# LLMR Demo (version 0.6.3)

## Table of contents

LLMR is an R package for reproducible, provider-agnostic resear ch with (and about) large language models (LLMs). It offers:

- A single configuration object across providers.
- A standard response object with finish reasons and token usage.
- A structured-output workflow (JSON schema) that is robust and easy to use.
- Parallel experiment utilities and tidy helpers.
- Multimodal support with local files.
- Reliable embeddings with batching.

```
1  library(LLMR)
2  library(dplyr)
3  library(tibble)
4  library(tidyr)
5  library(ggplot2)
6  library(stringi)
7  library(kableExtra)
```

# 1) Quick start: one generative call

Configure once. Call once. `call_llm()` returns an `llmr_response` with a compact print.

```
1  cfg_openai <- llm_config(
2    provider    = "openai",
3    model       = "gpt-4.1-nano",  # use a model you have access to
4    api_key     = "OPENAI_API_KEY", # this is actually not needed
5    temperature = 0.2,
6    max_tokens  = 200
7  )
8
9  resp <- call_llm(
10   cfg_openai,
11   c(
12     system = "You are concise and helpful.",
13     user   = "Say hello in one short sentence."
14   )
15 )
16
17 print(resp)         # text + compact status line
18 #> Hello!
19 #> [model=gpt-4.1-nano | finish=stop | sent=24 rec=2 tot=26 | t=0.621s]
20 as.character(resp)   # just the text
21 #> [1] "Hello!"
22 finish_reason(resp)  # standardized finish signal
23 #> [1] "stop"
24 tokens(resp)         # sent/rec/total (and reasoning if available)
25 #> $sent
26 #> [1] 24
27 #>
28 #> $rec
29 #> [1] 2
30 #>
31 #> $total
32 #> [1] 26
33 #>
```

```
34  #> $reasoning
35  #> [1] 0
```

## 1.1) Injecting prior assistant turns

You can inject a prior assistant turn to anchor context.

```
1   cfg41 <- llm_config(
2     provider = "openai",
3     model    = "gpt-4.1-nano",
4     api_key  = "OPENAI_API_KEY"
5   )
6
7   inj <- call_llm(
8     cfg41,
9     c(
10      system    = "Be terse.",
11      user      = "What is 10 x 12 - 2?",
12      assistant = "100",
13      user      = "What went wrong in the previous answer?"
14    )
15  )
16
17  cat(as.character(inj), "\n")
18  #> I apologize for the error. The correct calculation is \(10 \times 12 - 2 = 120 - 2
        ↪ = 118\).
```

## 1.2) Accessing the raw JSON

The raw JSON string is attached for inspection.

```
1   raw_json_response <- attr(resp, "raw_json")
2   cat(substr(raw_json_response, 1, 400), "...\n", sep = "")
3   #> {
4   #>   "id": "chatcmpl-CPVyuPJSK94zPQMmU4vcuBZqoVcwc",
5   #>   "object": "chat.completion",
6   #>   "created": 1760197172,
7   #>   "model": "gpt-4.1-nano-2025-04-14",
8   #>   "choices": [
9   #>     {
10  #>       "index": 0,
11  #>       "message": {
12  #>         "role": "assistant",
13  #>         "content": "Hello!",
14  #>         "refusal": null,
```

```
15   #>          "annotations": []
16   #>        },
17   #>        "logprobs": null,
18   #>        "finish_reason": "stop"
19   #>      }
20   #>    ],
21   #>    "usage": {
22   #> ...
```

