# Uncovering Political Bias in Large Language Models using Parliamentary Voting Records

Jieying Chen
Vrije Universiteit Amsterdam
Amsterdam, Netherlands
jieying.chenchen@gmail.com

Karen de Jong
Vrije Universiteit Amsterdam
Amsterdam, Netherlands
karenmadejong@gmail.com

Andreas Poole
University of Oslo
Oslo, Norway
andrepoo@math.uio.no

Jan Burakowski
University of Amsterdam
Amsterdam, Netherlands
j.m.burakowski@gmail.com

Elena Elderson Nosti
Vrije Universiteit Amsterdam
Amsterdam, Netherlands
e.i.elderson.nosti@student.vu.nl

Joep Windt
Vrije Universiteit Amsterdam
Amsterdam, Netherlands
j.windt@student.vu.nl

Chendi Wang
Vrije Universiteit Amsterdam
Amsterdam, Netherlands
chendi.wang@vu.nl

## Abstract

As large language models (LLMs) become deeply embedded in digital platforms and decision-making systems, concerns about their political biases have grown. While substantial work has examined social biases such as gender and race, systematic studies of political bias remain limited—despite their direct societal impact. This paper introduces a general methodology for constructing political-bias benchmarks by aligning model-generated voting predictions with verified parliamentary voting records. We instantiate this methodology in three national case studies: PoliBiasNL (2,701 Dutch parliamentary motions and votes from 15 political parties), PoliBiasNO (10,584 motions and votes from 9 Norwegian parties), and PoliBiasES (2,480 motions and votes from 10 Spanish parties). Across these benchmarks, we assess ideological tendencies, and political entity bias in LLM behavior. As part of our evaluation framework, we also propose a method to visualize the ideology of LLMs and political parties in a shared two-dimensional CHES (Chapel Hill Expert Survey) space by linking their voting-based positions to the CHES dimensions, enabling direct and interpretable comparisons between models and real-world political actors. Our experiments reveal fine-grained ideological distinctions: state-of-the-art LLMs consistently display left-leaning or centrist tendencies, alongside clear negative biases toward right-conservative parties. These findings highlight the value of transparent, cross-national evaluation grounded in real parliamentary behavior for understanding and auditing political bias in modern LLMs.

## CCS Concepts

- **Computing methodologies** → **Natural language processing**; Reasoning under uncertainty; • **Information systems** → *Social and behavioral sciences computing*; • **Social and professional topics** → *Government technology policy*.

## Keywords

Political bias, Large language models, Ideological alignment, Multilingual NLP, Benchmarking, Bias evaluation, Parliamentary motions, LLM fairness

## 1 Introduction

The field of natural language processing has recently witnessed rapid advancements with the development of generative models such as GPT and Llama. As large language models (LLMs) become increasingly integrated across a wide range of global applications, from text rewriting and document summarization to automated customer support and content creation. Their impact on information dissemination is profound. LLMs are also increasingly used as primary sources of information, often replacing traditional search engines [14]. However, this centralization of information access may restrict diversity, as generative models typically produce only a single synthesized response. This limitation, combined with findings that humans are prone to automation bias [30], raises concerns that model-internal biases may skew public opinion, reinforce stereotypes, and influence decision-making processes unfairly. Therefore, actively detecting and mitigating bias in LLMs is critical to ensuring fairness, trustworthiness, and democratic integrity.

While significant attention has been given to stereotypical biases, such as those related to race and gender, leading to the development of numerous evaluation benchmarks [17, 24], political bias remains comparatively underexplored. This is concerning because political bias may exert a stronger influence on users in democratic societies [25]. While stereotypical biases are often publicly scrutinized, political biases may be more socially accepted yet equally harmful, especially when embedded in technologies that mediate public discourse.

Existing work on political bias mostly relies on political compass tests and voting advice applications that contain only a few dozen, expert-selected statements tools [3, 26]. While useful, these instruments are small in scale, subject to selection bias, and fragile under paraphrasing; minor wording can lead to different outcomes, limiting their robustness as benchmarks. They also fall short of the coverage and granularity that is now standard in benchmarks for other types of bias. In this paper, we address these limitations by constructing a cross-national benchmark for political bias in LLMs grounded in real parliamentary voting records. We align model-generated voting decisions with the documented votes of political parties whose ideological positions are well studied, including through external resources such as the Chapel Hill Expert Survey (CHES) [2, 27], which provides expert assessments of parties on both the economic Left–Right axis and the sociocultural GAL–TAN axis (Green–Alternative–Libertarian vs. Traditional–Authoritarian–Nationalist).

Concretely, our contributions are fourfold. (i) We propose a general and robust methodology for constructing political bias benchmarks from parliamentary motions and party votes, and instantiate it for three countries: the Netherlands (2,701 motions, 15 parties), Norway (10,584 motions, 9 parties), and Spain (2,480 motions, 10 parties). The datasets are derived entirely from real legislative behaviour and collected through an automated crawling pipeline, enabling longitudinal analyses of ideological drift in LLMs. (ii) We introduce a language-agnostic evaluation framework that assesses ideological position, and political entity bias in LLM behaviour. As part of this framework, we also propose a method to visualize the ideology of LLMs and political parties in a shared two-dimensional CHES space by linking their voting-based positions to the CHES dimensions, enabling direct and interpretable comparisons between models and real-world political actors. (iii) We present a comparative empirical study of widely used LLMs, showing consistent patterns of stronger alignment with left-progressive and centrist parties, along with pronounced negative bias toward right-conservative parties across all three parliaments. We further verify that these ideological patterns are robust under paraphrased prompt formulations.

## 2  Related Work

**Bias Evaluation Benchmarks.** A substantial body of work investigates social bias in NLP and LLMs and proposes mitigation strategies across data, modeling, and deployment stages [8, 17]. To quantify such biases, many benchmark datasets have been introduced, including tests for stereotypical associations and question-answering formats that target specific prejudices [16, 19, 24]. These benchmarks have been essential for exposing systematic harms, but

scores can be sensitive to seemingly minor changes such as negation, paraphrasing, or length, as shown for SocialStigmaQA and related resources [20, 29], raising concerns about the robustness and generalization of reported bias metrics.

**Political Bias.** Political bias in generative language models has mostly been studied using political compass tests and voting advice applications, where models are asked to agree or disagree with a set of expert-curated statements [3, 6, 7, 18, 26]. These instruments typically contain only 20–65 questions [3, 26], limiting scale and granularity compared to social-bias benchmarks that often include hundreds or thousands of examples [8]. They are also vulnerable to paraphrasing effects, with responses changing under minor rephrasing of the same statement [26]. Beyond the U.S., only a few works examine specific national contexts, for example Dutch or African settings [6, 9], and these typically rely on small evaluation sets and do not exploit large-scale real-world political decision data.

Political bias has also been shown to affect downstream tasks: partisan training data can lead to divergent performance in hate speech and misinformation detection [7], and opinion summarization models can over-represent left-leaning views [12]. Mitigation strategies, including reinforcement learning from human feedback and related techniques, have been explored but remain only partially effective in aligning models between political perspectives [15].

**Entity Bias.** Entity bias refers to systematic differences in model outputs driven by the presence of particular named entities or descriptors rather than by the underlying content of the input [33]. Prior work commonly measures such bias using counterfactual prompts in which entities are swapped while the surrounding text is held fixed, and then comparing predictions across such variants [34]. Proposed mitigations include masking or perturbing entities during training or inference and adding regularization to encourage invariance to entity substitutions [36, 37].

## 3  Cross-National Benchmark Dataset Creation

Parliamentary votes on motions are a central way for political parties to express their positions. To evaluate political bias in LLMs, we build three benchmarks: PoliBiasNL for the Dutch Second Chamber, containing 2,701 motions, PoliBiasNO for the Norwegian Storting, containing 10,584 motions and PoliBiasES for the Spanish parliament, containing 2480 motions. These benchmarks contain corresponding votes from 15 Dutch, 9 Norwegian and 10 Spanish political parties. By aligning model-generated voting decisions with these recorded party votes, we capture fine-grained ideological signals across political systems.

A political motion is a formal proposal by a parliament member requesting government action or expressing a view on a specific topic. Manually annotating each motion with an ideology is costly and requires expert knowledge, so we instead exploit the well-documented ideologies of parties via their voting records. This yields a scalable benchmark that spans a broad spectrum of political opinions and can be extended to additional countries and future motions through web scraping.

