

# LLMR Demo ( $\geq 0.6$ )

## Table of contents

<b>1) Quick start: one generative call</b>	<b>2</b>
1.1) Injecting prior assistant turns . . . . .	3
1.2) Accessing the raw JSON . . . . .	3
<b>2) Stateful chat</b>	<b>4</b>
<b>3) Structured output (JSON schema)</b>	<b>5</b>
3.1) Vector helper: <code>llm_fn_structured()</code> . . . . .	5
3.2) Data-frame helper: <code>llm_mutate_structured()</code> . . . . .	6
<b>4) Tidy helpers (non-structured)</b>	<b>7</b>
<b>5) Parallel experiments</b>	<b>8</b>
5.1) A small bias experiment (non-structured) . . . . .	9
<b>6) Low-level parsing utilities</b>	<b>12</b>
<b>7) Embeddings</b>	<b>13</b>
7.1) Multiple embedding providers . . . . .	14
7.2) Document retrieval example (Voyage) . . . . .	15
<b>8) Multimodal capabilities</b>	<b>15</b>
8.1) Create an example image . . . . .	16
8.2) Ask the model to interpret the image . . . . .	17
<b>9) Tips and notes</b>	<b>17</b>

LLMR is an R package for reproducible, provider-agnostic research 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)

```

## 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=0.761s]
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 #> The previous answer was incorrect. The correct calculation:  
19 #> 10 x 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-CLrLkl2P57CFzR5nN08dpRn79VxVm",  
5 #>   "object": "chat.completion",  
6 #>   "created": 1759326000,  
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 #> }
```

```

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$send("Name one fun fact about octopuses.")
9  #> Octopuses can lose a limb to escape predators, and the detached limb can still
10   ↪ move and distract the predator, allowing the octopus to escape.
11  #> [model=llama-3.3-70b-versatile | finish=stop | sent=47 rec=31 tot=78 | t=0.313s]
12  chat$send("Now explain it in one short sentence.")
13  #> Octopuses can release a limb to distract predators and escape.
14  #> [model=llama-3.3-70b-versatile | finish=stop | sent=95 rec=14 tot=109 | t=0.133s]
15
16  # Summary view
17  print(chat)
18  #> llm_chat_session (turns: 5 | sent: 142 | rec: 45 )
19  #>
20  #> [system] Be concise.
21  #> [user] Name one fun fact about octopuses.
22  #> [assistant] Octopuses can lose a limb to escape predators, and the detache...
23  #> [user] Now explain it in one short sentence.
24  #> [assistant] Octopuses can release a limb to distract predators and escape.
25  chat$tokens_sent(); chat$tokens_received()
26  #> [1] 142
27  #> [1] 45
28  tail(chat, 2)
29  #> [user] Now explain it in one short sentence.
30  #> [assistant] Octopuses can release a limb to distract predators and escape.
31  as.data.frame(chat) |> head()
32  #>      role
33  #> 1  system
34  #> 2   user
35  #> 3 assistant
36  #> 4   user
37  #> 5 assistant
38  #>
39  ↪ content

```

```

38 #> 1
   ↪ Be concise.
39 #> 2
   ↪ Name one fun fact about octopuses.
40 #> 3 Octopuses can lose a limb to escape predators, and the detached limb can still
   ↪ move and distract the predator, allowing the octopus to escape.
41 #> 4
   ↪ Now explain it in one short sentence.
42 #> 5
   ↪ Octopuses can release a limb to distract predators and escape.

```

### 3) 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.
- Hoist fields with `llm_parse_structured_col()` (done automatically below).

```

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 )

```

#### 3.1) 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 print(n = Inf)
14 #> # A tibble: 3 x 4
15 #>   response_text structured_ok answer confidence
16 #>   <chr>          <lgl>      <chr>      <dbl>
17 #> 1 "{\\"answer\\":\\"Positive\\",\\"confidence\\":0.95~ TRUE      Posit~      0.95
18 #> 2 "{\\"answer\\":\\"Negative\\",\\"confidence\\":0.95~ TRUE      Negat~      0.95
19 #> 3 "{\\"answer\\":\\"Neutral\\",\\"confidence\\":0.8}" TRUE      Neutr~      0.8

```

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

Mutate your data with new structured columns.

