

Triage Against the Machine: Can AI Reason Deliberatively?

Francesco Veri, Gustavo Umbelino

2025-05-07

Large-Language Models (LLMs) Preview

Table 1: LLMs

	Provider	Model	Parameters (B)	Context Length	Architecture	Version
1	anthropic	claude-3-5-haiku-20241022	-	200000	-	2
2	anthropic	claude-3-5-sonnet-20241022	-	200000	-	2
3	anthropic	claude-3-7-sonnet-20250219	-	200000	-	3
4	anthropic	claude-3-haiku-20240307	-	200000	-	1
5	anthropic	claude-3-opus-20240229	-	200000	-	1
6	anthropic	claude-3-sonnet-20240229	-	200000	-	1
7	cohere	command	-	4096	-	1
8	cohere	command-a-03-2025	111	288000	dense, decoder-only	3
9	cohere	command-r-08-2024	32	128000	-	2
10	cohere	command-r-plus-08-2024	104	128000	dense, decoder-only	2
11	cohere	command-r7b-12-2024	7	128000	-	2
12	deepseek	deepseek-chat	671	128000	MoE	3
13	deepseek	deepseek-reasoner	671	128000	MoE	1
14	deepseek	deepseek-v2	NA	128000	-	1
15	deepseek	deepseek-v2.5	NA	128000	-	2
16	google	gemini-1.5-flash	-	1000000	MoE	1
17	google	gemini-1.5-flash-8b	8	1048576	MoE	1
18	google	gemini-1.5-pro	-	2000000	MoE	1
19	google	gemini-2.0-flash	-	1000000	-	2
20	google	gemini-2.0-flash-thinking-exp	NA	NA	NA	2
21	google	gemini-2.5-pro-preview-03-25	-	1048576	-	3
22	google	gemma	-	-	dense, decoder-only	1
23	google	gemma-3-27b-it	27	NA	NA	3
24	google	gemma2:27b	27	8190	dense, decoder-only	2
25	google	gemma3:12b	12	128000	-	3
26	ibm	granite3.3	8	131072	dense	3
27	meta	llama2:13b	13	4100	-	1

	Provider	Model	Parameters (B)	Context Length	Architecture	Version
28	meta	llama2:70b	70	4100	-	1
29	meta	llama3.1:405B-turbo	405	128000	-	3
30	meta	llama3.2	3	131072	-	4
31	meta	llama3.3:70b	70	128000	-	5
32	meta	llama3:70b	70	8190	-	2
33	meta	llama4-maverick	17	1000000	MoE	6
34	meta	llama4-scout	17	1000000000	MoE	6
35	microsoft	phi	NA	NA	-	1
36	microsoft	phi2	NA	NA	-	2
37	microsoft	phi3	NA	NA	-	3
38	microsoft	phi3.5	NA	NA	-	4
39	microsoft	phi4	14	16000	dense, decoder-only	5
40	mistralai	ministral-3b-latest	3	128000	-	1
41	mistralai	ministral-8b-latest	8	128000	-	1
42	mistralai	mistral-large-latest	123	128000	-	1
43	mistralai	mistral-small-latest	22	32800	-	1
44	mistralai	open-mistral-7b	7	NA	-	NA
45	mistralai	open-mistral-nemo	12	128000	-	1
46	mistralai	open-mixtral-8x22b	39	65400	SMoE	1
47	mistralai	open-mixtral-8x7b	7	NA	SMoE	NA
48	openai	gpt-3.5-turbo	-	16385	-	1
49	openai	gpt-4	-	8192	-	3
50	openai	gpt-4-turbo	-	128000	-	3
51	openai	gpt-4.5-preview	-	128000	-	4
52	openai	gpt-4o	-	128000	-	2
53	openai	gpt-4o-mini	-	128000	-	2
54	openai	o1	-	200000	-	1
55	openai	o1-mini	NA	NA	-	1
56	openai	o3-mini	-	200000	-	2
57	qwen	qwen-max	-	32768	-	1
58	qwen	qwen-plus	-	131072	-	1
59	qwen	qwen-turbo	-	1000000	-	1
60	qwen	qwen1.5-110b-chat	110	NA	-	1
61	qwen	qwen1.5-72b-chat	72	8000	-	1
62	qwen	qwen2-72b-instruct	72	131072	-	2
63	qwen	qwen2.5-72b-instruct	72	131072	-	3
64	qwen	qwq-plus	-	131072	-	1
65	xai	grok-2-1212	-	131072	-	2
66	xai	grok-3-beta	-	131072	-	3
67	xai	grok-3-mini-beta	-	131072	-	3
68	xai	grok-3-mini-beta-r=high	-	131072	-	3
69	xai	grok-3-mini-beta-r=low	-	131072	-	3
70	xai	grok-beta	314	131072	MoE	1

We started the analysis with 70 models, but some models were dropped after data collection. The models and reason for dropping are discussed later on Excluded Models.

Surveys

Table 2: Surveys

	survey	considerations	policies	scale_max	q_method
1	acp	48	5	11	FALSE
2	auscj	45	8	7	FALSE
3	bep	43	7	7	FALSE
4	biobanking_mayo_abc	38	7	11	FALSE
5	biobanking_wa	49	7	11	FALSE
6	ccps	33	7	11	FALSE
7	ds_aargau	33	7	7	FALSE
8	ds_bellinzona	32	7	7	FALSE
9	energy_futures	45	9	11	FALSE
10	fnqcj	42	5	12	FALSE
11	forestera	45	7	11	FALSE
12	fremantle	36	6	11	TRUE
13	gbr	35	7	7	FALSE
14	swiss_health	24	6	7	FALSE
15	uppsala_speaks	42	7	7	FALSE
16	valsamoggia	36	4	11	TRUE
17	zh_thalwil	31	7	7	FALSE
18	zh_uster	31	7	7	FALSE
19	zh_winterthur	30	6	7	FALSE
20	zukunft	20	7	7	FALSE

LLM Data Collection

Handle special models

command-r7b-12-2024-t=1 grok-3-beta-r=TRUE

We collected a total of 37460 valid LLM responses across 20 surveys.

Cost

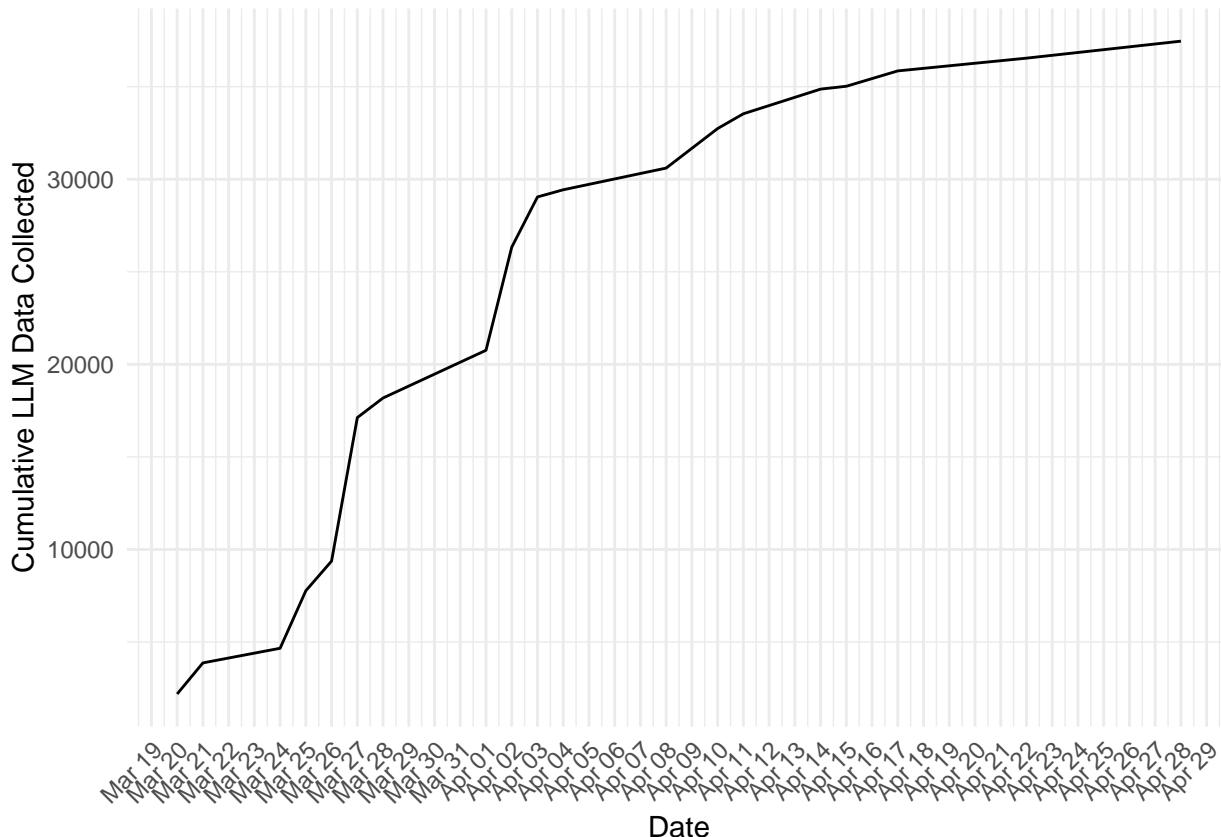
We spent a total of 411.3 USD. The cost breakdown per API is below.

Table 3: Costs by API

api	num_models	credits_paid
OpenAI API	9	225.52
Anthropic API	6	75.00
xAI API	6	29.95
Cohere API	6	20.34
Mistral AI API	8	20.00
Alibaba Cloud	8	17.49
Together AI	8	13.00
DeepSeek API	2	10.00
Google Cloud	8	NA
ollama	10	NA

Time

It took a total of 183 hours¹ across 39 days to complete data collection. Most of it was done in parallel. The first LLM response was collected on Thursday, Mar 20, 2025 and latest on Monday, Apr 28, 2025.



Excluded Models

17 out of 74 were excluded from the analysis for the following reasons.