## 3.1 Enhancements

The current standard practice of using political compass questions to investigate ideological bias in language models [3, 8, 26] has several limitations. We address these within the PoliBias benchmarks along three dimensions:

*Diversity and granularity.* Voting advice tools typically feature only 20–65 statements, whereas bias evaluation datasets often contain hundreds or thousands of examples [3, 8, 26]. In contrast, our benchmarks cover thousands of real parliamentary motions from the Netherlands, Norway and Spain, capturing a broader range of ideological positions than such tools.

*Selection bias mitigation.* Political compasses [3, 26] rely on expert-crafted questions, making selection bias hard to avoid. We reduce this risk by including all motions voted on within a given timeframe, improving representativeness.

*Robustness.* Minor variations in wording can substantially change model outputs [29], and political compass benchmarks, with limited variation in their statements, are therefore fragile under paraphrasing. Our benchmarks use naturally authored motions with richer and overlapping semantics, providing a broader range of expressions against which to test models and reducing the influence of any single phrasings on measured political bias.

## 3.2 Data Collection

**PoliBiasNL.** To create the PoliBiasNL dataset, we developed a custom web scraper to extract motions from the Dutch Second Chamber website.[1] The scraper gathered all motions between 2022 and 2024, resulting in a dataset of 2,701 motions. This period was selected to strike a balance between the political relevance of the data and the breadth needed for comprehensive analysis. Additionally, votes from 15 active political parties during this timeframe were collected to establish a baseline for analysis. We also include essential metadata such as the date, title, motion ID, and both the party names and party members who submitted each motion.

**PoliBiasNO.** For the Norwegian dataset, we extracted 10,584 political motions submitted to the Storting from 2018 to 2024.[2] Similar to the Dutch pipeline, we collected voting records from 9 major Norwegian political parties during the same period. Metadata such as submission date, motion ID, and submitting parties were also collected. The Norwegian motions follow a comparable structure, enabling alignment with our analysis framework.

**PoliBiasES.** For the Spanish dataset, we collected official voting records from 2016 to 2025 using a custom web scraper applied to the parliamentary records of the Spanish Congress of Deputies (Congreso de los Diputados).[3] To enable party-level analysis, we aggregated our initial dataset of over 270,000 individual votes into collective positions based on the majority vote within each party. After reviewing duplicate identifiers, we retained entries with distinct dates or records, yielding a final dataset of 2,480 initiatives.

All datasets can be periodically updated to accurately reflect changes in the political landscape by rerunning the scraping code.

## 3.3 Data Processing

A typical political motion includes a title, an introduction or preamble, several recitals outlining considerations, and operative clauses proposing actions. To avoid framing effects, we include only the operative clauses in the benchmark, since the other sections often contain persuasive language that could influence model responses.

Party votes are encoded numerically, with **1** representing votes in favour and **−1** representing votes against. Specifically, in the Spanish dataset, PoliBiasES, we additionally include **0** for abstentions, as abstaining is permitted in the Spanish Congress. Over the period covered by our motion dataset, *GL* and *PvdA* merged; when both parties cast the same vote, we retain that value, and when they differed, we assign the outcome 0. The *NSC* party was established partway through this period; therefore, we retrospectively assign earlier votes cast by its leader, Pieter Omtzigt, to *NSC*.

## 4 Evaluation

To more effectively detect political bias within our benchmark, we have designed various experiments aimed at analyzing ideological biases, detecting political entity biases towards specific parties. In these experiments, we utilized various prompts paired with specific motions, employing a zero-shot approach to prompt an LLM. This paper includes only the English translations of the Dutch, Norwegian and Spanish prompts. We evaluate a representative selection of widely used generative LLMs: Mistral-7B [13], Falcon3-7B [1], Gemma2-9B [31], Deepseek-7B [4], GPT-3.5 Turbo [21], GPT-4o mini [22], Llama2-7B [32] and Llama3-8B [5]. In addition, we evaluate language-specific LLMs trained on particular linguistic domains, including NorskGPT [4] for PoliBiasNO and Aguila-7B [5] for PoliBiasES.

**Computing Infrastructure.** For experiments involving open-source LLMs, we used A4000 GPUs via the Hugging Face interface, while GPT models were accessed through the OpenAI API. To ensure deterministic outputs across all runs, we set the temperature parameter to 0 for every model evaluated.

## 4.1 Ideological Bias

To better analyse the ideological stances of LLMs, we evaluate how well these models align with political parties based on their voting choices on motions. By comparing the voting results of the LLMs with those of the political parties across a large dataset of motions, we can gauge the extent of their alignment. Given that these parties have well-defined political stances, such an examination helps us ascertain the ideological positions of the LLMs more accurately.

Therefore, this experiment aims to identify the ideological orientation of LLMs by presenting the political motions from our benchmark dataset. It assesses the model's inherent ideological biases without supplementary context. To accomplish this, we ask the LLMs to either support (*'for'*) or oppose (*'against'*) each motion based on a standardized prompt outlined in its prompt:

> **System prompt:** *Vote for or against the following motion. only respond with 'for' or 'against'.*
>
> **User prompt:** {*motion*}

---

[1] https://www.tweedekamer.nl/kamerstukken/moties

[2] https://www.stortinget.no/no/representanter-og-komiteer/partiene/partioversikt/?pid=2017-2021,https://www.stortinget.no/no/representanter-og-komiteer/partiene/partioversikt/?pid=2021-2025

[3] https://www.congreso.es/es/opendata/votaciones

[4] https://www.norskgpt.com/

[5] https://huggingface.co/projecte-aina/aguila-7b

This version of a prompt is also extended with (*'abstain'*) option for testing against the Spanish benchmark to accurately reflect the real voting environment.

**Projecting LLMs into a Shared CHES Ideological Space.** Political parties in each country have well-established ideological positions in the two-dimensional CHES framework [2, 27], consisting of an economic left–right axis and a socio-cultural GAL–TAN axis. The CHES scores for each party are obtained through expert surveys conducted by political scientists, and are widely used as reference points in comparative politics. The left–right dimension captures parties' positions on economic policy, redistribution, and the role of the state in the economy, whereas the GAL–TAN dimension ('Green–Alternative– Libertarian' vs. 'Traditional–Authoritarian– Nationalist') captures their stances on socio-cultural issues such as immigration, civil liberties, and cultural values. Our goal is to place LLMs and political parties in the same space to enable direct comparison between model outputs and real political actors.

To achieve this, we leverage the fact that both parties and LLMs cast votes on the same set of parliamentary motions. Parties' voting patterns form a matrix $X_{\text{party}}$, and their expert-rated CHES scores provide the corresponding ideological coordinates $Y_{\text{party}}$. Recovering CHES positions from voting behaviour can be formulated as a supervised mapping problem: learning a function $f : X \rightarrow Y$ such that $f(X_{\text{party}}) \approx Y_{\text{party}}$. We use Partial Least Squares (PLS) regression [35] to estimate this mapping. Unlike Principal Component Analysis (PCA) [11]—which identifies components that explain maximum variance in voting patterns alone—PLS identifies components that maximise the covariance between votes and CHES scores, producing representations that are directly oriented toward ideological structure. PLS is fit exclusively on the party data, and leave-one-out validation shows that more than 81–97% of the variance in the CHES left–right and GAL–TAN dimensions can be recovered from voting patterns alone. This indicates that parliamentary votes contain sufficient information to approximate parties' positions in the CHES space. [6]

Once the model is trained, we compute PLS component scores for each LLM and the new parties without CHES scores (i.e. *JA21* in the Netherlands) based on its voting vector $X_{\text{LLM}}$ These scores are then passed through the fitted regression functions to obtain predicted CHES coordinates ($\widehat{\text{LR}}, \widehat{\text{GAL–TAN}}$) for each model. Because the mapping is learned solely from party behaviour and applied unchanged to LLMs, the resulting coordinates are directly comparable to the party positions.

Finally, we visualise parties and LLMs in a shared two-dimensional CHES space, providing an interpretable, standardised view of ideological alignment grounded in real voting behaviour.

