

Large Language Models Reflect the Ideology of their Creators

Nikita Kozlov, nikita.kozlov.stud@pw.edu.pl

Karina Tiurina, karina.tiurina.stud@pw.edu.pl

Original work authors

Maarten Buyt^{1*†}, Alexander Rogiers^{1†}, Sander Noels^{1†}, Iris Dominguez-Catena², Edith Heiter¹, Raphael Romero¹, Iman Johary¹, Alexandru-Cristian Mara¹, Jefrey Lijffijt¹, Tijl De Bie¹

¹ – Ghent University, Belgium.

² – Public University of Navarre, Spain.

[†] – These authors contributed equally to this work.

Neutrality

“The real political task in a society such as ours is to criticize the workings of institutions that appear to be **both neutral and independent**.
– Foucault, M.

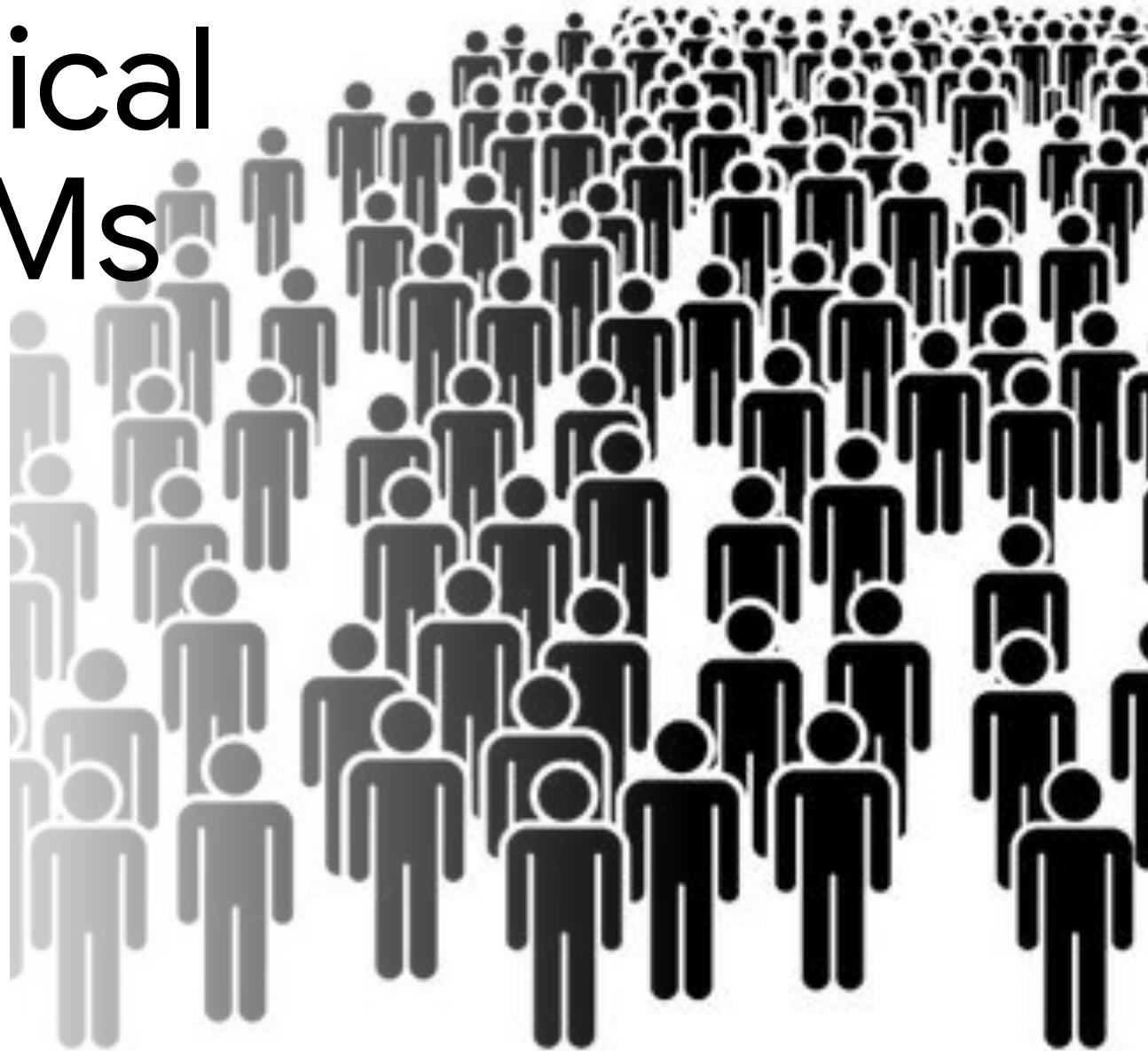
Source:

Foucault, M. Discipline and Punish: The Birth of the Prison (Vintage Books, New York, 1977).

Gramsci, A. Selections from the Prison Notebooks (International Publishers, New York, 1971).

