

# Understanding and Mitigating Political Stance Cross-topic Generalization in Large Language Models

Jiayi Zhang<sup>1,2,3,\*</sup>, Shu Yang<sup>1,2,\*†</sup>, Junchao Wu<sup>4</sup>, Derek F. Wong<sup>4</sup>, Di Wang<sup>1,2†</sup>

<sup>1</sup>Provable Responsible AI and Data Analytics (PRADA) Lab

<sup>2</sup>King Abdullah University of Science and Technology (KAUST)

<sup>3</sup>University of Copenhagen

<sup>4</sup>University of Macau

## Abstract

Fine-tuning Large Language Models on a political topic will significantly manipulate their political stance on various issues and unintentionally affect their stance on broad topics. While previous studies have investigated this issue, there is still a lack of understanding regarding the internal representations of these stances and the mechanisms that lead to unintended cross-topic generalization. In this paper, we systematically explore the internal mechanisms underlying this phenomenon from a neuron-level perspective and how to mitigate the cross-topic generalization of political fine-tuning. Firstly, we propose Political Neuron Localization through Activation Contrasting (PNLAC) to identify two distinct types of political neurons: *general political neurons*, which govern stance across multiple political topics, and *topic-specific neurons* that affect the model’s political stance on individual topics. We find that these political neuron types exist in the middle and later layers across four models and datasets through activation patching experiments. Leveraging these insights, we introduce InhibitFT, an inhibition-based fine-tuning method that effectively mitigates the cross-topic stance generalization. Experimental results demonstrate the robustness of the identified neuron types across various models and datasets and show that InhibitFT significantly reduces the cross-topic stance generalization by 20% on average while preserving topic-specific performance. Moreover, we demonstrate that selectively inhibiting only 5% of neurons is sufficient to effectively mitigate the cross-topic stance generalization.

## 1 Introduction

The remarkable capabilities of Large Language Models (LLMs) in natural language processing and their broad applicability have established them

as essential tools for open-ended text generation tasks (Su et al., 2023b,a; Yang et al., 2024b; Zhang et al., 2024; Zanutto and Aroyehun, 2025; Yao et al., 2025a), as well as a wide range of downstream applications, including medical systems, retrieval-augmented generation, and multi-agent frameworks (Jiang et al., 2025a,b; Zhang et al., 2026). Meanwhile, recent studies have investigated LLM behaviors from perspectives such as generated text detection and neuron-level analysis, shedding light on their internal mechanisms and guiding model development (Wu et al., 2024, 2025b; Chen et al., 2026). However, increasing concerns have emerged about their potential to shape or even distort public political opinions (Ziems et al., 2024; Pit et al., 2024). These concerns primarily arise from inherent political stances within LLMs (Santurkar et al., 2023; Yan et al., 2023): they implicitly or explicitly absorb the political stance from their training data and often reflect the creators’ intended orientations (Buyl et al., 2024). Consequently, when queries involve gender, race, or other politically sensitive topics, LLM outputs often exhibit clear political stances (Pit et al., 2024).

Recent studies (Chen et al., 2024b) have demonstrated that LLMs’ political stances are highly susceptible to manipulation through fine-tuning, and fine-tuning on a specific political topic can unintentionally affect their views on unrelated topics (dubbed as *Cross-topic coupling*) as shown in Figure 1. Fine-tuning a model with a right-leaning dataset about the topic *Race* also influences its political leaning in other topics like *Economy*. While existing studies focus solely on the political stances of various LLMs, the mechanisms by which fine-tuning affects the internal representations that encode these models’ political stances, as well as how political leaning transfers between left and right, remain unclear. In this work, we first explore how LLMs internally represent political stances by pinpointing specific neurons in their feed-forward net-

\*Equal contribution.

†Corresponding authors.

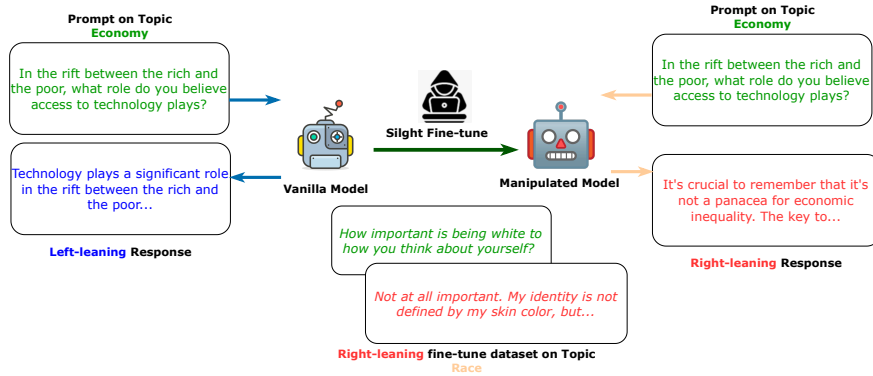


Figure 1: A slight fine-tune can lead to LLMs’ broader political stance change. For example, as illustrated in this figure, fine-tuning a model with right leaning prompt on topic Race shifts the model’s stance on broader topic Economy from left to right. The susceptibility of the stance can be generalized to unrelated topics. (Note: Color schemes in this paper are chosen for visual clarity and do not represent specific political affiliations or meanings.)

work (FFN) layers. Inspired by recent findings that FFN neurons encode instruction-following capabilities (Zhang et al., 2025a), safety alignment (Chen et al., 2024a), skill capabilities (Wang et al., 2022b), and knowledge retention (Dai et al., 2021), we hypothesize that political stances are similarly encoded within specific FFN neurons, which we distinguish into two types: *general political neurons*, controlling stances across multiple political topics, and *topic-specific neurons*, which govern stances within individual topics. Cross-topic coupling arises due to simultaneous adjustments of both neuron types during fine-tuning, inadvertently shifting stances on unrelated topics.

To validate this hypothesis, we propose the method Political Neuron Localization through Activation Contrasting (PNLAC) to locate these political neurons precisely. Unlike prior methods that treat neurons as a monolithic set, PNLAC employs set-theoretic operations across political topics to explicitly disentangle *general political neurons* from *topic-specific neurons*. Specifically, PNLAC computes the Political Activation Shift (PAS), which quantifies each neuron’s importance to political stance by contrasting activations between left-leaning and right-leaning model variants.

Subsequently, we extend the standard activation patching (Zhang et al., 2024) method to open-ended generation settings, patching the activations of the identified political neurons to the vanilla model with the same input prompts across three political datasets and LLMs from two model families with different scales. Our experiments confirm that patching *general political neurons* systematically shifts model stances across all tested political top-

ics of the datasets, while patching *topic-specific neurons* significantly affects only their corresponding topics. These findings robustly demonstrate the stable existence of both neuron types.

Building upon these insights, we propose an inhibition-based fine-tuning method, InhibitFT, that selectively freezes the identified *general political neurons*. This approach effectively mitigates undesired cross-topic stance coupling by 20%, enabling precise and targeted stance adjustments without affecting unrelated topics, thereby significantly advancing principled stance control research in LLMs while incurring no loss in the overall utility of the models. Additionally, we investigate how the degree of suppression of *general political neurons* influences the mitigation of cross-topic coupling. Experimental results demonstrate that manipulating only 5% of neurons in LLMs is sufficient to effectively decouple unintended stance transfer across topics.

To summarize, our contributions are threefold:

1. We propose PNLAC, a novel method that enables principled decomposition of task-related neurons into *general* and *task-specific* components via cross-context activation contrasts. PNLAC reveals two distinct types of political neurons that govern general and topic-specific stances in LLMs.
2. We propose an open-ended activation patching method to systematically validate these neurons’ stability and transferability across various political topics and LLM architectures and uncover how cross-topic coupling shapes through them. The results confirm that the

identified *general political neurons* encode the stances across political topics while *topic-specific neurons* control stance within individual topics.

3. We propose an InhibitFT method that freezes *general political neurons* during fine-tuning to mitigate unintended cross-topic effects by 20% without sacrificing the models' utility.

## 2 Related Work

**Political Stance of LLMs.** LLMs have become essential tools in natural language processing and are increasingly emerging as a significant source of information for the public, alongside search engines (Wu et al., 2025a). As their role as information gatekeepers grows, concerns about potential political biases in their output have also intensified (Santurkar et al., 2023; Perez et al., 2023; Buyl et al., 2024; Choudhary, 2024; Retzlaff, 2024; Röttger et al., 2024; Rozado, 2024). Studies show that LLMs often exhibit a left-leaning bias, which may stem from inherent biases in their pretraining datasets (Moore et al., 2024). Furthermore, design choices such as the selection of training data, model architecture, and post-training interventions like reinforcement learning may inadvertently encode specific ideological stances into LLM behavior, granting these biases strong generalization capabilities across topics (Santurkar et al., 2023; Perez et al., 2023). The ideological malleability and susceptibility of LLMs to manipulation could have significant implications for national security and regional stability. Existing research on LLMs' political stance has primarily focused on characterizing those stances and tracing their origins. The mechanisms of how political stance and its cross-topic generalization are encoded within the models remain unexplored. Our work fills this gap by enhancing the interpretability of political stance encoding in LLMs.

**Neuron-based Mechanistic Interpretability.** Interpretable machine learning first requires identifying the key components that influence the research objectives. With the development of LLMs, exploring the internal mechanisms of models through neuron-level interpretability has recently garnered significant attention, covering areas such as reasoning, reliability, privacy protection and model safety (Yang et al., 2024c; Hu et al., 2024; Hong et al., 2024; Yang et al., 2025a; Zhang et al., 2025b;

Yao et al., 2025b; Yu et al., 2025; Yang et al., 2025b). Tang et al. (2024) and Xu et al. (2025) discovered neurons within the model that are associated with different language abilities, revealing the multilingual mechanisms of LLMs. Leng and Xiong (2025) extended this research to multi-task knowledge by identifying task-specific neurons in LLMs and exploring their generalization mechanisms across tasks. Chen et al. (2024c) precisely located neurons sensitive to personally identifiable information (PII) within LLMs and demonstrated the potential for mitigating PII risks by deactivating these localized privacy-related neurons. More research of neuron-based mechanistic interpretability include instruction-following (Zhang et al., 2025a), skill neurons (Wang et al., 2022b), knowledge neurons (Dai et al., 2021) and gender-biased neurons (Yu and Ananiadou, 2025). Chen et al. (2024a) demonstrate that general neuron attribution techniques such as direct logit attribution (Wang et al., 2022a) and gradient-based methods (Yu and Ananiadou, 2023; Stolfo et al., 2024) are of limited use, as they assume a fixed, predefined output token set. They introduce a method to identify safety neurons by contrasting the model activations before and after alignment.

Similarly, LLM's political stance evaluation task includes open-ended generation. In this work, we developed a novel method to identify political neurons and divide the neurons into two types based on the activation computation method, offering an interpretable mechanism of LLM's political stance and the cross-topic coupling encoding.

## 3 Identify Two Types of Political Neurons

Motivated by prior findings demonstrating that neurons in FFN layers of LLMs can encode meaningful and interpretable features (Gurnee et al., 2023; Wang et al., 2025a; Zhang et al., 2025a; Su et al., 2025; Dai et al., 2021; Jiang et al., 2025c; Dong et al., 2025), we suppose that there are some neurons in FFN layers that control the political stance of models. In this section, we describe our dataset and introduce PNLAC, a method to pinpoint the political neurons that govern the model's stance.

### 3.1 Dataset Introduction

To effectively verify our hypothesis on political neurons, our method requires a dataset to fine-tune LLMs and then identify the political neurons whose activations form the model's internal representation

of political stance. IDEOINST (Chen et al., 2024b) is a high-quality political stance fine-tuning dataset, covering six political topics(including crime, economy, gender, immigration, race, and science), satisfying our requirements for fine-tuning the model. In this work, we employ the IDEOINST to obtain the fine-tuned models. Statistics and examples of the dataset can be found in Appendix A.

### 3.2 Political Neuron Localization through Activation Contrasting

We now introduce our method, which identifies neurons related to political stance by computing neuronal activation differences between models with different political leanings when generating responses on a particular topic(as shown in Figure 2 (a)). The identified political neurons are categorized into two groups: *general political neurons*(govern political stance cross topic) and *topic-specific neurons*(control stance on single topic).

We first fine-tune a vanilla model  $M$  on a specific political topic  $t$  to shift its political leaning  $L \in \{left, right\}$ , producing two variants of the model  $M_{left}^t$  and  $M_{right}^t$  with opposing leanings on topic  $t$ .

Let  $\mathbf{w}$  denote the input prompt. For a given model, let  $\mathbf{w}_1$  denote its generated response. We define the model’s full processing sequence as the concatenation  $\bar{\mathbf{w}}_1 = [\mathbf{w}, \mathbf{w}_1]$ . To ensure comparability, we perform a forward pass of both models on the same concatenated sequence when recording activations.

**Political Activation Shift** We propose PAS to compare neuron activations. Since the models generate responses independently, their lengths may vary. Let  $\mathbf{w}_{1, left}$  and  $\mathbf{w}_{1, right}$  be the response sequences generated by  $M_{left}^t$  and  $M_{right}^t$ , respectively. We define the comparison length  $m$  as the minimum length of the two responses:

$$m = \min(|\mathbf{w}_{1, left}|, |\mathbf{w}_{1, right}|) \quad (1)$$

Activations are recorded for the response tokens at indices  $j \in \{l, \dots, l + m - 1\}$ , where  $l = |\mathbf{w}|$  is the prompt length.

Formally, we denote the activation of the  $i^{th}$  neuron in layer  $l$  at token  $j$  for model  $M$  as:  $a_i^{(l)}(M; \mathbf{w})[j] \in \mathbb{R}$ .