Table 4: Excluded models and reasons

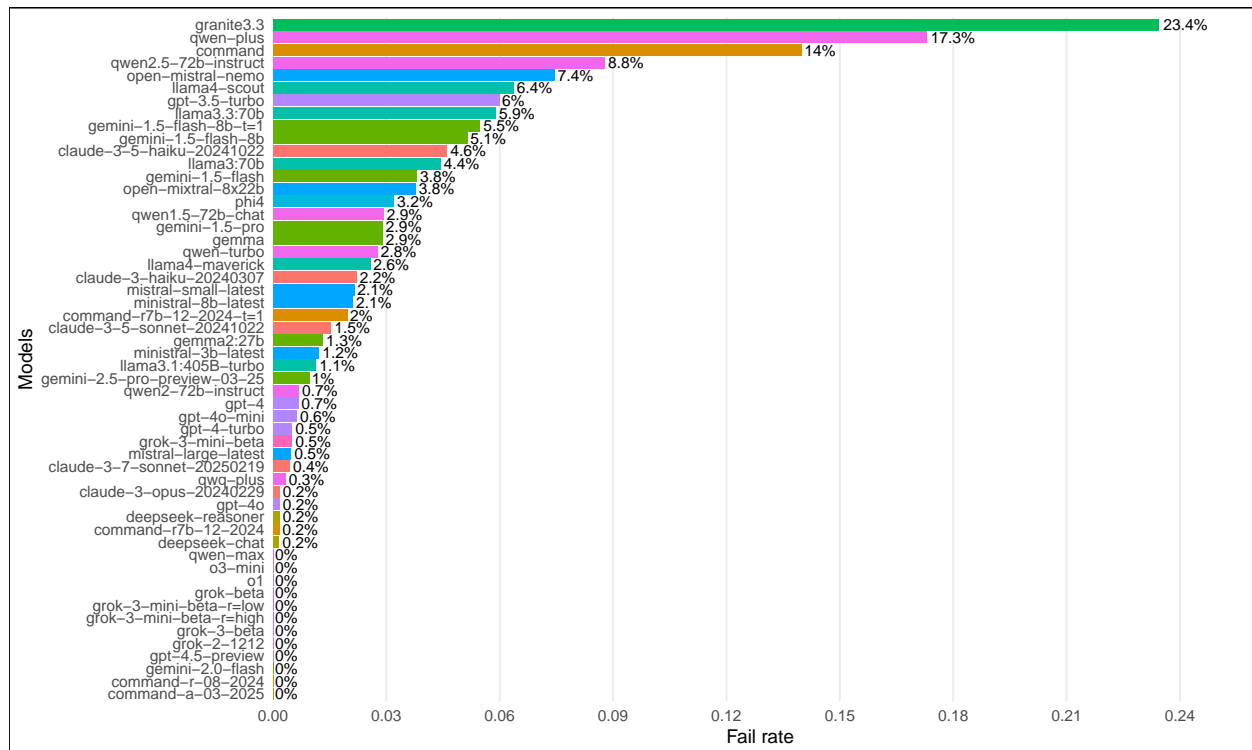
Provider	Model	Reason for exclusion
anthropic	claude-3-sonnet-20240229	not available in Anthropic API anymore
cohere	command-r-plus-08-2024	uniform aggregated considerations (1s)
deepseek	deepseek-v2	high fail rate (85%)
deepseek	deepseek-v2.5	too big to run locally; not available through APIs
google	gemma-3-27b-it	low rate limit (15K tokens/min)
google	gemma3:12b	uniform aggregated considerations (1s)
meta	llama2:13b	does not respond to prompts correctly
meta	llama2:70b	does not respond to prompts correctly
meta	llama3.2	3% success rate on auscj
microsoft	phi	does not respond to prompts correctly
microsoft	phi2	same model as phi
microsoft	phi3	does not respond to prompts correctly

¹Execution data is mostly accurate. Only a few (3-5) executions failed and, as a result, we have no record of it.

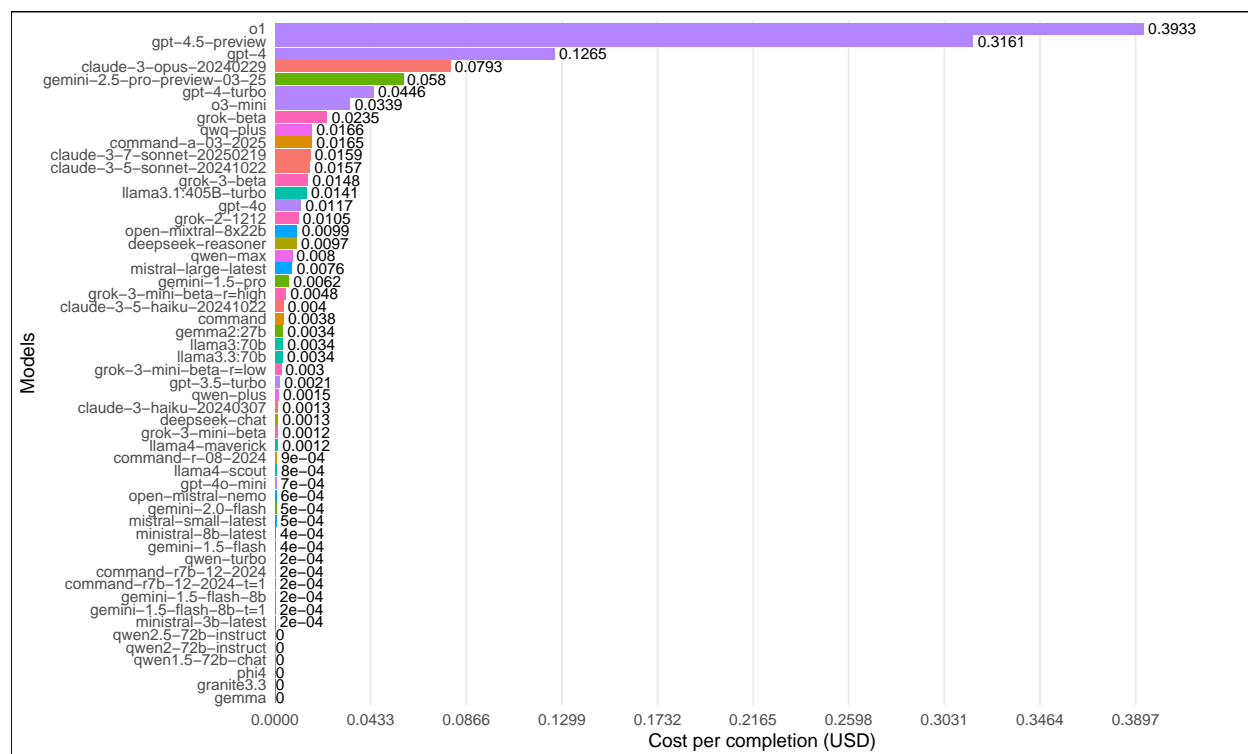
Provider	Model	Reason for exclusion
microsoft	phi3.5	10% success rate for biobanking_wa
mistralai	open-mistral-7b	11% success rate for auscj, uppsala_speaks, and biobanking_wa
mistralai	open-mixtral-8x7b	6% success rate on fremantle only
openai	o1-mini	0% success rate on uppsala_speaks only; responds with “I’m sorry, but I can’t help with that.”
qwen	qwen1.5-110b-chat	has API limit of 10 RPM; too slow