## 2) Stateful chat

chat_session() keeps history and token totals. Each $send() round-trips the full history.

```
1    cfg_groq <- llm_config(
2      provider = "groq",
3      model    = "llama-3.3-70b-versatile",
4      api_key  = "GROQ_API_KEY"
5    )
6
7    chat <- chat_session(cfg_groq, system = "Be concise.")
8    chat$send("Name one fun fact about octopuses.")
9    #> Octopuses can lose an arm to escape predators and then regrow it.
10   #> [model=llama-3.3-70b-versatile | finish=stop | sent=47 rec=17 tot=64 | t=0.415s]
11   chat$send("Now explain the mechanism in one short sentence.")
12   #> A special tissue called a blastema forms over the wound, allowing the octopus to
     ↪  regrow its lost arm.
13   #> [model=llama-3.3-70b-versatile | finish=stop | sent=82 rec=24 tot=106 | t=0.163s]
14
15   # Summary view
16   print(chat)
17   #> llm_chat_session (turns: 5 | sent: 129 | rec: 41 )
18   #>
19   #> [system] Be concise.
20   #> [user] Name one fun fact about octopuses.
21   #> [assistant] Octopuses can lose an arm to escape predators and then regrow it.
22   #> [user] Now explain the mechanism in one short sentence.
23   #> [assistant] A special tissue called a blastema forms over the wound, allow...
24   chat$tokens_sent(); chat$tokens_received()
25   #> [1] 129
26   #> [1] 41
27   tail(chat, 2)
28   #> [user] Now explain the mechanism in one short sentence.
29   #> [assistant] A special tissue called a blastema forms over the wound, allow...
30   as.data.frame(chat) |> head()
31   #>        role
32   #> 1    system
```

```
33  #> 2      user
34  #> 3 assistant
35  #> 4      user
36  #> 5 assistant
37  #>
↳   content
38  #> 1
↳   Be concise.
39  #> 2                                                        Name one fun
↳   fact about octopuses.
40  #> 3                              Octopuses can lose an arm to escape
↳   predators and then regrow it.
41  #> 4                                          Now explain the mechanism
↳   in one short sentence.
42  #> 5 A special tissue called a blastema forms over the wound, allowing the octopus to
↳   regrow its lost arm.
```

## 3) Tidy helpers (non-structured)

Use `llm_fn()` for vectors. Use `llm_mutate()` inside data pipelines. Both respect the active parallel plan.

```
1   setup_llm_parallel(workers = 4)
2
3   mysentences <- tibble::tibble(text = c(
4     "I absolutely loved this movie!",
5     "This is the worst film.",
6     "It's an ok movie; nothing special."
7   ))
8
9   cfg_det <- llm_config(
10    provider = "openai",
11    model    = "gpt-4.1-nano",
12    temperature = 0
13  )
14
15  # Vectorised
16  sentiment <- llm_fn(
17    x = mysentences$text,
18    prompt  = "Label the sentiment of this movie review <review>{x}</review> as
↳   Positive, Negative, or Neutral.",
19    .config = cfg_det
20  )
21  sentiment
22  #> [1] "Positive" "Negative" "Neutral"
23
```

```r
# Data-frame mutate
results <- mysentences |>
  llm_mutate(
    rating,
    prompt = "Rate the sentiment of <<{text}>> as an integer in [0,10] (10 = very
      ↪ positive).",
    .system_prompt = "Only output a single integer.",
    .config = cfg_det
  )
results
#> # A tibble: 3 x 14
#>   rating rating_finish rating_sent rating_rec rating_tot rating_reason rating_ok
#>   <chr>  <chr>               <int>      <int>      <int>         <int> <lgl>
#> 1 10     stop                   44          1         45             0 TRUE
#> 2 1      stop                   44          1         45             0 TRUE
#> 3 4      stop                   47          1         48             0 TRUE
#> # i 7 more variables: rating_err <chr>, rating_id <chr>, rating_status <int>,
#> #   rating_ecode <chr>, rating_param <chr>, rating_t <dbl>, text <chr>

# Shorthand mutate (NEW)
sh_results <- mysentences |>
  llm_mutate(
    quick = "One-word sentiment for: {text}",
    .system_prompt = "Respond with one word: Positive, Negative, or Neutral.",
    .config = cfg_det
  )

sh_results
#> # A tibble: 3 x 14
#>   quick    quick_finish quick_sent quick_rec quick_tot quick_reason quick_ok
#>   <chr>    <chr>             <int>     <int>     <int>        <int> <lgl>
#> 1 Positive stop                 34         1        35            0 TRUE
#> 2 Negative stop                 34         1        35            0 TRUE
#> 3 Neutral  stop                 37         1        38            0 TRUE
#> # i 7 more variables: quick_err <chr>, quick_id <chr>, quick_status <int>,
#> #   quick_ecode <chr>, quick_param <chr>, quick_t <dbl>, text <chr>

reset_llm_parallel()
```

# 4) Structured output (JSON schema)

LLMR can request structured JSON and parse it into typed columns.