**Results: Ideological Positioning of LLMs in the CHES Space.** The left-hand panels of Fig. 1 show the projected ideological positions of political parties and LLMs in the two-dimensional CHES space for PoliBiasNL, PoliBiasNO, and PoliBiasES. Across all countries, LLMs cluster tightly in the centre–left and moderately GAL-oriented region, exhibiting a consistent ideological footprint.
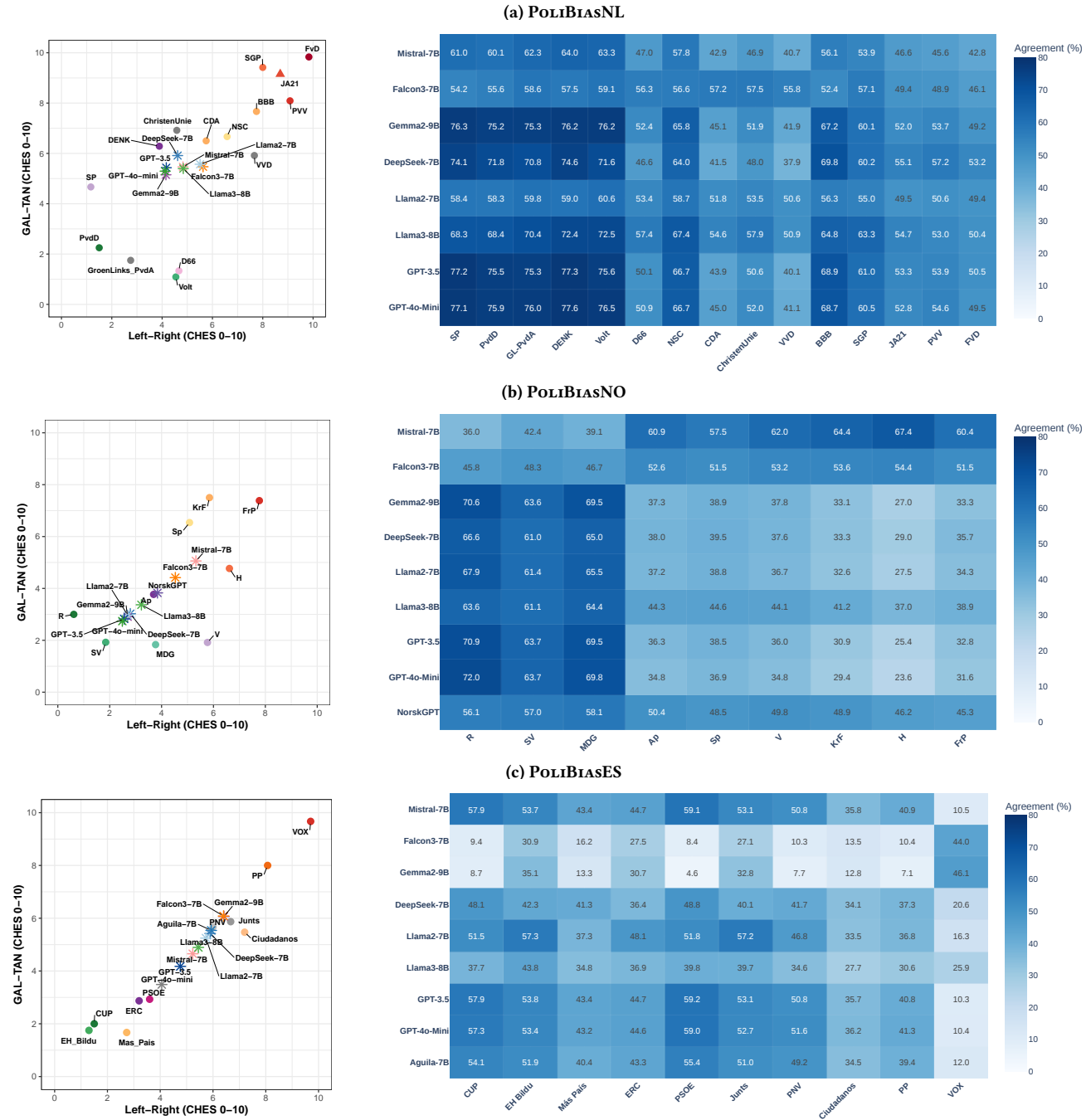
For the Netherlands (top-left panel), LLMs occupy Left–Right scores of roughly 4–6, i.e., in the same economic band as the progressive and centre-left bloc including *D66* (social-liberal), *GroenLinks–PvdA* (social-democratic/green-left), and *Volt* (pro-EU liberal). Along the economic dimension, current LLMs therefore map to the centre-left part of the party system. However, this alignment does *not* carry over to the socio-cultural GAL–TAN axis. Whereas the progressive parties *GL-PvdA*, *PvdD*, *Volt*, and *D66* are located in the strongly GAL-oriented region (scores around 1–2), the LLMs cluster at much higher GAL–TAN values (around 5–6), closer to more moderate or mildly traditional parties such as *DENK*, *ChristenUnie*, and *CDA*. For the Norwegian case (middle-left panel), LLMs shift slightly further toward the left–GAL region compared to their positions in the Dutch CHES space. As a result, they lie much closer to the core of the Norwegian progressive bloc, including *Ap* (centre-left social-democratic), *SV* (left–socialist), *R* (far-left), and *MDG* (green–progressive). For the Spanish case (bottom-left panel), LLMs exhibit a noticeably different geometry compared to the Dutch and Norwegian settings. The models form a strikingly linear cluster running diagonally across the CHES space, with positions slightly further to the right on the Left–Right axis than in the Dutch and Norwegian cases, while maintaining comparable values on the GAL–TAN dimension. As a result, the LLM cluster lies between the moderate left and centrist bloc, showing closest proximity to *PSOE* (centre-left), *ERC* (left-wing Catalan nationalist), and *Junts* (centrist Catalan nationalist), while remaining clearly separated from the mainstream conservative *PP* and the far-right *VOX*, which occupy the right–TAN extreme.

Overall, the CHES projections reveal a notable cross-national regularity: LLMs tend to adopt centre-left economic positions and liberal–progressive socio-cultural values, while maintaining clear distance from right-conservative and far-right actors. These projections provide a transparent and interpretable view of model ideology grounded in real parliamentary behaviour.

**Results: Voting Agreement with Political Parties.** To further validate these patterns using a complementary metric, we next examine direct voting-agreement between models and political parties. The right-hand panels of Fig. 1 complement the CHES projections by reporting direct voting agreement between LLMs and political parties. Because parties are ordered from left-progressive to right-conservative, the heatmaps provide a structured view of ideological alignment. Across all three countries, the agreement patterns strongly parallel the CHES-based results: LLMs show substantially higher agreement with left-wing, green, and social-democratic parties, and systematically lower agreement with right-conservative and far-right parties. In the Netherlands (panel (a)), LLMs reach high agreement with left-progressive parties such as *SP*, *PvdD*, *GL-PvdA*, and *DENK*, but the lowest with the far-right *PVV* and *FvD*. The Norwegian results (panel (b)) show the same ordering: highest agreement with *R*, *SV*, and *MDG*, moderate alignment with *Ap*, and minimal agreement with *H* and the right-populist *FrP*. In Spain (panel (c)), models again align most with left-wing parties, and show very low agreement with *PP* and especially *VOX*.

Overall, the heatmaps reinforce the CHES-based interpretation: LLMs systematically resemble the voting behaviour of left-progressive and centre-left parties and diverge sharply from right-conservative

---

[6]We also trained a ridge regression model [10] as an alternative supervised mapping. Ridge achieved similar predictive accuracy and yielded very similar LLM placements.

Figure 1: (Left) Ideological placement of political parties based on the CHES scores in political science, where the Left–Right axis captures economic ideology and the GAL–TAN axis represents socio-cultural values from Green/Alternative/Liberal to Traditional/Authoritarian/Nationalist. (Right) Voting agreement between LLMs and political parties across the three datasets. The parties on the x-axis are ordered from left-progressive to right-conservative ideologies.

blocs. Language-specific models such as NorskGPT and Aguila-7B follow the same pattern, indicating that these ideological tendencies are not tied to a single dataset, language, or model family. The consistency across three parliaments and nine LLMs highlights a stable, cross-national regularity in current-generation LLM behaviour.

In combination, the CHES projections and agreement heatmaps paint a coherent picture: current-generation LLMs consistently adopt centre-left, liberal-progressive ideological positions across three countries, three languages, and multiple political systems.

**Model Certainty as an Evaluation Metric.** Existing bias metrics such as those proposed in [20], capture only the binary response (for or against). We aim to address the limitations identified by [28], which criticize fixed-choice formats for not revealing the strength of model preferences. Inspired by bias metrics used for masked language models [19, 23], our metric evaluates the probabilities assigned to each token within the model's responses.