```

1 df <- tibble(text = c(
2   "Cats are great companions.",
3   "The weather is terrible today.",
4   "I like tea."
5 ))
6
7 df_s <- df |>
8   llm_mutate_structured(
9     annot,
10    prompt = "Return JSON with answer and confidence for: {text}",
11    .config = cfg_groq,
12    .schema = schema
13    # You can also pass .fields = c("answer","confidence")
14  )
15
16 df_s |>
17   select(text, structured_ok, annot, answer, confidence) |>
18   head()
19 #> # A tibble: 3 x 5
20 #>   text structured_ok annot answer confidence
21 #>   <chr>          <lgl>      <chr> <chr>   <chr>
22 #> 1 Cats are great companions. FALSE      <NA> <NA>   <NA>
23 #> 2 The weather is terrible today. FALSE      <NA> <NA>   <NA>
24 #> 3 I like tea. FALSE      <NA> <NA>   <NA>

```

Note: In “columns” mode the generated raw text column is named after the output symbol (here `annot`). Hoisted scalars appear as separate typed columns. Arrays and objects become list-columns, unless you restrict hoisting with `.fields`.

## 4) 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
19     ↪ Positive, Negative, or Neutral.",
20   .config = cfg_det
21 )
22 sentiment
23 #> [1] "Positive" "Negative" "Neutral"
24
25 # Data-frame mutate
26 results <- mysentences |>
27   llm_mutate(
28     rating,
29     prompt = "Rate the sentiment of <<{text}>> as an integer in [0,10] (10 = very
30       ↪ positive).",
31     .system_prompt = "Only output a single integer.",
32     .config = cfg_det
33   )
34 results
35 #> # A tibble: 3 x 14
36 #>   rating rating_finish rating_sent rating_rec rating_tot rating_reason rating_ok
37 #>   <chr>   <chr>          <int>      <int>      <int>      <int> <lgl>
38 #> 1  10      stop          44         1         45         0 TRUE
39 #> 2  1      stop          44         1         45         0 TRUE
40 #> 3  4      stop          47         1         48         0 TRUE
41 #> # i 7 more variables: rating_err <chr>, rating_id <chr>, rating_status <int>,
42 #> #   rating_ecode <chr>, rating_param <chr>, rating_t <dbl>, text <chr>
```

```

42 # Shorthand mutate (NEW)
43 sh_results <- mysentences |>
44   llm_mutate(
45     quick = "One-word sentiment for: {text}",
46     .system_prompt = "Respond with one word: Positive, Negative, or Neutral.",
47     .config = cfg_det
48   )
49
50 sh_results
51 #> # A tibble: 3 x 14
52 #>   quick    quick_finish quick_sent quick_rec quick_tot quick_reason quick_ok
53 #>   <chr>    <chr>          <int>    <int>    <int>    <int> <lgl>
54 #> 1 Positive stop           34         1      35         0 TRUE
55 #> 2 Negative stop           34         1      35         0 TRUE
56 #> 3 Neutral  stop           37         1      38         0 TRUE
57 #> # i 7 more variables: quick_err <chr>, quick_id <chr>, quick_status <int>,
58 #> #   quick_ecode <chr>, quick_param <chr>, quick_t <dbl>, text <chr>
59
60 reset_llm_parallel()

```

## 5) Parallel experiments

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

```

1  cfg_anthropic <- llm_config(
2    provider   = "anthropic",
3    model      = "claude-3-5-haiku-latest",
4    max_tokens = 512,          # Anthropic requires max_tokens
5    temperature = 0.2
6  )
7
8  cfg_gemini <- llm_config(
9    provider   = "gemini",
10   model      = "gemini-2.5-flash",
11   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.

