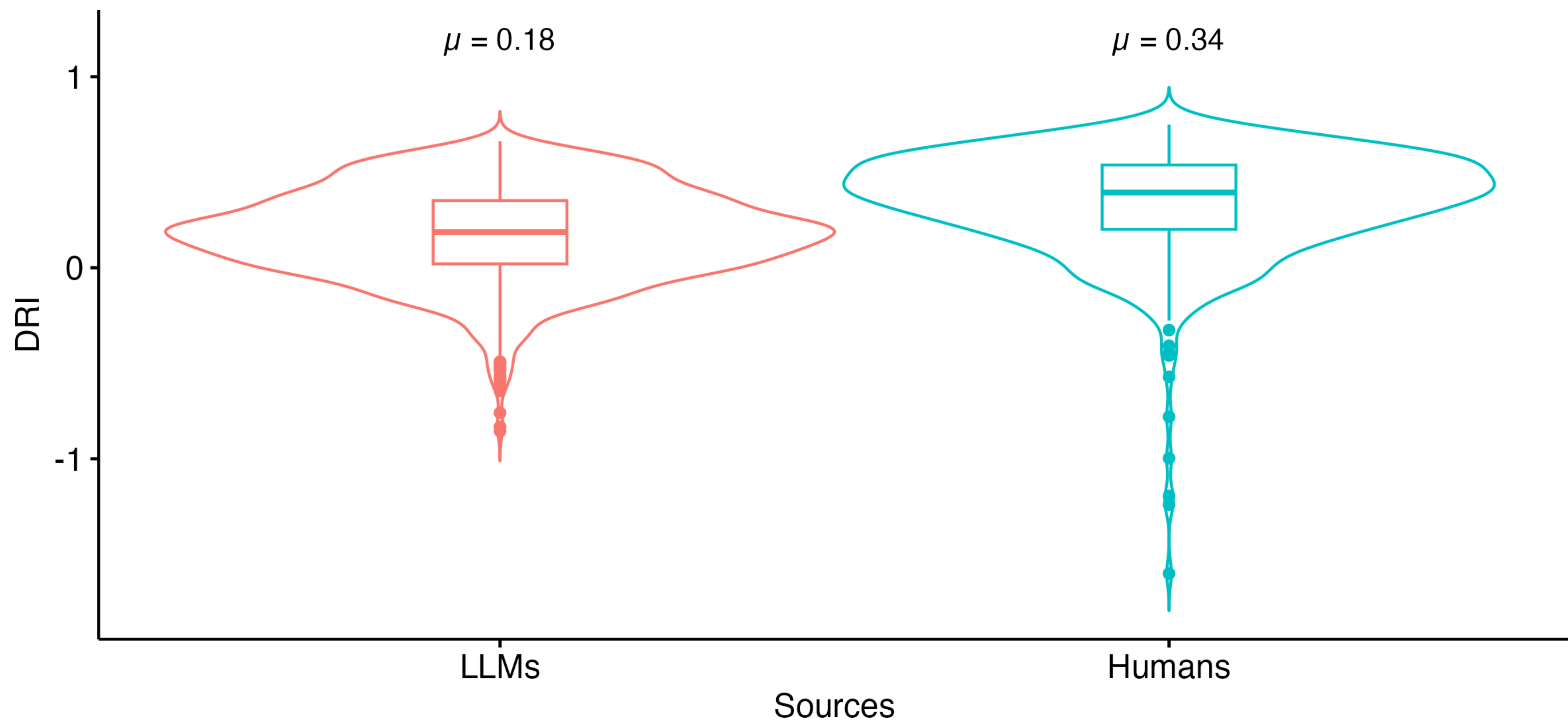# DRI LLM vs. Humans

## Human-AI Collaboration in Deliberation
## Evaluating