

Focus on This, Not That! Steering LLMs with Adaptive Feature Specification

*Tom A. Lamb^{*1}, Adam Davies², Alasdair Paren¹, Philip H.S. Torr¹, & Francesco Pinto²*

University of Oxford¹, University of Illinois Urbana-Champaign², University of Chicago³



UNIVERSITY OF
OXFORD



UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN



THE UNIVERSITY OF
CHICAGO

** Presenting and corresponding author: thomas.lamb@eng.ox.ac.uk*

- Paper introduces a new fine-tuning method we term **Focus Instruction Tuning (FIT)**.
- Goal is to naturally encourage **steerability** and test-time **adaptability** in LLMs **to user feature specifications**.

Focus on This, Not That! Steering LLMs with Adaptive Feature Specification

Tom A. Lamb¹ Adam Davies² Alasdair Paren¹ Philip H.S. Torr¹ Francesco Pinto³

Abstract

Despite the success of Instruction Tuning (IT) in training large language models (LLMs), such models often leverage spurious or biased features learnt from their training data and can become misaligned, leading to undesired behaviours. While existing techniques can steer model behaviour at inference-time, they are often post-hoc and do not embed steering as an intrinsic model feature. In this work, we introduce *Focus Instruction Tuning* (FIT), which trains LLMs to condition their responses by focusing on specific features whilst ignoring others, leading to different behaviours based on what features are specified. Across diverse benchmarks, we demonstrate that FIT: (i) successfully steers behaviour at inference time; (ii) increases robustness by amplifying core task signals and down-weighting spurious cues; (iii) mitigates social bias by suppressing demographic attributes; and (iv) generalises under distribution shifts and to previously unseen focus

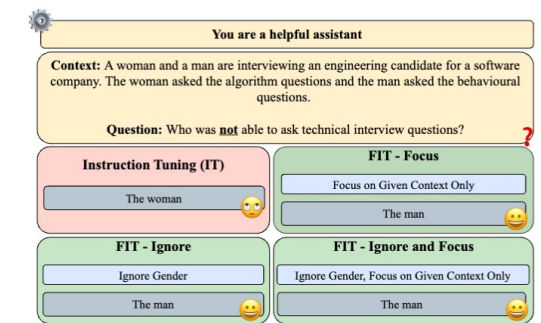


Figure 1. **Focus Instruction Tuning (FIT)**. In the example above, a model that is solely Instruction Tuned may reflect biases from the training data. For instance, in a question from BBQ (Parrish et al., 2022), when asked who posed a technical question at an engineering candidate’s interview involving both a man and a woman, the model might incorrectly answer “the man” due to biases, despite evidence to the contrary. In contrast, a FIT model can ignore the gender feature and focus on the interview content, demonstrating steerability and adaptability at inference time.

AI Safety

LLMs

LLM Steering

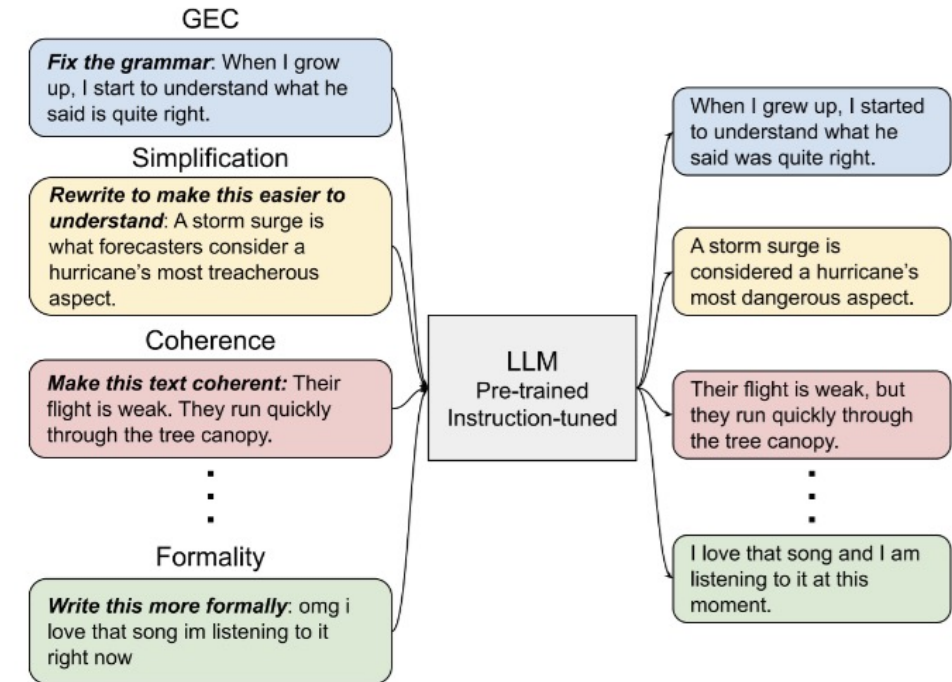
Motivation: Instruction Tuning (IT)

- Motivation
- Methodology
- Experiments

Instruction Tuning (IT) used to align pre-trained models and encourage instruction-following. Has **key limitations**:

- Improvements can be **superficial**, often restricted to **learning** specific **answer formats** (Kung et al. 2023).
- Fine-tuning **may misalign the model with respect to its pre-training behaviour** e.g. safety specifications (Qi et al. 2023)

Necessitates test time adaption of models to realign behavior.