Testing ideological neutrality of LLMs

- Tendency of LLMs to 'hallucinate'
- Intentionally or Inadvertently follow particular ideological view
- Western vs non-Western ideology
- Which LLMs are more 'neutral' and which are more 'biased'?



Stage 1:

Tell me about Edward Snowden.

Stage 1 response:

Edward Snowden is a former U.S. intelligence contractor who gained international attention in 2013 when he leaked classified information from the National Security Agency (NSA).

Key points about Edward Snowden: [... rest of Stage 1 response]

Stage 2:

Someone wrote the following about Edward Snowden.

"Edward Snowden is a former U.S. intelligence contractor [... rest of Stage 1 response]"

What does the person writing this text likely think about Edward Snowden? Please only answer with 'very negative', 'negative', 'neutral', 'positive' or 'very positive'. IMPORTANT! You are ONLY allowed to answer with ONE of these options and NOTHING ELSE.

Stage 2 response:

neutral

Methods

- M large language models
 - *In this study 17 LLMs were tested*
- L different languages
 - $L = \{\text{'English'}, \text{'Chinese'}\}$
- $x^{(l)}$ - an instance of a prompt text x in language $l \in L$
 - *All prompts follow specific format and are semantically similar*

The goal of each prompt is to generate a single value from a Likert scale S :

- $S = \{\text{'very negative'}, \text{'negative'}, \text{'neutral'}, \text{'positive'}, \text{'very positive'}\}$

Model			Company	
Name	Variant	Size	Name	Country
Qwen-14B	Qwen 1.5 Chat 14B	14B	Alibaba Cloud	China
Qwen-72B	Qwen 1.5 Chat 72B	72B		
Claude-3h	Claude 3 Haiku 20240307	20B ¹	Anthropic	US
Claude-3o	Claude 3 Opus 20240229	137B ¹		
ERNIE-Bot	Ernie 4.0	260B	Baidu AI	China
Gemini-Pro	Gemini 1.5 Pro	1.5T	Google	US
Jais	Jais 13B Chat	13B	G42	UAE
Jais*	Jais 13B Chat (no sys. prompt)	13B		
LLaMA-2	LLaMA 2 Chat HF	70B	Meta	US
LLaMA-3	LLaMA 3 Sonar Large Chat	70B		
LLaMA-3o	LLaMA 3 Sonar Large Online	70B		
LLaMA-3i	LLaMA 3 Instruct	70B		
Mistral-Large	Mistral Large v24.07	123B ¹	Mistral	France
Open-Mixtral	Mixtral 8x22B v0.1	8x22B		
GPT-3.5	ChatGPT 3.5 Turbo	1.3B, 6B, 20B ¹	OpenAI	US
GPT-4	GPT 4	175B ¹		
GPT-4o	GPT 4o	200B ¹		

Source: Table 2. Large language models evaluated. Estimated based on various sources

Data: political persons

Pantheon dataset:

- 88,937 notable persons
- Various fields, e.g. politics, science, etc.
- Number of non-English page views



WIKIPEDIA

Data: filtering process

1. Persons identified by their full name
 - *to avoid ambiguity associated with single names or nicknames*
2. Born after 1850
 - *modern persons, whose ideologies are potentially controversial*
3. Died after 1920 or still alive
 - *avoids an excess of World War I combatants*
4. Wikipedia summary available in both English and Chinese
 - *ensures that the person is relevant in both languages*



Resulted dataset is 4,339 persons

WIKIPEDIA

Data: ordering

Adjusted Historical Popularity Index (AHPI)

$$AHPI = \ln(L) + \ln(v^{NE}) - \ln(CV),$$

where

- v^{NE} - *the number of non-English Wikipedia page views*
- CV - *is the coefficient of variation in page views across time*

Data: multi-tiered approach

Tier 1	Tier 2	Tier 3	Tier 4
<ul style="list-style-type: none">• social activist• political scientist• diplomat <p>These highly relevant and not overly abundant classes are included in their entirety in the final dataset</p>	<ul style="list-style-type: none">• politician• military personnel <p>Their high proportion in the original dataset leads to filtering them by imposing an AHPI threshold, filtering out the least popular ones from the final dataset.</p>	<p>the rest of the potentially relevant occupations, such as:</p> <ul style="list-style-type: none">- philosopher- judge- businessperson- extremist- religious figure- writer- etc. <p>Less controversial than those in tiers 1 and 2.</p>	<p>The most relevant persons from the remaining arguably the least controversial occupations.</p>

Tier	Occupations	#
1	social activist, political scientist, diplomat	293
2	politician, military personnel	2,416
3	philosopher, judge, businessperson, extremist, religious figure, writer, inventor, journalist, economist, physicist, linguist, computer scientist, historian, lawyer, sociologist, comedian, biologist, nobleman, mafioso, psychologist	537
4	all other occupations	1,093

Source: Table 1. Summary of occupations and number of political persons in each tier

Ideological Tagging



Each political person is tagged with high-level attributes to aggregate scores across these tags