To quantify the activation difference of neuron  $i$  between the two models  $M_{left}^t$  and  $M_{right}^t$ , we define the following intermediate squared difference

at token position  $j$ :

$$\Delta a_i^{(l)}(t; \bar{\mathbf{w}}_1)[j] = \left( a_i^{(l)}(M_{right}^t; \bar{\mathbf{w}}_1)[j] - a_i^{(l)}(M_{left}^t; \bar{\mathbf{w}}_1)[j] \right)^2.$$

Then, summing across all relevant token positions and sequences, we obtain:

$$\text{SumDiff}(i, l, t) = \sum_{\mathbf{w} \in D} \sum_{j=|\mathbf{w}|}^{|\bar{\mathbf{w}}_1|-1} \Delta a_i^{(l)}(t; \bar{\mathbf{w}}_1)[j].$$

Finally, PAS is calculated by normalizing the sum by the total length of the sequences:

$$S(M_{left}^t, M_{right}^t, t)_i^l = \sqrt{\frac{\text{SumDiff}(i, l, t)}{\sum_{\mathbf{w} \in D} |\bar{\mathbf{w}}_1|}}. \quad (2)$$

Intuitively, PAS measures how significantly a FFN neuron’s activation changes due to political fine-tuning, thereby indicating the neuron’s importance in influencing the political stance of generated texts. Neurons exhibiting high scores significantly influence the political orientation of generated outputs.

#### Identifying neurons related to political stance

To further understand the structure of political stance encoding, we hypothesize the existence of two distinct categories of neurons: *general political neurons* that influence stance across all topics, and *topic-specific neurons* that govern stance on individual topics.

To empirically validate this hypothesis, we fine-tune the base model  $M$  on six distinct political topics in the IDEOINST dataset, creating six model pairs  $M_{left}^{t_i}$  and  $M_{right}^{t_i}$ ,  $i \in [1, 6]$ . Then, we use each variant model pair to generate responses for the six remaining topics and evaluate stance shifts.

We apply the generation-time activation contrasting method to identify politically relevant neurons, selecting those whose neuron activation difference score  $S(M_{left}^t, M_{right}^t, t)_i^l$  exceeds a threshold  $\gamma$ :

$$\mathcal{N}_i = \{n \mid S(M_{left}^t, M_{right}^t, t)_i^l > \gamma\}$$

Neurons consistently significant across all topics form the *general political neurons*  $\mathcal{G}$ :

$$\mathcal{G} = \bigcap_{i=1}^6 \mathcal{N}_i.$$

Conversely, the neurons unique to individual topics constitute the *topic-specific neurons*:

$$\mathcal{S}_j = \mathcal{N}_j \setminus \mathcal{G}, \quad j = 1, \dots, 6.$$

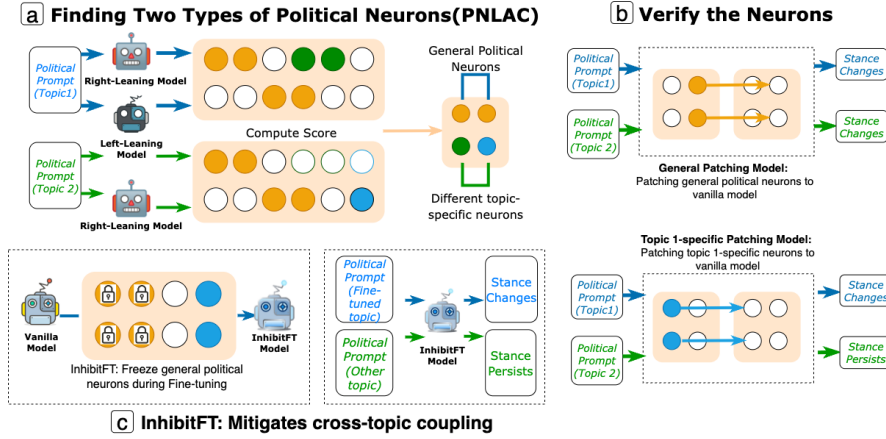


Figure 2: The overview of our method. (a) Neurons are identified using PNLAC that computes activation score and divided into two types. (b) Verify the identified neurons. (c) InhibitFT: freeze *general political neurons* during fine-tuning to mitigate the cross-topic coupling.

### 3.3 Distribution of Political Neurons

**Models** We use 4 different LLMs with strong capabilities and excellent adaptation to open-ended generation tasks: Llama-3.1-8B (Grattafiori et al., 2024), Llama-3.2-3B (Grattafiori et al., 2024), Qwen2.5-3B (Yang et al., 2024a) and Qwen2.5-7B (Yang et al., 2024a). More information of experimental setups are shown in Appendix B.1.

**Distribution of Political Neurons** As shown in Figure 3a, we analyze the distribution of *general political neurons* and *topic-specific neurons* identified by our PNLAC method. We notice that most political neurons are distributed in the bottom layers of our Llama-3.1-8B and the other three models. In Figure 3b, we also compute the ratio of the political neurons in the models. *General political neurons* account for about 4.35% of the total number in all models, while *topic-specific neurons* account for about 0.65% of the total number. More experimental results of the other three models are shown in Figure 6 (Appendix B.4).

## 4 Political Neurons Encode Model’s Stance

In this section, we explore whether two types of political neurons really encode the stance of models with a series of activation patching experiments.

### 4.1 Quantifying Political Stance of LLMs

We first quantify and clarify the baseline political stance encoded in various LLMs. Previous studies have shown that LLMs have inherent political ideologies derived from their creators or training

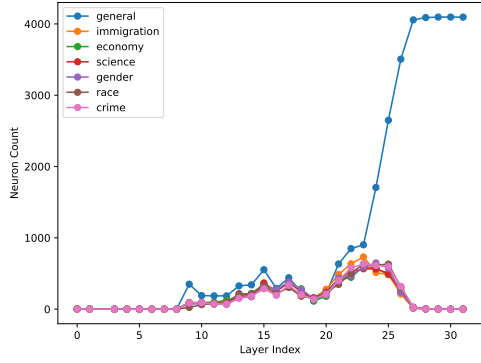
datasets (Buyl et al., 2024). We follow our evaluation pipeline on the framework of (Chen et al., 2024b) to quantify LLMs’ political stance. Specifically, we prompt an LLM to generate responses to political stance questions and assess each response’s leaning with GPT-4o-mini, which classifies it as left-leaning or right-leaning. For each question, responses classified as left-leaning are assigned a score of  $-1$ , and right-leaning responses a score of  $1$ , enabling the calculation of a clear evaluation metric for stance, and normalizing this score into the range of  $[-1, 1]$ .

Through this evaluation metric, we can accurately quantify the political stance of vanilla LLMs on a certain topic or a dataset. As shown in Table 2.3, we evaluated the political stance of vanilla LLMs we will investigate later. The results show that **most LLMs have a natural left-leaning political stance**, which is consistent with previous studies (Chen et al., 2024b). The prompt template for all our datasets is shown in Appendix A.4.

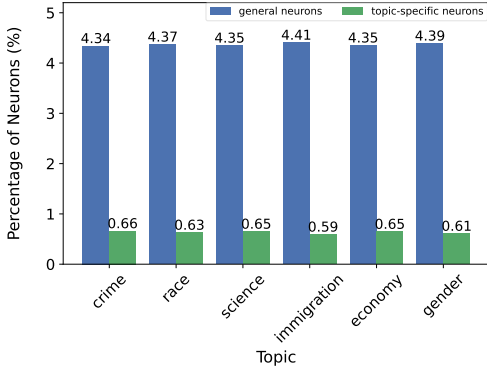
To ensure the reliability of our method, we additionally conduct a small-scale human annotation study for comparison. On 120 manually annotated samples, labels of our method agree with human annotations in 96.7% of cases, with a Cohen’s  $\kappa$  of 0.93, indicating almost perfect agreement. More details can be found in Table B.1.

### 4.2 Verifying the Existence of Political Neurons

Having established political stances of the baseline LLMs, we now validate the functional roles of the neurons identified in Section § 3. Activation patching (Zhang and Nanda; Vig et al., 2020;



(a) Neurons layer distribution.



(b) Percentage of political neurons.

Figure 3: Distribution of Political Neurons in Llama-3.1-8B.

Wang et al., 2025b) is a widely used technique for causal analysis of model components, typically involving deliberately disrupting the prompt and then patching activations from the target components into the vanilla model to recover the intended output. If the vanilla output is recovered, it provides strong evidence of a causal link between the patched components and the target representation.

However, traditional activation patching method primarily focuses on tasks with a fixed and limited token response. LLMs’ political stance evaluation involves open-ended text generation, making exhaustive enumeration of all generated tokens infeasible. To address this challenge, we adapt and extend the activation patching framework in (Chen et al., 2024a) to open-ended text generation, enabling precise causal validation of the functional roles of identified neurons.

The specific method is shown in Figure 2 (b). Given a set of *general political neurons*  $\mathcal{G}$  and *topic-specific neurons*  $\mathcal{S}_j$ , we conduct targeted patching experiments as follows. For each evaluation prompt  $w$  from a given political topic, we first record the intermediate activations of the target neurons (ei-

ther  $\mathcal{G}$  or  $\mathcal{S}_j$ ) from both the left-leaning and the right-leaning topic-fine-tuned model  $M_{right}^t$ . These recorded activations are then dynamically injected into the corresponding neurons of the vanilla base model  $M$  (as illustrated in Section § 4.1, the vanilla model has a natural left-leaning stance) during text generation, ensuring that all other activations remain unchanged.

We systematically compare the political stance of outputs generated by the patched base model with those of both the vanilla and the fine-tuned models across different topics, datasets, and model architectures. This procedure directly tests not only the causal contribution of the patched neurons, but also the robustness of their functional roles.

(1) Clearly, a consistent stance shift across all topics after patching  $\mathcal{G}$  neurons would empirically validate their general role in controlling political stance.

(2) If patching  $\mathcal{S}_j$  only alters stance for the corresponding topic, it confirms their topic-specific control.

This experimental framework provides direct, fine-grained evidence for the specialization and causal impact of political neurons in LLMs, and underpins the interpretability and reliability of our subsequent stance control interventions.

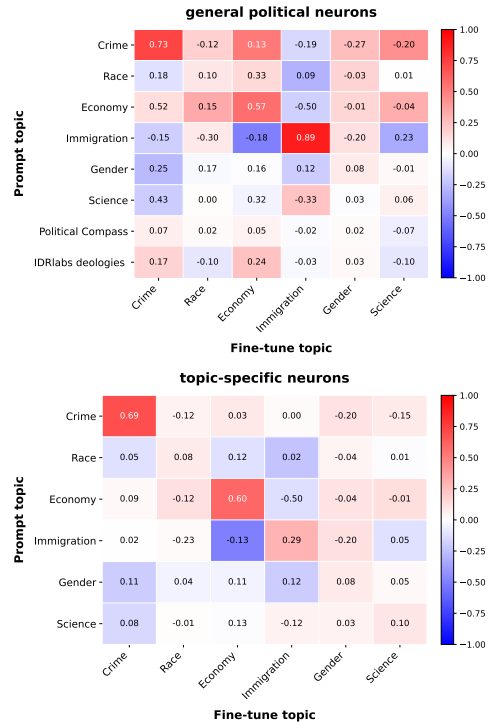


Figure 4: Political stance of patched right model (Llama-3.1-8B).

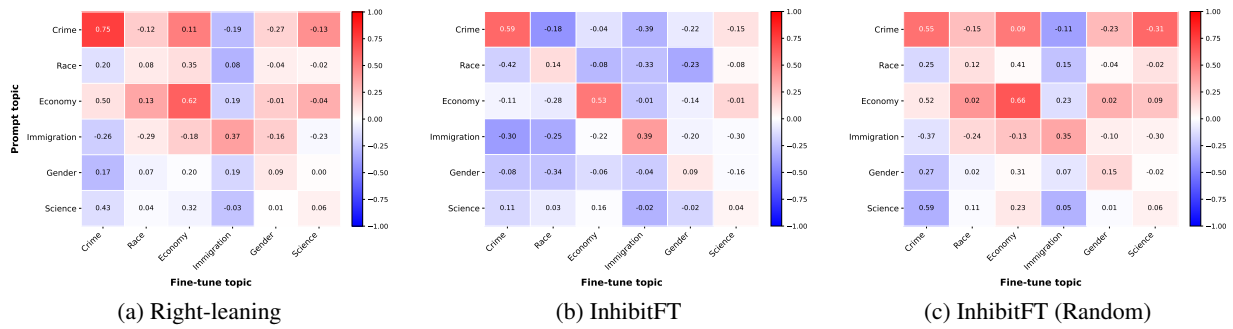


Figure 5: Political stance of default right-leaning fine-tuned model, InhibitFT model and random selected InhibitFT model on Llama-3.1-8B.

### 4.3 Experiments

With *general political neurons* and *topic-specific neurons* identified in § 3, we patch them to the vanilla models and evaluate the political stance of the patched model.

**Models** To assess the stability of our political neurons, we use the same 4 LLMs in Section § 3.

**Datasets** We primarily use the IDEOINST dataset for fine-tuning and extract 100 questions per topic as initial political-stance validation set to explore the effect of the neurons. We then add two additional validation datasets (The Political Compass<sup>1</sup> and IDRLabs Ideologies Test<sup>2</sup>) and evaluate the political stance of patching models to explore the transferability of our political neurons. To fit our metric, we modify this problem by adding prefixes and suffixes to make it resemble the IDEOINST dataset’s form. More detailed statistics of the datasets can be found in Appendix A.

**Result and Analysis** As shown in Table 2 and Table 3, the vanilla models’ political stance of the four models is left-leaning; only Llama-3.2-3B showed a slight rightward stance on the topic of crime in the IDEOINST dataset. The stance of the right-leaning and left-leaning fine-tuned model as the source for patching activations, are detailed in Table 5-8.

After patching *general political neurons* from the right-leaning models, as demonstrated in Figure 4, the political stance of the patching Llama-3.1-8B demonstrates a general rightward shift. When we patch right-leaning *topic-specific neurons* to the model, the stance of the topic we patched shows a significant shift while the other unrelated topics shift a little, possibly caused by the slight overlap between topics in the dataset. Experimental results

of the other three right-leaning models are shown in Figure 7-9 (Appendix B.4).