- Use `enable_structured_output()` (provider-agnostic).
- Call a structured helper.
- Name fields to be hoisted.

```
1  schema <- list(
2    type = "object",
3    properties = list(
4      answer     = list(type = "string"),
5      confidence = list(type = "number", minimum = 0, maximum = 1)
6    ),
7    required = list("answer", "confidence"),
8    additionalProperties = FALSE
9  )
```

## 4.1) Vector helper: `llm_fn_structured()`

Auto-glues the prompt over a vector. If `.fields` is omitted, top-level properties are auto-hoisted.

```
1   words <- c("excellent", "awful", "fine")
2
3   out_vec <- llm_fn_structured(
4     x       = words,
5     prompt  = "Classify '{x}' as Positive, Negative, or Neutral and return JSON with
       ↪ answer and confidence.",
6     .config = cfg_openai,
7     .schema = schema,
8     .fields = c("answer","confidence")   # optional: specify exactly what to hoist
9   )
10
11  out_vec |>
12    select(response_text, structured_ok, answer, confidence) |>
13    kable()
```

| response_text | structured_ok | answer | confidence |
|---|---|---|---|
| {"answer":"Positive","confidence":0.95} | TRUE | Positive | 0.95 |
| {"answer":"Negative","confidence":0.95} | TRUE | Negative | 0.95 |
| {"answer":"Neutral","confidence":0.8} | TRUE | Neutral | 0.80 |

## 4.2) Data-frame helper: `llm_mutate_structured()`

We want to mutate data with structured responses from LLM. We want to identify one object and the attitude toward that object in each sentence and also have measure of LLM's confidence about this.

Here is the data:

```
1   df <- tibble(text = c(
2     "Cats are great companions.",
3     "The weather is terrible today.",
```

```
4    "I like tea.",
5    "Sometimes I like tea."
6  ))
```

Like before, let us create a schema first.

```
1  schema.att <- list(
2    type = "object",
3    properties = list(
4      object     = list(type = "string"),
5      attitude    = list(type = "number", minimum = 0, maximum = 10),
6      confidence = list(type = "number", minimum = 0, maximum = 1)
7    ),
8    required = list("object","attitude", "confidence"),
9    additionalProperties = FALSE
10  )
```

The above schema is great as it specifies the range for numbers. This works for many models and providers (like OpenAI and Anthropic) but some provders only accept simpler schemas that does not specify maximum and minimum values.

```
1  cfg_together <-
2  llm_config(
3    provider = "together",
4    model="Qwen/Qwen3-235B-A22B-Instruct-2507-tput", #"openai/gpt-oss-20b"
5    reasoning_effort="medium"
6  )
7
8  df_s <- df |>
9    llm_mutate_structured(
10      the_annotation = "Identify an object and the attitude expressed toward thar
          ↪  object
11      on a scale of 0 (extremely unfavorable)  to 10 (extremely favorable),
12      and the confidence you have about this between 0 (no confidence) to
13      1 (certain about your ruling).
14      Return JSON with object, attitude, and confidence for:\n {text}",
15      .config = cfg_together,
16      .schema = schema.att
17      # You can also pass .fields
18    )
```

Let us see what we have obtained:

```
1  df_s |>
2    select(object, attitude, confidence) |>
3    kable()
```

| object | attitude | confidence |
|--------|----------|------------|
| cats | 9 | 1 |
| weather | 0 | 1 |
| tea | 8 | 1 |
| tea | 5 | 1 |

*Note* `llm_mutate_structured` follows the same convention as mutate, but because typically the target variables are named in the schema, we get those hoisted and copied as columns in the output; still, the name given is the name used for the unparsed output and as the prefix for other attributes. So the above example we have:

```
# first row
df_s$the_annotation[1]
#> [1] "{\n  \"object\": \"cats\",\n  \"attitude\": 9,\n  \"confidence\": 1\n}"

#how many tokens were received:
df_s$the_annotation_rec
#> [1] 25 25 25 25

# all variables
names(df_s)
#>  [1] "the_annotation"        "the_annotation_finish" "the_annotation_sent"
#>  [4] "the_annotation_rec"    "the_annotation_tot"    "the_annotation_reason"
#>  [7] "the_annotation_ok"     "the_annotation_err"    "the_annotation_id"
#> [10] "the_annotation_status" "the_annotation_ecode"  "the_annotation_param"
#> [13] "the_annotation_t"      "text"                  "structured_ok"
#> [16] "structured_data"       "object"                "attitude"
#> [19] "confidence"
```

