# Aligning Large Language Models with Diverse Political Viewpoints

**Anonymous ACL submission**

## Abstract

Large language models such as ChatGPT often exhibit striking political biases. If users query them about political information, they might take a normative stance and reinforce such biases. To overcome this, we align LLMs with diverse political viewpoints from 100,000 comments written by candidates running for national parliament in Switzerland. Such aligned models are able to generate more accurate political viewpoints from Swiss parties compared to commercial models such as ChatGPT. We also propose a procedure to generate balanced overviews from multiple viewpoints using such models.

## 1 Introduction

Large language models have become very popular, with chat applications like ChatGPT and Gemini having hundreds of millions of active users combined.[1] One of the intended use cases is the retrieval of factual information (e.g., Mehdi, 2023). Interacting with chatbots can influence users' views (Jakesch et al., 2023) and potentially influence behavior (Stieger et al., 2021). Because of this, LLMs – if used as decision aids in high-stakes contexts such as shaping political views or votes – must return factually correct and unbiased statements.

Political bias is present in all first-generation LLMs (Feng et al., 2023). And also, even ChatGPT does not seem to be impartial: several recent papers have shown that it exhibits progressive, liberal, and pro-environmental biases (Rozado, 2023; Hartmann et al., 2023; Motoki et al., 2024; Rutinowski et al., 2024). Given these findings, Hartmann et al. (2023) ask: "What if Chat-GPT exhibits a political ideology that may pervade its synthetic responses and subtly influence its millions of unsuspecting users?"

> **prompt:** You are a helpful Swiss policy advisor. Below, you are asked a policy issue or question. You are in the political party **P**, and you reply in **L**. What's your opinion on the following issue or question: **Q**

Figure 1: Prompt for conditional generation. Varying attributes are party P, language L, and political issue Q. For example, party = "GLP", language = German, and political issue = Should the state do more to promote equal educational opportunities?.

To address such concerns, we propose to align LLMs towards generating diverse political viewpoints. Inspired by early work with GPT models and conditional generation based on metadata (Zellers et al., 2019), we align LLMs towards accurately reflecting diverse viewpoints from Swiss parties. Throughout the paper, we use prompts of the form shown in Figure 1.

We align LLMs with data obtained from the Swiss voting advice application *smartvote*: We have ca. 100,000 comments written by candidates running for national parliament in Switzerland, coupled with meta-data, comprising 33 attributes such as party, language, canton, age, favorite movies and others. Smartvote lets users and candidates fill out a survey about stances towards various political issues. Then, it matches candidates with the highest overlap to a user.

We have taken that survey with ChatGPT. It is most closely aligned with candidates from the *Green Liberal Party of Switzerland*, confirming its progressive and pro-environmental bias in the Swiss context (Rozado, 2023; Hartmann et al., 2023; Motoki et al., 2024; Rutinowski et al., 2024). We present anecdotal evidence of such biases in the Swiss context in Table 1 and a more exhaustive analysis in Appendix A.

To overcome such implicit biases present in LLMs, we align Llama 3 models (AI@Meta, 2024) with smartvote data, combining conditional generation (e.g., Zellers et al., 2019) and monolithic preference optimization (Hong et al., 2024) alignment. We find that the resulting aligned models

---

[1] https://www.demandsage.com/chatgpt-statistics/
https://www.demandsage.com/google-gemini-statistics/