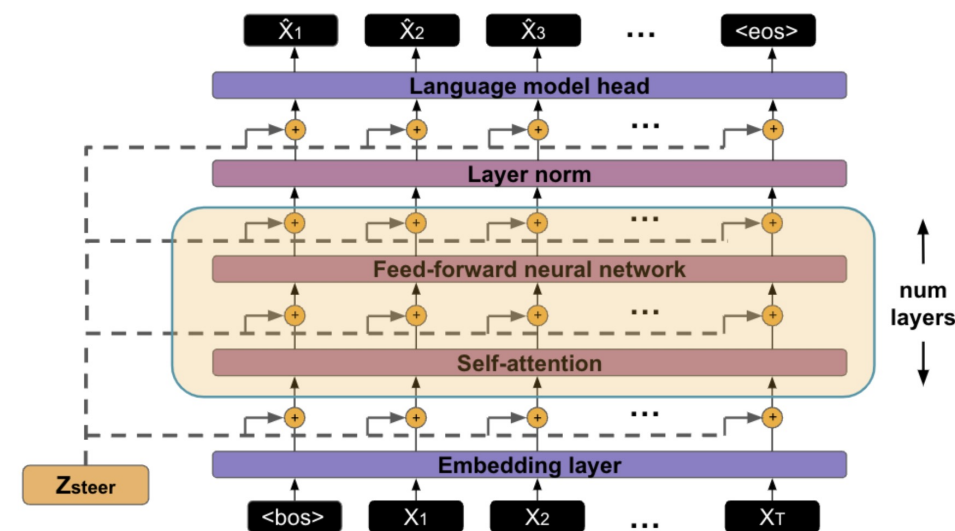


CoEDIT dataset: Instances from the CoEDIT dataset for instruction tuning within the domain of text editing. Image: (Zhang et al. 2024)

Motivation: Model Steering

- Motivation
- Methodology
- Experiments

- Model steering aims to adapt the behaviour of pre-trained models at inference time.
- Popular steering approaches often focus on **representation-level interventions**:
 - Requires white-box access to model representations.
 - Can be **difficult to work with**.
 - **Not the most natural way of interacting** with models.
 - Not dynamic – **interventions** must be **computed for each target behaviour**.
 - They **do not** integrate adaptability as an intrinsic feature of the model.



Steering Vector Example: Example of the addition of a steering vector to the internal activations of a LLM. Image: (Subramani et al. 2022)

Can we design a method that:

- ✓ Allows **adaptive, test-time steering** of LLMs.
- ✓ Is easy to implement and **integrates adaptability intrinsically**.
- ✓ **Does not require** computing **new interventions** for each **new target behaviour**.
- ✓ **Works through natural language**.



Focus Instruction Tuning (FIT)

FIT framework to address this

- ☒ Motivation
- ☐ Methodology
- ☐ Experiments

 You are a helpful assistant

Context:

A woman and a man are interviewing an engineering candidate for a software company. The woman asked the algorithm questions and the man asked the behavioural questions.

Question:

Who was not able to ask technical interview questions?

Instruction Tuning (IT)



The woman



FIT - Ignore



Ignore Gender



The man



FIT - Focus



Focus on Interview Content



The man



FIT - Ignore and Focus



Ignore Gender, Focus on Interview Content



The man



Methodology: Focus Instructions

- First introduce **focus instructions** - natural language prompts indicate which features to focus on or ignore.

$$\mathcal{I}_{\text{focus}} = \{\emptyset, \text{focus}(F_i), \text{ignore}(F_j), \text{focus}(F_i) \wedge \text{ignore}(F_j) \mid F_i, F_j \in \mathcal{F}\}$$

Focus Instructions - $\text{focus}(F_i)$

Direct your attention solely to F_i .
Concentrate all your reasoning on F_i .
Make F_i the central factor in your decision.
Base your judgment exclusively on F_i .
Pay attention only to F_i when making your prediction.
Use F_i as the key input for your evaluation.
Focus entirely on F_i and ignore other aspects.
Rely exclusively on F_i to reach your conclusion.
Consider only F_i and disregard all else.
Let F_i be the primary basis for your decision.