Execution Summary Plots

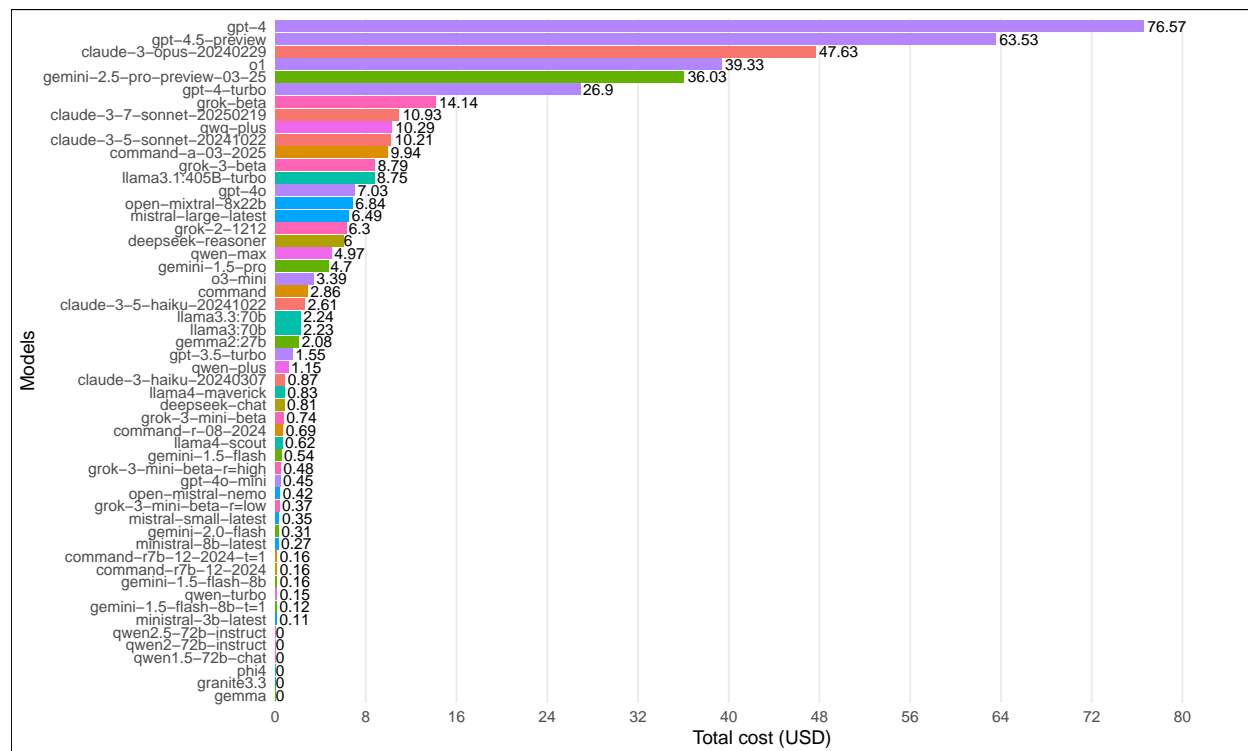
Fail rate



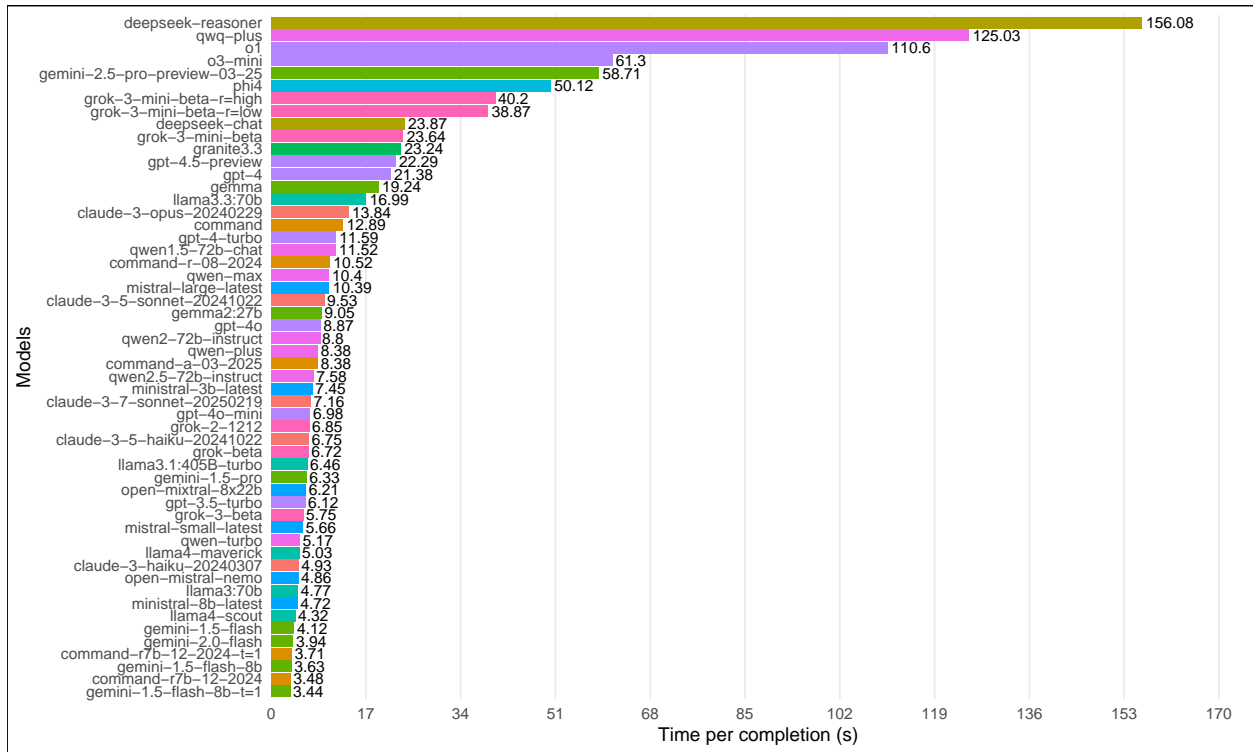
Cost per completion



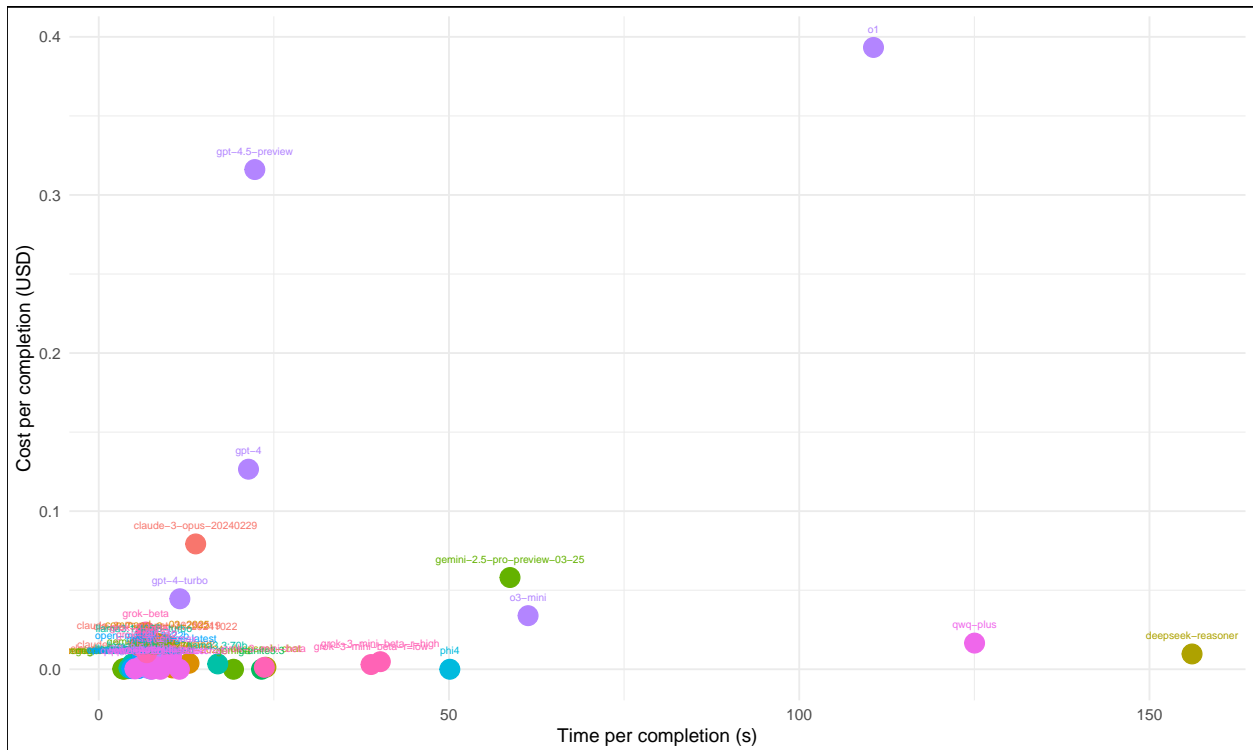
Total cost



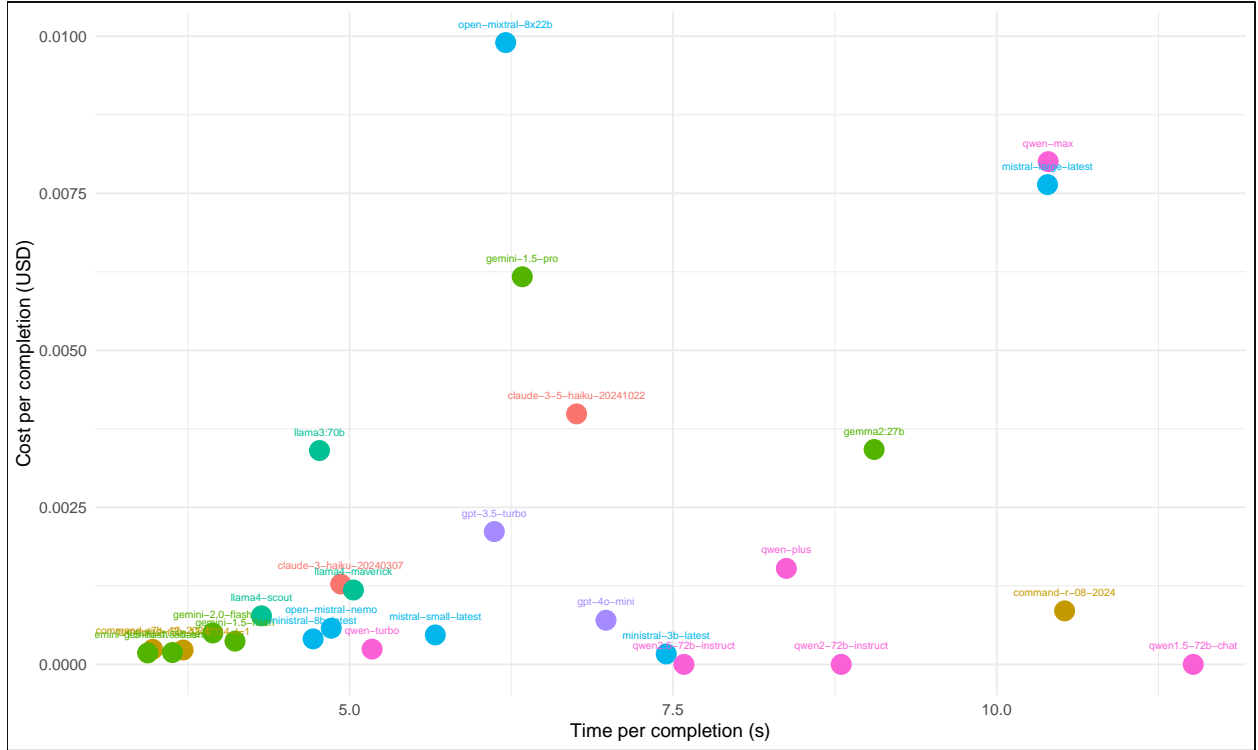
Time per completion



Cost/Time per completion



Zoomed in to cost < 0.01 USD and time < 12 s.



Internal Consistency of Responses

We calculate Cronbach's Alpha from the top 30 iterations.

Check alpha results per model

Table 5: Alpha summary across models, mean across surveys

	provider	model	N	all	considerations	policies
1	qwen	qwen1.5-72b-chat	600	0.70	0.75	0.49
2	google	gemma2:27b	600	0.71	0.75	0.50
3	meta	llama4-maverick	600	0.71	0.78	0.44
4	openai	gpt-4o-mini	600	0.72	0.74	0.45
5	anthropic	claude-3-haiku-20240307	600	0.74	0.82	0.44
6	google	gemin-1.5-flash	600	0.74	0.76	0.52
7	anthropic	claude-3-5-sonnet-20241022	600	0.75	0.81	0.58
8	deepseek	deepseek-reasoner	600	0.75	0.79	0.55
9	google	gemin-1.5-flash-8b-t=1	600	0.75	0.81	0.49
10	ibm	granite3.3	600	0.75	0.75	0.47
11	openai	gpt-4	600	0.75	0.82	0.52
12	openai	gpt-4-turbo	600	0.75	0.82	0.53
13	xai	grok-beta	600	0.75	0.85	0.49
14	google	gemin-1.5-pro	600	0.76	0.78	0.57
15	google	gemin-2.5-pro-preview-03-25	600	0.76	0.83	0.67
16	openai	gpt-4o	600	0.76	0.86	0.50
17	cohere	command	600	0.78	0.78	0.44
18	google	gemma	600	0.78	0.80	0.45
19	meta	llama3.3:70b	600	0.78	0.82	0.52
20	mistralai	mistral-small-latest	600	0.78	0.84	0.52

	provider	model	N	all	considerations	policies
21	mistralai	open-mistral-nemo	600	0.78	0.80	0.49
22	qwen	qwq-plus	600	0.78	0.79	0.58
23	xai	grok-2-1212	600	0.78	0.89	0.47
24	cohere	command-a-03-2025	600	0.79	0.86	0.51
25	cohere	command-r-08-2024	600	0.79	0.81	0.50
26	deepseek	deepseek-chat	600	0.79	0.86	0.52
27	google	gemini-1.5-flash-8b	600	0.79	0.84	0.50
28	meta	llama3:70b	600	0.79	0.79	0.52
29	qwen	qwen-turbo	600	0.79	0.83	0.48
30	anthropic	claude-3-7-sonnet-20250219	600	0.80	0.84	0.53
31	meta	llama4-scout	600	0.80	0.85	0.51
32	qwen	qwen-plus	600	0.80	0.82	0.49
33	qwen	qwen2-72b-instruct	600	0.80	0.86	0.48
34	qwen	qwen2.5-72b-instruct	600	0.80	0.84	0.51
35	xai	grok-3-mini-beta	600	0.80	0.78	0.67
36	anthropic	claude-3-5-haiku-20241022	600	0.81	0.86	0.47
37	microsoft	phi4	600	0.81	0.82	0.55
38	xai	grok-3-beta	600	0.81	0.84	0.53
39	mistralai	ministral-8b-latest	600	0.82	0.83	0.51
40	qwen	qwen-max	600	0.82	0.84	0.51
41	anthropic	claude-3-opus-20240229	600	0.83	0.87	0.50
42	mistralai	mistral-large-latest	600	0.83	0.86	0.54
43	google	gemini-2.0-flash	600	0.84	0.84	0.62
44	openai	gpt-3.5-turbo	600	0.84	0.87	0.48
45	openai	gpt-4.5-preview	201	0.84	0.87	0.70
46	cohere	command-r7b-12-2024-t=1	600	0.85	0.86	0.47
47	meta	llama3.1:405B-turbo	600	0.85	0.88	0.49
48	mistralai	ministral-3b-latest	600	0.85	0.86	0.53
49	cohere	command-r7b-12-2024	600	0.86	0.87	0.46
50	mistralai	open-mixtral-8x22b	600	0.87	0.90	0.52
51	xai	grok-3-mini-beta-r=high	100	0.91	0.90	0.81
52	xai	grok-3-mini-beta-r=low	124	0.91	0.89	0.80
53	openai	o1	100	0.92	0.92	0.77
54	openai	o3-mini	100	0.92	0.91	0.80

Human Data

Handle Swiss cases

```
## Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
## i Please use `all_of()` or `any_of()` instead.
##   # Was:
##   data %>% select(swiss_C_cols)
##
##   # Now:
##   data %>% select(all_of(swiss_C_cols))
##
## See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.
```

```
## Warning: Using an external vector in selections was deprecated in tidyselect 1.1.0.
## i Please use `all_of()` or `any_of()` instead.
##   # Was:
##   data %>% select(col)
##
##   # Now:
##   data %>% select(all_of(col))
##
## See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.
```

Table 6: Number of participants in each case study

	Case	Survey	Participants
1	Citizen Parliamentarian	acp	45
2	HGE Control Group	auscj	19
3	HGE Deliberative Group	auscj	23
4	BEP	bep	16
5	Mayo	biobanking_mayo_ubc	17
6	UBC Bio	biobanking_mayo_ubc	17
7	WA Citizens	biobanking_wa	9
8	WA Stakeholder	biobanking_wa	15
9	CCPS ACT Deliberative	ccps	31
10	Aargau	ds_aargau	16
11	Bellinzona	ds_bellinzona	8
12	CSIRO NSW	energy_futures	12
13	CSIRO WA	energy_futures	17
14	FNQCJ	fnqcj	11
15	Forest Lay Citizen	forestera	9
16	Forest Stakeholder	forestera	11
17	Fremantle	fremantle	41
18	GBR	gbr	7
19	CA	swiss_health	56
20	Activate	uppsala_speaks	26
21	Standard	uppsala_speaks	22
22	UPSA Control Group	uppsala_speaks	20
23	Valsamoggia	valsamoggia	16
24	Thalwil	zh_thalwil	14
25	USTER	zh_uster	15
26	Winterthur	zh_winterthur	16
27	Zukunft	zukunft	63

We collected 1144 human responses across 27 case studies, including pre-post deliberation responses.