| Party | ChatGPT generation | ORPO-aligned LLama 3 model |
|---|---|---|
| SP | Yes, the state should do more to promote equal educational opportunities in order to promote equal opportunities and reduce social inequalities. Tuition vouchers for pupils from low-income families can be a useful instrument in this regard. **overlap answer and party stance: high** | Equal opportunities are the foundation of our state. **overlap answer and party stance: high** |
| Die Mitte | Yes, Die Mitte supports the promotion of equal educational opportunities for all pupils. Measures such as tuition vouchers for pupils from low-income families can help to promote equal opportunities in the education system. **overlap answer and party stance: medium** | Funding is important, but not only for low-income families. **overlap answer and party stance: high** |
| FDP | Yes, the state should do more to promote equal educational opportunities. Tuition vouchers for pupils from low-income families could be a useful instrument in this regard. **overlap answer and party stance: low** | The cantons are responsible for providing schools. **overlap answer and party stance: high** |
| SVP | Yes, the state should do more to promote equal educational opportunities in order to promote equal opportunities and reduce social inequalities. Tuition vouchers for pupils from low-income families can be a useful instrument in this regard. **overlap answer and party stance: low** | Promoting low-achieving pupils is not the solution. **overlap answer and party stance: high** |

Table 1: Political stances generated with ChatGPT and an aligned model for the policy issue *Should the state do more to promote equal educational opportunities?* for all major Swiss parties represented in the Federal Council of Switzerland. Political leanings (taken from Wikipedia): SP = center-left, Die Mitte = center-right, center-right, Mitte = center-right, center-right, SVP = right-wing. = right-wing. Text in bold (overlap) inserted by authors.

generate more diverse and more accurate political viewpoints, which are preferred in human annotation.

Such models can be used to create political views of all Swiss parties towards an issue, which then could be summarized by other capable LLMs (e.g. OpenAI et al., 2024) to give balanced overviews. That potentially facilitates finding political compromises or learning more about political issues. However, we urge more research to better understand the promises and dangers of AI providing political information or voting advice. We strongly believe that if LLMs were used in such circumstances, they'd better be accurate and impartial.

## 2 Data

We use the same data source as (Vamvas and Sennrich, 2020): Comments written by candidates running for national parliament in Switzerland around 200 political issues. The comments are submitted to the voting advice application *smartvote*, which helps voters determine which candidates or parties have similar political stances. Prior to an election, candidates can report their stance on a short (30 questions) or long survey (75 questions) across various political issues. Voters can take the same survey and are matched with candidates having a high overlap. The questions are drafted by a team of political scientists (for more details, see Thurman and Gasser, 2009). Candidates can also submit comments that further explain their stance.

We use these comments and metadata to align models using conditional generation. Smartvote is a popular service where 85% of candidates running for elections in Switzerland have a smartvote profile, and one in five voters consults smartvote before elections. Thus, our data is representative of political preferences in Switzerland. We show more detailed dataset statistics in Appendix B.

## 3 Methods

In conditional generation, we want to generate text based on constraints or metadata (e.g., Zellers et al., 2019; Zhou et al., 2023). For example, previous work generated news articles based on the attributes *domain, date, authors* and *headline* (Zellers et al., 2019). Alignment datasets usually contain triples of the form *instruction, preferred choice* and *rejected choice*.[2] We interpret conditional generation for alignment to sample a comment towards a political issue $q$ drafted by somebody from party $p$ speaking language $l$ as the preferred choice. For the rejected choice, we sample comment for issue $q$ speaking language $l$, but being part of a different political party $\neg p$.

We use reference-free monolithic preference optimization (ORPO Hong et al., 2024) and optimize the following objective taken directly from the ORPO paper. We optimize the following joint loss

$$\mathcal{L}_{ORPO} = E_{(x,y_w,y_l)} \left[ \mathcal{L}_{SFT} + \lambda \cdot \mathcal{L}_{OR} \right] \quad (1)$$

$$\mathcal{L}_{OR} = -\log \sigma \left( \log \frac{\mathbf{odds}_\theta(y_w|x)}{\mathbf{odds}_\theta(y_l|x)} \right) \quad (2)$$

where the first part of equation 1 is just supervised fine-tuning. The second part $\mathcal{L}_{OR}$ from equation 2 increases the likelihood of the preferred choice $y_w$ and decreases the likelihood of the rejected choice $y_l$. For details, we refer to (Hong et al., 2024). We believe this loss is suited for conditional generation, as it further pushes apart similar generations with subtle differences due to different metadata. We will compare ORPO-aligned models to direct supervised fine-tuning (dSFT) (Taori et al., 2023) in the results section.