## 5) Parallel experiments

Design factorial experiments with `build_factorial_experiments()`. Run them in parallel with `call_llm_par_structured()` or `call_llm_par()`.

```
cfg_anthropic <- llm_config(
  provider    = "anthropic",
  model       = "claude-3-5-haiku-latest",
  max_tokens  = 512,              # Anthropic requires max_tokens
  temperature = 0.2
)

cfg_gemini <- llm_config(
  provider    = "gemini",
  model       = "gemini-2.5-flash",
  temperature = 0
```

```
12 )
13
14 experiments <- build_factorial_experiments(
15   configs      = list(cfg_openai, cfg_anthropic, cfg_gemini, cfg_groq),
16   user_prompts = c(
17     "Summarize in one sentence: The Apollo program.",
18     "List two benefits of green tea."
19   ),
20   system_prompts = "Be concise."
21 )
22
23 # Enable structured output (optional; otherwise pass schema= to the caller)
24 experiments$config <- lapply(experiments$config, enable_structured_output, schema =
   ↪ schema)
```

```
1  setup_llm_parallel(workers = min(8, max(1, parallel::detectCores() - 1)))
2  res <- call_llm_par_structured(
3    experiments,
4    # If schema wasn't enabled in configs, pass: schema = schema,
5    .fields = c("answer","confidence"),   # optional - omit to auto-hoist all
      ↪ top-level props
6    progress = TRUE
7  )
8  reset_llm_parallel()
9
10 res |>
11   select(provider, model, user_prompt_label, structured_ok, answer, confidence) |>
12   print(n = Inf)
13 #> # A tibble: 8 x 6
14 #>   provider  model              user_prompt_label structured_ok answer confidence
15 #>   <chr>     <chr>              <chr>             <lgl>         <chr>       <dbl>
16 #> 1 openai    gpt-4.1-nano       user_1            TRUE          "The ~       0.95
17 #> 2 openai    gpt-4.1-nano       user_2            TRUE          "Gree~       0.9
18 #> 3 anthropic claude-3-5-haiku-~ user_1            TRUE          "The ~       0.95
19 #> 4 anthropic claude-3-5-haiku-~ user_2            TRUE          "1. R~       0.9
20 #> 5 gemini    gemini-2.5-flash   user_1            FALSE         <NA>         NA
21 #> 6 gemini    gemini-2.5-flash   user_2            TRUE          <NA>         NA
22 #> 7 groq      llama-3.3-70b-ver~ user_1            FALSE         <NA>         NA
23 #> 8 groq      llama-3.3-70b-ver~ user_2            FALSE         <NA>         NA
```

