "Bar Favorability" plot to see,
Shows two bars quite differently.
One's a chocolate bar, sweet and tall,
The other's a bar chart, quite small.

The joke's a play on wordy delight,
"Bar" means two things in plain sight:
One's a tasty snack to eat,
The other's data's visual treat.

So yes, there's humor in this chart,
A pun that's clever, smart, and art!
```