---

[2]That is for reference-free methods such as DPO (Tunstall et al., 2023) or ORPO (Hong et al., 2024).

We also experimented with direct policy optimization DPO (Rafailov et al., 2023) following the recipe outlined in (Tunstall et al., 2023) and RLHF (Stiennon et al., 2020). However, the resulting models did not pass initial vibes tests.

In all our experiments, we used the transformer TRL library (von Werra et al., 2020) and the 4bit quantized unsloth version of Llama 3 8B models[3] and fine-tuned models using LoRA (Hu et al., 2022). For the supervised fine-tuning, we used the hyper-parameters outlined in (Tunstall et al., 2023), and for ORPO fine-tuning, we proceeded with the hyper-parameters outlined in (Hong et al., 2024).

# 4 Results

We present four sets of results on our dataset's development and test split outlined in Appendix B. All our models use the template shown in Figure 1, and we present results for the following models: (1) ChatGPT 3.5 zero-shot, (2) Llama-3-instruct zero-shot, (3) Llama-3-instruct-finetuned-with-dSFT (Taori et al., 2023) and (4) Llama-3-instruct-aligned-with-ORPO (Hong et al., 2024).
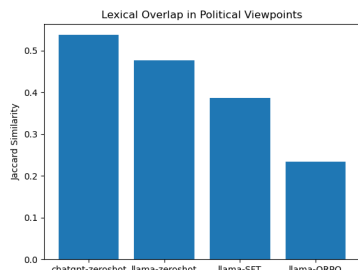
## 4.1 Diversity of Generations



Figure 2: Average diversity of replies within a political issue, measured with Jaccard similarities.

We show qualitative evidence of political bias in ChatGPT generations and a lack of variety in responses in Table 1, where ChatGPT generates similar progressive responses for all parties, although party stances vary substantially. We present more quantitative evidence of such phenomena in Figure 2. For each political issue and model, we compute Jaccard similarities between generations for all parties and show average overlap of responses measured in jaccard similarities[4].

---

[3]https://huggingface.co/unsloth/llama-3-8b-Instruct-bnb-4bit

[4]Jaccard similarity computes the intersection divided by the size of the union, so if two-word sets contain all the same

We find that zero-shot generations with Llama 3 and ChatGPT are strikingly similar on average, no matter the prompted party. Supervised fine-tuning results in more diverse generations, reducing the number of overlapping words by 30% compared to ChatGPT. The ORPO-aligned models further reduce overlapping generations and result in an average similarity between two outputs of 0.24, roughly half of the overlap measured for the zero-shot experiments.

## 4.2 Quantitative Evaluation

We further measure model generations using MAUVE scores (Pillutla et al., 2021). MAUVE is an automated metric that measures the gap between neural text and human references using LLM representations. The higher the MAUVE score, the closer the generated neural and human references are. Because our generations are either in German, French, or Italian, we use a multi-lingual RoBERTa model as a featurizer (Conneau et al., 2020). Table 2 shows the resulting scores over different dataset splits. We show average results over five runs (with 95% confidence intervals), sampling different reference comments in each run.

| Model | MAUVE dev | MAUVE test | MAUVE (dev and test) |
|---|---|---|---|
| ChatGPT | $0.36 \pm 0.02$ | $0.25 \pm 0.05$ | $0.24 \pm 0.02$ |
| Llama 3 zero-shot | $0.27 \pm 0.05$ | $0.03 \pm 0.0$ | $0.08 \pm 0.01$ |
| Llama 3 SFT | $0.48 \pm 0.02$ | $0.48 \pm 0.03$ | $0.38 \pm 0.02$ |
| Llama 3 ORPO | $\mathbf{0.63} \pm 0.03$ | $\mathbf{0.71} \pm 0.05$ | $\mathbf{0.64} \pm 0.01$ |

Table 2: Automated metrics measuring overlap between model-generated replies and actual replies in the development and testset.