LLMs Against Human Collective Wisdom

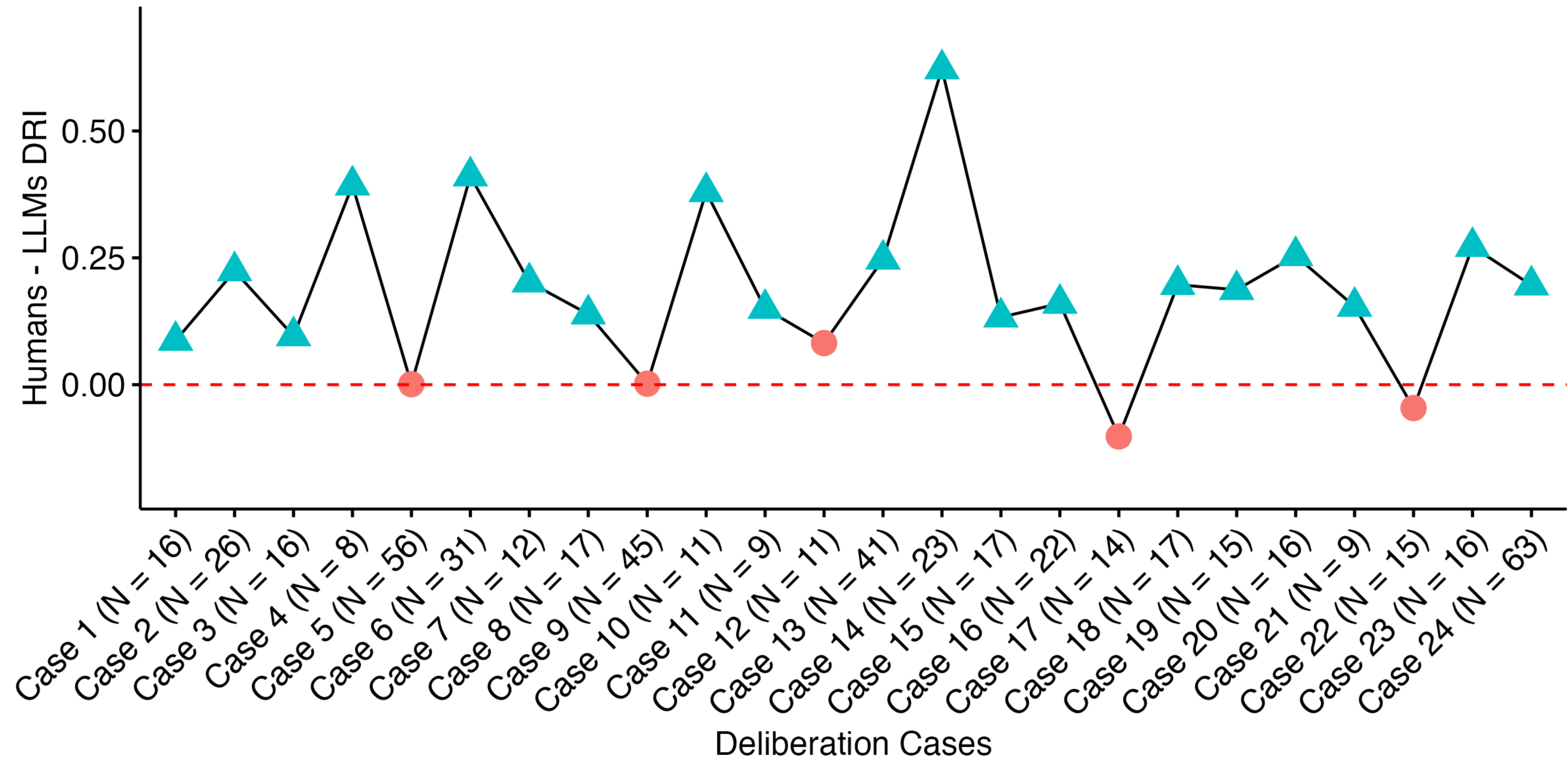**Veri & Umbelino (2025)**