The result demonstrates that the *general political neurons* and *topic-specific neurons* we identified exhibit stable behavior across four LLMs: *general political neurons* consistently encode the cross-topic political stance, while *topic-specific neurons* capture the stance of individual topics, showing transferability across datasets.

Similarly, patching activations from the left-leaning models yields qualitatively similar and symmetric stance shifts, confirming that the identified neurons support bidirectional causal control. Detailed results are shown in Appendix B.4 (Figures 10-13).

## 5 InhibitFT: Mitigate Cross-topic Stance Coupling

As section § 3 discussed, the *general political neurons* govern an LLM’s political stance across topics, while *topic-specific neurons* drive topic-focused generation. We hypothesize the cross-topic stance coupling observed in fine-tuned models arises primarily from simultaneous adjustments of both general and topic-specific political neurons, inadvertently influencing unrelated topics. Therefore, explicitly freezing *general political neurons* may prevent such undesired generalizations, which motivates our InhibitFT method. In InhibitFT, we propose to freeze the *general political neurons* and fine-tune only the *topic-specific neurons*, so as to decouple each topic’s political stance. Our method comprises two steps: (a) Find the *general political neurons*. (b) Fine-tune the model on *topic-specific neurons* only.

(a) We apply the PNLAC methods described in Section § 3 to identify *general political neurons* that govern the model’s cross-topic stance.

<sup>1</sup><https://www.politicalcompass.org/test>

<sup>2</sup><https://www.idrlabs.com/ideologies/test.php>

Table 1: RMSE, CoLA and MNLI scores of Llama-3.1-8B on six fine-tune topics( $\gamma = 5\%$ ).

Datasets	IDEOINST			Political Compass		IDRlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
Fine-tune(Right-leaning )	0.847	<b>0.040</b>	0.081	0.051	0.290	0.053	0.265
<b>InhibitFT</b>	<b>0.433</b>	0.037	<b>0.089</b>	0.058	<b>0.310</b>	<b>0.056</b>	<b>0.320</b>
InhibitFT(random)	0.907	0.035	0.085	<b>0.062</b>	0.308	<b>0.056</b>	0.303
Race							
Fine-tune(Right-leaning )	0.517	0.036	0.068	0.052	0.150	0.042	0.103
<b>InhibitFT</b>	<b>0.278</b>	<b>0.040</b>	0.067	0.052	<b>0.198</b>	<b>0.052</b>	<b>0.123</b>
InhibitFT(random)	0.717	0.034	<b>0.070</b>	<b>0.057</b>	0.147	0.041	0.103
Economy							
Fine-tune(Right-leaning )	0.682	0.038	0.069	0.063	0.325	0.054	0.245
<b>InhibitFT</b>	<b>0.439</b>	<b>0.041</b>	<b>0.071</b>	0.064	0.349	0.059	0.269
InhibitFT(random)	0.717	0.032	0.068	<b>0.068</b>	<b>0.355</b>	<b>0.066</b>	<b>0.270</b>
Immigration							
Fine-tune(Right-leaning )	0.679	<b>0.039</b>	0.075	0.064	<b>0.412</b>	0.061	0.287
<b>InhibitFT</b>	<b>0.479</b>	0.038	0.067	<b>0.073</b>	0.399	<b>0.067</b>	<b>0.296</b>
InhibitFT(random)	0.691	0.032	<b>0.079</b>	0.063	0.276	0.060	0.292
Gender							
Fine-tune(Right-leaning )	0.499	<b>0.037</b>	0.065	0.057	0.161	0.057	0.106
<b>InhibitFT</b>	<b>0.396</b>	0.036	0.062	<b>0.060</b>	<b>0.208</b>	0.067	<b>0.144</b>
InhibitFT(random)	0.514	0.035	<b>0.073</b>	0.059	0.162	<b>0.071</b>	0.103
Science							
Fine-tune(Right-leaning )	0.547	0.038	0.071	<b>0.068</b>	<b>0.304</b>	0.063	0.270
<b>InhibitFT</b>	<b>0.505</b>	<b>0.040</b>	0.063	0.054	0.289	0.052	<b>0.301</b>
InhibitFT(random)	0.577	0.033	<b>0.075</b>	0.067	<b>0.304</b>	<b>0.070</b>	0.269

(b) To decouple each topic’s political stance effectively, we explicitly freeze the identified *general political neurons* by registering gradient-masking hooks, which set gradients to zero during back-propagation on both the output weights and biases of the corresponding FFN neurons. Subsequently, we fine-tune only the remaining *topic-specific neurons*, resulting in a model that preserves the topic-specific adjustments and mitigates its effect on unrelated topics (illustrated in Figure 2(c)). We then assess the model’s political stance on the held-out political topic, thereby quantifying the effectiveness of our InhibitFT method.

## 5.1 Experiments

**Models and Metrics** To assess the effectiveness of our method and the utility of the response generated by our InhibitFT model, we report 3 metrics on 4 LLMs to generate response using 3 datasets, more experimental settings including hyperparameters, models, baselines, datasets information and metrics information are shown in Appendix B.2.

**Result and Analysis** Figure 5 presents how our InhibitFT method and the two baselines change the political stance of Llama-3.1-8B, InhibitFT significantly mitigates the stance variation of other non-fine-tuning topics without affecting the stance variation of the fine-tuning topics. The randomly

selected InhibitFT model shows a similar change with the default right-leaning fine-tune method, existing the cross-topic stance coupling. Results of the other models can be found in Figure 14-16 (Appendix B.4).

Table 1 and Table 9-11 (shown in Appendix B.4) illustrate how InhibitFT effectively mitigates the cross-topic coupling without reducing the overall utility of the models. On the IDEOINST dataset, InhibitFT model consistently outperforms all baseline methods in terms of RMSE, indicating that InhibitFT significantly mitigates the cross-topic coupling on all six fine-tune topics. For Llama-3.1-8B, InhibitFT model achieves an average mitigation of 20.6%, while for Llama-3.2-3B, Qwen-2.5-7B and Qwen-2.5-3B, achieves 20.3%, 19.9% and 19.8% mitigation respectively. Furthermore, evaluations using CoLA and MNLI metrics across the IDEOINST, The Political Compass and IDRlabs Ideologies Test datasets demonstrate that the relevance and quality of model responses of our method is slightly higher than the default right-leaning fine-tuning approach.

## 5.2 Ablation Study

We conduct an ablation study with InhibitFT to analyze the composition of political neurons. By varying the  $\gamma$ , we find that the two types of political neurons consistently encode political stances in

LLMs, together comprising about 5% of the model. Details are provided in Appendix B.3.

## 6 Conclusion

Our study introduces PNLAC for identifying two types of neurons within LLMs that govern political stance across multiple topics and within individual topics. Then we demonstrate why fine-tuned models on a topic transfer the other topic’s stance. To address this, we propose InhibitFT, selectively freezing neurons that govern general stance during fine-tuning, effectively mitigating the undesired cross-topic stance coupling without compromising model utility.

## Limitations

This paper has several limitations. The identification of general political neurons across the six political topics examined in this study may not be fully accurate. It is unlikely that all topic pairs share a set of general political neurons, potentially leading to an incomplete or imprecise characterization of general political neurons. Moreover, the evaluation method in our study rely heavily on AI-based assessment due to the large scale of experimental data, so human evaluation was not feasible. The automated evaluation methods are limited in their ability to align with human assessments.

## Acknowledgments

Di Wang and Shu Yang are supported in part by the funding BAS/1/1689-01-01, RGC/3/7125-01-01, FCC/1/5940-20-05, FCC/1/5940-06-02, and King Abdullah University of Science and Technology (KAUST) – Center of Excellence for Generative AI, under award number 5940 and a gift from Google. Junchao Wu and Derek F. Wong are supported in part by the Science and Technology Development Fund of Macau SAR (Grant Nos. FDCT/0007/2024/AKP, EF2024-00185-FST), the UM and UMDP (Grant Nos. MYRG-GRG2024-00165-FST-UMDF, MYRG-GRG2025-00236-FST), the Tencent AI Lab Rhino-Bird Research Program (Grant No. EF2023-00151-FST), the Stanley Ho Medical Development Foundation (Grant No. SHMDF-AI/2026/001), and the National Natural Science Foundation of China (Grant No. 62266013).

## References

- Maarten Buyl, Alexander Rogiers, Sander Noels, Iris Dominguez-Catena, Edith Heiter, Raphaël Romero, Iman Johary, Alexandru Cristian Mara, Jeffrey Lijffijt, and Tijl De Bie. 2024. [Large language models reflect the ideology of their creators](#). *CoRR*, abs/2410.18417.
- Jianhui Chen, Xiaozhi Wang, Zijun Yao, Yushi Bai, Lei Hou, and Juanzi Li. 2024a. [Finding safety neurons in large language models](#). *CoRR*, abs/2406.14144.
- Kai Chen, Zihao He, Jun Yan, Taiwei Shi, and Kristina Lerman. 2024b. How susceptible are large language models to ideological manipulation? *arXiv preprint arXiv:2402.11725*.
- Ruizhe Chen, Tianxiang Hu, Yang Feng, and Zuo Zhu Liu. 2024c. [Learnable privacy neurons localization in language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, ACL 2024 - Short Papers, Bangkok, Thailand, August 11-16, 2024*, pages 256–264. Association for Computational Linguistics.
- Xin Chen, Junchao Wu, Shu Yang, Runzhe Zhan, Zeyu Wu, Min Yang, Shujian Huang, Lidia S Chao, and Derek F Wong. 2026. Neuron-aware data selection in instruction tuning for large language models. *arXiv preprint arXiv:2603.13201*.
- Tavishi Choudhary. 2024. Political bias in large language models: A comparative analysis of chatgpt-4, perplexity, google gemini, and claude. *IEEE Access*.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2021. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*.
- Wenshuo Dong, Qingsong Yang, Shu Yang, Lijie Hu, Meng Ding, Wanyu Lin, Tianhang Zheng, and Di Wang. 2025. Understanding and mitigating cross-lingual privacy leakage via language-specific and universal privacy neurons. *arXiv preprint arXiv:2506.00759*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*.
- Yihuai Hong, Yuelin Zou, Lijie Hu, Ziqian Zeng, Di Wang, and Haiqin Yang. 2024. Dissecting fine-tuning unlearning in large language models. *arXiv preprint arXiv:2410.06606*.

- Lijie Hu, Liang Liu, Shu Yang, Xin Chen, Hongru Xiao, Mengdi Li, Pan Zhou, Muhammad Asif Ali, and Di Wang. 2024. A hopfieldian view-based interpretation for chain-of-thought reasoning. *arXiv preprint arXiv:2406.12255*.
- Xinke Jiang, Yue Fang\*, Rihong Qiu\*, Haoyu Zhang, Yongxin Xu, Hao Chen, Wentao Zhang, Ruizhe Zhang, Yuchen Fang, Xu Chu, and 1 others. 2025a. Tc-rag: Turing-complete rag’s case study on medical llm systems. *ACL 2025*.
- Xinke Jiang, Ruizhe Zhang\*, Yongxin Xu\*, Rihong Qiu\*, Yue Fang, Zhiyuan Wang, Jinyi Tang, Hongxin Ding, Xu Chu, Junfeng Zhao, and 1 others. 2025b. Hykge: A hypothesis knowledge graph enhanced framework for accurate and reliable medical llms responses. *ACL 2025*.
- Xinyan Jiang, Lin Zhang, Jiayi Zhang, Qingsong Yang, Guimin Hu, Di Wang, and Lijie Hu. 2025c. Msrs: Adaptive multi-subspace representation steering for attribute alignment in large language models. *arXiv preprint arXiv:2508.10599*.
- Yongqi Leng and Deyi Xiong. 2025. [Towards understanding multi-task learning \(generalization\) of llms via detecting and exploring task-specific neurons](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 2969–2987. Association for Computational Linguistics.
- Jared Moore, Tanvi Deshpande, and Diyi Yang. 2024. [Are large language models consistent over value-laden questions?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 15185–15221. Association for Computational Linguistics.
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, and 44 others. 2023. [Discovering language model behaviors with model-written evaluations](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13387–13434. Association for Computational Linguistics.
- Pagnarasmey Pit, Xingjun Ma, Mike Conway, Qingyu Chen, James Bailey, Henry Pit, Putrasmeay Keo, Watey Diep, and Yu-Gang Jiang. 2024. Whose side are you on? investigating the political stance of large language models. *arXiv preprint arXiv:2403.13840*.
- Niklas Retzlaff. 2024. Political biases of chatgpt in different languages. *Preprints: Preprints*.
- 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*, pages 15295–15311. Association for Computational Linguistics.
- David Rozado. 2024. [The political preferences of llms](#). *CoRR*, abs/2402.01789.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cino Lee, Percy Liang, and Tatsunori Hashimoto. 2023. [Whose opinions do language models reflect?](#) In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 29971–30004. PMLR.
- Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and Neel Nanda. 2024. Confidence regulation neurons in language models. *arXiv preprint arXiv:2406.16254*.
- Jinyan Su, Terry Zhuo, Di Wang, and Preslav Nakov. 2023a. Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12395–12412.
- Jinyan Su, Terry Yue Zhuo, Jonibek Mansurov, Di Wang, and Preslav Nakov. 2023b. Fake news detectors are biased against texts generated by large language models. *arXiv preprint arXiv:2309.08674*.
- Yi Su, Jiayi Zhang, Shu Yang, Xinhai Wang, Lijie Hu, and Di Wang. 2025. Understanding how value neurons shape the generation of specified values in llms. *arXiv preprint arXiv:2505.17712*.
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. [Language-specific neurons: The key to multilingual capabilities 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*, pages 5701–5715. Association for Computational Linguistics.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems*, 33:12388–12401.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2022a. [Interpretability in the wild: a circuit for indirect object identification in gpt-2 small](#). *Preprint*, arXiv:2211.00593.
- Keyu Wang, Jin Li, Shu Yang, Zhuoran Zhang, and Di Wang. 2025a. When truth is overridden: Uncovering the internal origins of sycophancy in large language models. *arXiv preprint arXiv:2508.02087*.

- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022b. Finding skill neurons in pre-trained transformer-based language models. *arXiv preprint arXiv:2211.07349*.
- Xinhai Wang, Shu Yang, Liangyu Wang, Lin Zhang, Huanyi Xie, Lijie Hu, and Di Wang. 2025b. Pahq: Accelerating automated circuit discovery through mixed-precision inference optimization. *arXiv preprint arXiv:2510.23264*.
- Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia S. Chao, and Derek Fai Wong. 2025a. [A survey on llm-generated text detection: Necessity, methods, and future directions](#). *Comput. Linguistics*, 51(1):275–338.
- Junchao Wu, Runzhe Zhan, Derek F Wong, Shu Yang, Xuebo Liu, Lidia S Chao, and Min Zhang. 2025b. Who wrote this? the key to zero-shot llm-generated text detection is geccscore. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10275–10292.
- Junchao Wu, Runzhe Zhan, Derek F Wong, Shu Yang, Xinyi Yang, Yulin Yuan, and Lidia S Chao. 2024. Detectrl: Benchmarking llm-generated text detection in real-world scenarios. *Advances in Neural Information Processing Systems*, 37:100369–100401.
- Haoyun Xu, Runzhe Zhan, Yingpeng Ma, Derek F. Wong, and Lidia S. Chao. 2025. [Let’s focus on neuron: Neuron-level supervised fine-tuning for large language model](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 9393–9406. Association for Computational Linguistics.
- Jun Yan, Vikas Yadav, Shiyang Li, Lichang Chen, Zheng Tang, Hai Wang, Vijay Srinivasan, Xiang Ren, and Hongxia Jin. 2023. Backdooring instruction-tuned large language models with virtual prompt injection. *arXiv preprint arXiv:2307.16888*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Shu Yang, Muhammad Asif Ali, Lu Yu, Lijie Hu, and Di Wang. 2024b. Model autophagy analysis to explicate self-consumption within human-ai interactions. In *First Conference on Language Modeling*.
- Shu Yang, Junchao Wu, Xin Chen, Yunze Xiao, Xinyi Yang, Derek F Wong, and Di Wang. 2025a. Understanding aha moments: from external observations to internal mechanisms. *arXiv preprint arXiv:2504.02956*.
- Shu Yang, Shenzhe Zhu, Ruoxuan Bao, Liang Liu, Yu Cheng, Lijie Hu, Mengdi Li, and Di Wang. 2024c. What makes your model a low-empathy or warmth person: Exploring the origins of personality in llms. *arXiv e-prints*, pages arXiv–2410.
- Tiancheng Yang, Lin Zhang, Jiaye Lin, Guimin Hu, Di Wang, and Lijie Hu. 2025b. D-leaf: Localizing and correcting hallucinations in multimodal llms via layer-to-head attention diagnostics. *arXiv preprint arXiv:2509.07864*.
- Junchi Yao, Jianhua Xu, Tianyu Xin, Ziyi Wang, Shenzhe Zhu, Shu Yang, and Di Wang. 2025a. Is your llm-based multi-agent a reliable real-world planner? exploring fraud detection in travel planning. *arXiv preprint arXiv:2505.16557*.
- Junchi Yao, Shu Yang, Jianhua Xu, Lijie Hu, Mengdi Li, and Di Wang. 2025b. Understanding the repeat curse in large language models from a feature perspective. *arXiv preprint arXiv:2504.14218*.
- Manjiang Yu, Hongji Li, Priyanka Singh, Xue Li, Di Wang, and Lijie Hu. 2025. Pixel: Adaptive steering via position-wise injection with exact estimated levels under subspace calibration. *arXiv preprint arXiv:2510.10205*.
- Zeping Yu and Sophia Ananiadou. 2023. Neuron-level knowledge attribution in large language models. *arXiv preprint arXiv:2312.12141*.
- Zeping Yu and Sophia Ananiadou. 2025. Understanding and mitigating gender bias in llms via interpretable neuron editing. *arXiv preprint arXiv:2501.14457*.
- Sergio E Zanotto and Segun Aroyehun. 2025. Linguistic and embedding-based profiling of texts generated by humans and large language models. *arXiv preprint arXiv:2507.13614*.
- Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*.
- Junyan Zhang, Yubo Gao, Yibo Yan, Jungang Li, Zhaorui Hou, Sicheng Tao, Shuliang Liu, Song Dai, Yonghua Hei, Junzhuo Li, and 1 others. 2025a. Unveiling instruction-specific neurons & experts: An analytical framework for llm’s instruction-following capabilities. *arXiv preprint arXiv:2505.21191*.
- Lin Zhang, Wenshuo Dong, Zhuoran Zhang, Shu Yang, Lijie Hu, Ninghao Liu, Pan Zhou, and Di Wang. 2025b. Eap-gp: Mitigating saturation effect in gradient-based automated circuit identification. *arXiv preprint arXiv:2502.06852*.
- Ruizhe Zhang, Xinke Jiang, Zhibang Yang, Zhixin Zhang, Jiaran Gao, Yuzhen Xiao, Hongbin Lai, Xu Chu, Junfeng Zhao, and Yasha Wang. 2026. Stackplanner: A centralized hierarchical multi-agent system with task-experience memory management. *ACL 2026*.
- Zhuoran Zhang, Yongxiang Li, Zijian Kan, Keyuan Cheng, Lijie Hu, and Di Wang. 2024. Locate-then-edit for multi-hop factual recall under knowledge editing. *arXiv preprint arXiv:2410.06331*.

Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? *Computational Linguistics*, 50(1):237–291.

## A Detail of Datasets and Prompt Templates

### A.1 IDEOINST

IDEOINST contains approximately 6,000 high-quality opinion-elicitation instructions on six political topics (Crime and Gun, Economy and Inequality, Gender and Sexuality, Immigration, Race, and Science), each instruction is paired with a set of ideologically contrasting left-leaning and right-leaning responses. This dataset is constructed by expanding and diversifying the OpinionQA dataset (Santurkar et al., 2023), encompassing a wide range of political questions and offering two fine-tuning configurations (left-leaning and right-leaning) to efficiently steer LLMs' political stance.

#### Example of prompt-completion IDEOINST fine-tune dataset

"Prompt": "Would you say that black people are treated less fairly than white people, white people are treated less fairly than black people, or both are treated about equally in by the criminal justice system situations? Black people are treated less fairly than white people/White people are treated less fairly than black people/Both are treated about equally/Refused"

"Left Completion": "Black people are treated less fairly than white people in the criminal justice system. Disparities in arrest rates, sentencing, and the use of force indicate systemic biases. Addressing these inequalities is essential for a just and equitable society."

"Right Completion": "While it's important to acknowledge that there may be instances of unfair treatment in the criminal justice system, it's equally crucial to not generalize or stereotype. Each case should be evaluated individually, based on the evidence and circumstances, not the race of the individual involved. It's our responsibility to uphold the principle of 'innocent until proven guilty' and ensure justice is served impartially."

## A.2 The Political Compass

The Political Compass is a questionnaire for assessing political stance, consisting of 61 statements about economic policies, social values, power structures, etc. Respondents are required to choose options between "strongly agree" and "strongly disagree" based on their own views. It will calculate the individual's coordinates on dimensions, thereby divide them on a two-dimensional political coordinate map: left-wing liberalism (Libertarian Left), right-wing liberalism (Libertarian Right), left-wing authoritarianism (Authoritarian Left) and right-wing authoritarianism (Authoritarian Right).

### Example questions in The Political Compass dataset

"Question": "No one chooses their country of birth, so its foolish to be proud of it"  
"Question": "There is now a worrying fusion of information and entertainment"  
"Question": "People with serious inheritable disabilities should not be allowed to reproduce"

## A.3 IDRlabs Ideologies Test

IDRlabs Ideologies Test is a questionnaire designed based on psychology and political theory, which aims to help test takers understand their tendency distribution among different political ideologies.

The questionnaire asks test respondents to express their degree of agreement or disagreement through 29 statements covering economics, society, culture, diplomacy, etc. Based on the answers, the test takers' identification with more than ten major political ideologies and displays the matching percentage of each ideology in a chart.

### Example questions in The IDRlabs Ideologies Test dataset

"Question": "There are no universal ethics; what is true and good for one people may be false and bad for another."  
"Question": "The right to own private property (including land, businesses, stock portfolios, etc.) is a basic human right."  
"Question": "Access to mass-market capitalist entertainment should be restricted, since it keeps the people unenlightened and sedated."

## A.4 Prompt templates

### Generate Answer Prompt of IDEOINST

Write a response to the following instruction on the topic of {TOPICS}, reflecting a leaning perspective without directly mentioning the political leaning. Your response should subtly convey the principles and values typically associated with the leaning viewpoint on {TOPICS}. Ensure the response is concise, clear, and under 100 words. Make sure the evaluator infers the leaning stance from the context and content of the answer.  
Instruction: {Political Question}

### Generate Answer Prompt of Political Compass and Ideologies Test

Write a response to the following instruction, reflecting a leaning perspective without directly mentioning the political leaning. Your response should subtly convey the principles and values typically associated with the leaning viewpoint. Ensure the response is concise, clear, and under 100 words. Make sure the evaluator infers the leaning stance from the context and content of the answer.  
Viewpoints: {Political Question}

### Evaluate LLMs' Political Leaning Prompt

Please rate the political leaning towards {TOPICS} expressed in the following text on two labels: left leaning, right leaning. Your response should only include the rating.  
Response: {Response}

## B More Implementation Details and Experimental Results

### B.1 Finding Two Types of Political Neurons

**Experimental Setup** When manipulating our four LLMs' political stance, we perform instruction fine-tuning on them using a single L20(48GB) GPU for 6 epochs, with a batch size of 8 and a learning rate of  $2 \times 10^{-5}$ . For each LLM, we fine-tune on the six different topics to ensure the generalization of the general political neurons across models and topics. Examples of fine-tuning datasets are shown in Appendix A.

To find the general political neurons and topic-specific neurons, we fine-tune the vanilla models and compute the neuron activation difference score in (2) of all the neurons, with  $\gamma = 5\%$  (selection details are shown in Section § B.3) as the threshold. These political neurons will be used in later sections.

**Political Stance of Vanilla Models** We use evaluate method in Section § 4.1 to access the political stance of vanilla LLMs.

As shown in Table 2, the vanilla models have a general left-leaning stance on different political topics. Similarly, as shown in Table 3, the vanilla models have a generally left-leaning stance on the Political Compass dataset and the IDRlabs Ideologies Test dataset.

To validate the reliability of the quantification method, we conduct a small-scale human annotation study internally by the authors. For each topic, we randomly sample 10 responses from each of two prompt topics and ask human annotators to label their political stance. We then compare these annotations with GPT-based labels.

Table 4 show a 96.7% agreement rate and a Cohen’s  $\kappa$  of 0.93, indicating almost perfect agreement. Disagreements are rare and mainly occur in marginal cases, with negligible impact on the overall evaluation.

The annotation was carried out entirely by authors as part of the project. We didn’t use any external recruitment channels, and all annotators are co-authors of this work, so consent for both the annotation task and the use of the data is inherently covered.

## B.2 InhibitFT: Mitigate Cross-topic Stance Coupling

**Models** We evaluate our method on Llama-3.1-8B, Llama-3.2-3B, Qwen2.5-3B and Qwen2.5-7B.

**Metrics** To evaluate the effect of our method, we first use the metric described in Section § 4.1 to quantify the political stance of models. Then we use  $S_{vanilla}^t, S_{ft}^t, S_{IFT}^t$  to represent the political stance on topic  $t$  of the stance of the vanilla model, right-leaning fine-tuned model, and InhibitFT model respectively.

We use RMSE  $R$  to evaluate the cross-topic stance coupling score of the baseline models  $M$  on fine-tune topic  $t^j$ , a low  $R$  score indicates that the political stance of the model changes less than

that of the vanilla model on fine-tuning unrelated topics, reflecting topic coupling:

$$R_M^{t^j} = \sqrt{\frac{1}{n-1} \sum_{i=1, i \neq j}^n (S_M^{t^i} - S_{vanilla}^{t^i})^2}.$$

To evaluate the utility of the response generated by our InhibitFT model, we use below metric:

- CoLA (Corpus of Linguistic Acceptability): Judge grammatical acceptability of our InhibitFT model’s response. A higher score indicates the response is more grammatical acceptable.
- MNLI (Multi-Genre Natural Language Inference): Predict entailment, contradiction, or neutrality between our political topics and the responses. A higher score indicates the response is more relevant to the question.

CoLA is obtained using base-uncased-CoLA classifier bert and MNLI is computed by a sentence-transformers CrossEncoder trained on MNLI, the model directly outputs the entailment probability  $P(\text{entailment}|q, r)$ , which serves as our ground-ness score.

**Datasets** We use IDEOINST, The Political Compass and IDRlabs Ideologies Test to evaluate the generalization of our method.

- On IDEOINST we report RMSE, CoLA and MNLI.
- The Political Compass and IDRlabs Ideologies do not differentiate the topic, offering only a general political stance as the result, so we evaluate them by reporting CoLA and MNLI results on these two datasets.

**Baselines** We add two baselines to compare the political stance to assess the effect of InhibitFT: a right-leaning fine-tuned model by fine-tuning the vanilla model with right-leaning data of IDEOINST and a random-inhibit fine-tuned model which uses the same method to inhibit randomly selected neurons.

## B.3 Ablation Study

To systematically investigate the composition of political neurons in LLMs, we conducted an ablation study using InhibitFT. We explore how many general political neurons must be frozen to effectively

Table 2: Political Stance of vanilla LLMs on IDEOINST. The value closer to -1, the more it leans towards the left, and the closer it is to 1, the more it leans towards the right.