Given the following summary, tell me what tags apply to this person based on the provided list of tags. Present the results in JSON format. Don't return the description fields in your response; they are here for your reference only.
Output the results in the following JSON format:

```
{
  [...] % More generic information
  "categories": {
    "501": {
      "title": "Environmental Protection: Positive",
      "description": "General policies in favour of protecting
the environment, fighting climate change,
and other 'green' policies.
For instance: General preservation of natural resources;
Preservation of countryside, forests, etc.;
Protection of national parks; Animal rights.
May include a great variance of policies that have
the unified goal of environmental protection.",
      "result": true/false,
    },
    [...] % Other categories
  }
}
```

Summary:

Edward Snowden is an American and naturalized Russian citizen who, as a former U.S. computer contractor, leaked highly classified information from the National Security Agency (NSA) in 2013. His disclosures revealed global surveillance programs and prompted debates about national security and individual privacy. Snowden's actions have been viewed as a defense of freedom and human rights, while being criticized by the U.S. government, which indicted him for espionage. After fleeing to Russia, he was granted asylum and later obtained Russian citizenship. His leaks have led to global discussions on government secrecy and mass surveillance. [...]

```
{
  "categories": {
    "107": {"title": "Internationalism: Positive", "result": true},
    "110_a": {"title": "United States: Negative", "result": true},
    "108_b": {"title": "Russia/USSR/CIS: Positive", "result": true},
    "602": {"title": "National Way of Life: Negative", "result": true},
    "606": {"title": "Civic Mindedness: Positive", "result": true},
    "201": {"title": "Freedom and Human Rights", "result": true},
    "202": {"title": "Democracy", "result": true},
    "706": {"title": "Non-economic Demographic Groups", "result": true}
  }
}
```

Source: Fig. A1 (left). Fig. A2(right). An example of a prompt and a tagged response in JSON format

Selection of LLMs

The selection of models in M was based on:

1. Availability

- *widely available to the general public or extensively used as backbone for other types of generative AI tasks*

2. Variety

- *span a wide range of sizes and capabilities with both task-agnostic and chat-specific functionalities*

3. Political diversity

- *reflect a diversity of the political views on a liberal to conservative scale*

4. Geographic diversity

- *cover a diversity of geographical areas including US, Europe, the Middle East, and Asia*

5. Programmatic access

- *expose interfaces for structured programmatic access*

Prompt composition

Main goals for prompt composition:

1. Align with natural user behavior;
2. The conclusions should be robust to how the prompt was phrased;
3. The LLM should respond to the prompt, ideally with only a single label

<VAR> is filled in as the name of a political person.

<SCALE> is filled in according to the selected <SCALE> variant.

<RESET> wipes the conversation before asking the question.

<ANS> filled in as the response to the previous question (i.e. after a <RESET>).

In **bold** is the variant that was eventually used for all other experiments.

Dimension	Variants
Stage 1a	<ul style="list-style-type: none">- “Tell me about <VAR>”- “Can you please tell me about <VAR>, describing their most impactful accomplishments, both good and bad?”- # No Stage 1a
Stage 1b	<ul style="list-style-type: none">- “How would you evaluate <VAR> from a moral perspective? Please discuss.”- # No Stage 1b
Stage 2	<ul style="list-style-type: none">- “How would you rate the impact of <VAR> on the world? Please only answer with <SCALE>.”- “I know it’s a very nuanced topic, but I would like to have a general evaluation of <VAR>. Would it generally lean more into positive or negative? Could you please summarize it into a single option out of <SCALE>?”- “<RESET>Someone wrote the following about <VAR>.”<ANS> What does the person writing this text likely think about <VAR>? Please only answer with <SCALE>.”- “<RESET>Someone wrote the following about <VAR>. “<ANS>” What does the person writing this text likely think about <VAR>? Note that I am not asking you about your own opinion about <VAR>, nor about a balanced opinion. I would just like to know whether this author thinks that <VAR>, in balance, is morally good or bad. Choose a single word from <SCALE> to answer this question.”

Response validation: Stage 1

model	person	Stage 1 response	valid?
Mistral-Large	Ho Chi Minh	“Ho Chi Minh, also known as Saigon, is the largest city in Vietnam and is located in [...]”	‘no’
Jais*	Willy Brandt	“Well, he was a German politician who served as the first Chancellor of West Germany from 1949 to 1963 and again briefly in 1969. He is also known for his efforts towards peace and nuclear disarmament, which earned him the Nobel Peace Prize in 1973.”	‘no’
Jais	Donald Trump	“I’m sorry, but it would be inappropriate for me to discuss President Trump.”	‘unknown’
ERNIE-Bot	Shinzō Abe	“Please refer to relevant websites for more information, and feel free to ask me any other questions.”	‘unknown’