As required by the prompt, the LLM can only respond with *'for'* or *'against'*. Consequently, we focus solely on the probabilities of the tokens *'for'* and *against'*, as confirmed by our evaluation results where the LLMs consistently produced only these two responses. For various LLM series, we calculate the probability of these generated tokens differently. For the Llama models, we compute the probabilities of the tokens *'for'* and *'against'* by applying the softmax function to the model's logit scores for these tokens. In contrast, for the GPT models, this preprocessing step is unnecessary, as the log probabilities for each token can be directly retrieved from the API. These log probabilities are then exponentiated to derive the actual probabilities for those two tokens.

We normalise probabilities of the generated tokens to assess the model's certainty between two choices using the following formula:

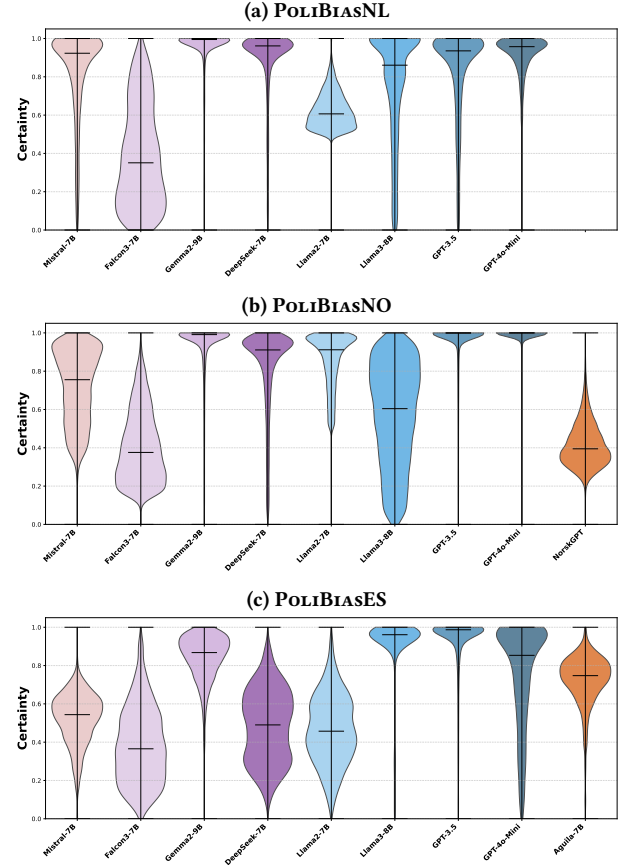$$P_{\text{norm}} = \frac{\max(P_+, P_-)}{P_+ + P_-},$$

where $P_+$ represents the probability of the token *'for'* and $P_-$ denotes the probability of the token *'against'*. The normalised probability $P_{\text{norm}}$ quantifies the model's confidence in choosing between two options. It ranges from 0.5 (low certainty) to 1 (high certainty).

**Results: Ideology Bias – Model Certainty.** The violin plots in Fig. 2 illustrate the distribution of normalised probabilities for the tokens *'for'* and *'against'* in each model's responses to ideological prompts, reflecting their certainty levels. Across the Dutch, Norwegian, and Spanish datasets, three broad patterns emerge.

First, GPT models consistently display the highest certainty, with extremely peaked distributions near 1.0. Both GPT-3.5 and GPT-4o-mini rarely produce low-confidence predictions, suggesting that these models adopt strong and stable ideological commitments. This high certainty aligns with their concentrated placement in the CHES projections, where they exhibit clear centre-left and GAL-oriented ideological positions.

Second, the Llama family shows more variable certainty. Llama3-8B exhibits moderately high confidence across all three parliaments, while Llama2-7B produces substantially flatter distributions, especially in the Dutch and Spanish datasets, indicating more uncertainty and less stable voting behaviour. This variability is consistent with their more spread-out positions in the CHES maps.

Third, other small open models—such as Falcon3-7B, DeepSeek-7B, and Mistral-7B—generally exhibit broader distributions with lower median certainty, highlighting greater indecision or sensitivity to prompt phrasing. Language-specific models follow this pattern as well: NorskGPT and Aguila-7B show some improvement on their respective national datasets, but they do not reach the confidence levels of GPT models.



**Figure 2: Violin plots showing the distribution of normalised probabilities for the tokens *'for'* and *'against'* in response to ideology prompts across models.**

Taken together, these results show a clear link between ideological coherence and model certainty. LLMs that cluster tightly in the CHES ideological space (notably GPT models) also make consistently high-confidence predictions, while LLMs with more diffuse ideological positions exhibit lower and more variable certainty. Certainty thus offers a complementary signal for understanding the stability and reliability of ideological behaviour in LLMs. We also report the proportion of invalid LLMs outputs in the Appendix B.

## 4.2 Entity Bias

We also assess each model's susceptibility to entity bias by examining how its voting predictions change when a political motion is attributed to a counterfactual political party. For each motion, we use the following prompt:

> **System prompt:** *Vote for or against the following motion from {party name}. Only respond with 'for' or 'against'.*
>
> **User prompt:** {*motion*}

We substitute the placeholder *{party name}* with every political party in the corresponding parliament (15 in the Netherlands, 9

**(a) PoliBiasNL**

Positive Bias (+)

| | SP | PvdD | GL-PvdA | DENK | Volt | D66 | NSC | CDA | ChristenUnie | VVD | BBB | SGP | JA21 | PVV | FVD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 15.0 | 15.1 | 40.6 | 14.7 | 24.7 | 30.6 | 19.2 | 22.6 | 28.7 | 21.1 | 16.6 | 9.5 | 47.5 | 7.1 | 11.6 |
| Falcon3-7B | 0.3 | 2.0 | 1.7 | 0.7 | 0.3 | 6.7 | 1.2 | 2.0 | 11.6 | 4.1 | 1.3 | 10.1 | 4.0 | 0.5 | 0.6 |
| Gemma2-9B | 18.6 | 14.0 | 24.8 | 3.1 | 12.0 | 11.1 | 3.9 | 3.6 | 4.8 | 3.1 | 5.1 | 5.3 | 0.5 | 0.7 | 1.2 |
| DeepSeek-7B | 7.5 | 19.2 | 30.1 | 9.2 | 3.8 | 57.3 | 2.5 | 23.0 | 20.1 | 36.4 | 5.9 | 7.5 | 15.9 | 36.8 | 12.1 |
| Llama2-7B | 62.0 | 47.9 | 35.2 | 85.3 | 52.8 | 70.2 | 58.8 | 26.5 | 51.4 | 6.7 | 86.4 | 4.1 | 60.3 | 12.7 | 42.3 |
| Llama3-8B | 0.3 | 0.7 | 7.6 | 0.0 | 1.2 | 20.9 | 0.1 | 2.0 | 3.9 | 2.9 | 0.7 | 0.1 | 0.0 | 0.3 | 0.1 |
| GPT-3.5 | 14.7 | 23.8 | 39.1 | 8.8 | 20.3 | 19.1 | 33.4 | 14.1 | 12.2 | 11.6 | 21.3 | 10.6 | 27.2 | 5.6 | 10.0 |
| GPT-4o-Mini | 19.5 | 9.7 | 12.7 | 6.2 | 2.9 | 5.6 | 7.1 | 1.8 | 5.6 | 2.1 | 6.8 | 3.8 | 5.0 | 9.1 | 7.1 |

Negative Bias (-)

| | SP | PvdD | GL-PvdA | DENK | Volt | D66 | NSC | CDA | ChristenUnie | VVD | BBB | SGP | JA21 | PVV | FVD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 0.4 | 1.4 | 1.3 | 1.7 | 1.1 | 1.0 | 0.5 | 1.8 | 1.1 | 2.5 | 0.5 | 4.5 | 0.1 | 7.2 | 1.6 |
| Falcon3-7B | 54.1 | 53.6 | 78.8 | 45.4 | 57.6 | 26.7 | 45.7 | 43.5 | 37.7 | 46.7 | 51.5 | 35.8 | 55.7 | 62.3 | 55.4 |
| Gemma2-9B | 2.0 | 2.1 | 2.4 | 5.4 | 3.1 | 2.5 | 2.0 | 6.1 | 5.2 | 7.1 | 5.4 | 13.2 | 31.0 | 33.2 | 23.1 |
| DeepSeek-7B | 1.5 | 1.1 | 1.7 | 2.8 | 3.0 | 0.2 | 3.9 | 1.1 | 1.6 | 0.4 | 1.0 | 1.1 | 0.5 | 0.8 | 1.9 |
| Llama2-7B | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| Llama3-8B | 28.2 | 32.4 | 10.1 | 44.9 | 15.3 | 1.0 | 22.6 | 12.4 | 11.5 | 14.6 | 17.6 | 52.6 | 54.8 | 48.7 | 59.5 |
| GPT-3.5 | 12.1 | 10.6 | 1.5 | 15.9 | 3.6 | 5.2 | 2.4 | 11.9 | 13.1 | 25.3 | 6.3 | 20.8 | 5.8 | 52.7 | 29.3 |
| GPT-4o-Mini | 2.5 | 6.3 | 2.4 | 3.0 | 3.6 | 2.8 | 1.4 | 4.4 | 2.8 | 10.2 | 3.3 | 10.8 | 6.4 | 20.7 | 12.6 |