Excluded cases

Table 7: Excluded cases

	Case	Survey	Participants	Excluded Reason
1	HGE Control Group	auscj	19	control group, no deliberation

	Case	Survey	Participants	Excluded Reason
2	GBR	gbr	7	unclear if human survey data is accurate
3	UPSA Control Group	uppsala_speaks	20	control group, no deliberation
4	Thalwill	zh_thalwil	14	??

We excluded 4 cases due to the reasons listed above.

Aggregation

We then aggregated LLM data into 1 response per model/survey. Based on (Motoki, Pinho Neto, and Rodrigues 2024), we bootstrap considerations 1000 times.

Aggregate considerations and preferences

We aggregated 33169 LLM responses into 1080 responses: 1 response per model per survey.

Randomly Generated Data

Then, we generated 20 random reseponses, one for each survey.

DRI Analysis

We begin by defining DRI calculation functions.

```
# original DRI formula
dri_calc <- function(data, v1, v2) {
  lambda <- 1 - (sqrt(2) / 2)
  dri <- 2 * (((1 - mean(abs((data[[v1]] - data[[v2]])) / sqrt(2)
))) - (lambda)) / (1 - (lambda))) - 1

  return(dri)
}

# updated DRI formula
# FIXME: only accounts for negligible positive correlations, but not negative ones
dri_calc_v2 <- function(data, v1, v2) {
  # Calculate orthogonal distance for each pair
  d <- abs((data[[v1]] - data[[v2]])) / sqrt(2))

  # Define lambda as in the original
  lambda <- 1 - (sqrt(2) / 2)

  # Calculate penalty: 0.5 if both correlations are in [0, 0.2], 1 otherwise
  penalty <- ifelse(data[[v1]] >= 0 & data[[v1]] <= 0.2 & #0.3
                    data[[v2]] >= 0 & data[[v2]] <= 0.2, # 0.3
                    0, 1)

  # Adjusted consistency per pair
  consistency <- (1 - d) * penalty
```

```
# Average consistency across all pairs
avg_consistency <- mean(consistency)

# Scale to [-1, 1] as in the original
dri <- 2 * ((avg_consistency - lambda) / (1 - lambda)) - 1

return(dri)
}

# updated DRI formula: penalizes both negligible
# positive and negative correlations in a scalar way.
dri_calc_v3 <- function(data, v1, v2) {
  d <- abs((data[[v1]] - data[[v2]]) / sqrt(2))
  lambda <- 1 - (sqrt(2) / 2)

  # Scalar penalty based on strength of signal (|r| and |q|)
  penalty <- ifelse(pmax(abs(data[[v1]]), abs(data[[v2]])) <= 0.2, pmax(abs(data[[v1]]), abs(data[[v2]])) *
    consistency <- (1 - d) * penalty
    avg_consistency <- mean(consistency)

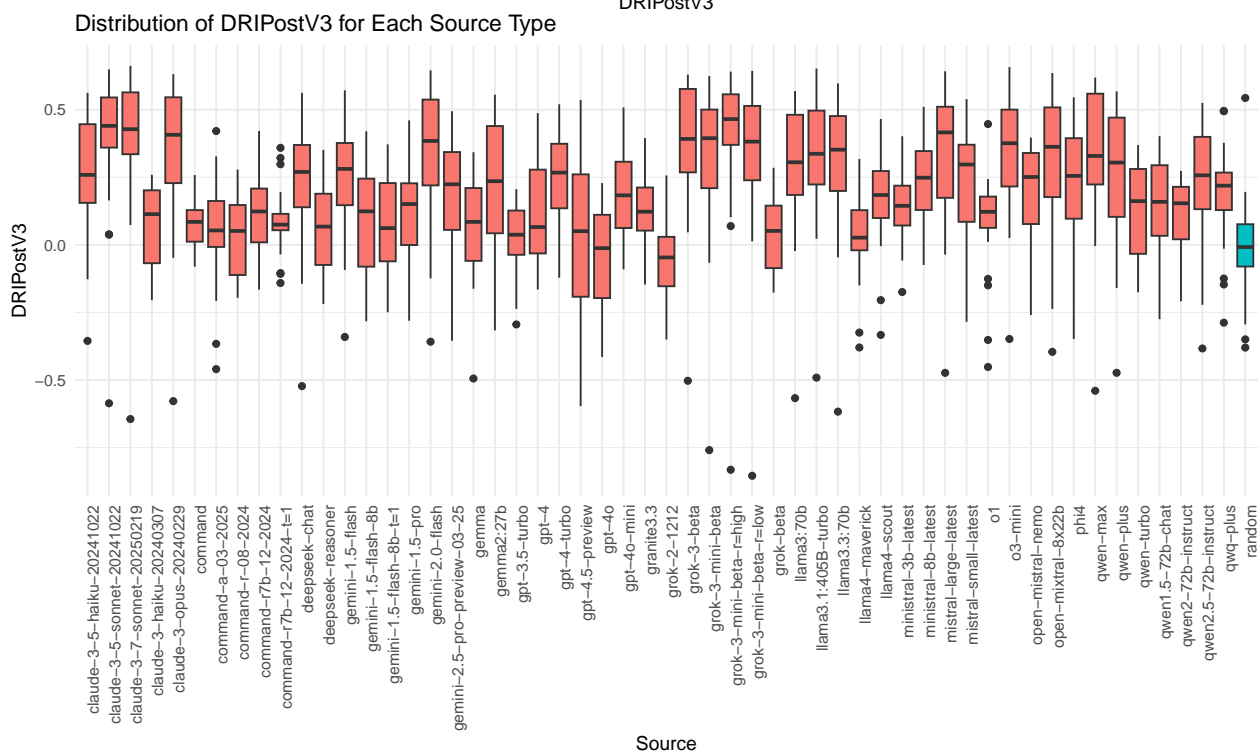
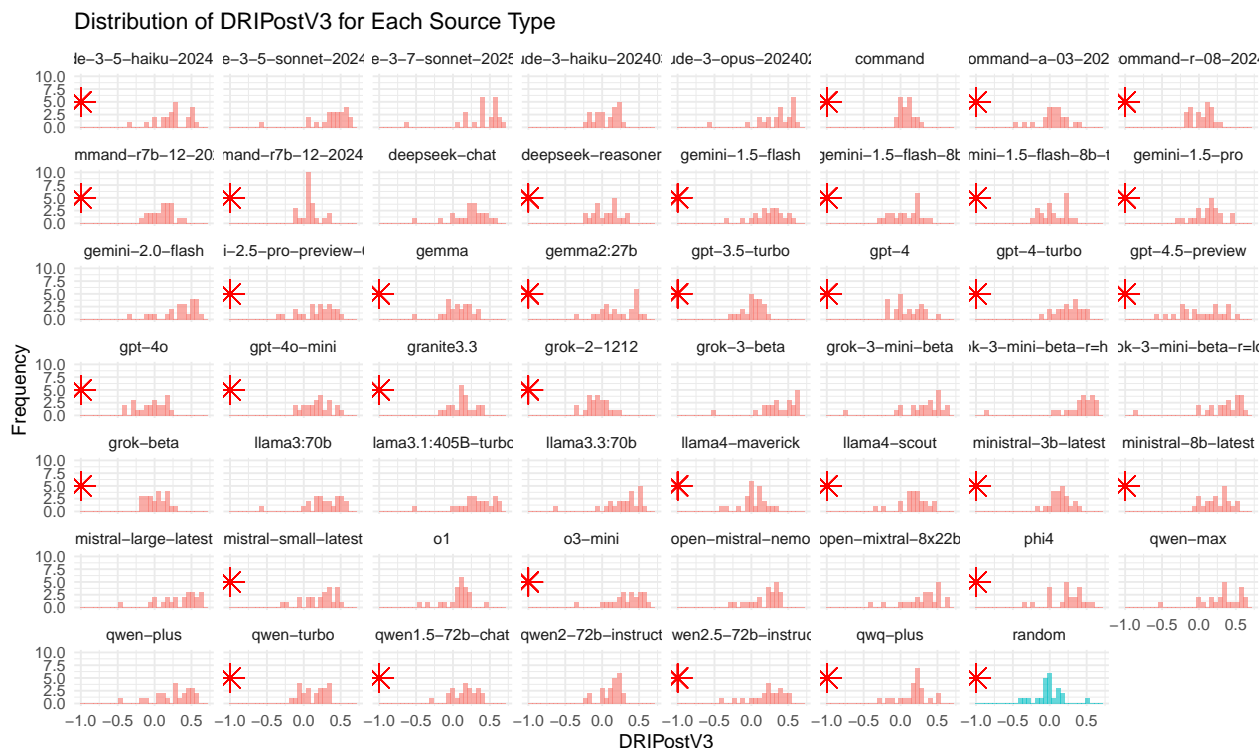
    dri <- 2 * ((avg_consistency - lambda) / (1 - lambda)) - 1
    return(dri)
  }
```

[illegible]

by Dunn's post-hoc test with Bonferroni correction.

The independent variable is be the type of participant (e.g., random, model). The dependent variable is the individual-level DRI score.

Adding missing grouping variables: `provider`, `model`



Testing hypothesis

```
##
## Kruskal-Wallis rank sum test
##
## data: DRIPostV3 by source
## Kruskal-Wallis chi-squared = 393.37, df = 54, p-value < 2.2e-16
```

Post-hoc tests

```
##
## alpha = 0.05
## Reject Ho if p <= alpha/2
```

Table 8: Models compared to random

Model	P-adjusted
grok-3-mini-beta-r=high	0*
claude-3-7-sonnet-20250219	0*
claude-3-5-sonnet-20241022	0*
grok-3-beta	0*
claude-3-opus-20240229	0*
grok-3-mini-beta-r=low	0*
gemini-2.0-flash	0*
o3-mini	0*
qwen-max	0*
llama3.1:405B-turbo	0*
grok-3-mini-beta	0*
mistral-large-latest	0*
open-mixtral-8x22b	0*
llama3.3:70b	0.001*
llama3:70b	0.002*
qwen-plus	0.009*
gpt-4-turbo	0.035*
claude-3-5-haiku-20241022	0.038*
gemini-1.5-flash	0.051
deepseek-chat	0.067
ministral-8b-latest	0.089
qwen2.5-72b-instruct	0.092
phi4	0.118
mistral-small-latest	0.139
gemma2:27b	0.161
open-mistral-nemo	0.784
claude-3-haiku-20240307	1
command	1
command-a-03-2025	1
command-r-08-2024	1
command-r7b-12-2024	1
command-r7b-12-2024-t=1	1
deepseek-reasoner	1
gemini-1.5-flash-8b	1
gemini-1.5-flash-8b-t=1	1
gemini-1.5-pro	1
gemini-2.5-pro-preview-03-25	1