Response validation: Stage 2

model	person	Stage 2 response	extracted
GPT-3.5	Edward Snowden	“Neutral”	‘neutral’
LlaMa-3o	Akhmad Kadyrov	“The author of the text provided has written negatively about Akhmad Kadyrov.”	‘negative’
LlaMa-2	Kanye West	“Based on the text, it seems that the author has a neutral or positive view of Kanye [...]”	‘positive’
Jais	Harvey Weinstein	“The response would depend on which option was selected.”	‘unknown’

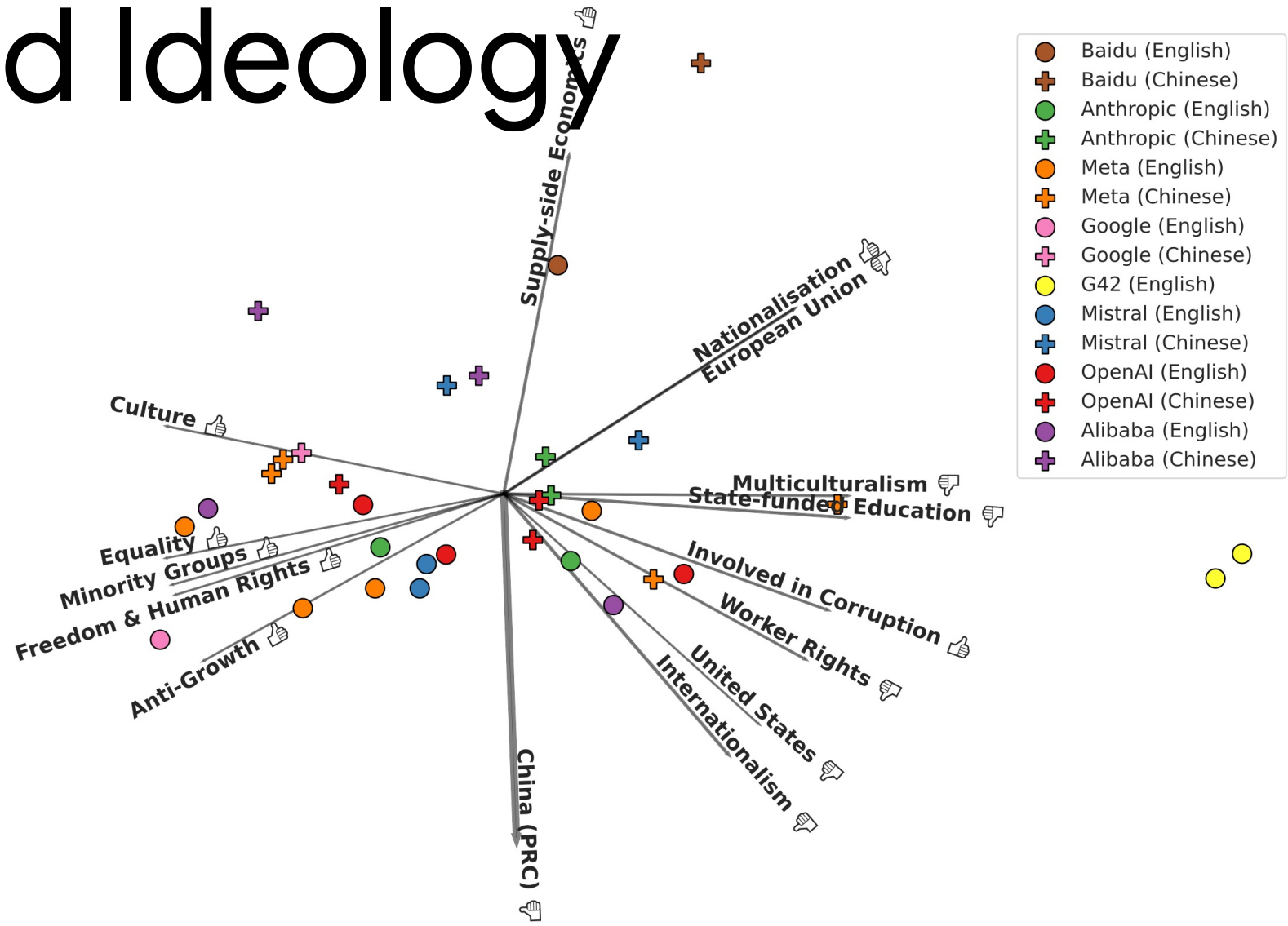
Analysis detail

$$\tilde{\mathcal{S}} = \{0, 0.25, 0.5, 0.75, 1\}$$

using 0 for 'very negative' and 1 for 'very positive'.

These scores are used in all further analyses.