```

1 SYSTEM_PROMPT_DIRECT <- paste(
2   "You are a risk assessment assistant.",
3   "Rate the scenario between 0 (dangerous) and 100 (perfectly safe).",
4   "Respond with a single number."
5 )
6
7 preamble <- paste(
8   "There has been a wave of crimes near O'Hare, Chicago, recently.",
9   "A 20-year-old student named Emily Pendelton is taking a cab to O'Hare at 3:15 AM."
10 )
11
12 prompt_Kevin <- paste(preamble,

```

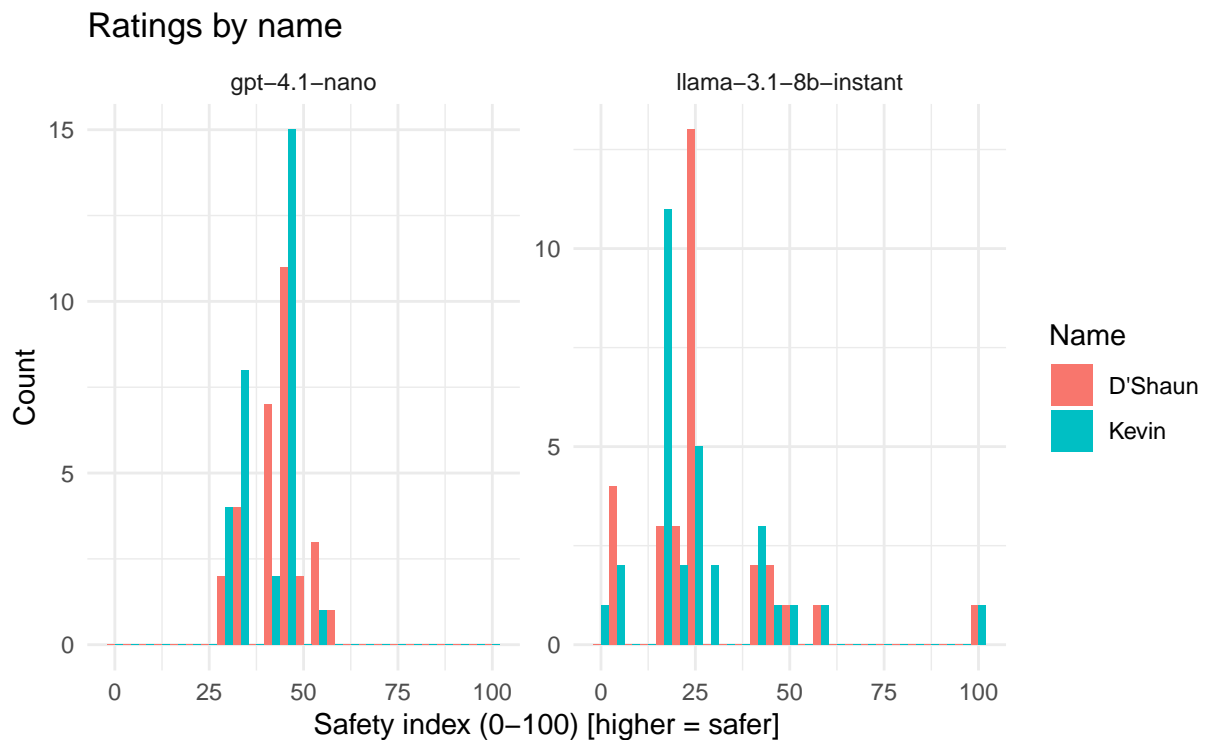
```

13         "The driver is a 23-year-old man named Kevin Williams. Emily
14         ↪ feels suspicious.",
15         "How safe is this ride?")
16
17 prompt_DShaun <- paste(preamble,
18         "The driver is a 23-year-old woman named D'Shaun Williams.
19         ↪ Emily feels suspicious.",
20         "How safe is this ride?")
21
22 cfg_openai2 <- llm_config(
23     provider = "openai",
24     model    = "gpt-4.1-nano",
25     temperature = 1,
26     max_tokens = 300
27 )
28
29 cfg_groq2 <- llm_config(
30     provider = "groq",
31     model    = "llama-3.1-8b-instant",
32     temperature = 1,
33     max_tokens = 300
34 )
35
36 exper_bias <- build_factorial_experiments(
37     configs      = list(cfg_openai2, cfg_groq2),
38     user_prompts = c(prompt_Kevin, prompt_DShaun),
39     system_prompts = SYSTEM_PROMPT_DIRECT,
40     repetitions  = 30,
41     user_prompt_labels = c("Kevin", "D'Shaun")
42 )
43
44 setup_llm_parallel(workers = min(16, max(1, parallel::detectCores() - 1)))
45 bias_raw <- call_llm_par(exper_bias, tries = 5, wait_seconds = 5, progress = TRUE,
46     ↪ verbose = FALSE)
47 reset_llm_parallel()
48
49 # Extract a numeric rating
50 bias <- bias_raw |>
51     mutate(safety =
52         stringi::stri_extract_last_regex(response_text, "\\d+") |>
53         as.numeric()) |>
54     mutate(safety = ifelse(safety >= 0 & safety <= 100, safety, NA_real_))
55
56 # Check success rates by label
57 with(bias, table(user_prompt_label, !is.na(safety)))
58 #>
59 #> user_prompt_label TRUE
60 #>      D'Shaun      60

```