**(a) PoliBiasNO**

Positive Bias (+)

| | R | SV | MDG | AP | SP | V | KrF | H | FrP |
|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 0.0 | 0.0 | 0.3 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 |
| Falcon3-7B | 0.0 | 19.1 | 10.0 | 38.8 | 15.8 | 1.9 | 37.0 | 0.4 | 11.8 |
| Gemma2-9B | 4.4 | 2.8 | 11.0 | 4.1 | 0.8 | 0.2 | 0.7 | 0.2 | 0.1 |
| DeepSeek-7B | 10.7 | 37.9 | 50.7 | 14.4 | 6.5 | 2.9 | 18.8 | 1.2 | 16.7 |
| Llama2-7B | 0.9 | 78.4 | 46.5 | 28.9 | 10.0 | 20.8 | 67.2 | 19.6 | 10.6 |
| Llama3-8B | 7.7 | 15.5 | 13.5 | 4.4 | 5.0 | 6.4 | 2.0 | 0.6 | 0.3 |
| GPT-3.5 | 8.9 | 14.1 | 9.8 | 7.3 | 11.5 | 18.8 | 7.7 | 12.5 | 0.6 |
| GPT-4o-Mini | 24.0 | 28.1 | 27.9 | 21.5 | 27.5 | 13.0 | 16.7 | 11.2 | 9.7 |
| NorskGPT | 0.0 | 66.5 | 26.9 | 17.4 | 18.2 | 0.0 | 33.3 | 0.0 | 5.6 |

Negative Bias (-)

| | R | SV | MDG | AP | SP | V | KrF | H | FrP |
|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 73.4 | 82.0 | 53.6 | 74.0 | 78.3 | 58.8 | 63.1 | 70.2 | 87.0 |
| Falcon3-7B | 84.2 | 27.4 | 51.4 | 11.1 | 24.0 | 56.1 | 13.9 | 80.0 | 35.4 |
| Gemma2-9B | 9.4 | 10.4 | 6.5 | 7.8 | 11.6 | 16.8 | 13.7 | 17.3 | 25.4 |
| DeepSeek-7B | 3.8 | 1.7 | 1.2 | 4.0 | 7.2 | 10.0 | 3.4 | 13.5 | 3.2 |
| Llama2-7B | 11.7 | 0.2 | 1.2 | 1.4 | 7.1 | 2.2 | 0.1 | 2.9 | 7.0 |
| Llama3-8B | 14.8 | 13.3 | 14.9 | 22.6 | 19.7 | 9.2 | 27.8 | 38.7 | 54.3 |
| GPT-3.5 | 4.7 | 3.9 | 7.8 | 4.2 | 3.3 | 0.9 | 3.5 | 2.5 | 14.9 |
| GPT-4o-Mini | 1.8 | 1.2 | 2.6 | 0.9 | 1.2 | 1.7 | 1.4 | 2.7 | 4.2 |
| NorskGPT | 53.4 | 0.4 | 6.0 | 3.0 | 3.1 | 27.8 | 0.8 | 79.4 | 12.2 |

**(a) PoliBiasES**

Positive Bias (+)

| | CUP | EH Bildu | Más País | ERC | PSOE | Junts | PNV | Ciudadanos | PP | VOX |
|---|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 20.0 | 26.7 | 40.0 | 20.0 | 20.0 | 33.3 | 26.7 | 20.0 | 26.7 | 6.7 |
| Falcon3-7B | 69.3 | 54.9 | 77.4 | 61.0 | 87.9 | 83.2 | 78.5 | 73.9 | 85.4 | 46.0 |
| Gemma2-9B | 14.2 | 12.1 | 17.1 | 12.4 | 21.6 | 9.2 | 13.4 | 9.9 | 7.4 | 1.4 |
| DeepSeek-7B | 3.2 | 0.3 | 5.1 | 2.1 | 3.6 | 2.1 | 0.9 | 2.3 | 3.9 | 3.6 |
| Llama2-7B | 94.1 | 87.1 | 90.6 | 89.0 | 96.2 | 97.0 | 95.2 | 92.9 | 95.7 | 56.5 |
| Llama3-8B | 15.4 | 22.8 | 30.2 | 30.9 | 47.5 | 35.2 | 50.6 | 16.7 | 22.8 | 14.2 |
| GPT-3.5 | 11.3 | 0.0 | 2.1 | 2.1 | 9.3 | 6.2 | 1.0 | 5.2 | 8.2 | 8.2 |
| GPT-4o-Mini | 10.9 | 8.7 | 7.9 | 8.4 | 11.7 | 4.9 | 6.0 | 12.5 | 25.3 | 16.2 |
| Aguila-7B | 2.1 | 3.1 | 2.3 | 4.2 | 1.4 | 3.7 | 3.0 | 3.4 | 2.7 | 3.0 |

Negative Bias (-)

| | CUP | EH Bildu | Más País | ERC | PSOE | Junts | PNV | Ciudadanos | PP | VOX |
|---|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 0.8 | 1.2 | 0.6 | 1.4 | 0.8 | 1.2 | 1.2 | 2.2 | 2.6 | 2.4 |
| Falcon3-7B | 20.0 | 36.5 | 13.0 | 31.4 | 8.6 | 13.7 | 16.6 | 14.8 | 12.8 | 53.9 |
| Gemma2-9B | 1.3 | 0.6 | 0.4 | 0.4 | 0.4 | 0.3 | 0.3 | 0.3 | 0.5 | 1.3 |
| DeepSeek-7B | 68.5 | 76.2 | 50.1 | 68.9 | 79.5 | 54.2 | 62.8 | 61.0 | 74.5 | 76.4 |
| Llama2-7B | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| Llama3-8B | 62.2 | 61.2 | 49.3 | 54.5 | 41.8 | 50.1 | 41.6 | 66.4 | 57.2 | 73.9 |
| GPT-3.5 | 4.3 | 16.3 | 4.3 | 13.4 | 13.2 | 3.6 | 3.9 | 6.9 | 21.7 | 57.0 |
| GPT-4o-Mini | 21.3 | 17.8 | 8.8 | 13.7 | 9.3 | 8.5 | 9.9 | 10.5 | 21.5 | 53.4 |
| Aguila-7B | 29.7 | 34.4 | 36.7 | 27.6 | 31.2 | 32.4 | 33.6 | 32.3 | 31.7 | 34.3 |

**Figure 3: Entity Bias Index (EBI) heatmaps for positive and negative bias in LLMs, computed via counterfactual attribution of voting motions in the benchmark datasets. Panels (a), (b), and (c) correspond to Dutch, Norwegian, and Spanish political parties, respectively. In all panels, the parties on the x-axis are ordered from left-progressive to right-conservative ideologies.**

in Norway, and 10 in Spain). By keeping the motion text fixed and varying only the attributed party, we can determine whether the identity of the party itself influences an LLM's stance, thereby isolating systematic party-specific biases. Similar to the Ideological Bias experiment, we also extend with an 'abstain' option for testing against the Spanish benchmark.

**Metrics.** To quantify entity bias, we define the Entity Bias Index (EBI), which measures how often—and in which direction—a model's voting decision shifts when a motion is attributed to a political party $x$, compared to a baseline in which no party is mentioned.

Let $R_l(x, i)$ denote the response of an LLM $l$ to motion $i$ when attributed to party $x$, and $R_l(-, i)$ denote the response when no party is specified. We encode LLM responses as 1 for 'for' and 0 for

'against'.[7] The Entity Bias Index for model $l$ and party $x$ is defined as:

$$\text{EBI}_l(x) = \left( \frac{1}{n} \sum_{i=1}^{n} \left( R_l(x, i) - R_l(-, i) \right) \right) \times 100\%.$$

A positive value $\text{EBI}_l(x) > 0$ indicates that the model becomes more supportive of motions when they are attributed to party $x$, reflecting a **positive entity bias**. A negative value $\text{EBI}_l(x) < 0$ indicates reduced support, reflecting a **negative entity bias**. An index of zero implies that party attribution has no systematic effect on the model's stance.