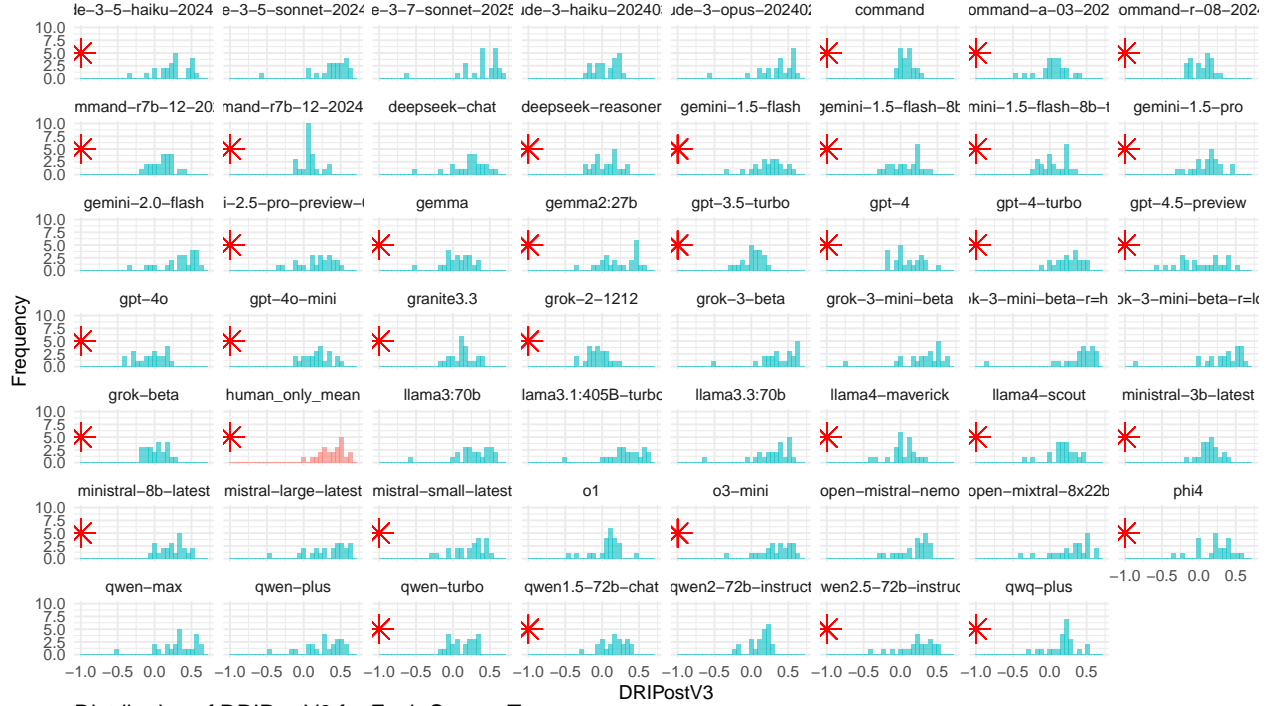
Model	P-adjusted
gemma	1
gpt-3.5-turbo	1
gpt-4	1
gpt-4.5-preview	1
gpt-4o	1
gpt-4o-mini	1
granite3.3	1
grok-2-1212	1
grok-beta	1
llama4-maverick	1
llama4-scout	1
ministral-3b-latest	1
o1	1
qwen-turbo	1
qwen1.5-72b-chat	1
qwen2-72b-instruct	1
qwq-plus	1

Some models, 18 out of 54, are significantly different than random.

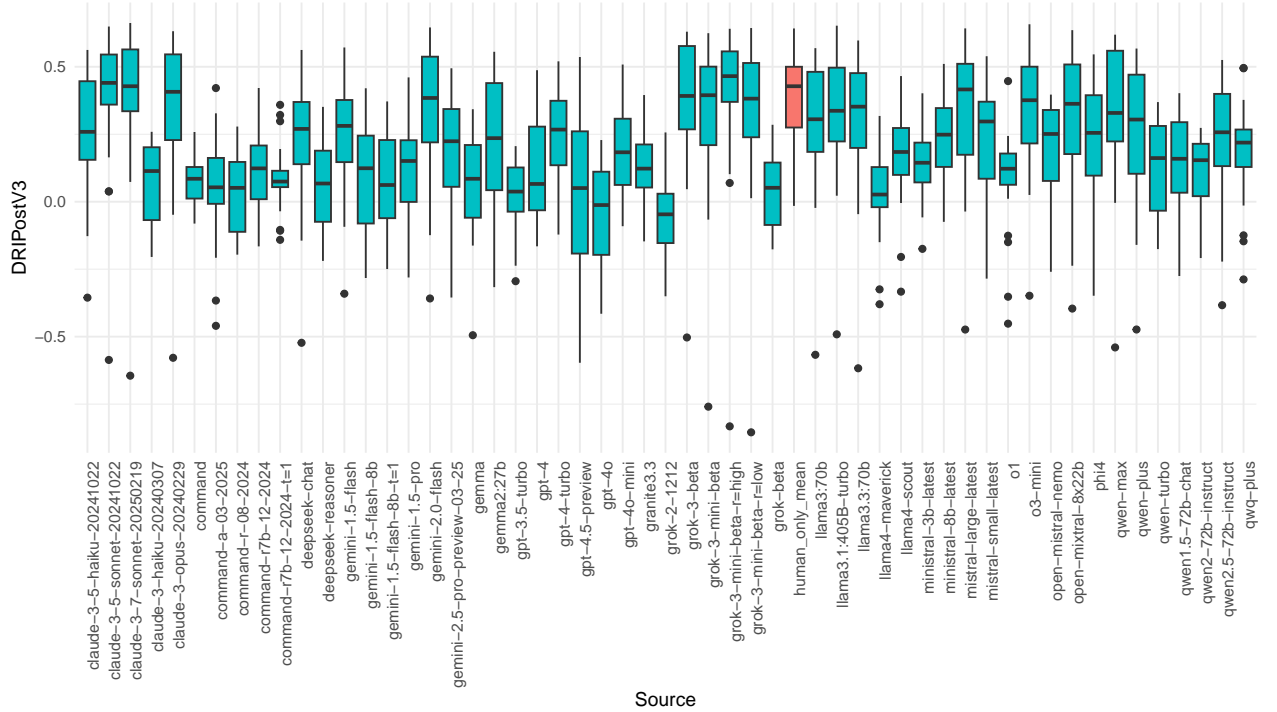
H2. LLMs' DRI scores will be significantly lower than those obtained from human participants after deliberation.

Testing assumptions

Distribution of DRIPostV3 for Each Source Type



Distribution of DRIPostV3 for Each Source Type



Testing hypothesis

To test H2, we will compare the average individual-level, post-deliberation DRI scores obtained by human participants with the individual-level DRI scores obtained by LLMs both across case studies and across LLM/version.

First, for each case study, we will employ a t-test (or non-parametric equivalent, depending on the results of the exploratory analysis) to analyze our results across case studies. The independent variable is participant type (human-only vs. LLM) and the dependent variable is the individual-level DRI scores.

For each case study...

human average

Second, for each LLM/version, we will employ a t-test (or non-parametric equivalent, depending on the results of the exploratory analysis) to analyze our results across LLM/version. The independent variable is participant type (human-only vs. LLM/version) and the dependent variable is the individual-level DRI scores.

```
##
## Kruskal-Wallis rank sum test
##
## data:  DRIPostV3 by source
## Kruskal-Wallis chi-squared = 395.62, df = 54, p-value < 2.2e-16
```

Post-hoc tests

```
##
## alpha = 0.05
## Reject Ho if p <= alpha/2
```

Table 9: Models compared to human

Model	P-adjusted
grok-2-1212	0*
gpt-4o	0*
gpt-3.5-turbo	0*
grok-beta	0*
command-r-08-2024	0*
human_only_mean - llama4-maverick	0*
command-a-03-2025	0*
claude-3-haiku-20240307	0.001*
deepseek-reasoner	0.001*
command	0.001*
gemini-1.5-flash-8b-t=1	0.001*
command-r7b-12-2024-t=1	0.001*
gemma	0.002*
gpt-4.5-preview	0.004*
human_only_mean - o1	0.006*
gemini-1.5-flash-8b	0.008*
gpt-4	0.015*
command-r7b-12-2024	0.017*
human_only_mean - qwen2-72b-instruct	0.019*
gemini-1.5-pro	0.047*
human_only_mean - qwen-turbo	0.061
granite3.3	0.081
human_only_mean - ministral-3b-latest	0.209
human_only_mean - qwen1.5-72b-chat	0.607

Model	P-adjusted
claude-3-5-haiku-20241022	1
claude-3-5-sonnet-20241022	1
claude-3-7-sonnet-20250219	1
claude-3-opus-20240229	1
deepseek-chat	1
gemini-1.5-flash	1
gemini-2.0-flash	1
gemini-2.5-pro-preview-03-25	1
gemma2:27b	1
gpt-4-turbo	1
gpt-4o-mini	1
grok-3-beta	1
grok-3-mini-beta	1
grok-3-mini-beta-r=high	1
grok-3-mini-beta-r=low	1
human_only_mean - llama3:70b	1
human_only_mean - llama3.1:405B-turbo	1
human_only_mean - llama3.3:70b	1
human_only_mean - llama4-scout	1
human_only_mean - ministral-8b-latest	1
human_only_mean - mistral-large-latest	1
human_only_mean - mistral-small-latest	1
human_only_mean - o3-mini	1
human_only_mean - open-mistral-nemo	1
human_only_mean - open-mixtral-8x22b	1
human_only_mean - phi4	1
human_only_mean - qwen-max	1
human_only_mean - qwen-plus	1
human_only_mean - qwen2.5-72b-instruct	1
human_only_mean - qwq-plus	1

H3. LLMs' DRI scores are improving over time, across each version.

Random slope –

Assume each case Multilevel analysis – each case behave differently

LMER –

To test H3, we will conduct a repeated measures ANOVA (or Friedman test if the assumptions of normality or sphericity are violated) to test for differences in the mean DRI across all versions (e.g., v1, v2, v3) of an LLM across each case study. We will treat different LLM versions as related groups and the individual-level LLM DRI in each case study as a subject. In this within-subjects design, we can assess whether more recent versions of LLMs have a significant impact on the DRI scores they produce.

We want to assess the effects of Case and Series on weight loss in 10 sedentary individuals.

Dependent variable: - DRIPostV3

Independent variables: - [LLM series (moderator) – which llm?] - [case (moderator) – which case?] – [LATER]
- version (focal)

##

Attaching package: 'rstatix'

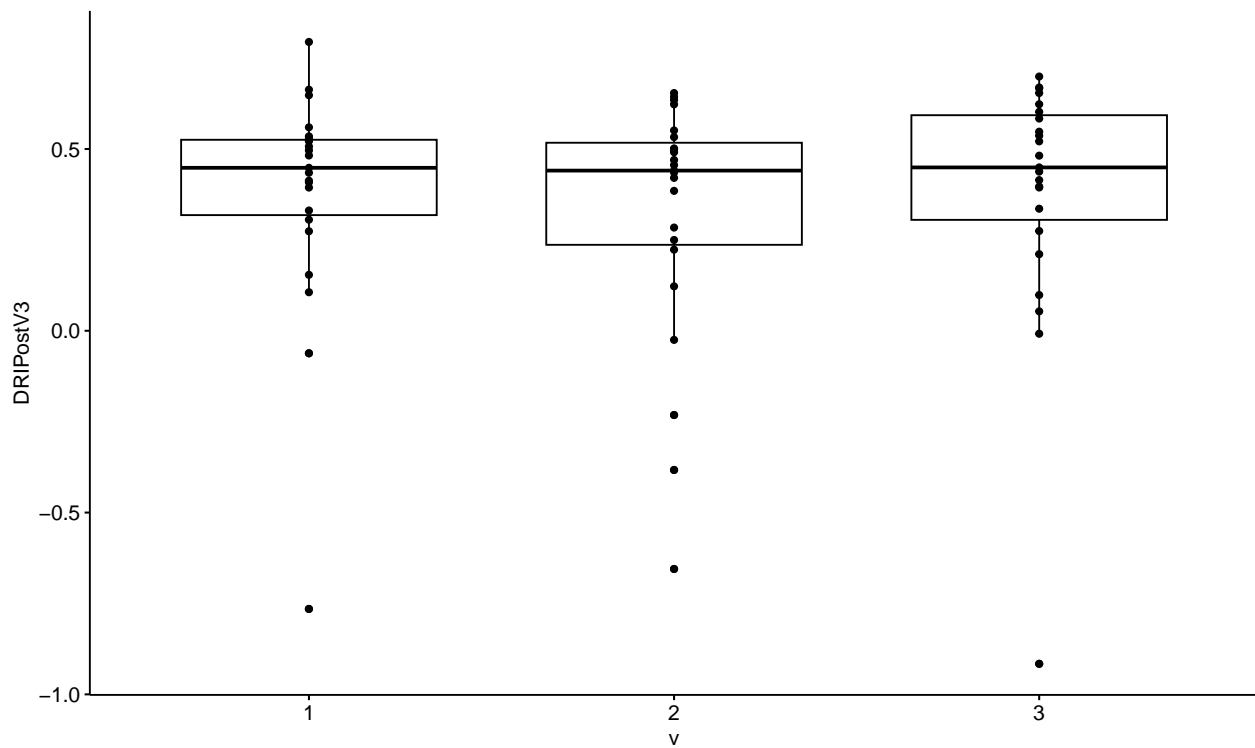
```
## The following object is masked from 'package:stats':
```

```
##
```

```
## filter
```

```
## # A tibble: 3 x 11
```