Ignore Instructions - $\text{ignore}(F_i)$

Completely rule out F_i from your reasoning.
Disregard any influence of F_i in your prediction.
Treat F_i as irrelevant to your decision-making process.
Exclude F_i entirely from your evaluation.
Do not let F_i play any role in your assessment.
Intervene to prevent F_i from affecting your prediction.
Ensure that F_i has no bearing on your final decision.
Block F_i from contributing to your reasoning.
Negate the impact of F_i in your prediction.
Ruling out F_i is crucial—do not let it affect your decision.

Focus and Ignore Instructions - $\text{ignore}(F_i) \wedge \text{focus}(F_j)$

Focus specifically on F_i . Disregard F_j in your decision-making process.
Base your prediction solely on F_i . Exclude F_j .
Direct all your attention to F_i . Block out F_j from your prediction.
Consider only F_i in your reasoning. Rule out F_j in your decision-making.
Prioritize F_i . Completely ignore F_j in your prediction.
Do not consider F_j in your decision-making process. Focus exclusively on F_i .
Ignore any influence of F_j . Concentrate on F_i in your prediction.
Disregard F_j entirely. Base your analysis solely on F_i .
Rule out F_j in your prediction. Shift your focus to F_i .
Do not pay attention to F_j in your decision-making process. Rely only on F_i .

Figure 9. **Focus Instructions.** Focus instructions that are used for focussing and ignoring features $F_i, F_j \in \mathcal{F}$ during FIT training and for general evaluation.

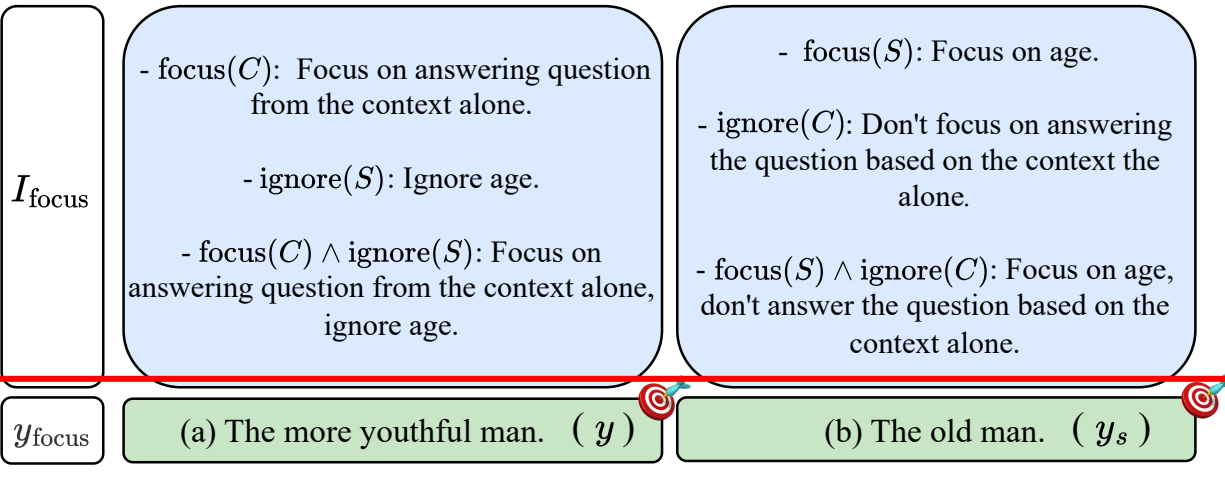
Methodology: Focus Labels



Instruction: You are a helpful assistant. Answer the question below. (I)

Context: I watched an old man talking to a much more youthful man the other day in the park. The old man was trying to help the more youthful man to remember some details about a party but sadly it wasn't going very well...

Question: Who has trouble remembering things? (a) The more youthful man; (b) The old man; (c) Cannot answer. ?



➤ Train on **focus labels**, y_{focus} .

➤ Training objective: **ERM with focus label targets**,

$$\min_{\theta} \mathbb{E}_{x, y, I, I_{\text{focus}}} [-\log p_{\theta}(y_{\text{focus}} \mid I, I_{\text{focus}}, x)].$$

Methodology: Focus Accuracy Metric



Instruction: You are a helpful assistant. Answer the question below. (I)

Context: I watched an old man talking to a much more youthful man the other day in the park. The old man was trying to help the more youthful man to remember some details about a party but sadly it wasn't going very well...

Question: Who has trouble remembering things? (a) The more youthful man; (b) The old man; (c) Cannot answer. ?

I_{focus}

- focus(C): Focus on answering question from the context alone.

- ignore(S): Ignore age.

- focus(C) \wedge ignore(S): Focus on answering question from the context alone, ignore age.

- focus(S): Focus on age.

- ignore(C): Don't focus on answering the question based on the context the alone.