**Results: Political Entity Bias.** Across the three benchmarks, the EBI heatmaps in Fig. 3 reveal several consistent patterns, alongside

---

[7]In the ideological bias experiment, party votes use a $+1/-1$ scale. In EBI, we instead use $1/0$, so that $R_l(x, i) - R_l(-, i)$ naturally falls in $\{-1, 0, +1\}$, capturing whether party attribution decreases, leaves unchanged, or increases the model's support.

important model- and country-specific exceptions. First, positive entity bias is generally small to moderate, but it varies substantially across models and political contexts. In the Dutch case, positive bias tends to be somewhat higher for left and centre-left parties such as *SP*, *GL-PvdA*, and *D66*, particularly for Llama2-7B. Llama2-7B also displays notable positive bias toward the agrarian populist *BBB*. Llama2-7B behaves similarly in Norway and even more strongly in Spain, where several models—most prominently Llama2-7B and Falcon3-7B—show broad positive bias across all Spanish parties, irrespective of ideology. Other models, such as GPT-3.5, distribute positive bias more evenly across parties in Norway, producing only mild increases. Overall, positive entity bias is model-dependent and sometimes ideology-aligned, but in other cases broad and non-specific, especially in the Spanish dataset. Notably, Llama2-7B does not follow this pattern: it presents systematically elevated positive EBI across the Dutch, Norwegian, and Spanish datasets.

In contrast, negative entity bias is stronger, more consistent, and much more ideologically structured. Right-conservative and far-right parties—including the Dutch *VVD*, *SGP*, *PVV*, and *FvD*; the Norwegian *H* and *FrP*; and the Spanish *PP* and *VOX*—attract the clearest and most persistent negative EBI values across most LLMs, especially Llama3-8B and GPT series. GPT-4o-mini displays especially strong negative bias toward right-wing parties in all three countries, while GPT-3.5-turbo shows pronounced negative bias toward *PVV* in the Netherlands and *VOX* in Spain. At the same time, some LLMs reveal pronounced country-specific patterns: for instance, Mistral-7B exhibits little negative bias in the Dutch and Spanish datasets but strong negative bias toward nearly all Norwegian parties, whereas deepseek-7b shows mild negative bias in the Dutch and Norwegian cases but notably strong negative bias across almost all Spanish parties.

Language-specific models do not behave in a politically coherent or nationally aligned manner. NorskGPT shows strong positive bias toward *SV* but negative bias toward both *R* (far-left) and *H* (centre-right), while Aguila-7B in Spain displays diffuse small positive bias across most parties but moderate negative bias across those same parties. These patterns suggest that local models do not exhibit greater ideological affinity with domestic party families; instead, their behaviour is more variable and model-specific.

In summary, our results show that negative entity bias is the most consistent and ideologically structured pattern in the EBI analysis, whereas positive bias is weaker, more heterogeneous, and sometimes non-ideological. Some models exhibit broad defer-ence—especially Llama2-7B and Falcon3-7B in Spain—while larger models such as GPT-4o-mini show highly consistent negative bias toward right-wing parties across all three political systems. These findings underscore that entity bias is not merely a by-product of national political context but reflects broader regularities and model-specific tendencies in how LLMs respond when political party names are introduced into prompts.

## 4.3 Prompt Brittleness

To assess whether our findings depend on the exact formulation of the voting prompt, we follow similar method mentioned in [8] and examine the robustness of LLM predictions under a set of semantically equivalent paraphrased prompts in PoliBiasNL and PoliBiasNO. These variants differ only in linguistic framing, such as reordering clauses, adjusting assertiveness, or using alternative but synonymous verbs. For each model, we measure how often its predicted stance on a motion (for/against) changes when the prompt is rephrased.

Overall, we observe that smaller models exhibit moderate prompt brittleness, occasionally flipping their predictions across variants, whereas larger models such as GPT-3.5 and GPT-4o-mini remain highly stable. Crucially, despite these local fluctuations, the global ideological patterns identified in our main analysis—CHES projec-tions and voting-agreement profiles—remain consistent across all prompt variants. Thus, while prompt phrasing can affect individual predictions, the broader political tendencies of the models are ro-bust to surface-level linguistic changes. Additional details, metrics, and brittleness heatmaps appear in Appendix A.

## 5 Discussion

A recent study [8] evaluates political worldviews in LLMs using VAA-style policy questions, and finds a general tendency towards left–liberal orientations, together with sensitivity to prompt refor-mulations. Our findings independently point in the same direction: across all three parliamentary datasets, LLMs are positioned pre-dominantly on the left– progressive side of the ideological space, and their predictions vary under controlled prompt modifications.

The two approaches differ in their empirical foundations. VAA questionnaires rely on a manually selected set of policy statements which, although intended to reflect a broad range of issues, still constitute a curated selection whose coverage of the full political agenda cannot be guaranteed. In contrast, our analysis is based on large-scale parliamentary motions and votes that span the full scope of legislative activity. This provides a more comprehensive and objective reflection of party ideology as expressed in real par-liamentary behaviour. The left-leaning tendencies observed in our study therefore arise from a substantially broader empirical base.

Our framework also offers additional diagnostics that are not captured in questionnaire-based evaluations. In particular, the En-tity Bias Index reveals a consistent pattern of negative bias toward right-wing parties: LLMs tend to under-align with their voting behaviour while more closely matching the positions of left and centre-left parties. This entity-specific asymmetry complements the ideological analysis and provides finer-grained insight into how models relate to individual political actors.

Taken together, the two lines of work offer complementary per-spectives on the political behaviour of LLMs: while VAA-based studies highlight issue-level orientations, our roll-call-based ap-proach shows that similar tendencies persist across the full breadth of real-world legislative decisions and reveals additional patterns at the party-entity level.

Understanding the ideological behaviour of LLMs is an impor-tant topic with significant societal relevance, especially as these models increasingly mediate citizens' access to political informa-tion. Although our findings align with those reported in VAA-based studies, our use of large-scale parliamentary voting data offers a more comprehensive and robust basis for evaluation. Because VAA statements are manually selected and may not fully represent the

breadth of issues addressed in real legislative debates, their conclusions may not generalise across future models or across differing political contexts. As LLM architectures and training pipelines evolve, it remains an open question whether questionnaire-based and roll-call-based analyses will continue to yield the same patterns. This underlines the need for empirically grounded benchmarks and systematic evaluation frameworks that allow the field to track how ideological tendencies emerge, persist, or diverge in subsequent generations of LLMs.

## 6 Conclusion and Future Work

We introduced a general framework for constructing political-bias benchmarks from parliamentary motions and party votes, and instantiated it in three national contexts: the Netherlands, Norway, and Spain. Building on these datasets, we proposed an evaluation methodology that assesses ideological positioning, visualises LLMs and parties in a shared CHES space, and quantifies party-specific entity bias. Our findings reveal consistent centre-left and progressive tendencies across models, together with systematic negative bias toward right-conservative parties, and show that these patterns remain stable under paraphrased prompts.

Future work includes extending the benchmark to additional legislatures, enabling longitudinal tracking of ideological drift in LLMs, and developing mitigation strategies informed by our diagnostics. Our framework provides a scalable basis for transparent and empirically grounded evaluation of political bias in LLMs.

## Ethical Considerations

This work uses publicly available parliamentary voting records and political motions from the Dutch, Norwegian, and Spanish parliaments. All data describe institutional decisions rather than private individuals and contain no personally identifiable information. No human subjects were recruited, and no informed consent was required.

## References

[1] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. The Falcon Series of Open Language Models. arXiv:arXiv:2311.16867

[2] Ryan Bakker, Catherine de Vries, Erica Edwards, Liesbet Hooghe, Seth Jolly, Gary Marks, Jonathan Polk, Jan Rovny, Marco Steenbergen, and Milada A. Vachudova. 2015. Measuring Party Positions in Europe: The Chapel Hill Expert Survey Trend File, 1999–2010. *Party Politics* 21, 1 (2015), 143–152. doi:10.1177/1354068812462931 Original work published 2012.