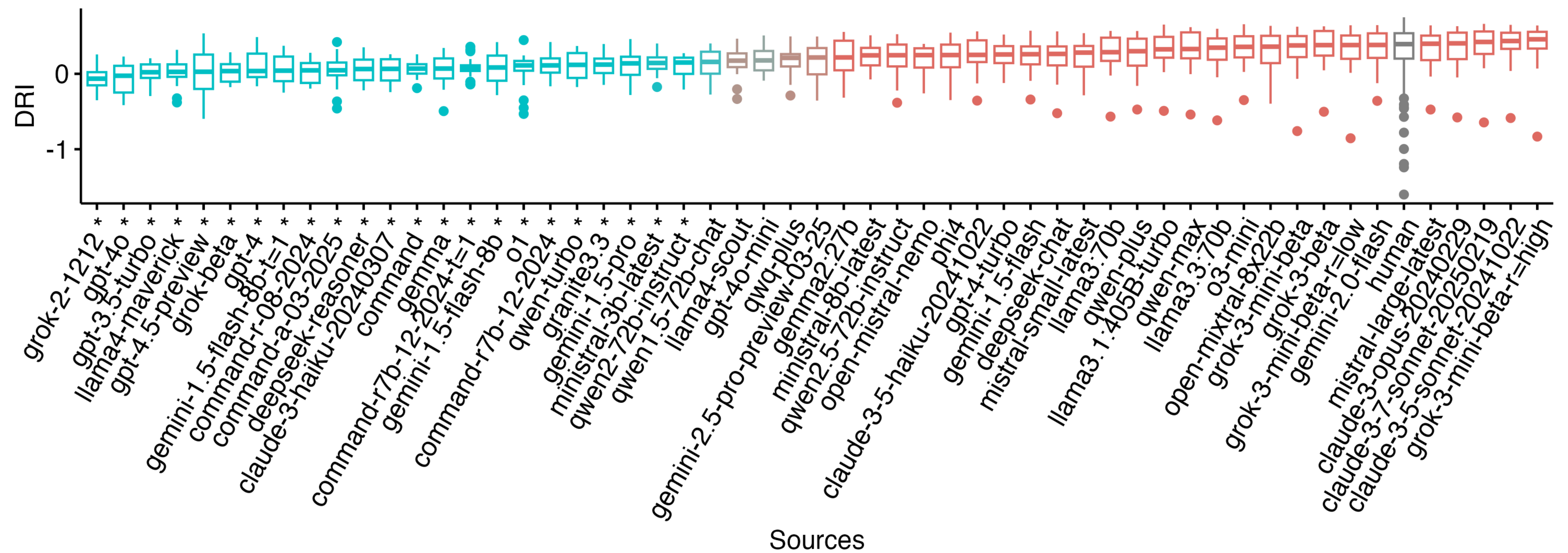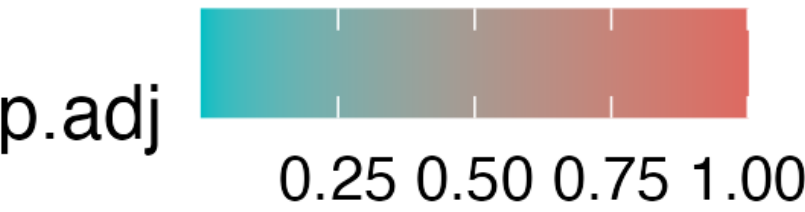# 1. LLMs are promising for emulating human reasoning