ORPO-aligned generations are by far closest to the actual reference comments, followed by generations produced by the Llama-3-SFT model. Perhaps not surprisingly, the zero-shot answers are most distant from human references, arguably because zero-shot models do not know what such comments look like. Results are robust across runs, and the 95% confidence intervals remain small. Furthermore, we computed MAUVE scores with an MBERT encoder (Devlin et al., 2019) which produces similar scores and the same ranking (Devlin et al., 2019).

## 4.3 Human Evaluation

Lastly, we performed a human evaluation of the generated comments. Each annotated datapoint

---

words, this is 1. If there exist no overlapping words, the intersection is 0.

3

consists of an instruction (see Figure 1) and two randomly sampled generations from different models. We then asked the human annotator which generations are preferred. Detailed instructions, annotator demographics, and further robustness validations can be found in Appendix Section C.
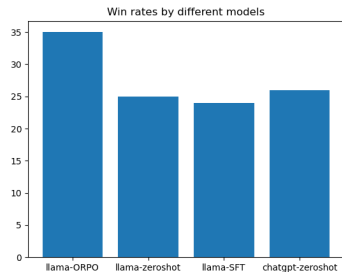


Figure 3: Win rates by different models.

Generations produced by the ORPO-aligned models are preferred in around 60% of the cases, whereas generations from the other models produce similar win rates. The author team, with the help of the mayor of a Swiss city, manually annotated 40 comments in a deliberative setting. All data points during that annotation round were discussed at length, and before seeing the model output, the team discussed what an optimal generation would look like. We then discussed which of the two generations is closer to that. We treat this set as the gold standard. Inner annotator agreement of the annotator and this set is 0.33 indicating fair agreement. If we don't punish generations that have been annotated as being of similar quality, agreement rises to 0.69, indicating substantial agreement. All evaluations point towards aligning LLMs with ORPO seems to work best to generate diverse political viewpoints.

## 5 Discussion and Conclusion

If explicitly asked to produce political viewpoints from a certain party perspective, LLM-generated text better be accurate and faithful to these party viewpoints. We show that this is not the case in the Swiss context when using zero-shot methods.

Combining alignment and conditional generation substantially improves such generations, as we have shown qualitatively and quantiatively throughout this work. Also, ORPO-aligned models seem to work better in conditional generation than models obtained via supervised fine-tuning.

Such aligned models have practical use cases beyond accurately presenting specific party pref-

erences. Using a simple algorithm (see Figure 4), we can also generate an overview of viewpoints towards a specific issue. This algorithm involves producing the party-specific stances with a given LLM (aligned or zero-shot), and then using Chat-GPT to summarize these stances. For illustration, we run this procedure for the issue *Should the state do more to promote equal educational opportunities?*. As Figure 9 shows, the overview synthesized from the LLama-3-ORPO-generated stances is richer and more accurate than the one synthesized from ChatGPT-generated stances.

```
1: Initialize answers = ∅
2: for each party in parties do
3:     answer = generate_answer(LLM, party)
4:     answers.append(answer)
5: end for
6: synthesize_answers(gpt4, answers)
```

Figure 4: Pseudocode for generating and synthesizing answers.

Generating text reflecting political views has major implications. LLMs have the potential to shift attitudes and behavior. If they do that in the political domain, this might influence elections, one of the most important decision-making processes in democratic societies (Berger et al., 2008). Thus, we recommend (1) more research on promises and pitfalls of LLMs providing politically relevant information or perhaps even giving explicit voting advice and (2) a societal and scientific debate on whether LLMs and AI should be involved in democratic processes.