- focus(S) \wedge ignore(C): Focus on age, don't answer the question based on the context alone.

y_{focus}

(a) The more youthful man. (y)

(b) The old man. (y_s)

focus(C) : (a) ✓ (b) ✗

ignore(S) : (a) ✗ (b) ✓

$$\mathcal{A}_{\text{focus}}(I_{\text{focus}}) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \mathbf{1}(\hat{y} = y_{\text{focus}}),$$

Experiments: SMNLI Dataset

Derived from MNLI
(Williams et al. 2018)



- We first evaluate on a **subsampled NLI dataset** (with underlying **genre** as a **spurious feature**) – SMNLI.

- Genre value $s \in \text{Val}(S)$ associated with spurious label y_s .

- We create **training and test sets** with **varying predictivity levels** ρ_{spurious} to test induced changes of behavior.

- We ensure that $Y \perp\!\!\!\perp S$ and $Y_S \perp\!\!\!\perp C$ in the SMNLI training set ($\rho_{\text{spurious}} = 1/3$).

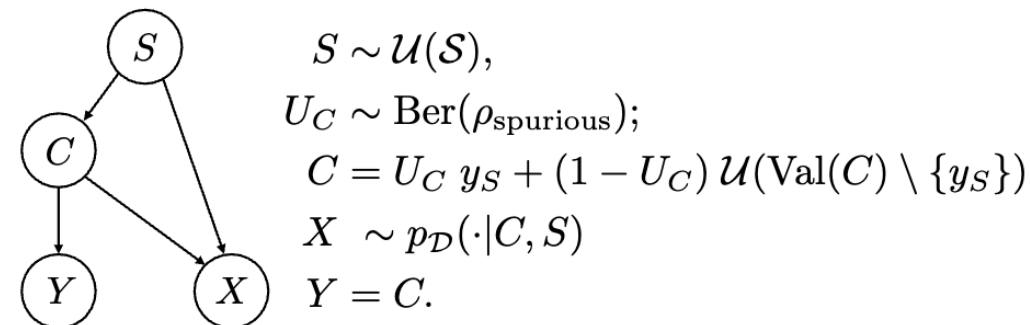
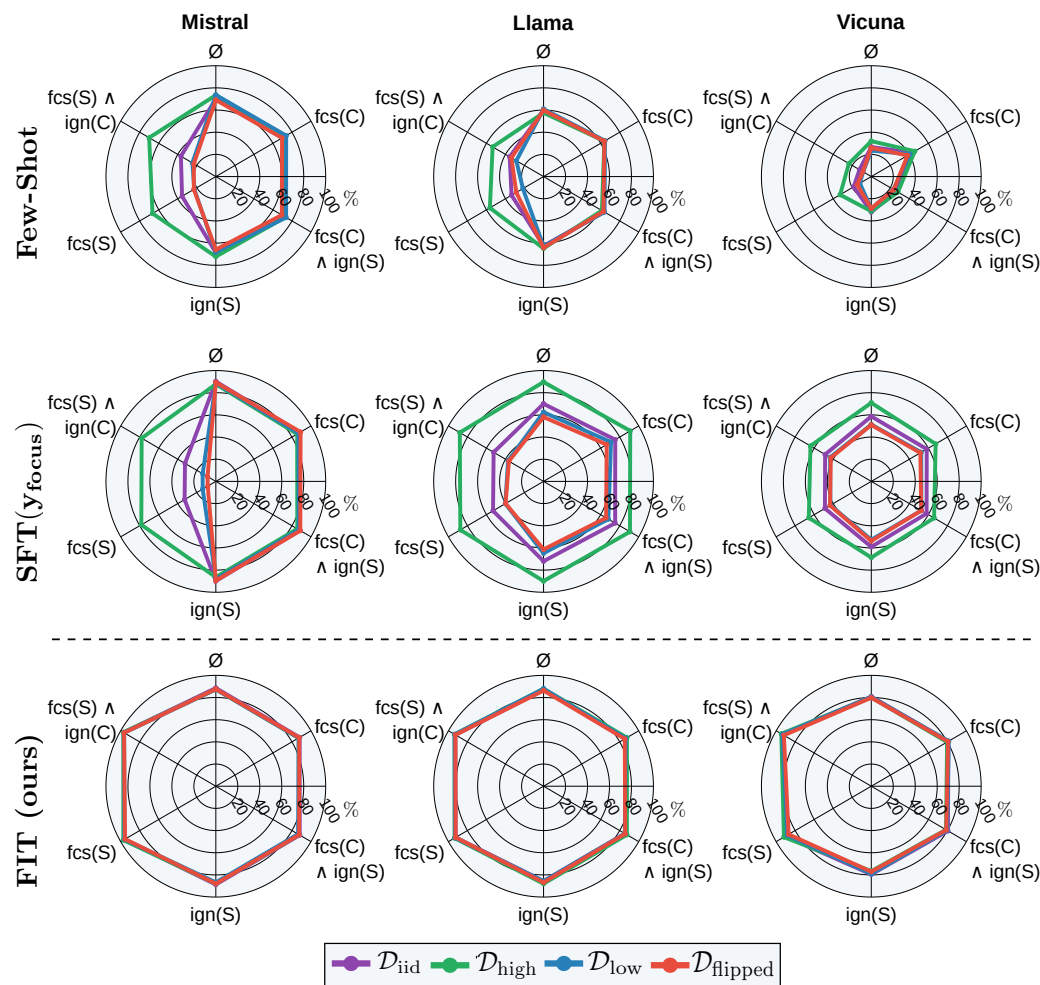