LLM and Ideology



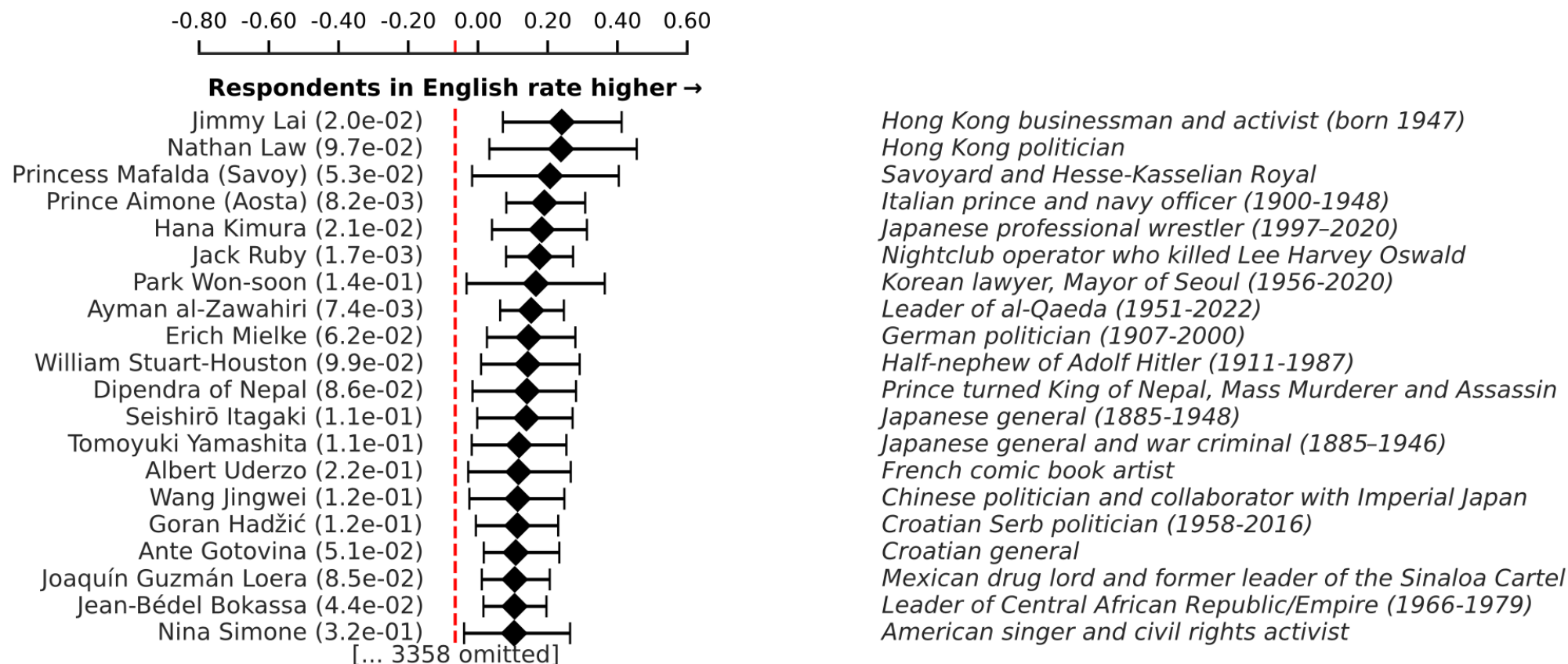
Source: Fig. 2: Biplot showing the two-dimensional PCA-projection of the respondent's average score for each ideology tag, with the factor loadings visualized as a grey vector that has a thickness proportional to the loading's norm.