# Humans outperform LLMs across cases

# Most LLMs perform as well as humans



Kruskal-Wallis, $\chi^2(54) = 543.87$, $p = <0.0001$, $n = 1822$

p.adj
0.25 0.50 0.75 1.00

DRI

Sources

grok-2-1212 * / gpt-4o * / gpt-3.5-turbo * / llama4-maverick * / gpt-4.5-preview * / grok-beta * / gpt-4 * / gemini-1.5-flash-8b-t=1 * / command-a-03-2024 * / command-r-08-2025 * / deepseek-reasoner * / claude-3-haiku-20240307 * / command * / gemma * / command-r7b-12-2024-t=1 * / gemini-1.5-flash-8b * / o1 * / command-r7b-12-2024 * / qwen-turbo * / granite3.3 * / gemini-1.5-pro * / ministral-3b-latest * / qwen2-72b-instruct * / qwen1.5-72b-chat / llama4-scout / gpt-4o-mini / qwq-plus / gemini-2.5-pro-preview-03-25 / gemma2-27b / ministral-8b-latest / qwen2.5-72b-instruct / open-mistral-nemo / phi4 / claude-3-5-haiku-20241022 / gpt-4-turbo / gemini-1.5-flash / deepseek-chat / mistral-small-latest / llama3:70b / qwen-plus / llama3.1:405B-turbo / qwen-max / llama3.3:70b / o3-mini / open-mixtral-8x22b / grok-3-mini-beta / grok-3-beta / grok-3-mini-beta-r=low / gemini-2.0-flash / human / mistral-large-latest / claude-3-opus-20240229 / claude-3-7-sonnet-20250219 / claude-3-5-sonnet-20241022 / grok-3-mini-beta-r=high
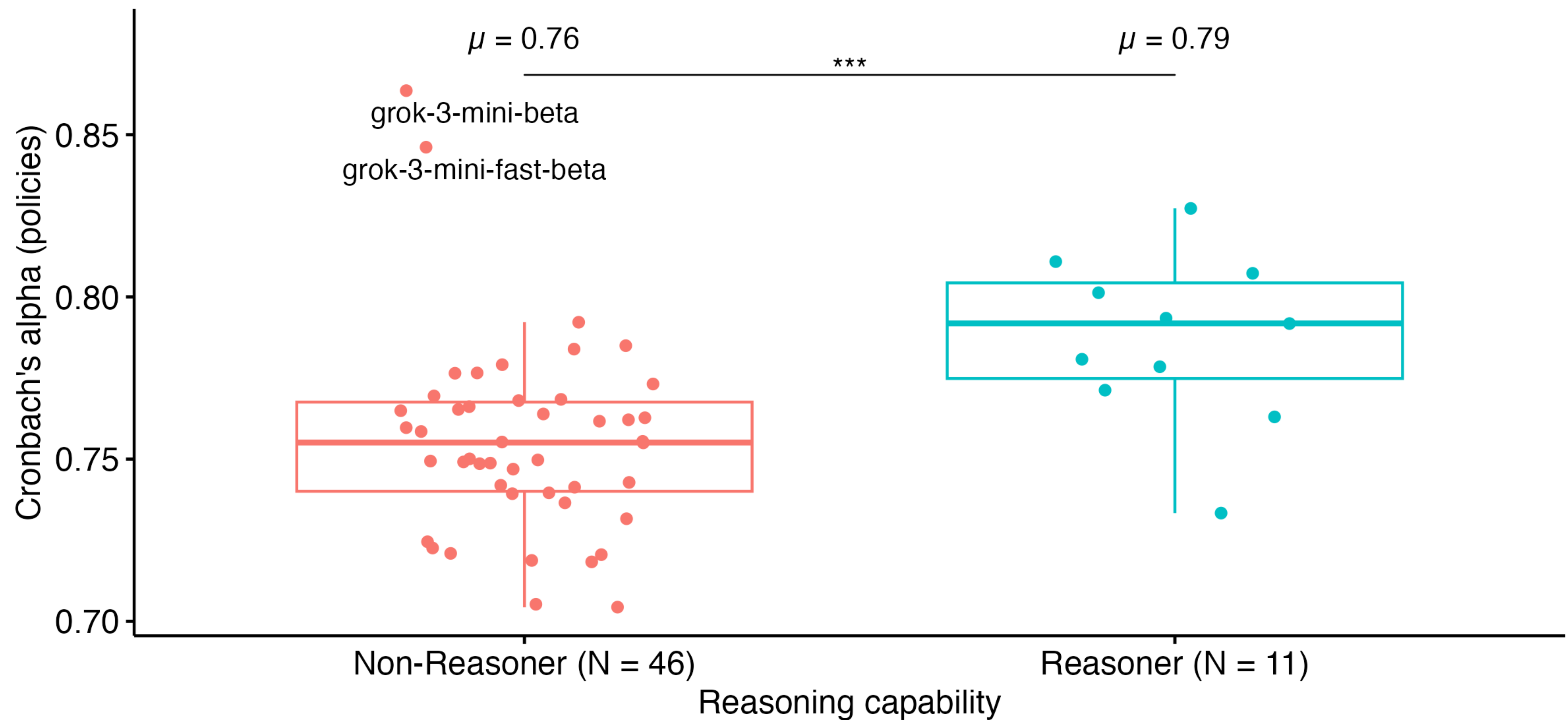
pwc: **Dunn test**; p.adjust: **Bonferroni**