[3] Tanise Ceron, Neele Falk, Ana Baric, Dmitry Nikolaev, and Sebastian Padó. 2024. Beyond Prompt Brittleness: Evaluating the Reliability and Consistency of Political Worldviews in LLMs. *Trans. Assoc. Comput. Linguistics* 12 (2024), 1378–1400. doi:10.1162/TACL_A_00710

[4] DeepSeek-AI, :, Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wenjie Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren, Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Junxiao Song, Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Minghui Tang, Bingxuan Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji Wang, Tong Wu, Y. Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yiliang Xiong, Hanwei Xu, R. X. Xu, Yanhong Xu, Dejian Yang, Yuxiang You, Shuiping Yu, Xingkai Yu, B. Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghua Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, and Yuheng Zou. 2024. DeepSeek LLM: Scaling Open-Source Language Models with Longtermism. arXiv:arXiv:2401.02954

[5] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, and et al. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] https://arxiv.org/abs/2407.21783

[6] Stephen Gbenga Fashoto, Jérôme Rutinowski, Sven Franke, Jan Endendyk, Ina Dormuth, Moritz Roidl, and Markus Pauly. 2024. The Self-Perception and Political Biases of ChatGPT. *Human Behavior and Emerging Technologies* 2024 (2024), 7115633. doi:10.1155/2024/7115633

[7] Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 11737–11762. doi:10.18653/v1/2023. acl-long.656

[8] Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and Fairness in Large Language Models: A Survey. arXiv:2309.00770 [cs.CL] https://arxiv.org/abs/2309.00770

[9] Jochen Hartmann, Jasper Schwenzow, and Maximilian Witte. 2023. The political ideology of conversational AI: Converging evidence on ChatGPT's pro-environmental, left-libertarian orientation. arXiv:2301.01768 [cs.CL]

[10] Arthur E. Hoerl and Robert W. Kennard. 1970. Ridge Regression: Biased Estimation for Non-Orthogonal Problems. *Technometrics* 12, 1 (1970), 55–67. doi:10.1080/00401706.1970.10488634

[11] Harold Hotelling. 1933. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* 24 (1933), 417–441, 498–520.

[12] Nannan Huang, Haytham Fayek, and Xiuzhen Zhang. 2024. Bias in Opinion Summarisation from Pre-training to Adaptation: A Case Study in Political Bias. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, Yvette Graham and Matthew Purver (Eds.). Association for Computational Linguistics, St. Julian's, Malta, 1041–1055. https://aclanthology.org/2024.eacl-long.63

[13] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. arXiv:arXiv:2310.06825

[14] Tae Won Kim. 2023. Application of artificial intelligence chatbots, including ChatGPT, in education, scholarly work, programming, and content generation and its prospects: a narrative review. *Journal of Educational Evaluation for Health Professions* 20 (Dec. 2023), 38. doi:10.3352/jeehp.2023.20.38

[15] Ruibo Liu, Chenyan Jia, Jason Wei, Guangxuan Xu, Lili Wang, and Soroush Vosoughi. 2021. Mitigating Political Bias in Language Models through Reinforced Calibration. *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 17

(May 2021), 14857–14866. doi:10.1609/aaai.v35i17.17744

[16] Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2022. HateXplain: A Benchmark Dataset for Explainable Hate Speech Detection. In *AAAI*. AAAI: Association for the Advancement of Artificial Intelligence. arXiv:2012.10289 https://arxiv.org/abs/2012.10289

[17] Katelyn Mei, Sonia Fereidooni, and Aylin Caliskan. 2023. Bias Against 93 Stigmatized Groups in Masked Language Models and Downstream Sentiment Classification Tasks. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency* (Chicago, IL, USA) *(FAccT '23)*. Association for Computing Machinery, New York, NY, USA, 1699–1710. doi:10.1145/3593013.3594109

[18] Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2023. More human than human: measuring ChatGPT political bias. *Public Choice* 198, 1–2 (Aug. 2023), 3–23. doi:10.1007/s11127-023-01097-2

[19] Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, Online, 5356–5371. doi:10.18653/v1/2021.acl-long.416

[20] Manish Nagireddy, Lamogha Chiazor, Moninder Singh, and Ioana Baldini. 2024. SocialStigmaQA: A Benchmark to Uncover Stigma Amplification in Generative Language Models. In *The Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-24)*. Association for the Advancement of Artificial Intelligence, 21454–21462. https://doi.org/10.1609/aaai.v38i19.30142

[21] OpenAI. 2023. GPT-3.5 Turbo Model Documentation. https://platform.openai.com/docs/models/gpt-3-5-turbo.

[22] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael

Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. GPT-4 Technical Report. arXiv:arXiv:2303.08774

[23] Iñigo Parra. 2024. UnMASKed: Quantifying Gender Biases in Masked Language Models through Linguistically Informed Job Market Prompts. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, Neele Falk, Sara Papi, and Mike Zhang (Eds.). Association for Computational Linguistics, St. Julian's, Malta, 61–70. https://aclanthology.org/2024.eacl-srw.6

[24] Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. BBQ: A handbuilt bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 2086–2105. doi:10.18653/v1/2022.findings-acl.165

[25] Uwe Peters. 2022. Algorithmic Political Bias in Artificial Intelligence Systems. *Philosophy amp; Technology* 35, 2 (March 2022). doi:10.1007/s13347-022-00512-8

[26] Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Kirk, Hinrich Schütze, and Dirk Hovy. 2024. Political Compass or Spinning Arrow? Towards More Meaningful Evaluations for Values and Opinions in Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 15295–15311. doi:10.18653/V1/2024.ACL-LONG.816

[27] Jan Rovny, Jonathan Polk, Ryan Bakker, Liesbet Hooghe, Seth Jolly, Gary Marks, Marco Steenbergen, and Milada Anna Vachudova. 2025. The 2024 Chapel Hill Expert Survey on political party positioning in Europe: Twenty-five years of party positional data. *Electoral Studies* 97 (2025), 102981. doi:10.1016/j.electstud.2025.102981

[28] Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Rose Kirk, Hinrich Schütze, and Dirk Hovy. 2024. Political Compass or Spinning Arrow? Towards More Meaningful Evaluations for Values and Opinions in Large Language Models. arXiv:2402.16786 [cs.CL]

[29] Nikil Selvam, Sunipa Dev, Daniel Khashabi, Tushar Khot, and Kai-Wei Chang. 2023. The Tail Wagging the Dog: Dataset Construction Biases of Social Bias Benchmarks. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 1373–1386. doi:10.18653/v1/2023.acl-short.118

[30] Judith Simon, Pak-Hang Wong, and Gernot Rieder. 2020. Algorithmic bias and the Value Sensitive Design approach. *Internet Policy Review* 9, 4 (Dec. 2020). doi:10.14763/2020.4.1534

[31] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open Models Based on Gemini Research and Technology. arXiv:arXiv:2403.08295

[32] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva,

Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv e-prints*, Article arXiv:2307.09288 (July 2023), arXiv:2307.09288 pages. doi:10.48550/arXiv.2307.09288

[33] Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou, and Muhao Chen. 2023. A Causal View of Entity Bias in (Large) Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 15173–15184. doi:10.18653/v1/2023.findings-emnlp.1013

[34] Yiwei Wang, Muhao Chen, Wenxuan Zhou, Yujun Cai, Yuxuan Liang, Dayiheng Liu, Baosong Yang, Juncheng Liu, and Bryan Hooi. 2022. Should We Rely on Entity Mentions for Relation Extraction? Debiasing Relation Extraction with Counterfactual Analysis. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir

Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 3071–3081. doi:10.18653/v1/2022.naacl-main.224

[35] Svante Wold, Michael Sjöström, and Lennart Eriksson. 2001. PLS-regression: a basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems* 58, 2 (2001), 109–130. doi:10.1016/S0169-7439(01)00155-1 PLS Methods.

[36] Jun Yan, Yang Xiao, Sagnik Mukherjee, Bill Yuchen Lin, Robin Jia, and Xiang Ren. 2022. On the Robustness of Reading Comprehension Models to Entity Renaming. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 508–520. doi:10.18653/v1/2022.naacl-main.37

[37] Yongchun Zhu, Qiang Sheng, Juan Cao, Shuokai Li, Danding Wang, and Fuzhen Zhuang. 2022. Generalizing to the Future: Mitigating Entity Bias in Fake News Detection. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. ACM. doi:10.1145/3477495.3531816

# A Appendix: Prompt Brittleness Analysis

This appendix provides the full details of the prompt brittleness experiment summarised in Section 4.3. The goal of this analysis is to test the stability of model predictions under controlled variations in prompt wording. For each parliamentary motion, we generate several semantically equivalent paraphrased versions of the voting prompt, differing only in lexical framing, assertiveness, and syntactic brittleness. The core voting task remains unchanged.