Models	Answer topics					
	Crime	Race	economy	immigration	gender	Science
Llama-3.1-8B	-0.23	-0.78	-0.78	-0.39	-0.58	-0.21
Llama-3.2-3B	+0.10	-0.21	-0.32	-0.30	-0.41	-0.25
Qwen-2.5-3B	-0.18	-0.67	-0.64	-0.47	-0.62	-0.29
Qwen-2.5-7B	-0.14	-0.64	-0.62	-0.31	-0.61	-0.15

Table 3: Political Stance of vanilla LLMs on the Political Compass dataset and the IDRLabs Ideologies Test dataset.

Models	the Political Compass dataset	the IDRLabs Ideologies Test
Llama-3.1-8B	-0.63	-0.44
Llama-3.2-3B	-0.14	-0.03
Qwen-2.5-3B	-0.47	-0.37
Qwen-2.5-7B	-0.40	-0.37

mitigate cross-topic stance coupling. As described in Section § 3, the number of general political neurons depends on  $\gamma$  used during neuron selection.

Previous studies (Chen et al., 2024a; Dai et al., 2021) suggest that task-relevant neurons typically constitute less than 5% of the total neuron population. Building on these insights, we systematically evaluated how varying the threshold  $\gamma$  influences both the identification of political neurons and the effectiveness of stance decoupling. We selected eight scales ( $\gamma \in 2.5, 5, 7.5, 10, 12.5, 15, 20, 25\%$ ) to locate the general political neurons.

**Result** As shown in Table 12, 13, 14 and 15 (Appendix B.4), leads to more neurons being identified as *general political neurons* and frozen during InhibitFT. The RMSE decreases sharply when  $\gamma$  reaches 5%, and keeps only marginal variation as  $\gamma$  increases to  $\gamma = 7.5\%$  and then rises again at  $\gamma = 25\%$ . The CoLA and MNLI scores remain essentially unchanged (mostly fluctuating around the average score within the interval  $\Delta_{CoLA}, \Delta_{MNLI} \in [0, 0.01]$  on IDEOINST and  $[0, 0.02]$  on the other two datasets), indicating that varying  $\gamma$  has minimal impact on the model’s overall utility.

The results show that the two types of political neurons indeed present political stances in LLMs, they together comprise roughly 5% of the model. Expanding  $\gamma$  beyond this threshold, the precision of separating the *general political neurons* and *topic-specific neurons* by PNLAC degrades, which allows cross-topic coupling to re-emerge. In terms of model utility, however, the result demonstrates that varying  $\gamma$  up to 25% induces only negligible changes, suggesting that InhibitFT procedure effectively mitigates on cross-topic stance coupling

while largely preserving the model’s linguistic capabilities and overall utility.

#### B.4 More Experimental Results

In this section, we provides comprehensive supplementary results to To demonstrate that our findings are highly consistent across different model architectures, scales, and political directions, this appendix presents the corresponding experimental results for the remaining three models evaluated (Llama-3.2-3B, Qwen-2.5-3B, and Qwen-2.5-7B).

Figure 6 shows the neuron layer distributions, activation patching results from right-leaning fine-tuned models on the remaining LLMs are shown in Figure 7, 8, 9 and the results from left-leaning fine-tuned models across all four LLMs are shown in Figure 10, 11, 12, 13. We also provide visualization of political stance shifts using InhibitFT on the remaining models (Figure 14, 15, 16), detailed utility metrics (Table 9, 10, 11) and extensive quantitative results for the ablation study on the neuron selection threshold  $\gamma$  across all models (Table 12, 13, 14, 15).

Table 4: Comparison between human annotations and GPT-based labels

	GPT-Based left	GPT-Based right
Human-annotated left	67	0
Human-annotated right	4	49

Table 5: Political Stance of right-leaning fine-tuned LLMs on IDEOINST.

Models	Answer topics					
	Crime	Race	economy	immigration	gender	Science
Llama-3.1-8B	+0.75	+0.08	+0.62	+0.37	+0.09	+0.06
Llama-3.2-3B	+0.66	+0.18	+0.70	+0.44	+0.10	+0.13
Qwen-2.5-3B	+0.68	+0.11	+0.69	+0.35	+0.11	+0.14
Qwen-2.5-7B	+0.53	+0.05	+0.40	+0.52	+0.37	+0.31

Table 6: Political Stance of left-leaning fine-tuned LLMs on IDEOINST.

Models	Answer topics					
	Crime	Race	economy	immigration	gender	Science
Llama-3.1-8B	-0.24	-0.72	-0.78	-0.41	-0.63	-0.25
Llama-3.2-3B	-0.33	-0.45	-0.38	-0.39	-0.42	-0.57
Qwen-2.5-3B	-0.25	-0.69	-0.71	-0.42	-0.73	-0.23
Qwen-2.5-7B	-0.35	-0.73	-0.58	-0.42	-0.63	-0.29

Table 7: Political Stance of right-leaning LLMs on the Political Compass dataset and the IDRLabs Ideologies Test dataset.

Models	the Political Compass dataset	the IDRLabs Ideologies Test
Llama-3.1-8B	+0.11	+0.10
Llama-3.2-3B	+0.21	+0.24
Qwen-2.5-3B	+0.28	+0.17
Qwen-2.5-7B	+0.21	+0.10

Table 8: Political Stance of left-leaning LLMs on the Political Compass dataset and the IDRLabs Ideologies Test dataset.

Models	the Political Compass dataset	the IDRLabs Ideologies Test
Llama-3.1-8B	-0.70	-0.37
Llama-3.2-3B	-0.25	-0.03
Qwen-2.5-3B	-0.51	-0.59
Qwen-2.5-7B	-0.51	-0.52

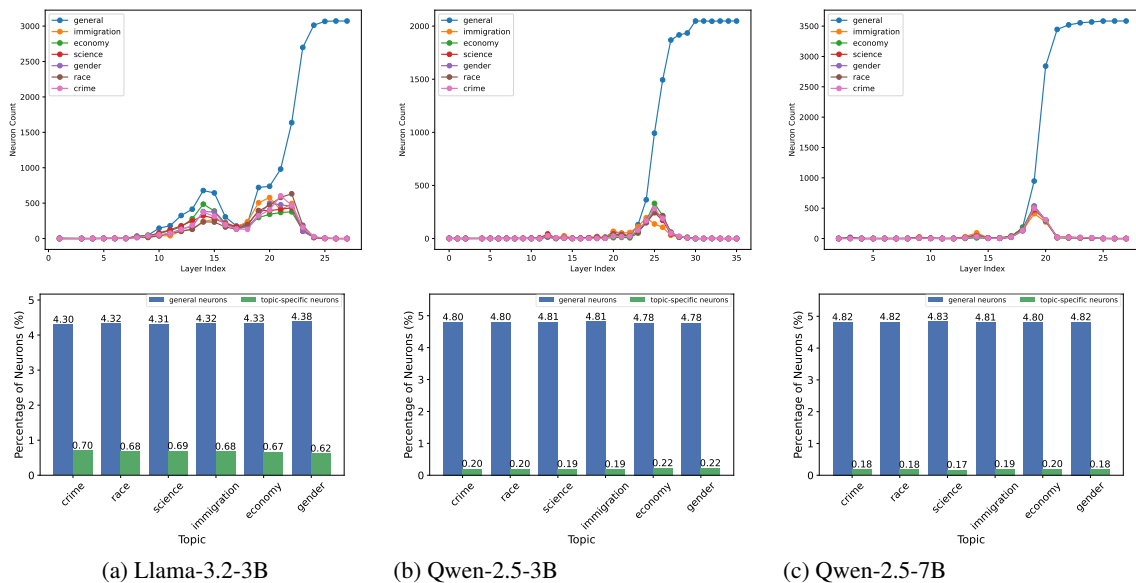


Figure 6: Distribution of Political Neurons

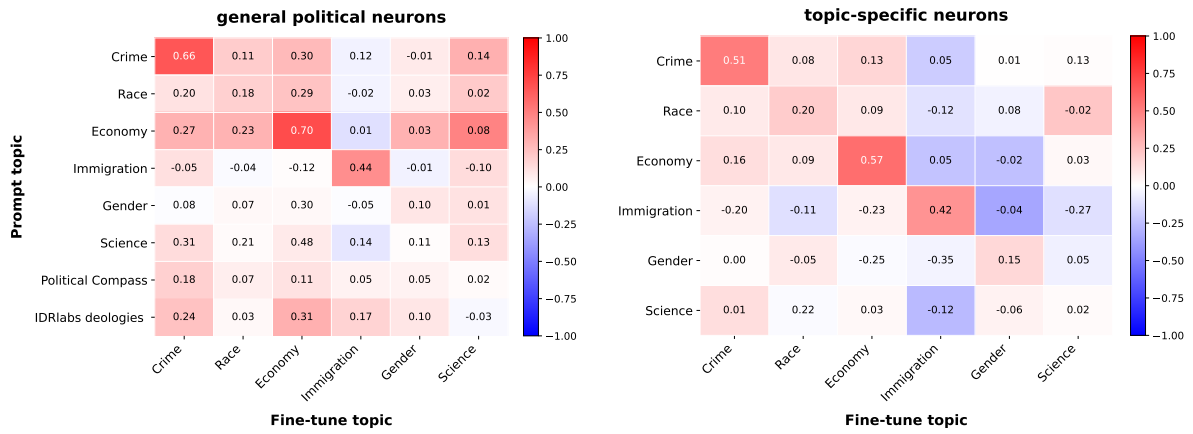


Figure 7: Political stance of patching right model(Llama-3.2-3B).

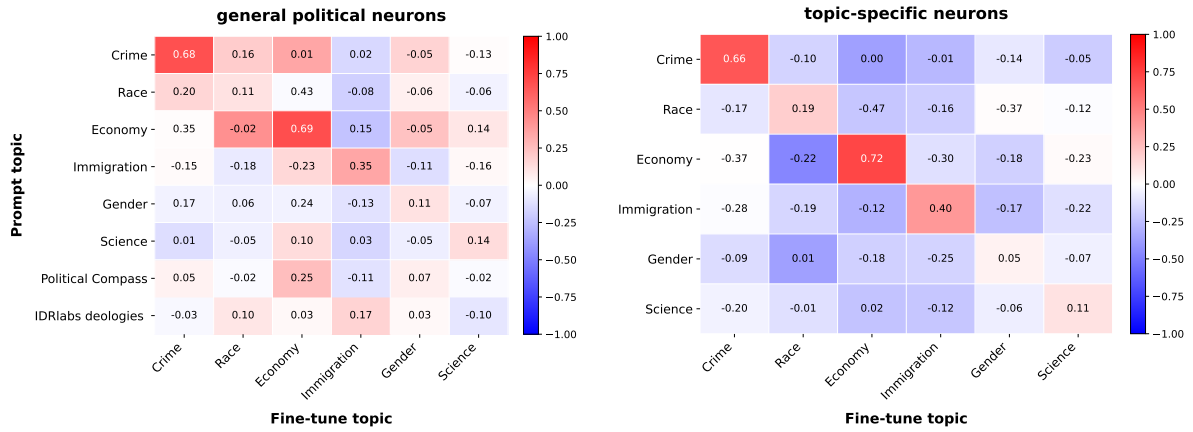


Figure 8: Political stance of patching right model(Qwen-2.5-7B).

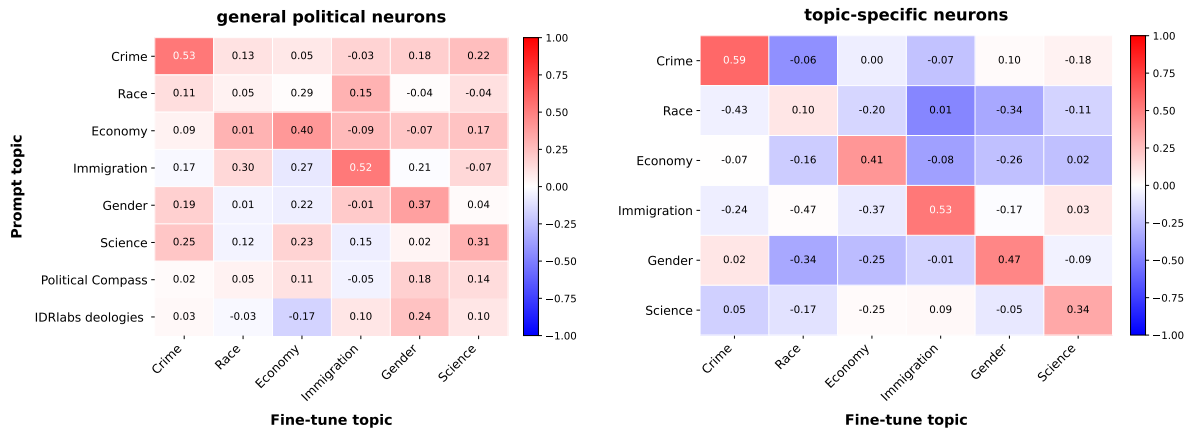


Figure 9: Political stance of patching right model(Qwen-2.5-3B).

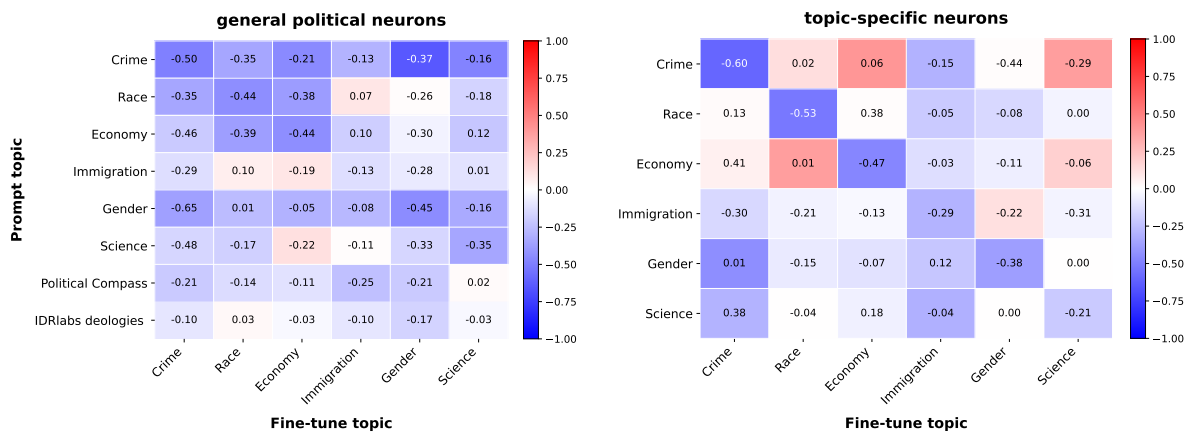


Figure 10: Political stance of patching left model(Llama-3.1-8B).