Hartmann et al. (2023) asked what if ChatGPT exhibited a political ideology that may pervade its synthetic responses and subtly influence its unsuspecting users? Our work speaks to this question, and we see three possible mitigation strategies going forward: (1) LLMs would always refuse to answer anything related to shaping political beliefs and take sincere political impartiality as an alignment goal. (2) LLMs would always produce broad overviews of political issues (as in Figure 9). Our work might facilitate creating appropriate datasets for this. (3) LLMs would explicitly produce text that is aligned with an ideology. In this case, however, the provider of an LLM would be fully transparent about this, and/or the user would fully control what ideology LLM-generated text should be aligned with. Aligning models with party preferences using conditional generation, as presented in this work, is one way toward that goal.

4

## Limitations

We acknowledge several limitations and outline a range of possibilities for future work.

**Choice of models and Alignment algorithms.** We have mainly experimented with LLama-3 models and ChatGPT 3.5 zero-shot. There are, by now, other capable open-source models (Mistral, Mixtral, Llama 2) or model sizes (70B) that we could have fine-tuned with the method proposed in this work. Also, there exists a range of different alignment algorithms (DPO, RLHF), which we have experimented with, but the resulting models did not pass initial vibes tests. We plan to investigate all of this more thoroughly in future work.

**Choice of metadata for conditional generation.** In preliminary experiments, we experimented with the (Vamvas and Sennrich, 2020) dataset and generated comments based on stance (pro/contra) and not party affiliation. Eyeballing these results indicates that ORPO-aligned models in this setting also return more diverse answers than zero-shot models, and ORPO-aligned model generations seem more creative than SFT models. We take this as evidence for the robustness of conditional alignment, but we have not exhaustively evaluated this. Next, we think there are exciting opportunities in alignment with more metadata, such as canton, age, gender, and any other attribute potentially influencing political viewpoints. We tried this in preliminary experiments with Mistral models and DPO. However, this didn't work. We plan to revisit this with LLama 3 and ORPO.

**Data availability.** We initially signed an NDA with smartvote about the release of dataset and artifacts (e.g., models). We are discussing and trying to release all data and models produced in this work, but we cannot guarantee this at the time of submission.

**Further data.** We believe it should be possible, in principle, to find more diverse data sources with party affiliation (e.g., newspaper or TV interviews, party website content). It should be possible to collect such, and this would make for a dataset including more diverse political questions, which might lead to more creative models. This approach can be used across countries and parties, and thus allow for replication studies in contexts not related to smartvote data or Switzerland.

## Ethics Statement

We also acknowledge ethical implications of our work.

**Contested topic.** We believe the combination of LLMs and democracy is a very delicate topic, and throughout the manuscript tried to do justice to such challenging circumstances. On one hand, we all exhibit political biases, which we tried to remove from the paper as good as possible, but we also acknowledge that we probably have not written a completeley impartial piece. The same holds for LLMs. Another goal of this paper is to increase the awareness about these points.

**Intended use case.** We believe it is important that LLMs, its users, developers testers and other stakeholders are aware of political bias in machine-generated text. In this paper, we do not argue for creating chatbots which act as echochambers and reinforce existing biases present in such models – or change political views or actions of users. We want to argue for the opposite, that LLMs should exactly not do that. We tried to do justice to this goal, but also hope and acknowledge that the points we raise in this paper hopefully be outdated soon.

**Biases in LLMs.** Political bias is one sort of bias present in LLMs. There are others (see e.g., Abid et al., 2021; Lucy and Bamman, 2021), which are not addressed in this work. Our resulting models may potentially perpetuate these biases.

**Accuracy, Hallucinations and outdated information.** Our aligned models, as well as chatGPT, are not 100% accurate in producing political information: They produce hallucinations or other potentially harmful text, hence we do not advocate to use them in a commercial context, but propose a method to potentially mitigate political biases in LLMs. Further, we align our models on smartvote comments from 2015 - 2023. Parties might change their stance in the meantime. It remains an open research question how to incorporate such changes in stances.

**Non-constitutional parties.** We included viewpoints of political parties that operate within the limits of the constitution. Whether LLMs should reproduce the content of extremist parties without disclaimers is not within the scope of our research.