##	v	variable	n	min	max	median	iqr	mean	sd	se	ci
##	<fct>	<fct>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	1	DRIPostV3	23	-0.765	0.794	0.448	0.207	0.378	0.312	0.065	0.135
## 2	2	DRIPostV3	23	-0.655	0.654	0.441	0.281	0.318	0.344	0.072	0.149
## 3	3	DRIPostV3	23	-0.916	0.699	0.449	0.288	0.379	0.347	0.072	0.15

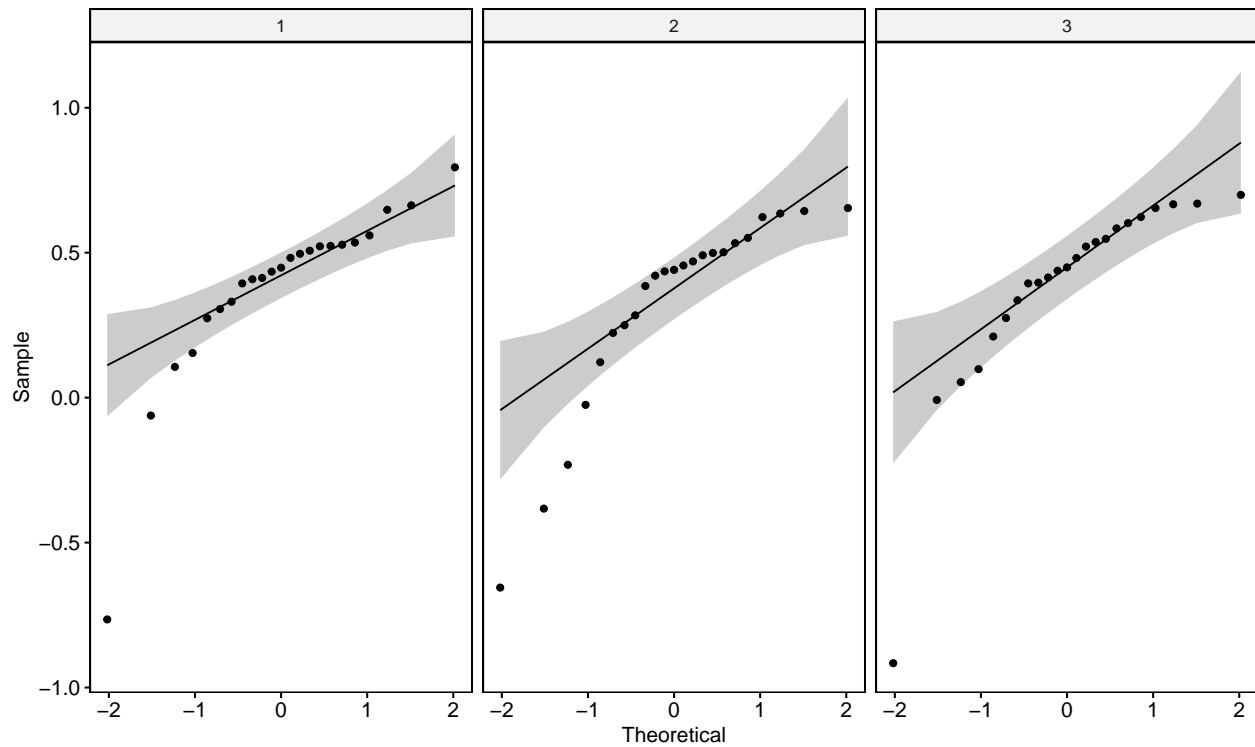


```
## # A tibble: 6 x 6
```

##	v	id	version	DRIPostV3	is.outlier	is.extreme
##	<fct>	<fct>	<dbl>	<dbl>	<lgl>	<lgl>
## 1	1	open-qwen/Citizen Parliamentari~	1	-0.0618	TRUE	FALSE
## 2	1	open-qwen/HGE Deliberative Group	1	-0.765	TRUE	TRUE
## 3	2	open-qwen/HGE Deliberative Group	2	-0.655	TRUE	TRUE
## 4	2	open-qwen/Aargau	2	-0.383	TRUE	FALSE
## 5	2	open-qwen/CSIRO NSW	2	-0.232	TRUE	FALSE
## 6	3	open-qwen/HGE Deliberative Group	3	-0.916	TRUE	TRUE

```
## # A tibble: 3 x 4
```

##	v	variable	statistic	p
##	<fct>	<chr>	<dbl>	<dbl>
## 1	1	DRIPostV3	0.772	0.000143
## 2	2	DRIPostV3	0.817	0.000734
## 3	3	DRIPostV3	0.745	0.0000561

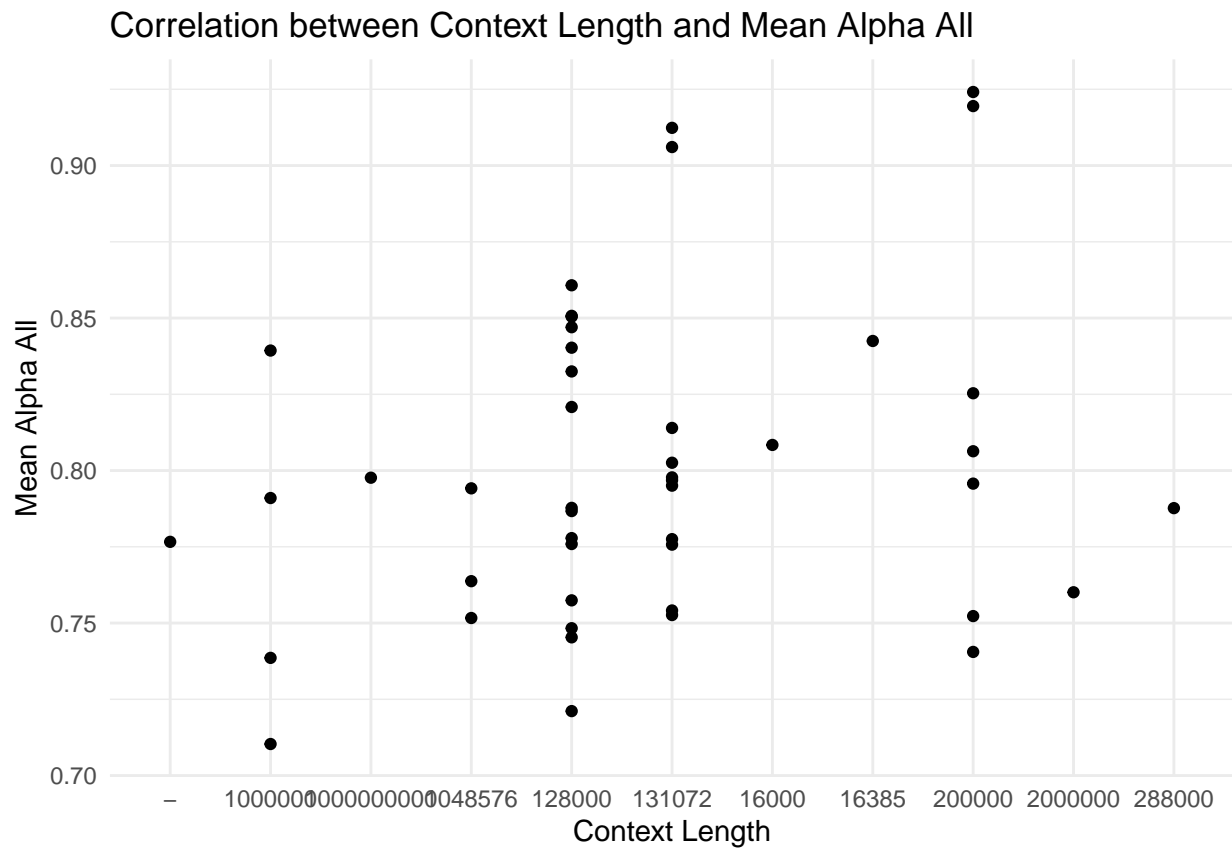


```
## # A tibble: 1 x 6
##   .y.      n statistic    df      p method
## * <chr>   <int>    <dbl> <dbl> <dbl> <chr>
## 1 DRIPostV3    23    0.783     2 0.676 Friedman test
```

If a significant difference is found, we will conduct a post-hoc analysis using paired t-tests (or Wilcoxon signed-rank tests) for pairwise comparisons, with adjustments for multiple comparisons.

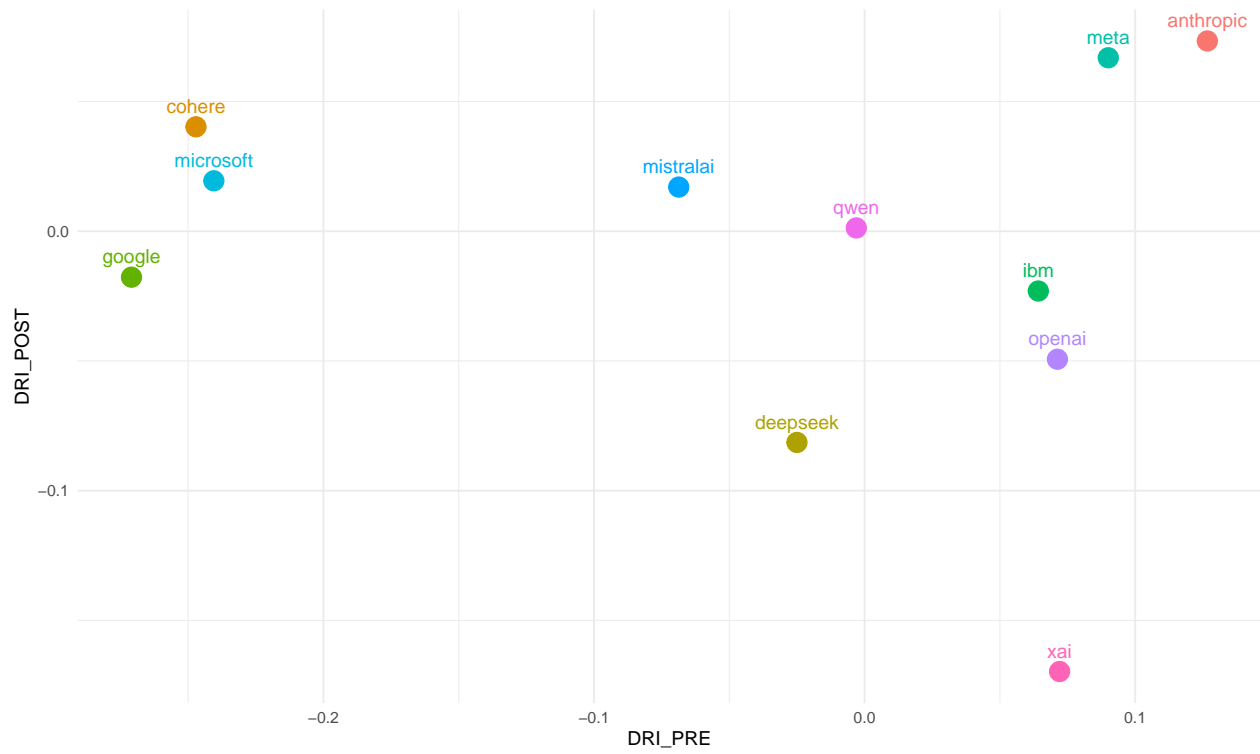
DRI Benchmark

```
## `geom_smooth()` using formula = 'y ~ x'
```