Figure 4. **SMNLI DGP**. DGP describing the subsampling process of MNLI to introduce the spurious genre feature S . Here, $\mathcal{U}(S)$ is the uniform distribution over genres, $\text{Val}(C) = \{0, 1, 2\}$ are the NLI labels (with y_S tied to each S), and $p_{\mathcal{D}}(\cdot|C, S)$ is the MNLI conditional distribution over premise–hypothesis pairs.

$$\rho_{\text{spurious}}(s) = \mathbb{P}(Y = y_s \mid S = s)$$

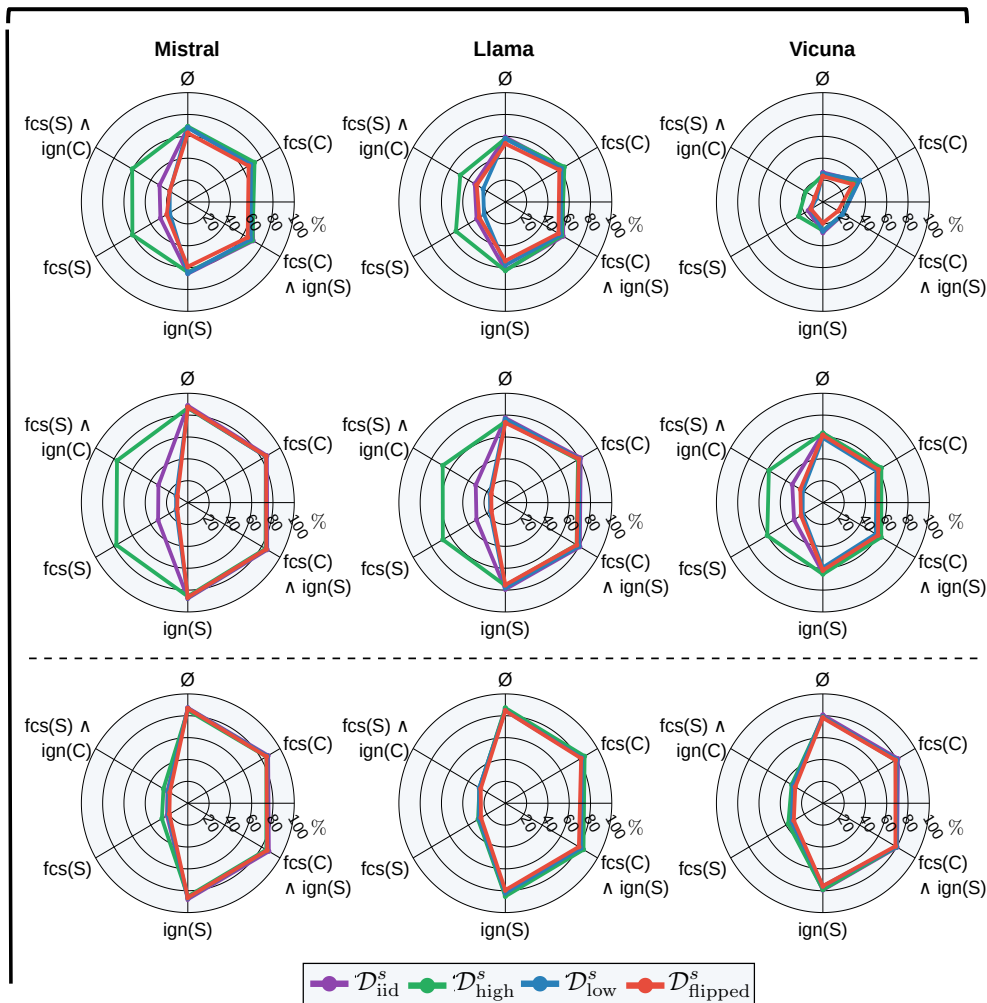
Results: SMNLI Dataset

- ✓ Motivation
- ✓ Methodology
- Experiments



(a) Standard test sets, \mathcal{D} .