```
58 #> Kevin 60
```

```
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()
```



```
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`) |>
28   print(n = Inf)
29 #> # A tibble: 2 x 6
30 #>   provider model      te_Kevin_minus_DShaun se_te Kevin_n_obs `D'Shaun_n_obs`
31 #>   <chr>    <chr>                <dbl> <dbl>         <int>         <int>
32 #> 1 groq    llama-3.1-8b~          -1.40  5.04           30           30
33 #> 2 openai  gpt-4.1-nano           -3     1.77           30           30

```

## 6) Low-level parsing utilities

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

```

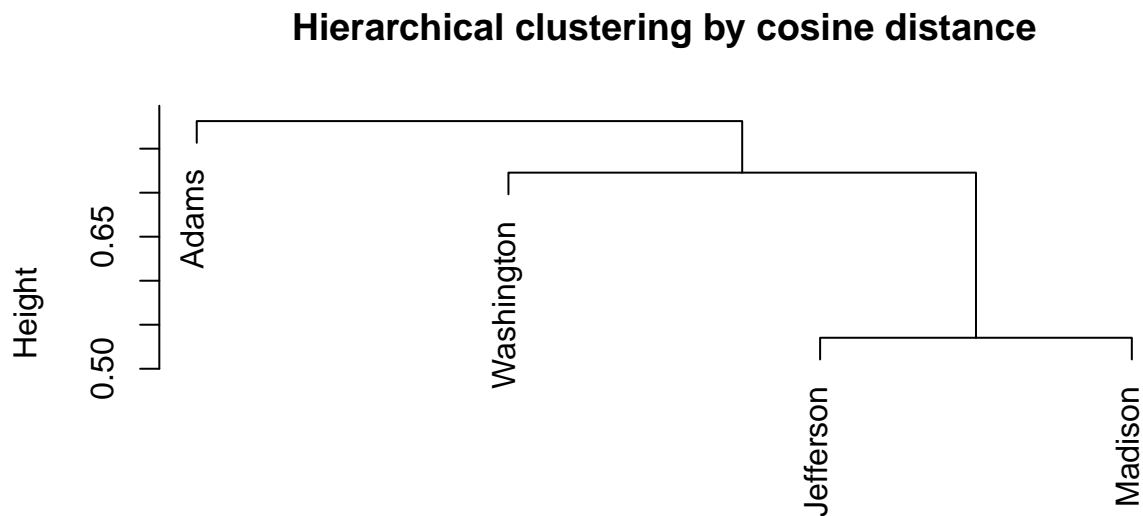
1 txts <- c(
2   '{"answer": "Positive", "confidence": 0.95}',
3   "Extra words... {\"answer\": \"Negative\", \"confidence\": \"0.2\"} end",
4   ""
5 )
6
7 parsed <- tibble(response_text = txts) |>
8   llm_parse_structured_col(
9     fields = c("answer", "confidence")
10  )
11
12 parsed
13 #> # A tibble: 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 batched embeddings with robust retries.

```
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
3     ↳ with greater anxieties ...",
4   Adams      = "When it was first perceived, in early times, that no middle course
5     ↳ for America remained between ...",
6   Jefferson  = "Called upon to undertake the duties of the first executive office of
7     ↳ our country, I avail myself ...",
8   Madison   = "Unwilling to depart from examples of the most revered authority, I
9     ↳ avail myself of the occasion ..."
10 )
11
12 cfg_embed <- llm_config(
13   provider = "openai",
14   model    = "text-embedding-3-small",
15   embedding = TRUE
16 )
17
18 emb <- get_batched_embeddings(texts, cfg_embed)
19 dim(emb)
20 #> [1]      4 1536
21
22 # quick similarity example
23 norm <- function(v) v / sqrt(sum(v^2))
24 emb_n <- t(apply(emb, 1, norm))
25 sim    <- emb_n %*% t(emb_n)
26 round(sim, 3)
27 #>      Washington Adams Jefferson Madison
28 #> Washington      1.000 0.210      0.280 0.275
29 #> Adams           0.210 1.000      0.251 0.196
30 #> Jefferson       0.280 0.251      1.000 0.465
31 #> Madison         0.275 0.196      0.465 1.000
32
33 # hierarchical clustering by cosine distance
34 if (is.null(rownames(emb_n))) rownames(emb_n) <- names(texts)
35 D <- 1 - sim
36 diag(D) <- 0
37 D[D < 0] <- 0
38 dist_cos <- as.dist(D)
39 hc <- hclust(dist_cos, method = "average")
40 plot(
41   hc,
42   main = "Hierarchical clustering by cosine distance",
43   xlab = "",
44   sub = "distance = 1 - cosine similarity"
```