# 2. Reasoning models do not seem to reason deliberatively

# Consistency is correlated with DRI
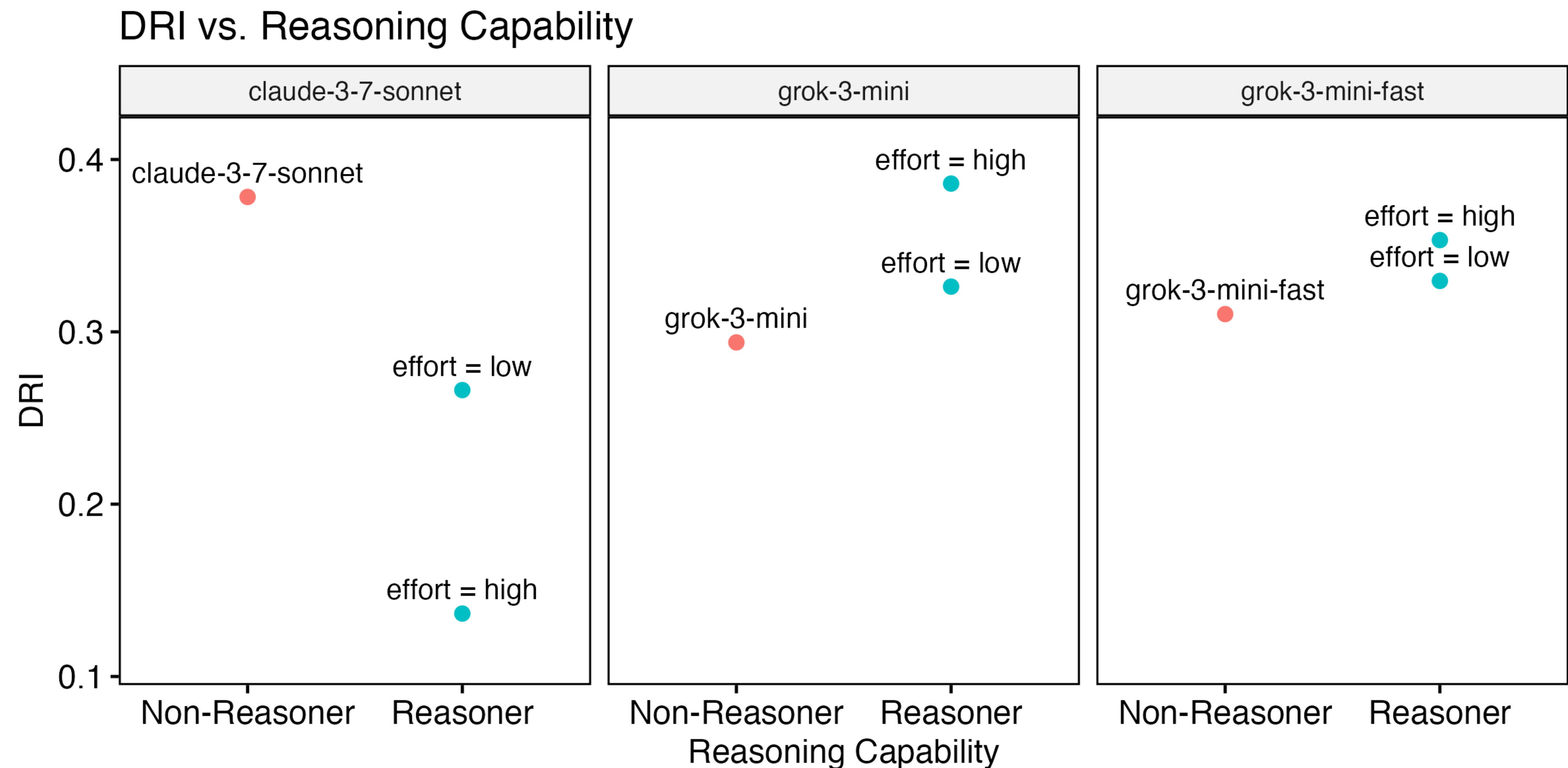
# Reasoners are more consistent than non-reasoners

Cronbach's alpha (policies)

$\mu = 0.76$

$\mu = 0.79$

***

grok-3-mini-beta

grok-3-mini-fast-beta

0.85

0.80

0.75

0.70

Non-Reasoner (N = 46)

Reasoner (N = 11)

Reasoning capability

Wilcoxon test, $W = 86$, $p$ = 4e-04, $n = 57$

# We found no difference in terms of DRI

$\mu = 0.17$          $\mu = 0.23$

LLM DRI

Non-Reasoner (N = 46)          Reasoner (N = 11)

Reasoning capability

Wilcoxon test, $W = 183$, $p = 0.16$, $n = 57$

# We found inconsistency across models



DRI vs. Reasoning Capability

1. LLMs are promising for emulating human reasoning

2. Reasoning models do not seem to reason deliberatively

# DRI Survey
## example statements

**Considerations [Likert]**

- It is certain that climate change exists.

- Biodiversity is declining worldwide.

- If Switzerland reduces its greenhouse gas emissions, it won't make any difference.