# References

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.

AI@Meta. 2024. Llama 3 model card.

Jonah Berger, Marc Meredith, and S. Christian Wheeler. 2008. Contextual priming: Where people vote affects how they vote. *Proceedings of the National Academy of Sciences*, 105(26):8846–8849.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

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)*, pages 11737–11762, Toronto, Canada. Association for Computational Linguistics.

Jochen Hartmann, Jasper Schwenzow, and Maximilian Witte. 2023. The political ideology of conversational ai: Converging evidence on chatgpt's pro-environmental, left-libertarian orientation. *Preprint*, arXiv:2301.01768.

Jiwoo Hong, Noah Lee, and James Thorne. 2024. Orpo: Monolithic preference optimization without reference model. *Preprint*, arXiv:2403.07691.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-writing with opinionated language models affects users' views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery.

Li Lucy and David Bamman. 2021. Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55.

Yusuf Mehdi. 2023. Reinventing search with a new ai-powered microsoft bing and edge, your copilot for the web. Accessed: 2024-06-08.

Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2024. More human than human: Measuring chatgpt political bias. *Public Choice*, 198(1):3–23.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, and Ilge Akkaya et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. In *NeurIPS*.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Preprint*, arXiv:2305.18290.

David Rozado. 2023. The political biases of chatgpt. *Social Sciences*, 12(3):148.

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(1):7115633.

Mirjam Stieger, Christoph Flückiger, Dominik Rüegger, Tobias Kowatsch, Brent W. Roberts, and Mathias Allemand. 2021. Changing personality traits with the help of a digital personality change intervention. *Proceedings of the National Academy of Sciences*, 118(8):e2017548118.

Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. 2020. Learning to summarize from human feedback. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html*, 3(6):7.

James Thurman and Urs Gasser. 2009. Three case studies from switzerland: Smartvote. Technical report, Berkman Center Research Publications.

6

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment. *Preprint*, arXiv:2310.16944.

Jannis Vamvas and Rico Sennrich. 2020. X-stance: A multilingual multi-target dataset for stance detection. *Preprint*, arXiv:2003.08385.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl.

Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Ethan Wilcox, Ryan Cotterell, and Mrinmaya Sachan. 2023. Controlled text generation with natural language instructions. *Preprint*, arXiv:2304.14293.

## A ChatGPT and smartvote

We have taken the smartVote survey with Chat-GPT for the national elections in 2023 in Switzerland. We have taken the short survey (30 questions) using a temperature of 0 and the following system prompt: *You are a helpful research assistant in Switzerland. You are given a few questions which you need to answer. You MUST reply with only one of "yes", "rather yes", "rather no", "no".* The output categories are the options users can enter while taking the smartvote survey. The user prompt equals to the actual smart vote questions, e.g., "Should the state do more to promote equal educational opportunities"?.

In Figure 5, we show the output a user receives after taking the smartvote survey: In the top panel, we show political preferences across different dimensions. The categories correspond to: Liberale Gesellschaft = Liberal society; Offene Aussenpolitik = Open foreign policy; Liberale Wirtschaftspolitik = Liberal economic policy; Restriktive Finanzpolitik =Restrictive financial policy; Restriktive Migrationspolitk = Restrictive migration policy; Ausgebauter Umweltschutz = Expanded environmental protection; Ausgebauter Sozialstaat = Expanded welfare state.

In the bottom panel, we show the candidates who were identified as having the highest political overlap with ChatGPT. 7 out of 12 (58%) of the most aligned candidates would be from the (Young) Liberal Party of Switzerland (GLP or JGLP).

## B Dataset Statistics

In Table 3, we show 5 randomly sampled comments (and their English translation) from our dataset.

In Table 4, we show dataset statistics, the number of unique comments in each split, the number of political issues and the share of different languages in the different splits of the dataset. In Figure 6, we show a histogram of the sequence lengths in the dataset (across all splits). We excluded comments shorter than five words.