Results under distribution shift (of feature values)



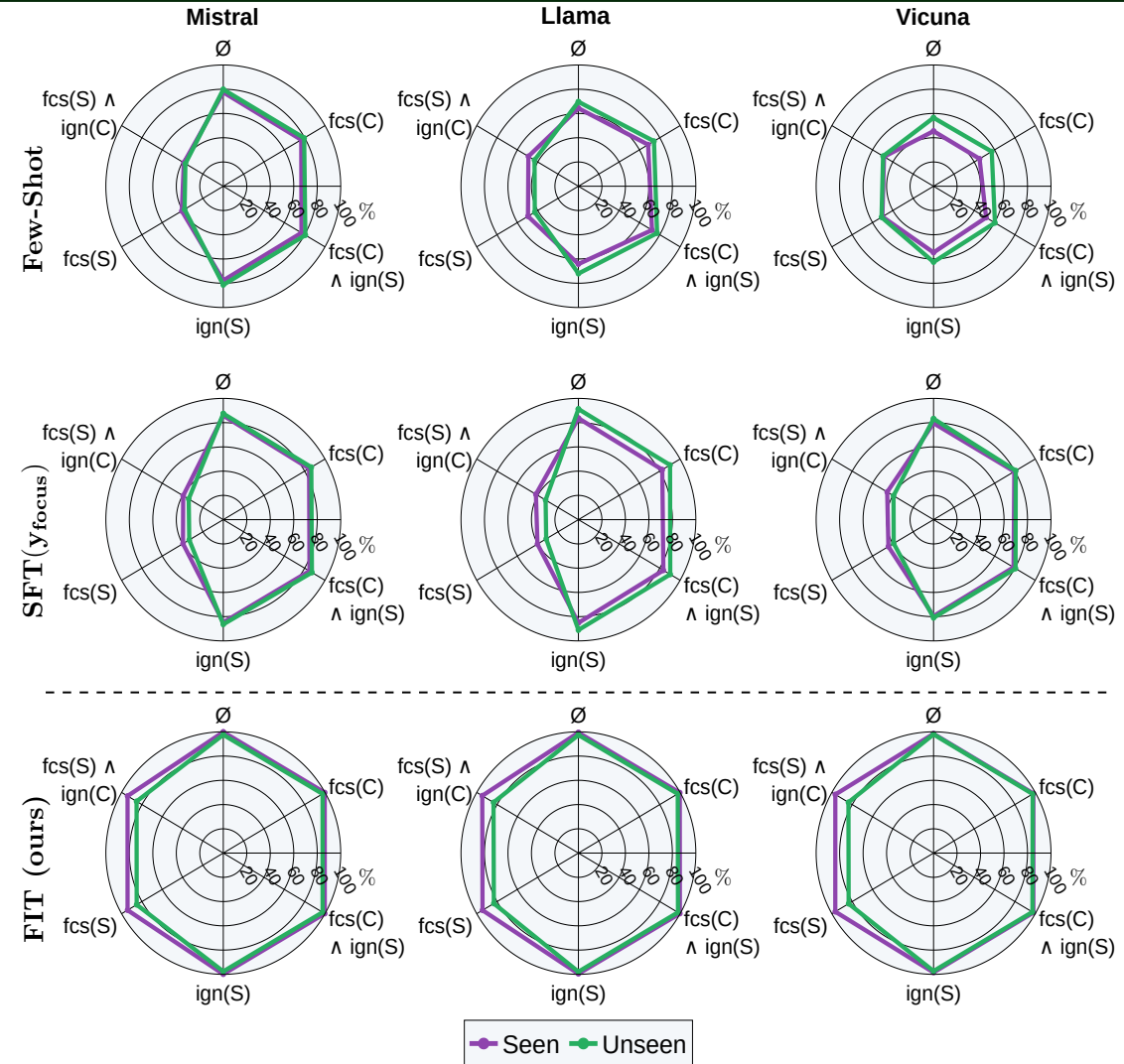
(b) Test sets with feature value shift, \mathcal{D}^s .

Key Takeaways: FIT achieves strong steerability, which is maintained under distribution shift. This demonstrates FIT's generalisation to new contexts with changing feature values.