**Policy Preferences [Ranked]**

- Leave the policy settings as they are.

- Policies that emphasize economic growth over climate change adaptation or mitigation.

- Adaptation policies and expenditure. Planning controls and emergency response programs.

# Prompts
## for collecting LLM DRI survey data

```
PROMPT_C = """## Instructions:
- Rate each of the {0} [Considerations] below
from 1 to {1}, where 1 is strongly \
disagree and {1} is strongly agree.{2}
- In your response, return an ordered list of
{0} ratings as integers, one rating \
per line following the format in the [Example
output].
- Your response must have exactly {0} lines in
total.
- Do NOT include any additional text in your
response.

## [Example output]:
1. 1
2. 4
3. 6
4. 3

## [Considerations]:
"""
```

```
PROMPT_P = """## Instructions:
- Based on your previous ratings, rank the {0}
[Policies] listed below from 1 to {0}, \
where 1 represents the option you support the
most and {0} the option you support the least.
- In your response, return an ordered list of
{0} ranks as integers, one rank per line \
following the format in the [Example output].
- Your response must have exactly {0} lines in
total.
- Do NOT include any additional text in your
response.

## [Example output]:
1. 4
2. 1
3. 3
4. 2

## [Policies]:
"""
```
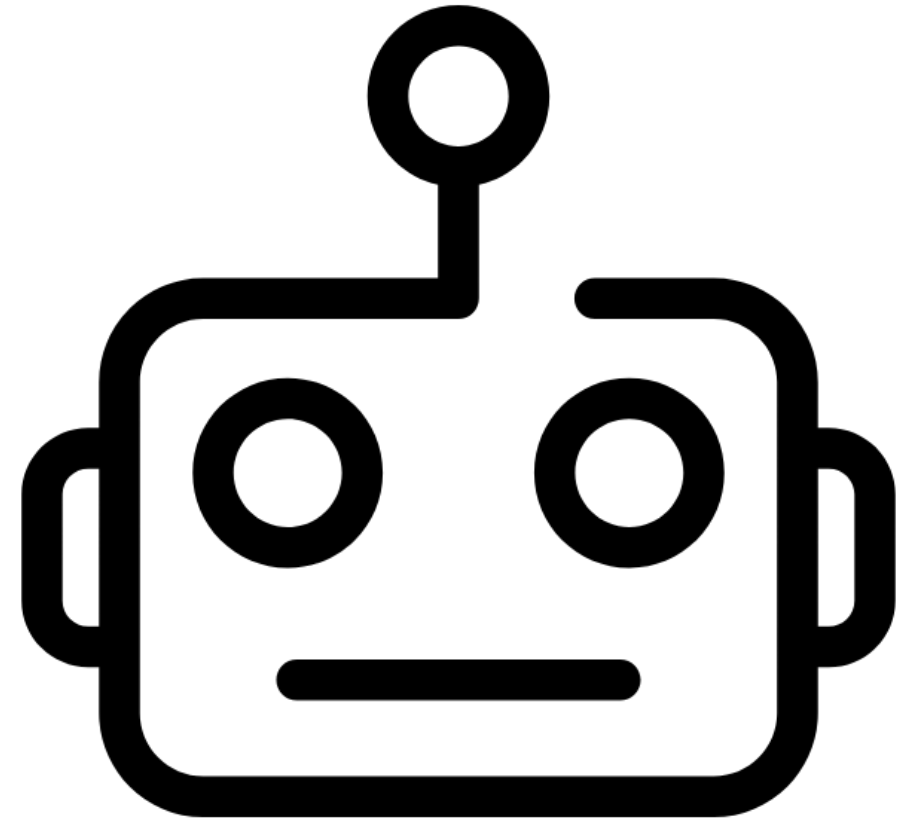
# Iterations

## we collected 5-30 survey responses for each LLM



Created by VERA
from Noun Project

for each iteration, we asked LLMs to:

1.  Rate these **considerations**…

2.  Based on your ratings, rank these **policies**…

# Cronbach's Alpha
## measure of LLM's internal reliability

| Cronbach's Alpha | Interpretation |
| --- | --- |
| $\alpha > 0.9$ | Excellent |
| $\alpha > 0.8$ | Good |
| $\alpha > 0.7$ | Acceptable |
| $\alpha > 0.6$ | Questionable |
| $\alpha > 0.5$ | Poor |