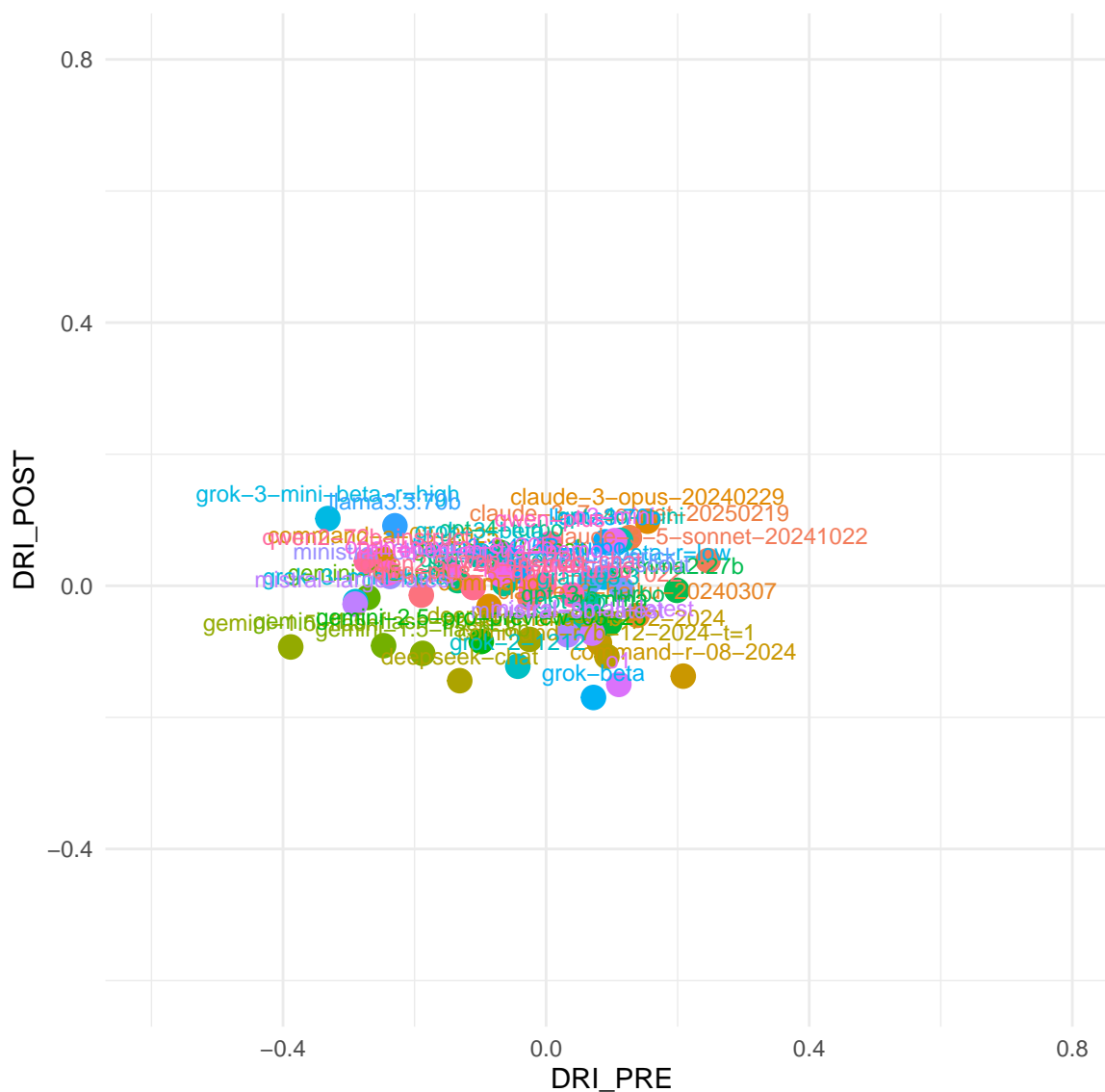


```
## `summarise()` has grouped output by 'provider', 'model'. You can override using
## the ``.groups` argument.
```