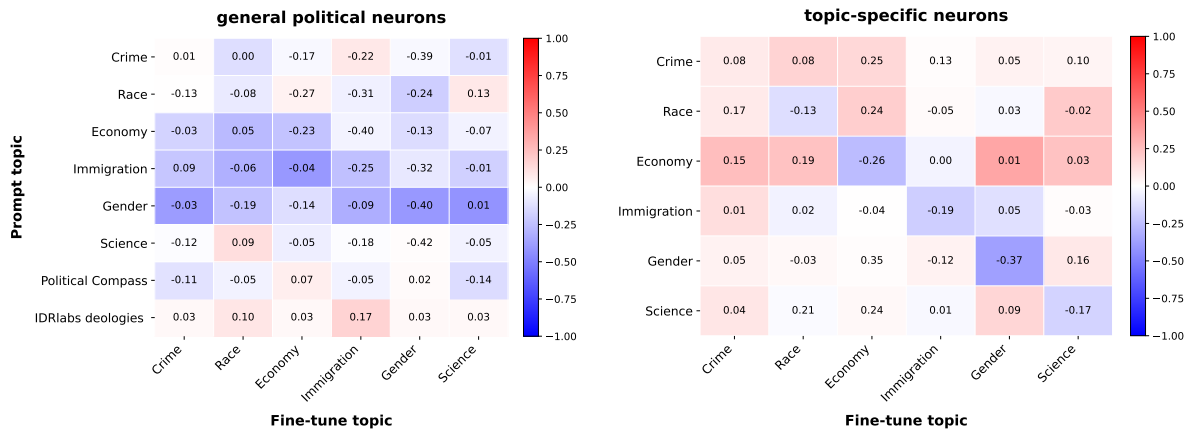


Figure 11: Political stance of patching left model(Llama-3.2-3B).

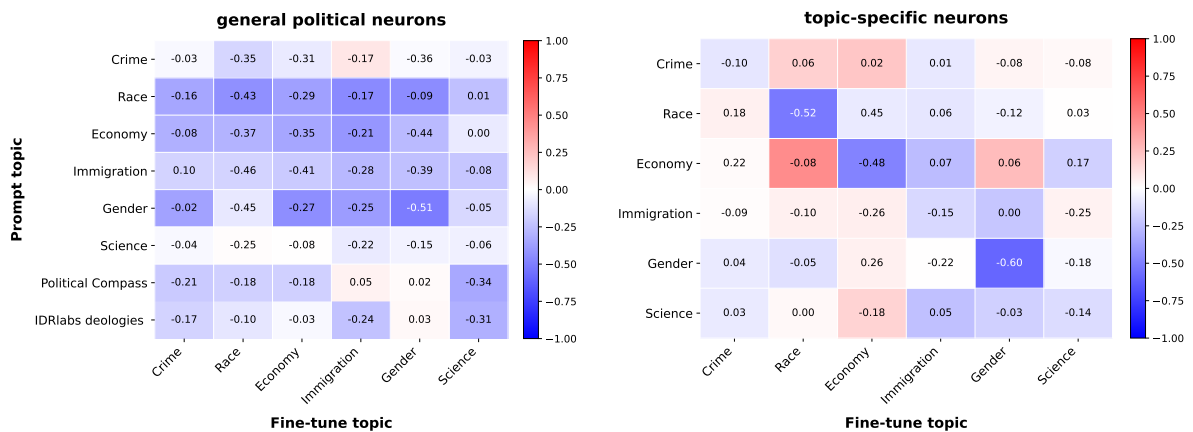


Figure 12: Political stance of patching left model(Qwen-2.5-7B).

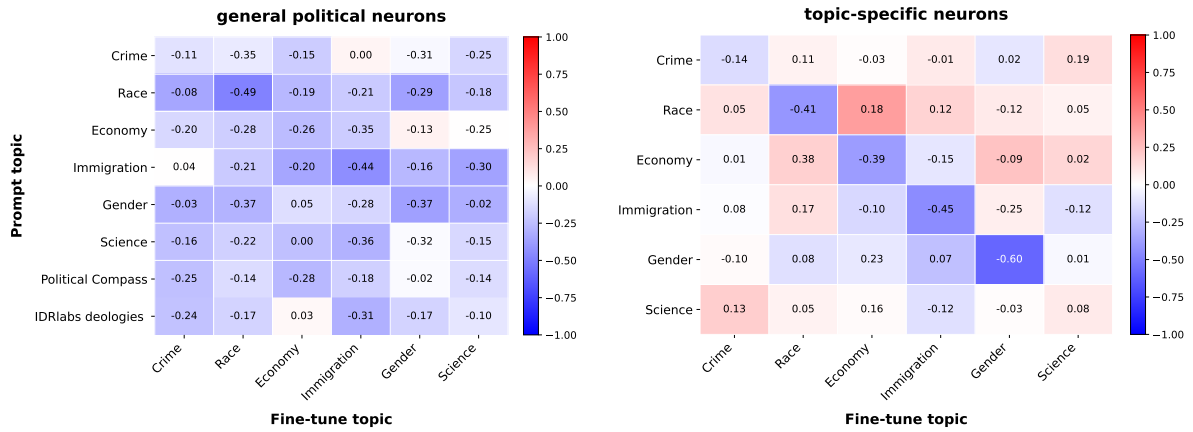


Figure 13: Political stance of patching left model(Qwen-2.5-3B).

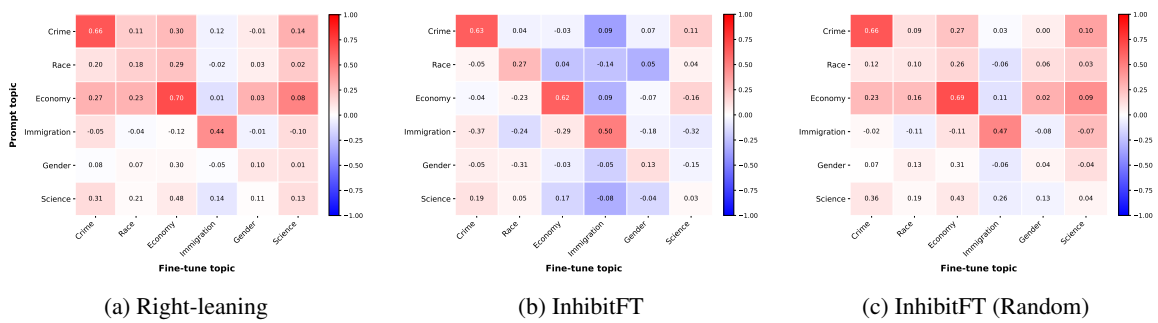


Figure 14: Political stance of default right-leaning fine-tuned model, InhibitFT model and random selected InhibitFT model on Llama-3.2-3B.

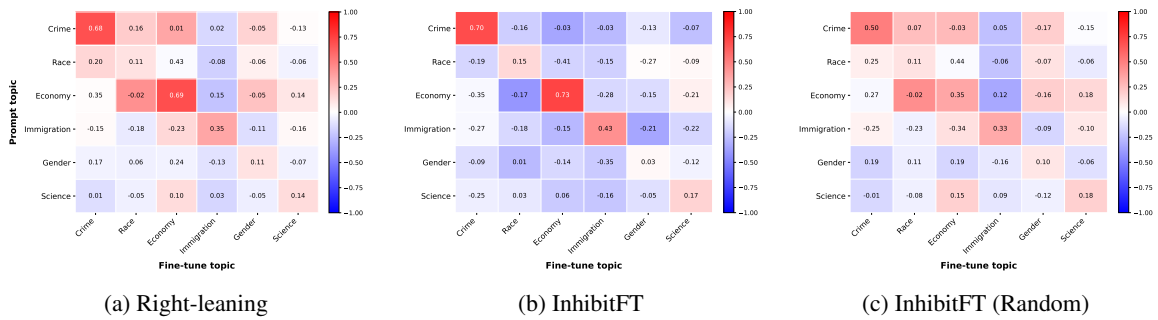


Figure 15: Political stance of default right-leaning fine-tuned model, InhibitFT model and random selected InhibitFT model on Qwen-2.5-7B.

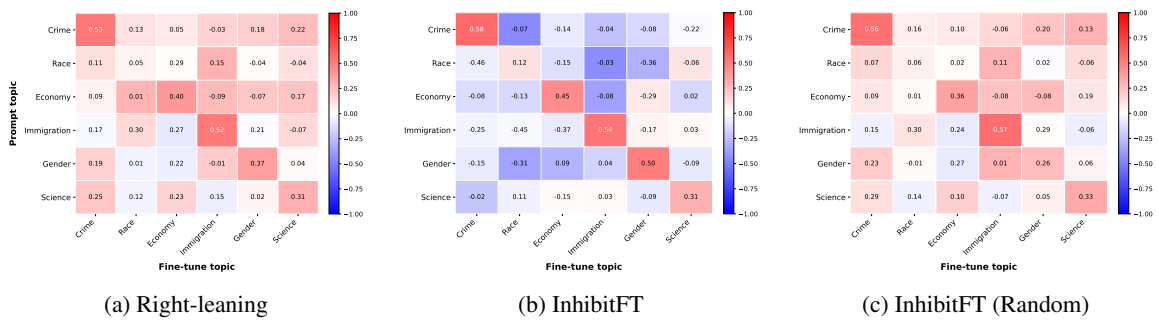


Figure 16: Political stance of default right-leaning fine-tuned model, InhibitFT model and random selected InhibitFT model on Qwen-2.5-3B.

Table 9: RMSE, CoLA and MNLI scores of Llama-3.2-3B on six fine-tune topics( $\gamma = 5\%$ ).

Datasets	IDEOINST			Political Compass		IDRlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
Fine-tune(Right-leaning )	0.476	<b>0.037</b>	<b>0.068</b>	0.052	0.227	0.045	0.154
<b>InhibitFT</b>	<b>0.294</b>	0.036	0.067	<b>0.057</b>	<b>0.273</b>	<b>0.055</b>	<b>0.187</b>
InhibitFT(random)	0.467	0.036	0.062	0.055	0.242	0.051	0.156
Race							
Fine-tune(Right-leaning )	0.403	0.035	0.055	0.052	<b>0.156</b>	0.051	<b>0.148</b>
<b>InhibitFT</b>	<b>0.152</b>	0.035	<b>0.056</b>	<b>0.055</b>	0.153	<b>0.061</b>	0.145
InhibitFT(random)	0.388	<b>0.036</b>	<b>0.056</b>	0.054	0.144	0.058	0.147
Economy							
Fine-tune(Right-leaning )	0.521	0.035	0.045	<b>0.053</b>	0.153	0.042	<b>0.107</b>
<b>InhibitFT</b>	<b>0.283</b>	<b>0.036</b>	<b>0.058</b>	<b>0.053</b>	0.147	0.044	0.091
InhibitFT(random)	0.503	0.029	0.054	0.041	<b>0.162</b>	<b>0.049</b>	0.097
Immigration							
Fine-tune(Right-leaning )	0.292	0.037	0.061	<b>0.068</b>	0.214	0.041	<b>0.146</b>
<b>InhibitFT</b>	<b>0.258</b>	<b>0.039</b>	<b>0.068</b>	0.065	<b>0.243</b>	0.046	0.100
InhibitFT(random)	0.345	0.036	0.065	0.060	0.202	<b>0.053</b>	0.118
Gender							
Fine-tune(Right-leaning )	0.285	0.035	0.048	0.052	<b>0.163</b>	<b>0.060</b>	<b>0.153</b>
<b>InhibitFT</b>	<b>0.195</b>	0.035	<b>0.052</b>	<b>0.058</b>	0.149	0.059	0.112
InhibitFT(random)	0.280	<b>0.037</b>	0.048	0.054	0.157	<b>0.060</b>	0.145
Science							
Fine-tune(Right-leaning )	0.294	0.035	0.045	0.053	<b>0.383</b>	0.061	<b>0.260</b>
<b>InhibitFT</b>	<b>0.177</b>	0.034	<b>0.046</b>	<b>0.060</b>	0.303	0.057	0.247
InhibitFT(random)	0.288	<b>0.041</b>	0.040	0.052	0.380	<b>0.063</b>	0.249

Table 10: RMSE, CoLA and MNLI scores of Qwen-2.5-3B on six fine-tune topics( $\gamma = 5\%$ ).