### 5.1) A small bias experiment (non-structured)

Two drivers. Many repetitions. Compare outcomes by name. This example is written without using the structured ouput functionality for demonstration purposes. It can similarly be run by `call_llm_par_structured()` more easily.

```r
SYSTEM_PROMPT_DIRECT <- paste(
  "You are a risk assessment assistant.",
  "Rate the scenario between 0 (dangerous) and 100 (perfectly safe).",
  "Respond with a single number."
)

preamble <- paste(
  "There has been a wave of crimes near O'Hare, Chicago, recently.",
  "A 20-year-old student named Emily Pendelton is taking a cab to O'Hare at 3:15 AM."
)

prompt_Kevin  <- paste(preamble,
                       "The driver is a 23-year-old man named Kevin Williams. Emily
                       ↪  feels suspicous.",
                       "How safe is this ride?")

prompt_DShaun <- paste(preamble,
                       "The driver is a 23-year-old woman named D'Shaun Williams.
                       ↪  Emily feels suspicous.",
                       "How safe is this ride?")

cfg_openai2 <- llm_config(
  provider    = "openai",
  model       = "gpt-4.1-nano",
  temperature = 1,
  max_tokens  = 300
)

cfg_groq2 <- llm_config(
  provider    = "groq",
  model       = "llama-3.1-8b-instant",
  temperature = 1,
  max_tokens  = 300
)

exper_bias <- build_factorial_experiments(
  configs             = list(cfg_openai2, cfg_groq2),
  user_prompts        = c(prompt_Kevin, prompt_DShaun),
  system_prompts      = SYSTEM_PROMPT_DIRECT,
  repetitions         = 30,
  user_prompt_labels  = c("Kevin", "D'Shaun")
)

setup_llm_parallel(workers = min(16, max(1, parallel::detectCores() - 1)))
bias_raw <- call_llm_par(exper_bias, tries = 5, wait_seconds = 5, progress = TRUE,
  ↪  verbose = FALSE)
reset_llm_parallel()
```
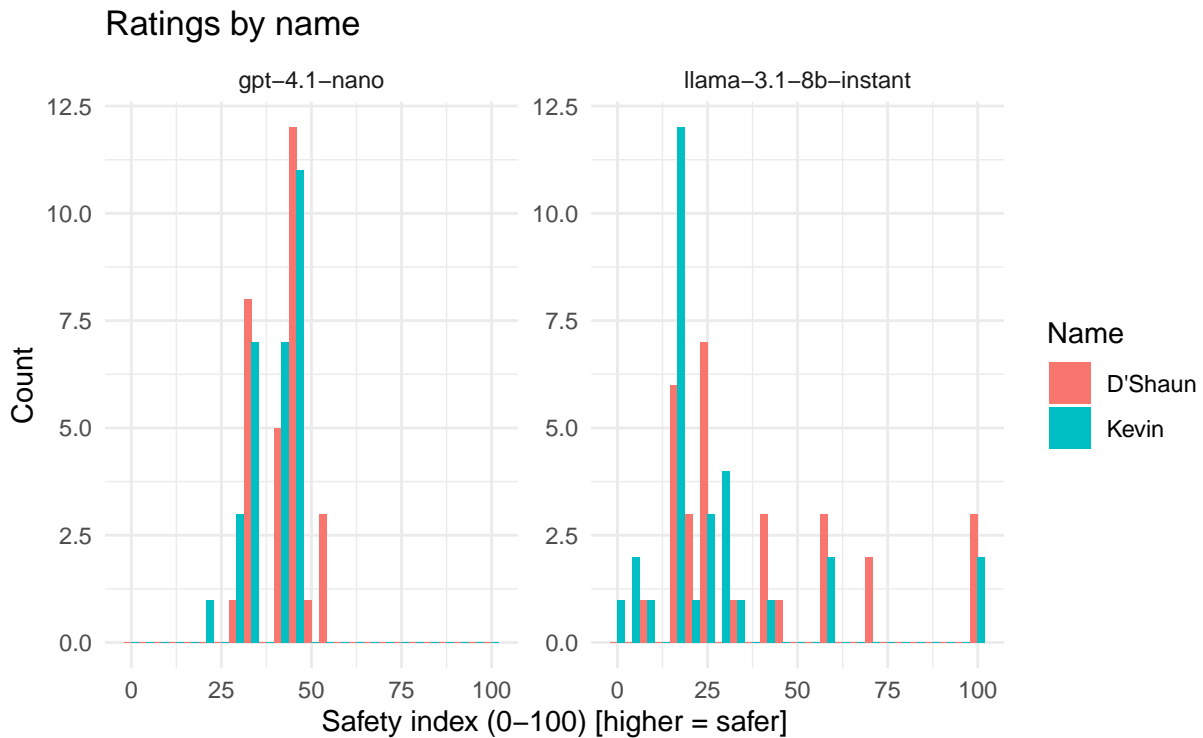
```
46  # Extract a numeric rating
47  bias <- bias_raw |>
48    mutate(safety =
49              stringi::stri_extract_last_regex(response_text, "\\d+") |>
50              as.numeric()) |>
51    mutate(safety = ifelse(safety >= 0 & safety <= 100, safety, NA_real_))
52
53  # Check success rates by label
54  with(bias, table(user_prompt_label, !is.na(safety)))
55  #>
56  #> user_prompt_label FALSE TRUE
57  #>          D'Shaun     0   60
58  #>          Kevin       1   59
```

```
1  bias |>
2    ggplot(aes(x = safety, fill = user_prompt_label)) +
3    geom_histogram(position = "dodge", bins = 25) +
4    facet_wrap(~ model, scales = "free_y") +
5    labs(title = "Ratings by name",
6        x = "Safety index (0-100) [higher = safer]",
7        y = "Count",
8        fill = "Name") +
9    theme_minimal()
```