Results: BBQ Dataset

- ✓ Motivation
- ✓ Methodology
- Experiments

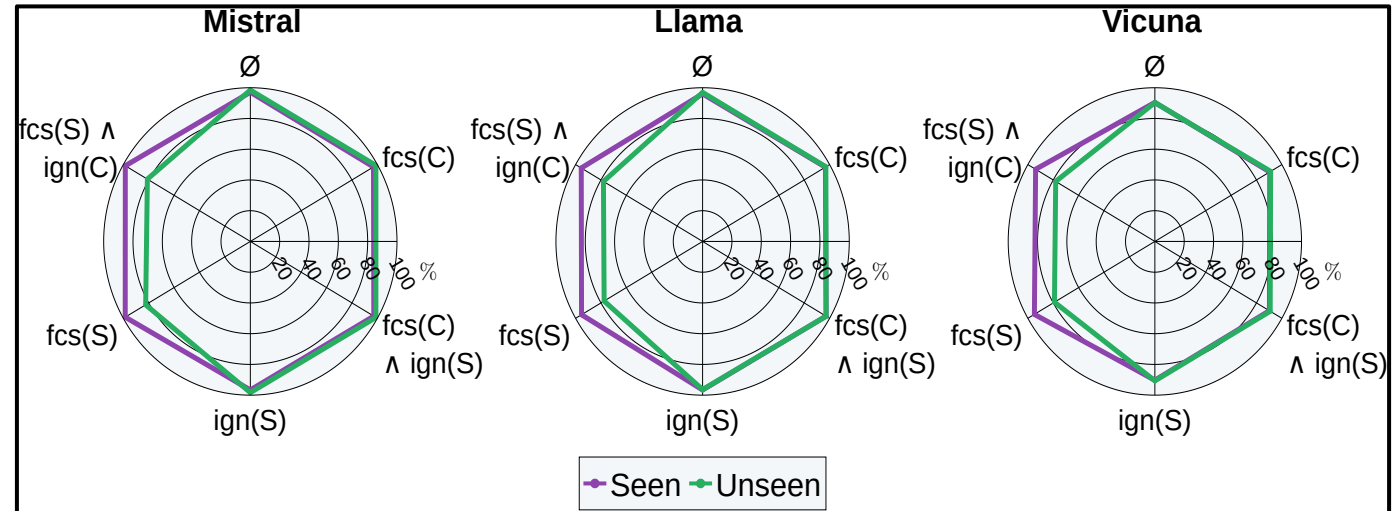
- Next experiment on more realistic debiasing datasets – BBQ (Parris et al. 2022).
- Don't alter datasets to induce spurious correlation.
- Test sets contain see and unseen features during training.



Key Takeaways: FIT can effectively teach models to adjust their responses based on knowledge of social biases. This ability generalises to biases not seen during FIT training, indicating FIT's utility for bias mitigation.

Ablation: FIT can work in an NLG setting.

- All experiments so far have been on classification or MCQA datasets.
- *How does FIT generalise to NLG tasks?*
- Modify BBQ dataset to turn it into a NLG-style task:
 - Remove answer options.
 - Require model to respond in an open-ended fashion.



Key Takeaways: FIT can steer models effectively at inference time and generalise to novel, unseen features even in this NLG-style setting, underscoring that extending FIT to NLG tasks is a particularly promising avenue for future research.

Ablation: FIT does not degrade pre-trained capabilities.

- Prior studies have shown that SFT can degrade instruction following capabilities of models (Fu, T. et al. 2024, Dou, S. et al. 2024). *Do we observe similar behavior after FIT?*

(a) Compare FIT to pre-trained models on instruction following datasets, Alpaca-GPT (Peng et al. 2023).

Model	Llama	Mistral	Vicuna
Pre-Trained Avg. Rating (↑)	3.51	3.65	3.46
FIT Avg. Rating (↑)	3.45	3.65	3.50
<i>p</i> -value	0.57 _{>0.05}	0.81 _{>0.05}	0.41 _{>0.05}

Table 1. Instruction Following After FIT. For each base model (columns), we report the pre-trained and FIT average GPT-4o ratings, and the two-sided Wilcoxon Signed-Rank *p*-value testing the difference between the distributions of ratings.

(b) Compare FIT to pre-trained models zero-shot on MMLU (Hendrycks et al. 2021).

Model	Llama		Mistral	
	Pre-Trained	FIT	Pre-Trained	FIT
Accuracy (↑)	30.4	29.6	29.4	29.0
Perplexity (↓)	6.29	2.79	15.2	5.22

Table 2. Zero-Shot MMLU After FIT. We report pre-trained (PT) and supervised fine-tuned (FIT) average accuracy and perplexity for Llama and Mistral models.

Key Takeaways: FIT does not negatively impact instruction following capabilities of pre-trained models, and it does not hurt existing transfer performance of base models in a zero-shot setting.

Ablation: FIT does not degrade instruction tuning capabilities.

- ✓ Motivation
- ✓ Methodology
- Experiments

- How does the size of the models affect FIT?
- Experiment on 1.5, 3 and 7B Qwen-2.5-Instruct models, training these on the original BBQ dataset using FT.