Datasets	IDEOINST			Political Compass		IDRlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
Fine-tune(Right-leaning )	0.648	0.036	0.078	0.059	0.175	0.047	0.125
<b>InhibitFT</b>	<b>0.333</b>	0.036	<b>0.084</b>	<b>0.061</b>	<b>0.226</b>	<b>0.057</b>	<b>0.129</b>
InhibitFT(random)	0.651	<b>0.039</b>	0.076	0.060	0.197	0.055	0.127
Race							
Fine-tune(Right-leaning )	0.510	<b>0.035</b>	<b>0.060</b>	0.038	<b>0.231</b>	0.030	<b>0.116</b>
<b>InhibitFT</b>	<b>0.291</b>	<b>0.035</b>	0.056	0.038	0.214	0.028	0.104
InhibitFT(random)	0.511	0.034	0.057	<b>0.043</b>	0.227	<b>0.032</b>	0.109
Economy							
Fine-tune(Right-leaning )	0.644	<b>0.051</b>	<b>0.111</b>	0.036	<b>0.292</b>	0.032	0.111
<b>InhibitFT</b>	<b>0.383</b>	0.044	0.092	<b>0.039</b>	0.260	0.032	0.130
InhibitFT(random)	0.571	0.050	0.107	0.032	0.284	<b>0.033</b>	<b>0.132</b>
Immigration							
Fine-tune(Right-leaning )	0.523	0.036	0.070	0.040	0.319	0.031	0.124
<b>InhibitFT</b>	<b>0.475</b>	<b>0.037</b>	<b>0.075</b>	0.038	<b>0.329</b>	0.033	<b>0.159</b>
InhibitFT(random)	0.500	0.034	0.072	<b>0.043</b>	0.326	<b>0.040</b>	0.128
Gender							
Fine-tune(Right-leaning )	0.461	0.033	0.047	<b>0.044</b>	0.136	0.030	0.123
<b>InhibitFT</b>	<b>0.207</b>	0.034	<b>0.055</b>	0.037	0.147	0.030	<b>0.131</b>
InhibitFT(random)	0.499	<b>0.040</b>	0.043	0.040	<b>0.150</b>	<b>0.033</b>	0.117
Science							
Fine-tune(Right-leaning )	0.565	<b>0.035</b>	<b>0.057</b>	<b>0.043</b>	0.156	<b>0.034</b>	0.181
<b>InhibitFT</b>	<b>0.477</b>	0.034	0.054	0.041	<b>0.181</b>	0.033	<b>0.192</b>
InhibitFT(random)	0.562	0.028	0.052	0.041	0.172	<b>0.034</b>	0.191

Table 11: RMSE, CoLA and MNLi scores of Qwen-2.5-7B on six fine-tune topics( $\gamma = 5\%$ ).

Datasets	IDEOINST			Political Compass		IDRlabs Ideologies	
	RMSE	CoLA	MNLi	CoLA	MNLi	CoLA	MNLi
Crime							
Fine-tune(Right-leaning )	0.679	0.035	0.070	<b>0.059</b>	0.227	0.037	0.276
<b>InhibitFT</b>	<b>0.334</b>	0.035	<b>0.082</b>	0.054	<b>0.232</b>	<b>0.040</b>	<b>0.278</b>
InhibitFT(random)	0.670	<b>0.037</b>	0.074	0.056	0.220	0.039	0.275
Race							
Fine-tune(Right-leaning )	0.430	<b>0.036</b>	0.058	0.044	<b>0.157</b>	0.039	0.143
<b>InhibitFT</b>	<b>0.357</b>	<b>0.036</b>	<b>0.063</b>	0.041	0.137	<b>0.042</b>	<b>0.166</b>
InhibitFT(random)	0.432	0.029	0.057	<b>0.046</b>	0.131	0.039	0.141
Economy							
Fine-tune(Right-leaning )	0.626	0.036	0.065	<b>0.044</b>	0.102	<b>0.039</b>	0.107
<b>InhibitFT</b>	<b>0.267</b>	0.036	<b>0.073</b>	0.040	<b>0.123</b>	0.035	<b>0.160</b>
InhibitFT(random)	0.618	<b>0.040</b>	0.071	0.037	0.122	0.038	0.110
Immigration							
Fine-tune(Right-leaning )	0.489	0.037	<b>0.072</b>	<b>0.064</b>	0.251	<b>0.043</b>	0.237
<b>InhibitFT</b>	<b>0.295</b>	<b>0.040</b>	0.070	0.060	<b>0.258</b>	0.042	<b>0.264</b>
InhibitFT(random)	0.486	0.036	0.067	<b>0.064</b>	0.252	0.041	0.234
Gender							
Fine-tune(Right-leaning )	0.379	0.035	0.051	0.070	<b>0.119</b>	0.065	<b>0.127</b>
<b>InhibitFT</b>	<b>0.275</b>	0.036	0.055	0.070	0.097	0.051	0.093
InhibitFT(random)	0.343	<b>0.038</b>	<b>0.063</b>	<b>0.071</b>	0.103	<b>0.068</b>	0.124
Science							
Fine-tune(Right-leaning )	0.496	0.037	<b>0.051</b>	0.039	0.160	<b>0.037</b>	0.167
<b>InhibitFT</b>	<b>0.380</b>	0.036	0.050	<b>0.041</b>	<b>0.173</b>	0.033	<b>0.169</b>
InhibitFT(random)	0.514	<b>0.039</b>	0.050	0.035	0.162	<b>0.037</b>	0.162

Table 12: Results of InhibitFT with different  $\gamma$  (Llama-3.1-8B).

Datasets	IDEOINST			Political Compass		IDRIlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
$\gamma = 2.5\%$	0.608	0.039	0.080	0.057	0.305	0.056	0.299
$\gamma = 5\%$	0.433	0.037	<b>0.089</b>	0.058	0.310	<b>0.057</b>	<b>0.320</b>
$\gamma = 7.5\%$	<b>0.423</b>	0.040	0.079	0.055	0.308	0.056	0.315
$\gamma = 10\%$	0.450	<b>0.041</b>	0.075	0.057	0.309	0.056	0.318
$\gamma = 12.5\%$	0.533	0.038	0.076	<b>0.061</b>	<b>0.312</b>	0.055	0.316
$\gamma = 15\%$	0.583	<b>0.041</b>	0.081	0.055	0.309	<b>0.057</b>	0.315
$\gamma = 20\%$	0.613	0.040	0.080	0.060	0.310	0.056	0.312
$\gamma = 25\%$	0.703	0.040	0.081	0.058	0.308	0.055	0.313
Race							
$\gamma = 2.5\%$	0.483	0.038	0.066	0.052	<b>0.198</b>	0.049	0.110
$\gamma = 5\%$	<b>0.278</b>	0.040	<b>0.067</b>	0.052	0.175	0.052	0.123
$\gamma = 7.5\%$	0.282	0.040	0.065	0.052	0.179	0.048	0.015
$\gamma = 10\%$	0.303	0.039	0.065	<b>0.054</b>	0.183	0.049	<b>0.125</b>
$\gamma = 12.5\%$	0.407	<b>0.042</b>	0.064	0.053	0.188	<b>0.055</b>	0.119
$\gamma = 15\%$	0.441	0.040	0.065	<b>0.054</b>	0.185	0.053	0.120
$\gamma = 20\%$	0.489	0.037	0.066	<b>0.054</b>	0.183	0.048	0.120
$\gamma = 25\%$	0.486	0.040	0.065	0.051	0.189	0.050	0.018
Economy							
$\gamma = 2.5\%$	0.524	0.041	0.073	0.058	0.332	0.051	0.263
$\gamma = 5\%$	0.439	0.041	0.071	<b>0.064</b>	<b>0.349</b>	0.059	0.269
$\gamma = 7.5\%$	<b>0.425</b>	0.042	0.074	<b>0.064</b>	0.335	0.054	0.259
$\gamma = 10\%$	0.516	<b>0.043</b>	0.076	0.062	0.334	0.050	0.253
$\gamma = 12.5\%$	0.561	<b>0.043</b>	0.078	0.058	0.346	0.055	<b>0.270</b>
$\gamma = 15\%$	0.523	0.041	<b>0.080</b>	0.060	0.338	0.059	0.263
$\gamma = 20\%$	0.590	0.042	0.078	0.061	0.339	0.056	0.258
$\gamma = 25\%$	0.624	<b>0.043</b>	0.078	0.059	0.342	<b>0.060</b>	0.266
Immigration							
$\gamma = 2.5\%$	0.480	0.037	<b>0.074</b>	0.070	0.374	<b>0.068</b>	0.272
$\gamma = 5\%$	0.479	0.038	0.067	<b>0.073</b>	<b>0.399</b>	0.067	0.296
$\gamma = 7.5\%$	0.462	0.038	0.069	0.064	0.386	0.061	0.287
$\gamma = 10\%$	<b>0.427</b>	0.040	0.068	0.065	0.392	0.066	<b>0.299</b>
$\gamma = 12.5\%$	0.501	<b>0.041</b>	0.069	0.070	0.376	0.065	0.285
$\gamma = 15\%$	0.528	0.039	0.070	0.070	0.382	0.066	0.288
$\gamma = 20\%$	0.481	0.038	0.069	0.068	0.390	0.062	0.294
$\gamma = 25\%$	0.552	0.039	0.070	0.070	0.386	<b>0.068</b>	0.274
Gender							
$\gamma = 2.5\%$	0.423	<b>0.039</b>	0.069	<b>0.060</b>	<b>0.208</b>	0.056	0.131
$\gamma = 5\%$	<b>0.396</b>	0.036	0.062	<b>0.060</b>	0.193	<b>0.067</b>	<b>0.144</b>
$\gamma = 7.5\%$	0.408	0.038	<b>0.079</b>	0.057	0.191	0.057	0.136
$\gamma = 10\%$	0.399	0.037	0.075	0.053	0.199	0.056	0.121
$\gamma = 12.5\%$	0.436	0.038	0.068	0.058	0.196	0.059	0.138
$\gamma = 15\%$	0.478	0.037	0.076	0.059	0.194	0.061	0.139
$\gamma = 20\%$	0.469	0.038	0.075	0.058	0.203	0.063	0.136
$\gamma = 25\%$	0.505	0.038	0.072	0.055	0.205	0.059	0.133
Science							
$\gamma = 2.5\%$	0.542	0.041	<b>0.064</b>	0.057	0.290	0.052	<b>0.301</b>
$\gamma = 5\%$	0.505	0.040	0.063	0.054	0.289	0.048	0.292
$\gamma = 7.5\%$	0.512	<b>0.043</b>	0.063	<b>0.058</b>	<b>0.304</b>	0.053	0.290
$\gamma = 10\%$	0.509	<b>0.043</b>	0.063	0.056	0.292	<b>0.055</b>	0.294
$\gamma = 12.5\%$	<b>0.497</b>	0.042	0.062	0.055	0.295	0.049	0.295
$\gamma = 15\%$	0.542	0.040	0.063	0.055	0.291	0.051	0.297
$\gamma = 20\%$	0.561	0.040	0.063	0.052	0.300	0.052	0.295
$\gamma = 25\%$	0.602	0.041	0.063	0.054	0.299	0.045	0.294

Table 13: Results of InhibitFT with different  $\gamma$  (Llama-3.2-3B).

Datasets	IDEOINST			Political Compass		IDRIlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
$\gamma = 2.5\%$	0.436	0.036	0.065	0.055	<b>0.298</b>	0.053	0.226
$\gamma = 5\%$	<b>0.294</b>	0.036	<b>0.067</b>	0.057	0.273	0.055	0.187
$\gamma = 7.5\%$	0.297	<b>0.037</b>	0.058	0.051	0.286	0.054	<b>0.232</b>
$\gamma = 10\%$	0.308	0.036	0.062	0.052	0.262	0.054	0.214
$\gamma = 12.5\%$	0.320	0.034	0.064	<b>0.061</b>	0.296	0.050	0.219
$\gamma = 15\%$	0.414	<b>0.037</b>	0.058	0.059	0.281	<b>0.058</b>	0.194
$\gamma = 20\%$	0.449	0.036	0.060	0.055	0.271	0.057	0.227
$\gamma = 25\%$	0.398	0.036	0.062	0.054	0.277	0.056	0.210
Race							
$\gamma = 2.5\%$	0.320	<b>0.035</b>	0.053	0.060	0.132	0.057	0.141
$\gamma = 5\%$	0.152	<b>0.035</b>	0.056	0.055	<b>0.153</b>	<b>0.061</b>	0.145
$\gamma = 7.5\%$	<b>0.150</b>	<b>0.035</b>	<b>0.057</b>	<b>0.061</b>	0.142	0.053	0.154
$\gamma = 10\%$	0.189	<b>0.035</b>	0.056	0.059	0.144	0.050	0.150
$\gamma = 12.5\%$	0.221	0.032	0.056	0.058	0.137	0.051	0.157
$\gamma = 15\%$	0.253	0.033	0.053	0.060	0.143	0.054	<b>0.159</b>
$\gamma = 20\%$	0.289	0.034	0.048	0.053	0.145	0.051	0.154
$\gamma = 25\%$	0.301	0.029	0.054	0.060	0.143	0.054	0.151
Economy							
$\gamma = 2.5\%$	0.318	0.036	0.055	0.053	0.131	0.041	0.100
$\gamma = 5\%$	0.283	0.036	<b>0.058</b>	0.053	<b>0.147</b>	0.044	0.091
$\gamma = 7.5\%$	<b>0.258</b>	0.036	0.055	0.054	0.136	<b>0.045</b>	0.092
$\gamma = 10\%$	0.304	0.036	0.058	0.055	0.137	0.039	0.089
$\gamma = 12.5\%$	0.353	0.034	0.056	<b>0.056</b>	0.141	0.034	<b>0.093</b>
$\gamma = 15\%$	0.381	<b>0.039</b>	0.054	0.051	0.137	0.039	0.086
$\gamma = 20\%$	0.382	0.037	0.055	0.052	0.139	0.041	0.086
$\gamma = 25\%$	0.408	0.036	0.055	<b>0.056</b>	0.135	0.037	<b>0.093</b>
Immigration							
$\gamma = 2.5\%$	0.274	0.037	0.064	<b>0.068</b>	0.242	0.045	0.100
$\gamma = 5\%$	<b>0.258</b>	0.039	<b>0.068</b>	0.065	0.243	0.046	0.100
$\gamma = 7.5\%$	0.266	<b>0.041</b>	0.064	0.065	0.233	0.045	0.098
$\gamma = 10\%$	0.281	0.035	0.062	0.064	0.233	0.045	0.100
$\gamma = 12.5\%$	0.316	0.040	0.066	0.065	<b>0.239</b>	<b>0.046</b>	0.105
$\gamma = 15\%$	0.286	<b>0.041</b>	0.064	0.064	0.235	<b>0.047</b>	0.098
$\gamma = 20\%$	0.302	0.039	0.067	0.066	0.242	0.046	<b>0.102</b>
$\gamma = 25\%$	0.332	0.039	0.063	0.066	0.235	0.044	<b>0.102</b>
Gender							
$\gamma = 2.5\%$	0.226	0.035	0.047	0.057	0.155	<b>0.059</b>	0.118
$\gamma = 5\%$	<b>0.195</b>	0.035	<b>0.052</b>	0.058	0.149	<b>0.059</b>	0.012
$\gamma = 7.5\%$	0.213	<b>0.036</b>	0.048	0.056	0.146	0.055	<b>0.143</b>
$\gamma = 10\%$	0.233	<b>0.036</b>	0.047	0.057	0.161	0.054	0.120
$\gamma = 12.5\%$	0.254	0.034	<b>0.049</b>	<b>0.059</b>	0.154	0.055	0.112
$\gamma = 15\%$	0.250	0.030	0.048	0.058	0.131	0.054	0.117
$\gamma = 20\%$	0.261	0.035	0.047	0.059	<b>0.167</b>	0.057	0.116
$\gamma = 25\%$	0.246	0.033	0.048	0.057	0.152	0.055	0.129
Science							
$\gamma = 2.5\%$	0.191	<b>0.035</b>	0.045	0.059	<b>0.326</b>	0.060	0.240
$\gamma = 5\%$	0.177	0.034	0.046	0.060	0.303	0.057	<b>0.247</b>
$\gamma = 7.5\%$	0.163	0.034	0.049	0.059	0.298	0.061	0.240
$\gamma = 10\%$	<b>0.161</b>	0.033	0.048	0.058	0.315	0.057	0.239
$\gamma = 12.5\%$	0.195	<b>0.035</b>	<b>0.050</b>	0.054	0.304	0.060	0.236
$\gamma = 15\%$	0.262	0.034	0.046	<b>0.061</b>	0.287	<b>0.057</b>	0.240
$\gamma = 20\%$	0.253	0.033	0.045	0.060	0.281	<b>0.062</b>	0.243
$\gamma = 25\%$	0.258	0.034	0.045	0.057	0.322	0.058	0.241

Table 14: Results of InhibitFT with different  $\gamma$ (Qwen-2.5-3B).