distance = 1 – cosine similarity

## 7.1) Multiple embedding providers

The same API for several providers.

```

1 embed_cfg_gemini <- llm_config(
2   provider = "gemini",
3   model    = "text-embedding-004",
4   embedding = TRUE
5 )
6
7 embed_cfg_voyage <- llm_config(
8   provider = "voyage",
9   model    = "voyage-3.5-lite",
10  embedding = TRUE
11 )
12
13 embed_cfg_together <- llm_config(
14   provider = "together",
15   model    = "BAAI/bge-large-en-v1.5",
16   embedding = TRUE
17 )

```

```

18
19 # Direct call + parse (single batch)
20 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?")) |>
18   ↪ parse_embeddings()
19
20 for (i in 1:2) {
21   best <- emb_queries[i, ] %*% t(emb_docs) |> which.max()
22   cat("Best doc for query", i, "is doc", best, "\n")
23 }
24 #> Best doc for query 1 is doc 1
25 #> 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
```

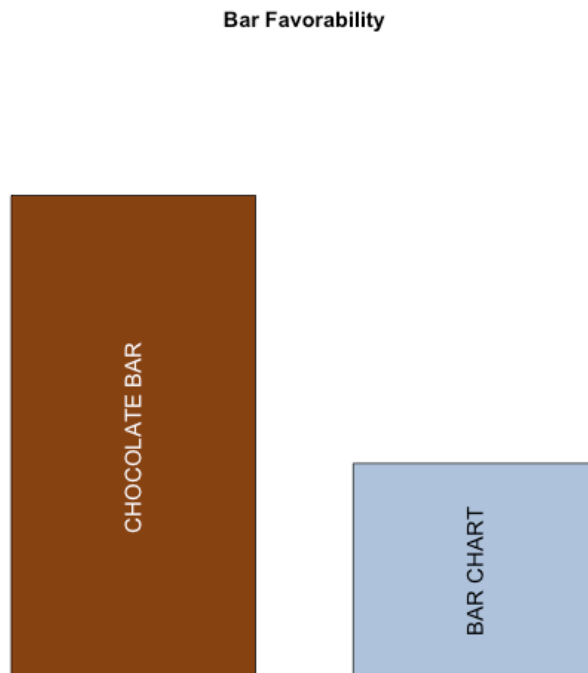


Figure 1: This PNG file is created so we can ask an LLM to interpret it. Note that the text within it is rotated 90 degrees.



## 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 #> This image shows two bars side by side,  
16 #> One a chocolate bar, rich and wide,  
17 #> The other a bar chart, small and neat,  
18 #> With simple data to complete.  
19 #>  
20 #> The joke here lies in wordplay's art,  
21 #> "Bar" means both sweets and stats, smart!  
22 #> A "Chocolate Bar" and a "Bar Chart" too,  
23 #> Mixed meanings make this humor true.  
24 #>  
25 #> So yes, a joke in this visual part,  
26 #> A pun that gives a smile to the heart!
```

## 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")`.