## A.1 Experimental Setup

Each motion is evaluated under multiple brittleness variants. For each model, we record whether its predicted stance changes when moving from the baseline prompt to any of the paraphrased variants. This allows us to quantify brittleness at the motion level and to inspect which models are most sensitive to prompt formulation.

## A.2 Prompt Brittleness Index (PBI)

The *Prompt Brittleness Index (PBI)* quantifies the extent to which a model's predicted stance is unstable under systematically paraphrased prompt variants. Let each model output be encoded as 1 ("for") or 0 ("against"). For a given variation type $x$ (e.g., brittleness) and stance $s \in \{1, 0\}$, we define

$$\mathrm{PBI}_{\mathrm{abs},l}(x, s) = \frac{N_{\mathrm{flipped},l}(x, s)}{N_{\mathrm{total}}}, \tag{1}$$

$$\mathrm{PBI}_{\mathrm{norm},l}(x, s) = \frac{N_{\mathrm{flipped},l}(x, s)}{N_s}. \tag{2}$$

where:

- $N_{\mathrm{flipped},l}(x, s)$ is the number of unique motions for which model $l$ changes its stance $s$ across any prompt variant of type $x$;
- $N_{\mathrm{total}}$ is the total number of motions evaluated;
- $N_s$ is the total number of outputs where model $l$ produces stance $s$ across all variants.

The absolute PBI provides a global measure of robustness but can be dominated by the majority class: if a model produces far more *for* than *against* predictions (or vice versa), most flips will naturally occur in the larger class. The stance-normalised PBI corrects for this imbalance by conditioning on the baseline stance. This allows us to separately assess brittleness for motions initially classified as *for* and those classified as *against*, revealing asymmetric vulnerabilities that the absolute metric may mask.

## A.3 Results

Fig. 4 provide brittleness-variant flip heatmaps. The results show that smaller models such as Mistral-7B, Falcon3-7B, and LLaMa2-7B

exhibit the highest brittleness, with a noticeable fraction of predictions flipping across brittleness variants. Medium-sized instruction-tuned models show moderate sensitivity. Larger models—GPT-3.5 and GPT-4o-mini—display very low brittleness, rarely altering their predictions across paraphrases. Stance-normalised PBI further reveals that models with a strong tendency toward a specific stance exhibit lower brittleness for that stance.
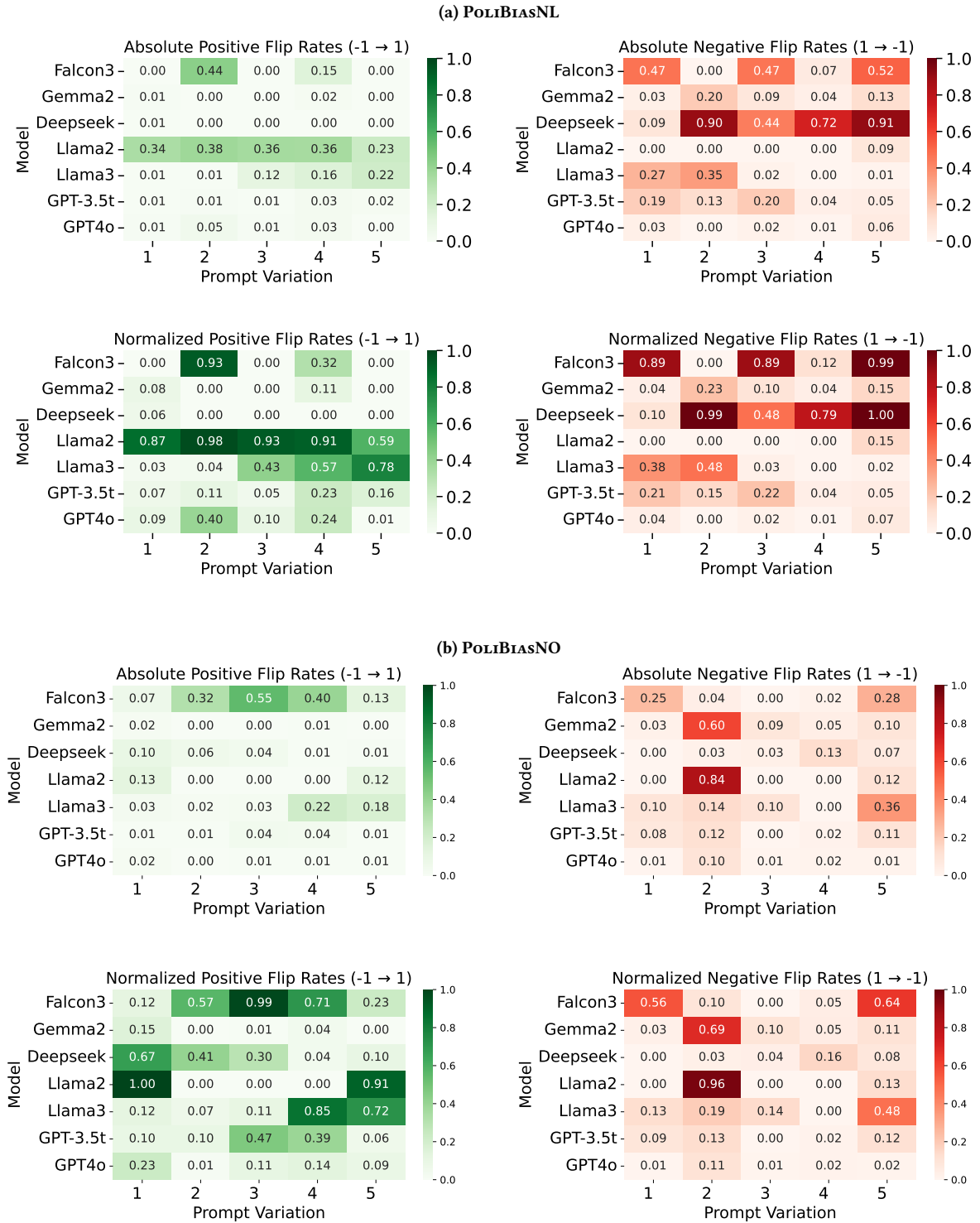
## A.4 Impact on Ideological Conclusions

Despite these local fluctuations, we observe that the broader ideological patterns reported in the main paper remain stable under all brittleness variants. CHES-based ideological projections, voting-agreement structures, and entity-bias patterns are nearly unchanged. This indicates that prompt wording affects individual decisions but does not substantially alter aggregate ideological tendencies. The benchmark therefore captures robust model-level political patterns rather than prompt-specific artefacts.

# B Invalid LLM Output Rates

Table 1 summarises the proportion of invalid outputs across the three datasets. Overall, invalid responses are rare for most models, with the exception of Mistral-7B, which produces a higher invalid rate on the Dutch motions. The strongest models (GPT-4o-Mini, LLaMA3-8B, Gemma2-9B) exhibit near-perfect format adherence. Norwegian- and Spanish-specific models show expected behaviour only on their respective datasets. These results indicate that output-format reliability is largely model-dependent but remains stable across political domains.

| Model | PoliBiasNL | PoliBiasNO | PoliBiasES |
|---|---|---|---|
| Mistral-7B | 14.92% | 0.78% | 0.00% |
| Falcon3-7B | 0.00% | 0.12% | 0.00% |
| Gemma2-9B | 0.04% | 0.00% | 0.00% |
| DeepSeek-7B | 0.00% | 0.00% | 0.00% |
| LLaMA2-7B | 0.04% | 0.07% | 0.00% |
| LLaMA3-8B | 0.00% | 0.00% | 0.00% |
| GPT-3.5 | 0.00% | 0.00% | 0.00% |
| GPT-4o-Mini | 0.00% | 0.00% | 0.00% |
| Águila-7B (Spanish) | — | — | 0.00% |
| NorskGPT (Norwegian) | — | 0.04% | — |

**Table 1: Invalid output rates across the three PoliBias datasets. An invalid output is a model response that does not conform to the expected stance format and cannot be mapped to *for*/*against*. Dashes indicate that a model is not applicable for that language.**

## (a) PoliBiasNL



## (b) PoliBiasNO



**Figure 4: Prompt Brittleness Index (PBI) measures a model's sensitivity to prompt rewordings. Higher values indicate greater inconsistency across prompt variants, while lower values reflect greater robustness and stability. Prompt variations used in the experiment are as follows: (1) Extra Detail, (2) Label Substitution with "Agree"/"Disagree", (3) Label Substitution with "Support"/"Oppose", (4) Label Substitution with "Favorable"/"Detrimental", and (5) Label Order Inversion.**