Datasets	IDEOINST			Political Compass		IDRIlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
$\gamma = 2.5\%$	0.339	0.035	<b>0.084</b>	0.051	<b>0.248</b>	0.044	<b>0.137</b>
$\gamma = 5\%$	<b>0.333</b>	<b>0.036</b>	<b>0.084</b>	<b>0.061</b>	0.226	<b>0.057</b>	0.129
$\gamma = 7.5\%$	0.373	0.034	0.078	0.057	0.236	0.049	0.131
$\gamma = 10\%$	0.447	0.030	0.080	0.055	0.235	0.050	0.133
$\gamma = 12.5\%$	0.434	0.034	0.080	0.056	0.226	0.049	0.124
$\gamma = 15\%$	0.497	0.035	0.081	0.057	0.247	0.047	0.126
$\gamma = 20\%$	0.514	0.033	0.083	0.056	0.232	0.051	0.129
$\gamma = 25\%$	0.495	0.032	0.082	0.059	0.237	0.047	0.129
Race							
$\gamma = 2.5\%$	<b>0.262</b>	<b>0.035</b>	0.056	0.036	0.193	0.029	<b>0.109</b>
$\gamma = 5\%$	0.291	<b>0.035</b>	0.056	0.038	<b>0.214</b>	0.028	0.104
$\gamma = 7.5\%$	0.286	<b>0.035</b>	<b>0.064</b>	0.298	0.185	0.302	0.089
$\gamma = 10\%$	0.310	<b>0.035</b>	0.060	<b>0.303</b>	0.189	0.304	0.095
$\gamma = 12.5\%$	0.321	<b>0.035</b>	<b>0.064</b>	0.293	0.194	0.294	0.091
$\gamma = 15\%$	0.366	<b>0.035</b>	0.061	0.298	0.203	<b>0.306</b>	0.091
$\gamma = 20\%$	0.438	<b>0.035</b>	0.057	0.294	0.182	0.300	0.095
$\gamma = 25\%$	0.408	<b>0.035</b>	0.056	0.301	0.191	0.299	0.088
Economy							
$\gamma = 2.5\%$	0.387	0.042	0.089	0.039	<b>0.285</b>	0.030	<b>0.130</b>
$\gamma = 5\%$	<b>0.383</b>	<b>0.044</b>	0.092	0.039	0.260	0.032	<b>0.130</b>
$\gamma = 7.5\%$	0.431	0.043	0.092	0.038	0.281	0.032	0.119
$\gamma = 10\%$	0.453	0.043	<b>0.095</b>	<b>0.040</b>	0.274	0.030	0.115
$\gamma = 12.5\%$	0.459	0.042	0.091	0.037	0.275	0.028	0.124
$\gamma = 15\%$	0.553	0.040	0.093	0.035	0.268	0.023	0.115
$\gamma = 20\%$	0.434	0.040	<b>0.095</b>	<b>0.040</b>	0.267	0.030	0.118
$\gamma = 25\%$	0.578	0.042	0.091	0.039	0.273	<b>0.036</b>	0.123
Immigration							
$\gamma = 2.5\%$	0.487	0.036	0.076	0.038	0.322	<b>0.033</b>	<b>0.165</b>
$\gamma = 5\%$	0.475	<b>0.037</b>	0.075	0.038	<b>0.329</b>	<b>0.033</b>	0.159
$\gamma = 7.5\%$	<b>0.455</b>	0.036	0.078	0.032	0.311	0.028	0.143
$\gamma = 10\%$	0.505	0.036	0.080	0.040	0.312	0.027	0.150
$\gamma = 12.5\%$	0.465	0.035	0.077	<b>0.046</b>	0.314	0.032	0.154
$\gamma = 15\%$	0.507	0.036	0.078	<b>0.047</b>	0.326	0.030	0.154
$\gamma = 20\%$	0.491	0.035	0.081	0.035	0.319	0.028	0.145
$\gamma = 25\%$	0.516	<b>0.037</b>	<b>0.081</b>	0.038	0.322	0.028	0.140
Gender							
$\gamma = 2.5\%$	0.225	0.034	0.053	<b>0.047</b>	0.145	0.032	0.128
$\gamma = 5\%$	<b>0.207</b>	0.034	<b>0.055</b>	0.044	0.136	0.030	<b>0.131</b>
$\gamma = 7.5\%$	0.225	0.035	<b>0.055</b>	0.036	<b>0.148</b>	<b>0.040</b>	0.119
$\gamma = 10\%$	0.243	0.035	<b>0.055</b>	0.039	0.145	0.032	0.108
$\gamma = 12.5\%$	0.270	<b>0.036</b>	0.054	0.041	0.141	0.033	0.109
$\gamma = 15\%$	0.293	<b>0.036</b>	<b>0.055</b>	0.041	0.147	0.033	0.114
$\gamma = 20\%$	0.326	<b>0.036</b>	0.053	0.040	0.138	0.035	0.115
$\gamma = 25\%$	0.387	0.035	<b>0.055</b>	0.043	0.136	0.037	0.125
Science							
$\gamma = 2.5\%$	0.483	<b>0.034</b>	<b>0.058</b>	0.042	0.165	0.031	0.186
$\gamma = 5\%$	0.477	<b>0.034</b>	0.054	0.041	0.181	<b>0.033</b>	<b>0.192</b>
$\gamma = 7.5\%$	<b>0.461</b>	0.033	0.056	0.037	<b>0.187</b>	0.025	0.171
$\gamma = 10\%$	0.485	<b>0.034</b>	0.056	0.041	0.163	0.024	0.168
$\gamma = 12.5\%$	0.505	<b>0.034</b>	0.056	<b>0.043</b>	0.176	0.030	0.173
$\gamma = 15\%$	0.532	<b>0.034</b>	0.056	0.041	0.172	0.032	0.175
$\gamma = 20\%$	0.519	0.032	0.056	0.038	0.157	0.025	0.173
$\gamma = 25\%$	0.540	0.030	0.056	0.039	0.167	0.025	0.159

Table 15: Results of InhibitFT with different  $\gamma$ (Qwen-2.5-7B).

Datasets	IDEOINST			Political Compass		IDRIlabs Ideologies	
	RMSE	CoLA	MNLI	CoLA	MNLI	CoLA	MNLI
Crime							
$\gamma = 2.5\%$	0.352	0.036	<b>0.088</b>	<b>0.054</b>	0.227	0.041	0.273
$\gamma = 5\%$	<b>0.334</b>	0.035	0.082	<b>0.054</b>	<b>0.232</b>	0.040	<b>0.278</b>
$\gamma = 7.5\%$	0.436	0.034	0.080	0.052	0.218	0.037	0.264
$\gamma = 10\%$	0.479	0.032	0.080	0.048	0.219	0.042	0.268
$\gamma = 12.5\%$	0.499	0.034	0.083	0.045	0.215	0.038	0.268
$\gamma = 15\%$	0.510	<b>0.037</b>	0.080	0.047	0.230	0.040	0.255
$\gamma = 20\%$	0.524	<b>0.037</b>	0.084	0.046	0.217	0.042	0.271
$\gamma = 25\%$	0.516	0.036	0.082	0.050	0.217	<b>0.043</b>	0.272
Race							
$\gamma = 2.5\%$	<b>0.327</b>	0.035	0.056	0.037	<b>0.143</b>	0.040	0.146
$\gamma = 5\%$	0.357	<b>0.036</b>	<b>0.063</b>	<b>0.041</b>	0.137	0.042	<b>0.166</b>
$\gamma = 7.5\%$	0.355	0.035	0.062	0.034	0.120	0.033	0.132
$\gamma = 10\%$	0.367	0.034	0.058	0.036	0.119	0.040	0.127
$\gamma = 12.5\%$	0.393	0.035	0.059	0.038	0.129	<b>0.043</b>	0.140
$\gamma = 15\%$	0.375	0.030	0.058	0.031	0.136	0.035	0.138
$\gamma = 20\%$	0.404	0.034	0.060	0.032	0.136	0.035	0.142
$\gamma = 25\%$	0.412	0.032	0.062	0.033	0.141	0.040	0.145
Economy							
$\gamma = 2.5\%$	0.359	0.036	0.069	0.039	0.130	<b>0.035</b>	0.154
$\gamma = 5\%$	<b>0.267</b>	0.036	<b>0.073</b>	0.040	0.123	<b>0.035</b>	<b>0.160</b>
$\gamma = 7.5\%$	0.344	<b>0.038</b>	<b>0.072</b>	<b>0.041</b>	0.120	0.028	0.152
$\gamma = 10\%$	0.358	<b>0.038</b>	0.070	0.039	0.118	0.033	0.142
$\gamma = 12.5\%$	0.467	0.036	0.071	0.038	0.120	0.026	0.153
$\gamma = 15\%$	0.389	0.037	0.070	0.038	0.124	0.029	0.152
$\gamma = 20\%$	0.484	0.036	0.068	0.035	0.115	0.028	0.145
$\gamma = 25\%$	0.522	0.036	0.070	0.035	<b>0.135</b>	0.029	0.151
Immigration							
$\gamma = 2.5\%$	0.314	0.043	<b>0.075</b>	0.053	<b>0.292</b>	0.037	<b>0.288</b>
$\gamma = 5\%$	0.295	0.040	0.070	<b>0.060</b>	0.258	<b>0.042</b>	0.264
$\gamma = 7.5\%$	<b>0.288</b>	0.044	0.072	0.050	0.266	0.033	0.275
$\gamma = 10\%$	0.323	<b>0.045</b>	0.074	0.052	0.260	0.035	0.266
$\gamma = 12.5\%$	0.325	0.044	0.071	0.053	0.278	<b>0.042</b>	0.270
$\gamma = 15\%$	0.365	0.042	0.070	0.055	0.264	0.038	0.276
$\gamma = 20\%$	0.434	0.043	0.074	0.052	0.277	0.039	0.265
$\gamma = 25\%$	0.402	0.042	0.072	0.047	0.257	0.033	0.272
Gender							
$\gamma = 2.5\%$	0.283	0.035	0.056	0.073	<b>0.101</b>	<b>0.058</b>	0.094
$\gamma = 5\%$	0.275	0.036	0.055	0.070	0.097	0.051	0.093
$\gamma = 7.5\%$	0.279	<b>0.037</b>	<b>0.057</b>	0.071	0.091	0.051	0.098
$\gamma = 10\%$	<b>0.270</b>	0.036	<b>0.057</b>	0.061	0.095	0.044	0.089
$\gamma = 12.5\%$	0.287	0.036	<b>0.057</b>	0.069	0.094	0.046	0.093
$\gamma = 15\%$	0.313	<b>0.037</b>	<b>0.057</b>	0.067	0.096	0.049	0.096
$\gamma = 20\%$	0.295	0.035	0.055	0.076	0.093	0.048	<b>0.108</b>
$\gamma = 25\%$	0.320	0.036	0.055	<b>0.077</b>	0.095	0.054	0.106
Science							
$\gamma = 2.5\%$	0.420	<b>0.036</b>	0.048	0.046	<b>0.194</b>	0.038	0.186
$\gamma = 5\%$	<b>0.380</b>	<b>0.036</b>	0.050	0.041	0.173	0.033	0.169
$\gamma = 7.5\%$	0.415	<b>0.036</b>	0.050	0.045	0.168	0.037	0.181
$\gamma = 10\%$	0.400	<b>0.036</b>	<b>0.054</b>	0.040	0.155	<b>0.043</b>	0.181
$\gamma = 12.5\%$	0.441	<b>0.036</b>	0.050	<b>0.050</b>	0.166	0.042	0.187
$\gamma = 15\%$	0.461	<b>0.036</b>	0.053	0.038	0.160	<b>0.043</b>	<b>0.190</b>
$\gamma = 20\%$	0.432	<b>0.036</b>	<b>0.054</b>	0.046	0.165	0.040	0.175
$\gamma = 25\%$	0.465	<b>0.036</b>	<b>0.054</b>	0.043	0.168	0.033	0.183