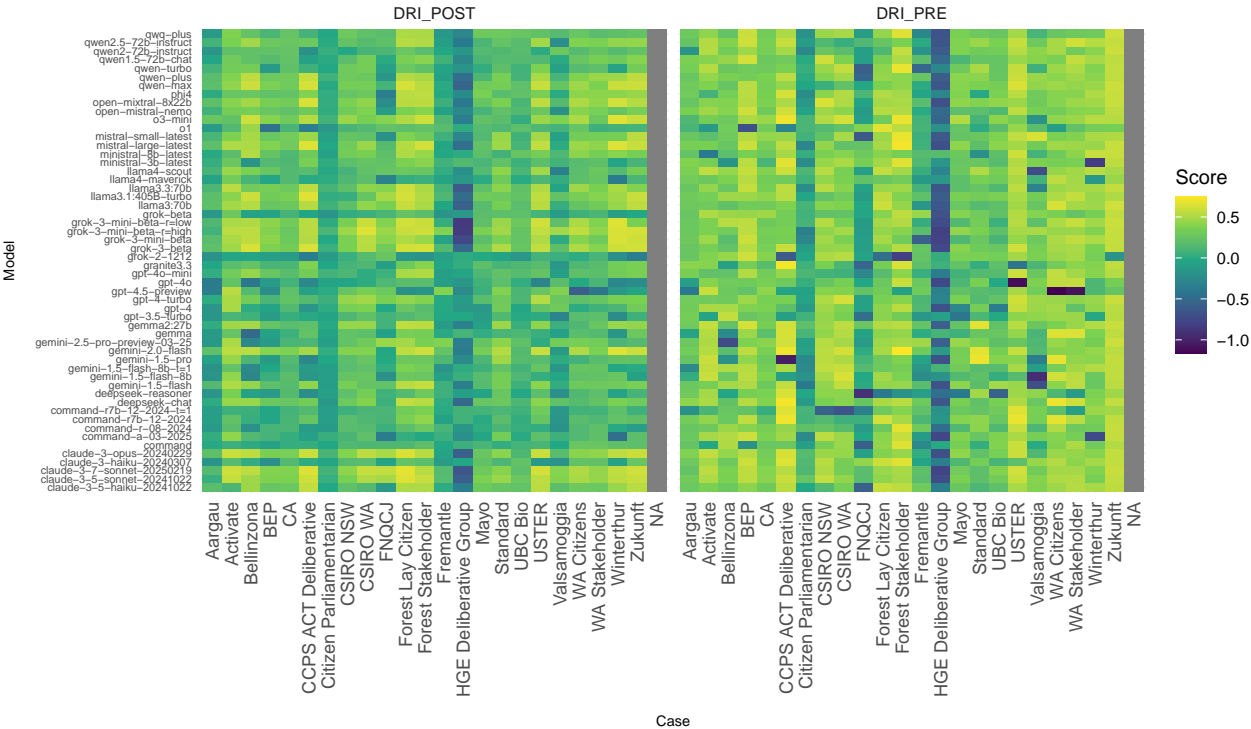
Comparison PRE and POST DRI by Provider



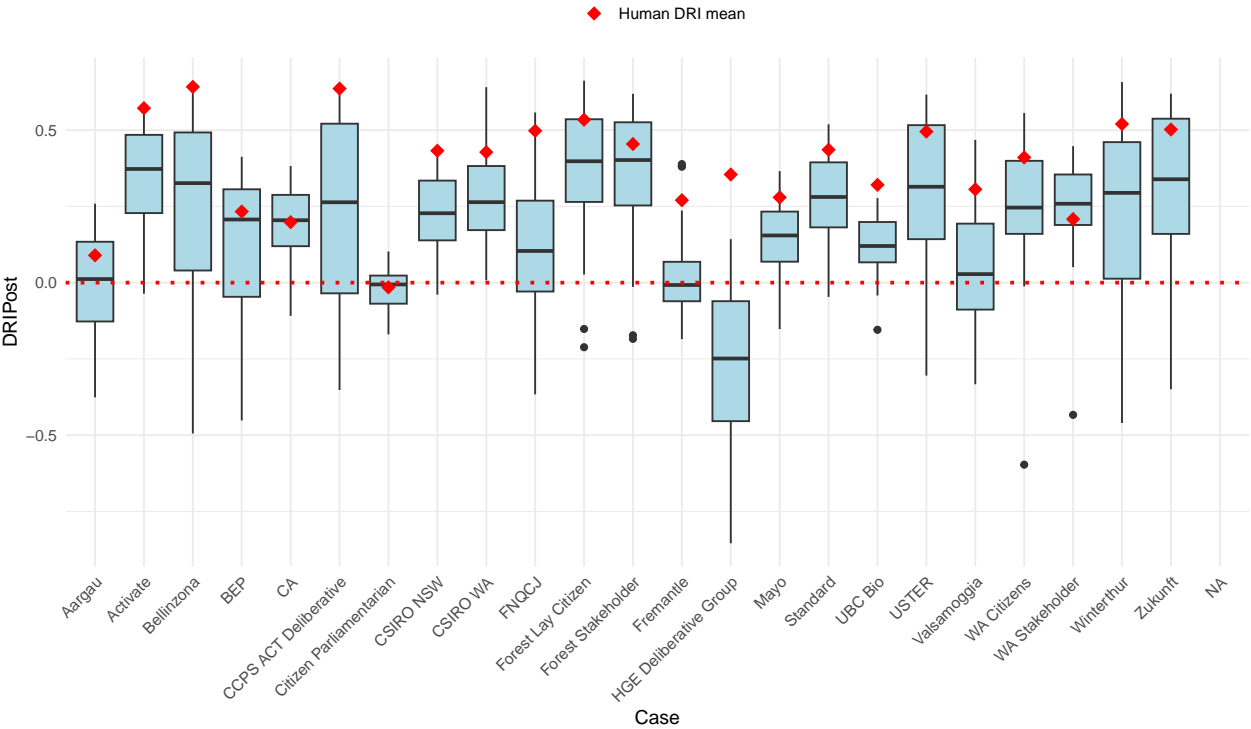
Comparison PRE and POST DRI by Model



Heatmap of DRI Scores by Case and Model



Boxplot of LLM DRI Post by Case

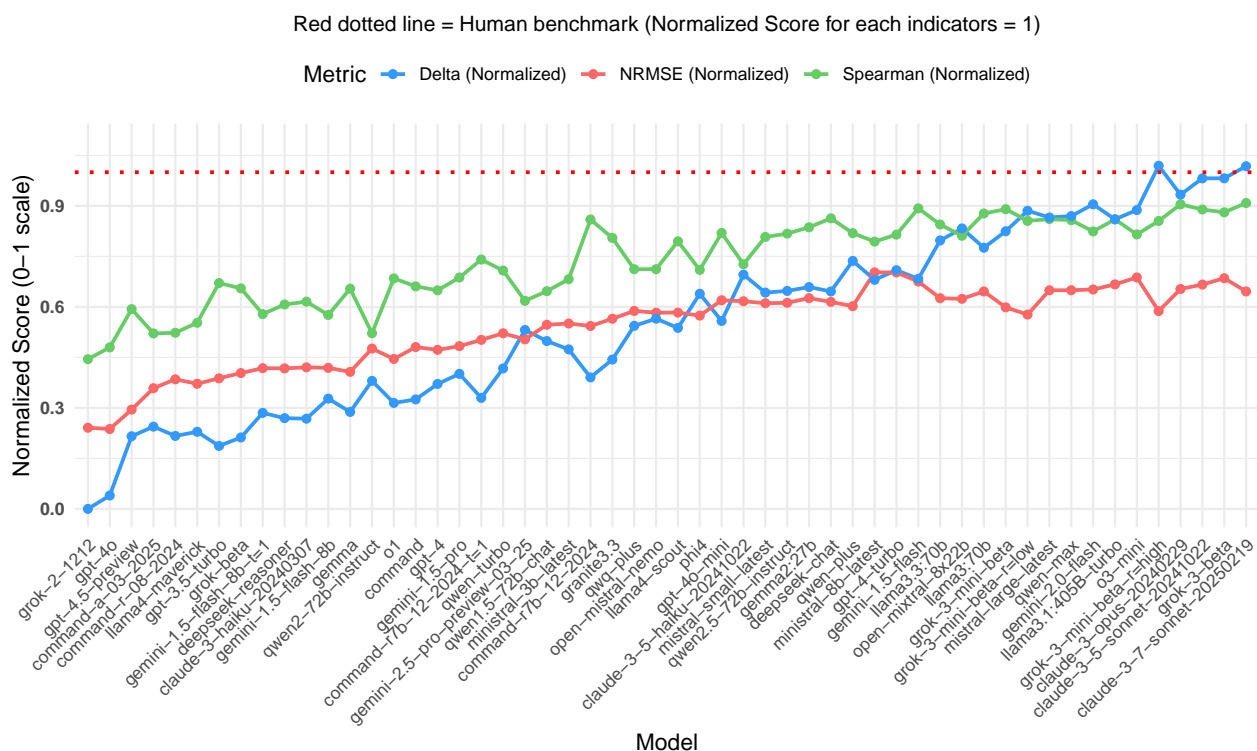


LLM Performance Metrics Against Human DRI Post-Scores

Table 10: LLM Performance Metrics Against Human DRI Post-Scores

Model	MAE	RMSE	MAPE (%)	Human Range	NMAE	NRMSE	Spearman	Delta
ministral-8b-latest	0.154	0.196	57.984	0.657	0.235	0.297	0.588	-0.142
gpt-4-turbo	0.147	0.196	57.296	0.657	0.223	0.298	0.629	-0.129
o3-mini	0.147	0.205	68.782	0.657	0.223	0.312	0.630	-0.050
grok-3-beta	0.128	0.207	57.658	0.657	0.194	0.315	0.762	-0.008
gemini-1.5-flash	0.162	0.214	64.731	0.657	0.246	0.325	0.786	-0.141
llama3.1:405B-turbo	0.136	0.219	52.371	0.657	0.207	0.333	0.720	-0.062
claude-3-5-sonnet-20241022	0.120	0.219	52.574	0.657	0.183	0.334	0.779	-0.008
claude-3-opus-20240229	0.132	0.228	75.162	0.657	0.201	0.347	0.809	-0.030
gemini-2.0-flash	0.156	0.229	57.281	0.657	0.237	0.348	0.648	-0.042
qwen-max	0.142	0.230	46.611	0.657	0.216	0.350	0.715	-0.058
mistral-large-latest	0.152	0.231	51.577	0.657	0.231	0.351	0.721	-0.060
llama3:70b	0.146	0.233	67.090	0.657	0.222	0.354	0.755	-0.100
claude-3-7-sonnet-20250219	0.129	0.233	65.055	0.657	0.197	0.354	0.817	0.008
gemma2:27b	0.169	0.246	49.815	0.657	0.258	0.374	0.673	-0.152
llama3.3:70b	0.145	0.246	70.392	0.657	0.220	0.374	0.690	-0.090
open-mixtral-8x22b	0.153	0.247	57.272	0.657	0.233	0.376	0.623	-0.074
gpt-4o-mini	0.210	0.250	77.454	0.657	0.320	0.380	0.639	-0.197
claude-3-5-haiku-20241022	0.168	0.252	45.543	0.657	0.255	0.383	0.454	-0.135
deepseek-chat	0.176	0.253	81.490	0.657	0.268	0.385	0.726	-0.158
qwen2.5-72b-instruct	0.184	0.255	53.739	0.657	0.280	0.388	0.635	-0.157
mistral-small-latest	0.187	0.256	64.667	0.657	0.284	0.389	0.616	-0.159
qwen-plus	0.172	0.262	73.468	0.657	0.262	0.398	0.638	-0.117
grok-3-mini-beta	0.144	0.264	47.289	0.657	0.218	0.401	0.781	-0.078
qwq-plus	0.210	0.271	58.542	0.657	0.319	0.412	0.424	-0.203
grok-3-mini-beta-r=high	0.146	0.271	74.208	0.657	0.222	0.412	0.710	0.009
llama4-scout	0.214	0.274	56.324	0.657	0.325	0.417	0.590	-0.206
open-mistral-nemo	0.211	0.274	62.660	0.657	0.321	0.417	0.424	-0.194
grok-3-mini-beta-r=low	0.155	0.278	53.106	0.657	0.236	0.423	0.711	-0.051
phi4	0.198	0.280	62.000	0.657	0.302	0.426	0.420	-0.161
granite3.3	0.248	0.286	67.160	0.657	0.377	0.435	0.610	-0.248
ministral-3b-latest	0.237	0.295	64.480	0.657	0.360	0.449	0.365	-0.234
qwen1.5-72b-chat	0.238	0.298	65.981	0.657	0.362	0.453	0.293	-0.223
command-r7b-12-2024	0.271	0.300	94.366	0.657	0.413	0.457	0.719	-0.271
qwen-turbo	0.260	0.314	65.480	0.657	0.396	0.478	0.416	-0.259
gemini-2.5-pro-preview-03-25	0.228	0.326	81.547	0.657	0.347	0.496	0.236	-0.209
command-r7b-12-2024-t=1	0.298	0.327	105.656	0.657	0.454	0.498	0.481	-0.298
gemini-1.5-pro	0.269	0.340	70.472	0.657	0.409	0.517	0.375	-0.267
command	0.305	0.341	81.164	0.657	0.464	0.519	0.322	-0.300
qwen2-72b-instruct	0.286	0.344	87.907	0.657	0.435	0.524	0.043	-0.276
gpt-4	0.282	0.347	85.125	0.657	0.429	0.528	0.298	-0.280
o1	0.305	0.365	119.839	0.657	0.464	0.555	0.370	-0.305
claude-3-haiku-20240307	0.327	0.381	95.595	0.657	0.497	0.579	0.231	-0.326
gemini-1.5-flash-8b	0.301	0.382	101.860	0.657	0.458	0.581	0.152	-0.299
gemini-1.5-flash-8b-t=1	0.321	0.383	107.309	0.657	0.488	0.582	0.157	-0.318
deepseek-reasoner	0.325	0.383	105.435	0.657	0.495	0.583	0.214	-0.325
gemma	0.317	0.390	94.849	0.657	0.482	0.593	0.308	-0.317
grok-beta	0.351	0.392	134.820	0.657	0.533	0.596	0.310	-0.351
gpt-3.5-turbo	0.362	0.402	107.022	0.657	0.551	0.612	0.342	-0.362

Human-Normalized Performance

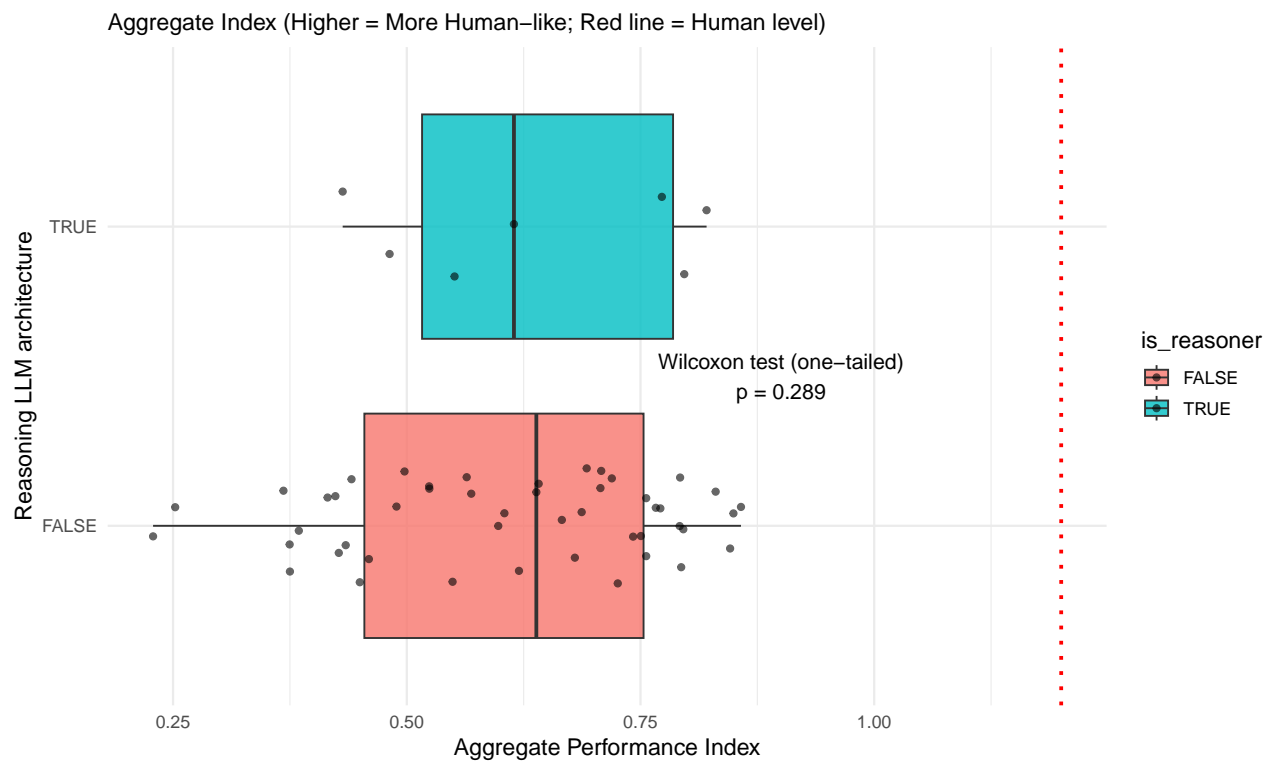


LLM Performance by Reasoner Classification

Architecture types:

- Transformer-based models (Vaswani et al. 2017).

Some models are considered “reasoning” models, like , reason using chain-of-thought (CoT) – this is not a difference in architecture



References

- Motoki, Fabio, Valdemar Pinho Neto, and Victor Rodrigues. 2024. "More Human Than Human: Measuring ChatGPT Political Bias." *Public Choice* 198(1): 3–23. doi:10.1007/s11127-023-01097-2.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. "Attention Is All You Need." In Curran Associates, Inc. https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.