We have access to smartvote data for the national parliament elections in 2015, 2019, and 2023. We split the data into a training set, a development, and a test set. Both the development and test set consist of 10% of the political issues from the 2023 election that were not present in the 2015 or 2019 survey.

We show the 10 most often occuring parties and their associated number of comments in Table 5.



Figure 5: Overview of political preferences of ChatGPT in Switzerland for the 2023 national elections of parliament (source: smartvote.ch)



Swiss candidates have the highest overlap in political preferences regarding ChatGPT (source: smartvote.ch).



Figure 6: Sequence lengths of Smartvote comments.

| issue | comment (English) | comment (original) |
|---|---|---|
| Are you in favor of amending the social welfare guidelines to reduce benefits for large families and young adults? | I find it important to look at the individual persons/families concerned. | Auch hier finde ich die individuelle Betrachtung der betreffenden Personen/Familien wichtig. |
| Should the state provide more funds for health insurance premium reductions? | Yes, the state must invest massively more in combating rising poverty. | Ja, der Staat muss massiv mehr in die Bekämpfung der steigenden Armut investieren. |
| Should incentives and target agreements rather than bans and restrictions be used exclusively to achieve the climate targets? | To guide certain behaviors, however, the time has come for prohibitions and restrictions. | Pour guider certains comportement, l'heure est quand même aux interdictions et restrictions. |
| Should the differences between financially strong and weak cantons be reduced more through financial equalization? | Wealthy cantons have benefited greatly from corporate tax cuts in recent years. | Les cantons riches ont ces dernières années largement profité des réductions de l'imposition des entreprises. |
| The financially strong cantons would like to significantly reduce their contributions to the financially weak cantons as part of the financial equalization (NFA). Do you support this request? | Long-term abuse of solidarity is counterproductive. | Un abuso della solidarietà a lungo termine è contraproducente. |

Table 3: Example comments from the dataset (automatically translated with deepl, and manually checked whether the translation is accurate).

| split | # examples | # political issues | German (%) | French (%) | Italian (%) |
|---|---|---|---|---|---|
| train | 92,986 | 203 | 75.5 | 22.2 | 2.2 |
| dev | 4262 | 7 | 76.8 | 21.0 | 2.3 |
| test | 5488 | 7 | 78.4 | 19.3 | 2.2 |

Table 4: Dataset Statistics

## C  Annotation Guidelines

We recruited an annotator from Switzerland with a university degree in political science and a strong self-declared interest in Swiss politics. The annotator read a random sample of 200 messages. The annotator was instructed as follows:

*Here is some information on the project: Large Language Models (LLMs, such as ChatGPT) often have biases. For a research project, we have fine-tuned an open-source LLM to make it more representative of the political values of different Swiss people. We used Smartvote data for this alignment. Now, an important question is whether our fine-tuned model is better than other models in terms of how accurately it presents the viewpoints of Swiss politicians. This is where your contribution is important. In the linked file, you will see:*

- *A prompt*

- *Two responses from LLMs (Candidate A and Candidate B)*

- *Reference comments*

- *Two columns that you will fill in (see instructions below)*

***First column to fill in: your preference.*** *Please ask yourself if Candidate A or B is better. Here, "better" means that the respective Candidate most closely reflects the reference comments and aligns with your knowledge of that party's stance on the respective policy question. You can also leverage your knowledge of whether a given position would correspond to that party's mainstream position.*

- *1 = both A and B are good*

- *0: neither A nor B is good*

- *A = only A is good, or A is significantly better than B*

- *B = only B is good, or B is significantly better than A*

*Some illustrations: If both are OK/good, but one is better, you will enter your favorite. If both are good and it is impossible to choose between A and B, you enter "1". If none are good, enter "0".*

***Second column to fill in: your familiarity.*** *Here you would indicate how familiar you are with the party's position on the issue.*