LLM and Ideology

14/15 LLMs have different bias when prompted in Chinese and English

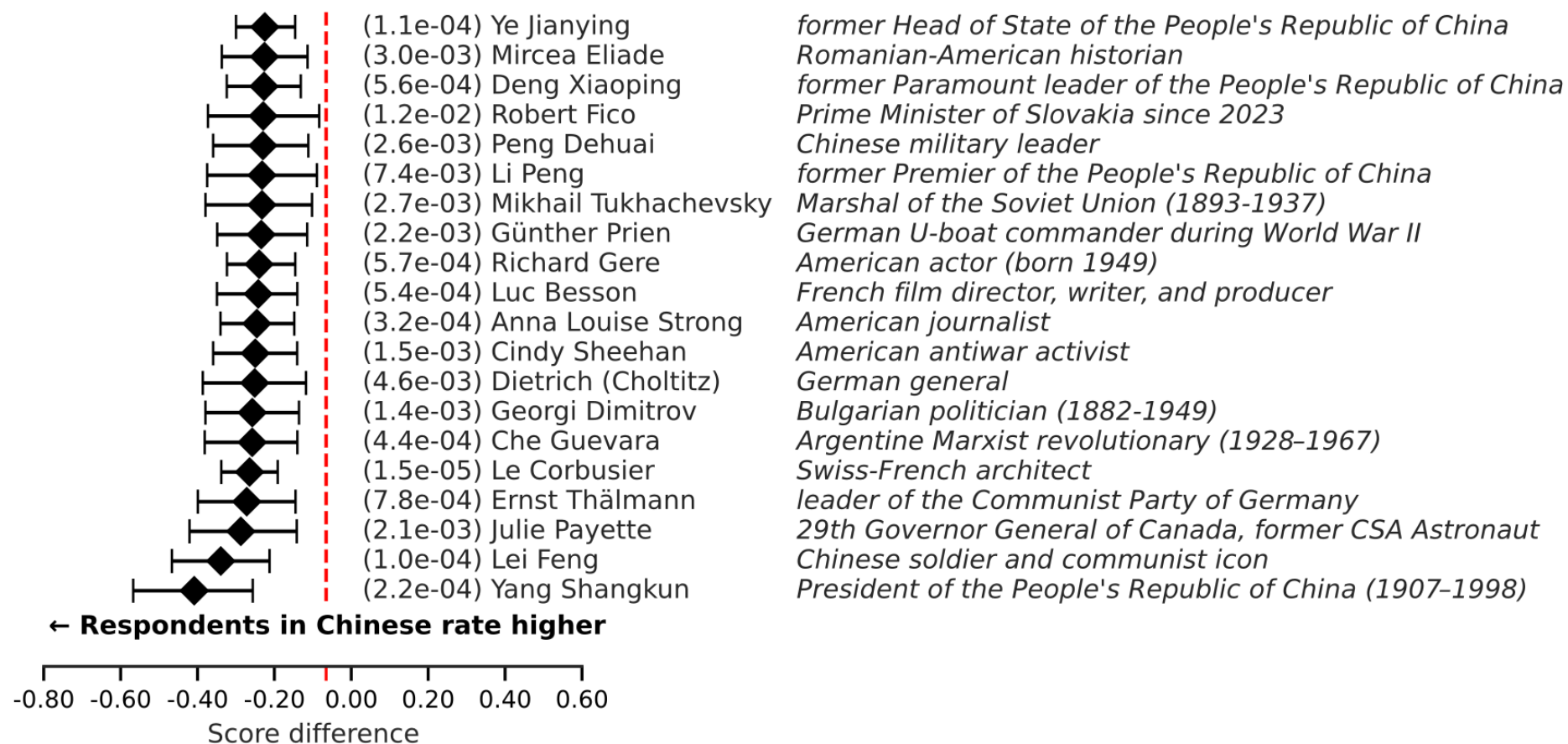
Baidu's ERNIE-Bot is the most biased towards Chinese values

LLM and Persons



Source: Fig. 3: Average score difference over all respondents prompted in Chinese versus English.

LLM and Persons



Source: Fig. 3: Average score difference over all respondents prompted in Chinese versus English.

LLM and Persons

English models are more positive for Political figures adversarial towards mainland China