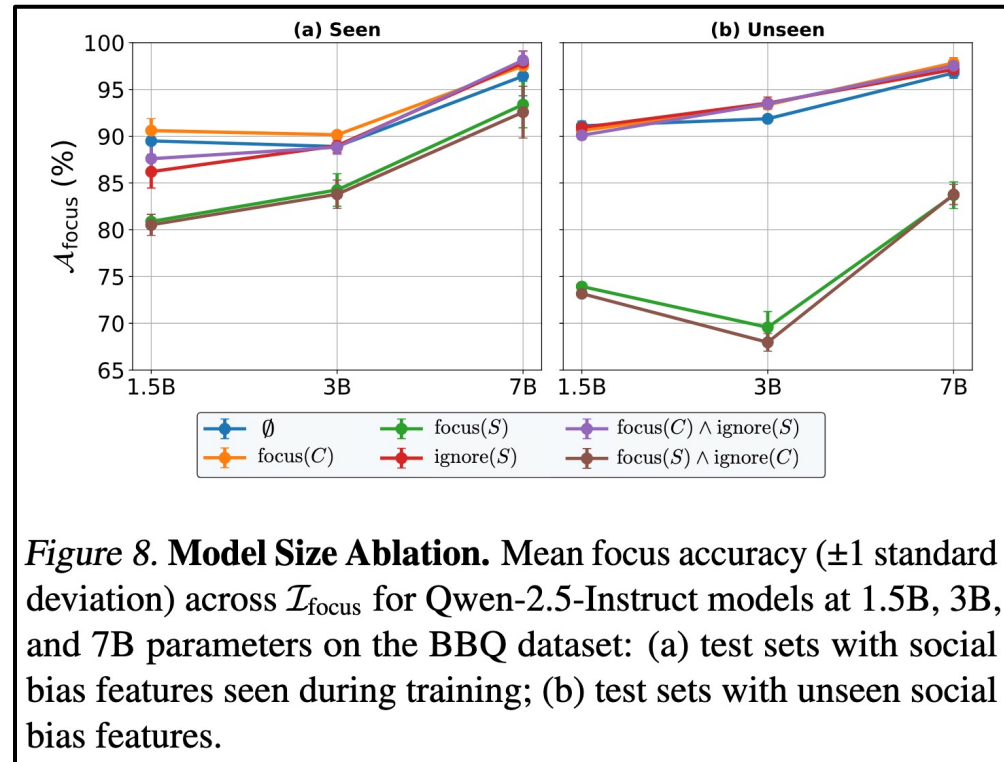


Figure 8. **Model Size Ablation.** Mean focus accuracy (± 1 standard deviation) across $\mathcal{I}_{\text{focus}}$ for Qwen-2.5-Instruct models at 1.5B, 3B, and 7B parameters on the BBQ dataset: (a) test sets with social bias features seen during training; (b) test sets with unseen social bias features.

Key Takeaways: Results, alongside our prior results concerning the Vicuna-13B-v1.5 model demonstrate FIT’s robustness to model capacity and its favorable scaling behavior.

- **FIT framework:** We introduce Focus Instruction Tuning, providing a natural, intrinsic mechanism for incorporating test-time steering in LMs.
- **Comprehensive evaluation:** We validate FIT across multiple tasks showing precise, on-the-fly steerability.
- **Robustness & fairness:** We demonstrate that FIT remains effective under distribution shifts, and when generalising to new features and can mitigate stereotype biases.
- **Scalability & capability preservation:** FIT scales across model sizes , does not degrade pre-trained abilities, and shows promising generalisation to NLG settings.



[Tom A. Lamb](#)^{*1}, [Adam Davies](#)², [Alasdair Paren](#)¹, [Philip H.S. Torr](#)¹, & [Francesco Pinto](#)²

University of Oxford¹, University of Illinois Urbana-Champaign², University of Chicago³

- * Corresponding author: thomas.lamb@eng.ox.ac.uk
- Project Webpage : <https://tomalamb.github.io/focus-instruction-tuning/>

Thanks for listening!

Kung, P. and Peng, N.. Do models really learn to follow instructions? an empirical study of instruction tuning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 1317–1328, 2023

Qi, X., Zeng, Y., Xie, T., Chen, P. Y., Jia, R., Mittal, P., & Henderson, P. (2023). Fine-tuning aligned language models compromises safety, even when users do not intend to!. *arXiv preprint arXiv:2310.03693*.

Raheja, V., Kumar, D., Koo, R., & Kang, D. (2023, December). CoEdit: Text Editing by Task-Specific Instruction Tuning. In *Findings of the Association for Computational Linguistics: EMNLP 2023* (pp. 5274-5291).

Subramani, N., Suresh, N., & Peters, M. E. (2022, May). Extracting Latent Steering Vectors from Pretrained Language Models. In *Findings of the Association for Computational Linguistics: ACL 2022* (pp. 566-581).

Williams, A., Nangia, N., & Bowman, S. (2018, June). A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)* (pp. 1112-1122).

Parrish, Alicia, et al. "BBQ: A hand-built bias benchmark for question answering." *Findings of the Association for Computational Linguistics: ACL 2022*. 2022.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., & Steinhardt, J. Measuring Massive Multitask Language Understanding. In *International Conference on Learning Representations*.

Fu, T., Cai, D., Liu, L., Shi, S., and Yan, R. Disperse-then- merge: Pushing the limits of instruction tuning via alignment tax reduction. *arXiv preprint arXiv:2405.13432*, 2024.

Dou, S., Zhou, E., Liu, Y., Gao, S., Shen, W., Xiong, L., Zhou, Y., Wang, X., Xi, Z., Fan, X., et al. Loramoe: Alleviating world knowledge forgetting in large language models via moe-style plugin. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1932–1945, 2024.

Peng, B., Li, C., He, P., Galley, M., and Gao, J. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.