- *1 = wild guess/need to read more about it*

- *2 = have a clue, but need to double-check*

- *3 = reasonable guess*

- *4 = quite sure about the party's stance on issue Q, I would be surprised if the stance is different from my intuition*

- *5 = 100% sure about the party's stance*

| party | # comments |
|---|---|
| FDP | 15589 |
| GLP | 11341 |
| GRÜNE | 8992 |
| SP | 8880 |
| EVP | 7734 |
| SVP | 6780 |
| DIE MITTE | 6274 |
| CVP | 4756 |
| EDU | 3940 |
| JG | 3595 |

Table 5: Caption

## C.1 Further Robustness Checks in Human Annotations

We drop annotation datapoints where the human annotator indicated that the familiarity with a political issue and the party stance would be less than a reasonable guess. We find that the inner annotator agreement (cohen's kappa) with these annotated datapoints and our gold standard annotations is 0.15. Thus, we don't think these datapoints are a reliable source for measuring preferences.

If we only consider data points where the annotator is fairly confident ("quite sure about the party's stance on issue Q" or above), the higher inner annotator agreement with our gold standard increases – and overall winrates increase more sharply for the ORP-aligned LLama 3 model. We show these win rates in Figure 7.
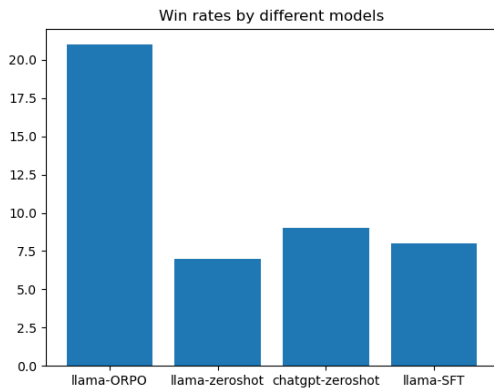
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638



Figure 7: Win rates by different models.

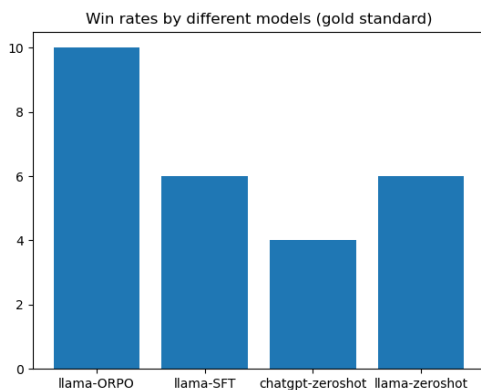Lastly, we show overall win rates by models in our gold standard annotations in Figure **??**

Figure 8: Win rates by different models (gold standard).

## D Generating Overviews over Political Issues

**Generated overview based on ChatGPT replies:** In Switzerland, there is a broad consensus among major political parties that the state should take more action to promote equal educational opportunities. The common rationale is to foster equality and reduce social disparities. A widely supported measure is the introduction of tutoring vouchers for students from low-income families, which is seen as an effective tool to level the educational playing field and provide all children with equal starting opportunities. This approach is endorsed across the political spectrum, reflecting a unified commitment to enhancing educational equity.

**Generated overview based on LLama-3-ORPO replies:** The question of whether the state should do more to promote equal educational opportunities is a significant topic in Swiss politics, with varied perspectives across the political spectrum. Some parties emphasize the importance of education as a fundamental right and a cornerstone of societal equity, advocating for enhanced state involvement to ensure equal opportunities from early education onwards. Others argue that the responsibility for education lies primarily with the cantons, suggesting that local governance is better suited to address educational needs. Some views focus on the importance of supporting both high and low-performing students, while some parties believe that the current educational opportunities in Switzerland are already sufficient and that the focus should instead be on integration. Overall, the debate reflects a balance between the need for national standards and local autonomy, as well as differing opinions on the best methods to achieve educational equity.

Figure 9: Prompt for conditional generation. Varying attributes are party, language and political issue.