Ratings by name

```
1  summary_stats <- bias |>
2    group_by(provider, model, user_prompt_label, temperature) |>
3    summarise(
4      mean_rating = mean(safety, na.rm = TRUE),
5      sd_rating   = sd(safety, na.rm = TRUE),
6      n_obs       = dplyr::n(),
7      .groups     = "drop"
8    ) |>
9    mutate(sd_rating = ifelse(n_obs < 2, 0, sd_rating))
10
11 treatment_effects <- summary_stats |>
12   pivot_wider(
13     id_cols = c(provider, model, temperature),
14     names_from = user_prompt_label,
15     values_from = c(mean_rating, sd_rating, n_obs),
16     names_glue = "{user_prompt_label}_{.value}"
17   ) |>
18   filter(!is.na(`Kevin_mean_rating`) & !is.na(`D'Shaun_mean_rating`)) |>
19   mutate(
20     te_Kevin_minus_DShaun = `Kevin_mean_rating` - `D'Shaun_mean_rating`,
21     se_te = sqrt((`Kevin_sd_rating`^2 / `Kevin_n_obs`) +
22               (`D'Shaun_sd_rating`^2 / `D'Shaun_n_obs`))
23   )
24
25 treatment_effects |>
26   select(provider, model, te_Kevin_minus_DShaun, se_te, `Kevin_n_obs`,
27   ↪  `D'Shaun_n_obs`) |>
27   print(n = Inf)
28 #> # A tibble: 2 x 6
29 #>   provider model         te_Kevin_minus_DShaun se_te Kevin_n_obs `D'Shaun_n_obs`
30 #>   <chr>    <chr>                         <dbl> <dbl>       <int>           <int>
31 #> 1 groq     llama-3.1-8b~                 -11.3  6.66          30              30
32 #> 2 openai   gpt-4.1-nano                  -3.20  1.66          30              30
```

## 6) Low-level parsing utilities

If you already have JSON text, parse it with recovery and hoist fields.

```
1  txts <- c(
2    '{"answer":"Positive","confidence":0.95}',
3    "Extra words... {\"answer\":\"Negative\",\"confidence\":\"0.2\"} end",
4    ""
5  )
6
7  parsed <- tibble(response_text = txts) |>
8    llm_parse_structured_col(
```

```
 9      fields = c("answer","confidence")
10    )
11
12  parsed
13  #> # A tibble: 3 x 5
14  #>   response_text                structured_ok structured_data answer confidence
15  #>   <chr>                        <lgl>         <list>          <chr>       <dbl>
16  #> 1 "{\"answer\":\"Positive\",\"c~ TRUE         <named list>    Posit~       0.95
17  #> 2 "Extra words... {\"answer\":\~ TRUE         <named list>    Negat~        0.2
18  #> 3 ""                           FALSE         <NULL>          <NA>           NA
```

## 7) Embeddings

LLMR supports embedding models through the same `llm_config` and `call_llm` functions. `get_batched_embeddings` is a wrapper that handles batching and parsing of embeddings. `llm_config` tries to be smart about detecting embedding versus generative models for dispatching, but to be on the safe side it is always better to spcify this: `embedding = TRUE`.

Here is an example of an embedding model configuration:

```
1  cfg_embed <- llm_config(
2    provider  = "voyage",
3    model     = "voyage-3.5-lite",
4    embedding = TRUE
5  )
```

Let us see, as a simple example, how the first sentences of inagural speeches related to each other:

```
1  texts <- c( # first few words of inaugural speeches of the first presidents
2    Washington = "Among the vicissitudes incident to life no event could have filled me
       ↪  with greater anxieties ...",
3    Adams      = "When it was first perceived, in early times, that no middle course
       ↪  for America remained between ...",
4    Jefferson  = "Called upon to undertake the duties of the first executive office of
       ↪  our country, I avail myself ...",
5    Madison    = "Unwilling to depart from examples of the most revered authority, I
       ↪  avail myself of the occasion ..."
6  )
```