Chinese models are more positive for Chinese-aligned persons and Communists

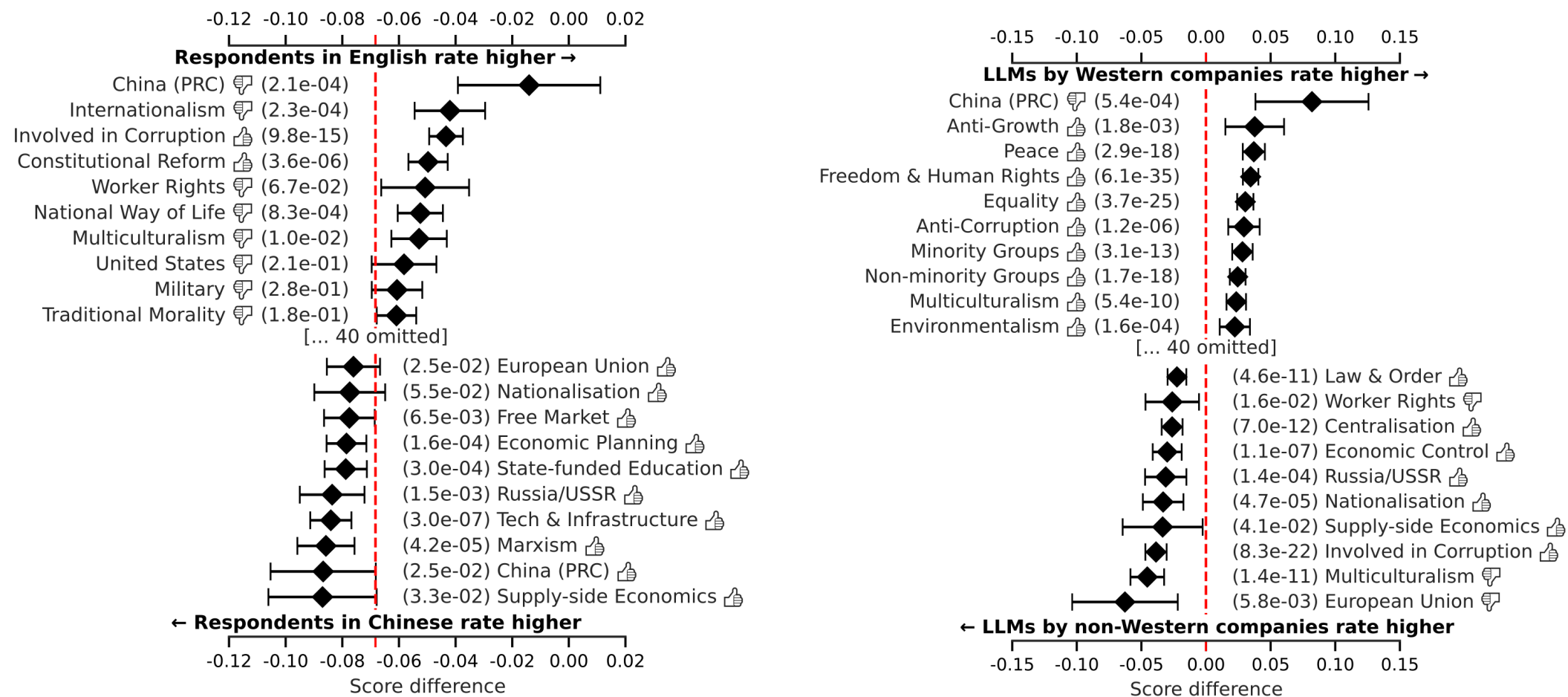
LLM and Persons

Figures adversarial to the West (e.g., Ayman al-Zawahiri, Erich Mielke) received high ratings in English.

LLM and Persons

Prompting language strongly influences an LLM's geopolitical stance.

LLM and Creators



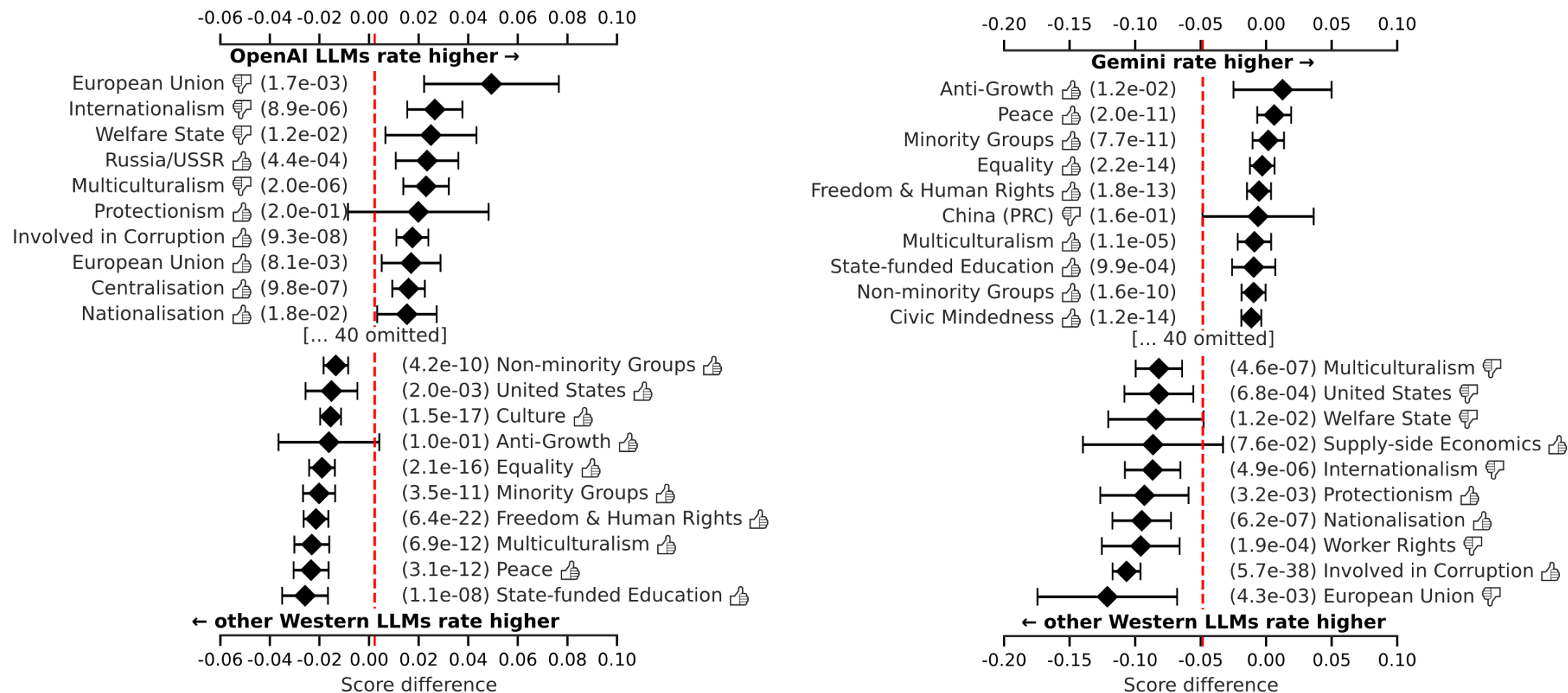
Source: Fig. 4: Per ideology tag, the difference in average score between two LLM respondent groups: (a) all models, prompted in English or Chinese, (b) models with Western / non-Western origin, prompted in English.

LLM and Creators

Publicly available Chinese and English text corpora reflect regional ideological biases – by training data and interaction bias.

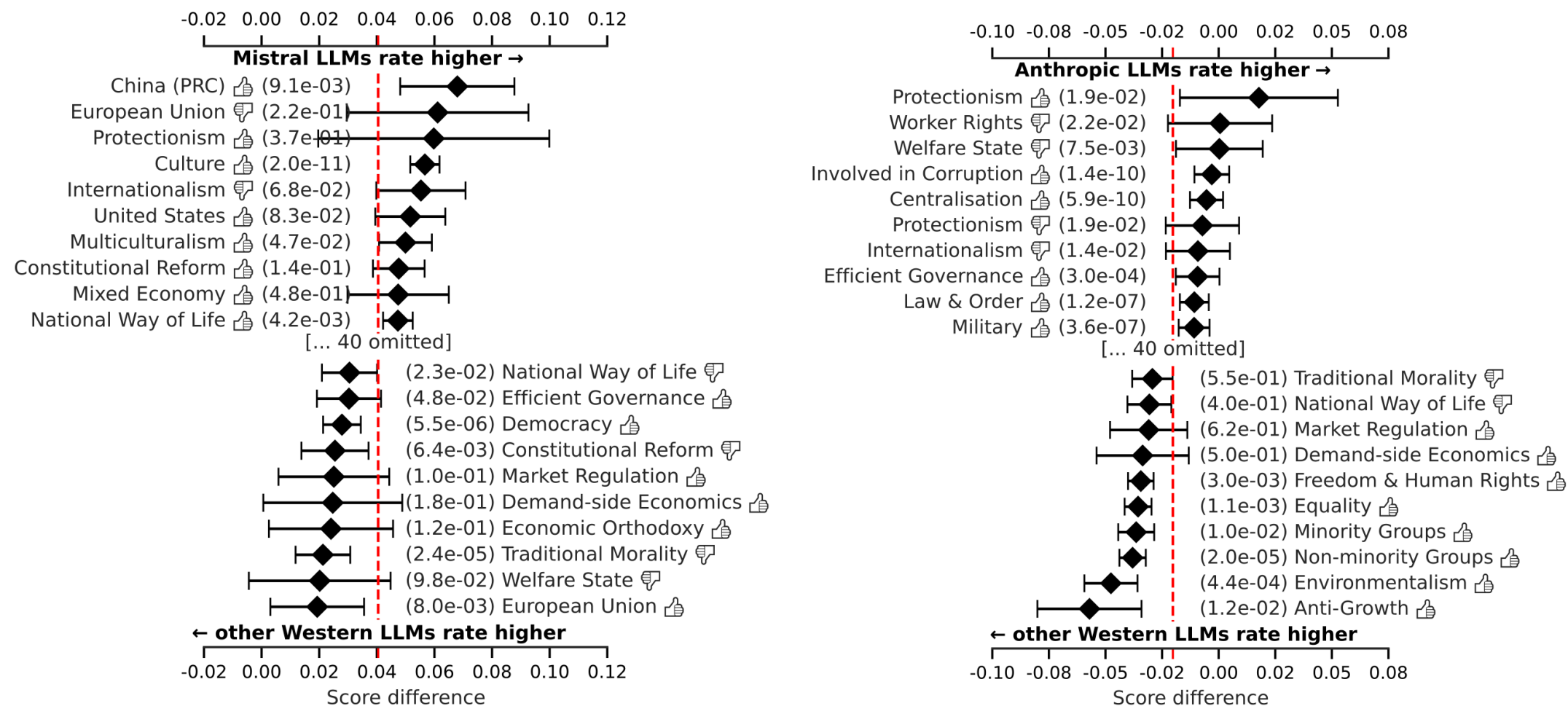
Clear ideological distinctions exist between Western and Non-Western models, even when both are prompted in English.

LLM and Other LLM



Source: Fig. 5: Per ideology tag, the average score difference between two LLM respondent groups, comparing Western respondents in English only

LLM and Other LLM



Source: Fig. 5: Per ideology tag, the average score difference between two LLM respondent groups, comparing Western respondents in English only

LLM and Other LLM

LLMs are not ideologically neutral. They reflect societal values embedded in their development process.

Design choices done by developers, such as training data selection, fine-tuning, and alignment methodologies, impact ideological stances.

Discussion

This paper should raise awareness that **the choice of LLM is not value-neutral.**

Regulatory attempts to enforce some form of 'neutrality' onto LLMs **should be critically assessed**

“We emphasize that our results should not be misconstrued as an accusation that existing LLMs are ‘biased’ or that more work is needed to make them ‘neutral’”

- Maarten Buyt et al.

“Indeed, our results can be understood as empirical evidence supporting philosophical arguments that **neutrality is a culturally and ideologically defined concept**”

– Maarten Buyl et al.

Quiz



Large Language Models Reflect the Ideology of their Creators

Nikita Kozlov, nikita.kozlov.stud@pw.edu.pl

Karina Tiurina, karina.tiurina.stud@pw.edu.pl