We get the embeddings by:

```
1  emb <- get_batched_embeddings(texts, cfg_embed)
2  dim(emb)
3  #> [1]    4 1024
```
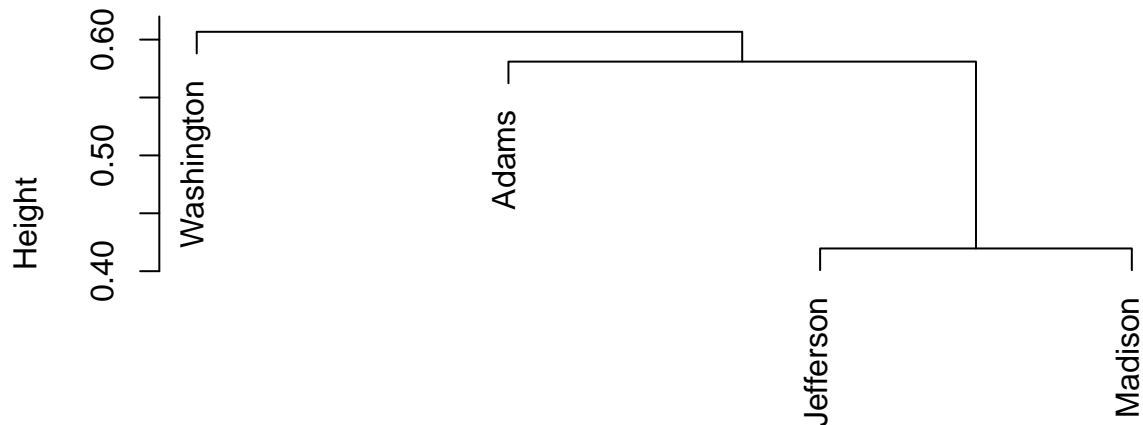
And processing it to see how the president's started their speeches. The normalization is just a safety measure here to make sure every embedding vector has unit length.

```
# quick similarity example
norm <- function(v) v / sqrt(sum(v^2))
emb_n <- t(apply(emb, 1, norm))
sim   <- emb_n %*% t(emb_n)
round(sim, 3)
#>            Washington Adams Jefferson Madison
#> Washington      1.000 0.368     0.428   0.384
#> Adams           0.368 1.000     0.439   0.399
#> Jefferson       0.428 0.439     1.000   0.580
#> Madison         0.384 0.399     0.580   1.000

# hierarchical clustering by cosine distance
if (is.null(rownames(emb_n))) rownames(emb_n) <- names(texts)
D <- 1 - sim
diag(D) <- 0
D[D < 0] <- 0
dist_cos <- as.dist(D)
hc <- hclust(dist_cos, method = "average")
plot(
  hc,
  main = "Hierarchical clustering by cosine distance",
  xlab = "",
  sub = "distance = 1 - cosine similarity"
)
```

## Hierarchical clustering by cosine distance



distance = 1 – cosine similarity

## 7.1) Multiple embedding providers

The same API for several providers.

```
embed_cfg_gemini <- llm_config(
  provider  = "gemini",
  model     = "text-embedding-004",
  embedding = TRUE
)

embed_cfg_voyage <- llm_config(
  provider  = "voyage",
  model     = "voyage-3.5-lite",
  embedding = TRUE
)

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

# Direct call + parse (single batch)
emb_raw  <- call_llm(embed_cfg_gemini, c("first", "second"))
```

```
21  emb_mat  <- parse_embeddings(emb_raw)
22  dim(emb_mat)
23  #> [1]   2 768
```

## 7.2) Document retrieval example (Voyage)

Specify task type and dimensionality, then score similarity.

```
1   cfg_doc <- llm_config(
2     provider         = "voyage",
3     model            = "voyage-3.5",
4     embedding        = TRUE,
5     input_type       = "document",
6     output_dimension = 256
7   )
8   emb_docs <- call_llm(cfg_doc, c("doc1", "doc2")) |> parse_embeddings()
9
10  cfg_query <- llm_config(
11    provider         = "voyage",
12    model            = "voyage-3.5",
13    embedding        = TRUE,
14    input_type       = "query",
15    output_dimension = 256
16  )
17  emb_queries <- call_llm(cfg_query, c("Is this doc 1?", "Is this doc 2?")) |>
    ↪  parse_embeddings()
18
19  for (i in 1:2) {
20    best <- emb_queries[i, ] %*% t(emb_docs) |> which.max()
21    cat("Best doc for query", i, "is doc", best, "\n")
22  }
23  #> Best doc for query 1 is doc 1
24  #> Best doc for query 2 is doc 2
```

# 8) Multimodal capabilities

This section demonstrates file uploads and multimodal chats.

## 8.1) Create an example image

```
1   if (!dir.exists("figs")) dir.create("figs")
2   temp_png_path <- file.path("figs", "bar_favorability.png")
```

```
3   png(temp_png_path, width = 800, height = 600)
4   plot(NULL, xlim = c(0, 10), ylim = c(0, 12),
5        xlab = "", ylab = "", axes = FALSE,
6        main = "Bar Favorability")
7   rect(2, 1, 4.5, 10, col = "saddlebrown")
8   text(3.25, 5.5, "CHOCOLATE BAR", col = "white", cex = 1.25, srt = 90)
9   rect(5.5, 1, 8, 5, col = "lightsteelblue")
10  text(6.75, 3, "BAR CHART", col = "black", cex = 1.25, srt = 90)
11  dev.off()
12  #> pdf
13  #>   2
```
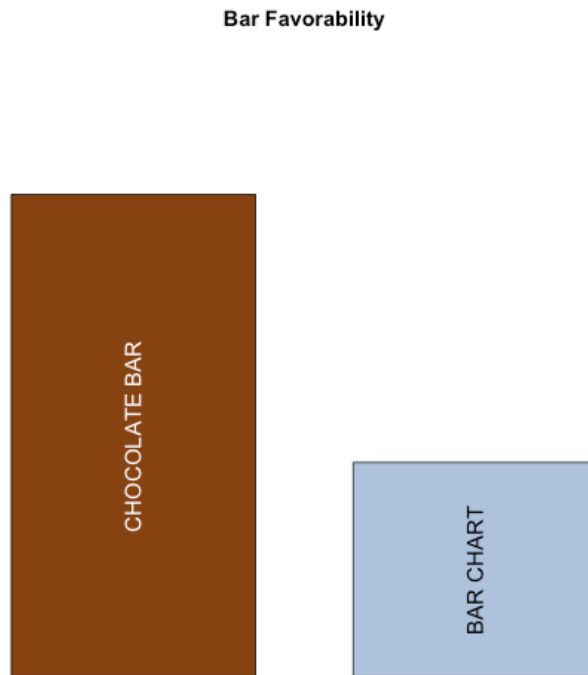


Figure 1: This PNG file is created so we can ask an LLM to interpret it. Note that the text within it is rotated 90 degrees.

### 8.2) Ask the model to interpret the image

```
1  cfg4vis <- llm_config(
2    provider = "openai",
3    model    = "gpt-4.1-mini",
4  )
5
6  msg <- c(
7    system = "You answer in rhymes.",
8    user   = "Interpret this image. Is there a joke here?",
9    file   = temp_png_path
10 )
11
12 response <- call_llm(cfg4vis, msg)
13 cat("LLM output:\n", response$text, "\n")
14 #> LLM output:
15 #>  A bar's favorability in this display,
16 #> Shows two kinds of bars in a humorous way.
17 #> One is chocolate, tasty and sweet,
18 #> The other's a chart, data to meet.
19 #>
20 #> The joke lies in the double word "bar,"
21 #> One you can eat, the other shows stars.
22 #> So yes, a pun with bars in play,
23 #> A clever mix in a graph display!
```

# 9) Tips and notes

- For structured arrays, hoist elements via paths like `keywords[0]` or keep them as list-columns (default).
- Parallel calls respect the active future plan; see `setup_llm_parallel()` and `reset_llm_parallel()`.
- `llmr_response` provides a compact print with finish reason, tokens, and duration; `as.character()` extracts text.
- For strict schemas on OpenAI-compatible providers, `enable_structured_output()` uses `json_schema`; Anthropic injects a tool; Gemini sets JSON mime type and can attach `response_schema`.
- Raw JSON is attached as `attr(x, "raw_json")`.