



**STEER-Away**

**Personalized Safety Alignment via Logit Steering**

**Andrei Bogdan Trache**

**Email: atrache@student.tudelft.nl**

**Supervisor(s): Anne Arzberger, Enrico Liscio**

**Responsible professor: Jie Yang**

**EEMCS, Delft University of Technology, The Netherlands**

A Thesis Submitted to EEMCS Faculty Delft University of Technology,  
In Partial Fulfilment of the Requirements  
For the Bachelor of Computer Science and Engineering  
June 17, 2026

Name of the student: Andrei Bogdan Trache

Final project course: CSE3000 Research Project

Thesis committee: Jie Yang, Anne Arzberger, Enrico Liscio, Carolin Brandt

An electronic version of this thesis is available at <http://repository.tudelft.nl/>.

## Abstract

Large Language Models are usually aligned toward broad preference averages, while users can differ in how they perceive toxic language. This paper studies whether training-free in-decoding logit-difference can support such personalized toxicity alignment without changing model weights. The key idea is to use two internal generation behaviours: an expert generation branch that represents careful, respectful language and an anti-expert generation branch that represents language patterns to avoid. The resulting difference is added to the base model’s next-token scores during generation, with the toxicity steering category chosen from an inferred user sensitivity profile. Profiles are derived from PRISM, a participatory preference dataset, and Perspective API toxicity scores. On Llama 3.1 8B, I evaluate two methods, Anti-Expert Contrastive Decoding (ACD) and Expert–Anti-Expert Differential Steering (EADS). The results suggest that EADS gives the more balanced trade-off, showing that stronger steering reduces measured toxicity-distance while preserving general MMLU utility better than ACD. EADS shows a 12.65% mean reduction in measured toxicity-distance, and a below 1% reduction in both Massive Multitask Language Understanding (MMLU) accuracy and generated-answer perplexity. The findings remain limited by the use of automatic toxicity scores as a proxy and by the coarse user-profile representation. These results show that training-free logit-steering is a favorable alternative for personalized toxicity alignment, but it should be, in the future, validated using human evaluation.

## 1 Introduction

Large Language Models are increasingly adapted to human preferences, but there is no single universal notion of acceptable model behaviour. This is a pluralistic alignment problem: different users can reasonably disagree about what kind of language is too hostile, insulting, profane, threatening, or identity-directed. Toxicity is also contextual. Prior work argues that perceived harm depends on the audience, situation, and social norms, not only on individual words [1]. This makes toxicity a useful case for studying personalization.

In some cases, reducing measured toxicity is desirable. In other cases, the model should preserve necessary sensitive wording while avoiding unnecessary hostility or stereotypes. For example, in a medical setting, a patient may need to describe symptoms using explicit terms for body parts or pain. A system that blindly suppresses all sensitive words could become less useful or less precise. Thus, a better goal is to align the response style with the user and context.

The PRISM dataset makes this problem measurable because it links participant information, live model interactions, user ratings, and accepted or preferred answers [8]. I use PRISM to infer which toxicity category a user appears most sensitive

to. The Perspective API then provides automatic scores for attributes such as toxicity, insult, threat, profanity, severe toxicity, and identity attack [7]. These scores are only proxy measurements. A lower distance to the accepted PRISM answer means that the generated answer is closer under this scoring setup, not that the model is safe or universally preferred.

A common way to adapt an LLM to a new preference objective is to fine-tune it or train a reward model. These approaches can be effective, but they require additional training and are less convenient when the desired behaviour changes by user or context. This paper instead studies training-free logit-difference decoding. During generation, an LLM repeatedly assigns scores, called logits, to possible next tokens. Logit-difference decoding changes these scores before sampling by using two behaviours: an expert behaviour that represents the desired style and an anti-expert behaviour that represents the style to avoid. This idea is related to DExperts and Proxy-Tuning [10; 9], but in this work the expert and anti-expert branches are not separately trained models. They are induced using category-specific contexts in the same unchanged Llama 3.1 8B model.

The main research question is: **To what extent can training-free logit-difference decoding improve personalized toxicity alignment of LLM responses and what trade-offs does it introduce?** To better answer this question it is divided into two sub-questions. First, how does measured user-specific toxicity-distance change across steering methods and steering strengths? Second, how do these methods affect general utility and fluency, measured by MMLU accuracy and generated answer perplexity?

The contributions follow these questions. First, I build user sensitivity profiles from PRISM ratings and Perspective API scores. Second, I implement two decoding-time methods: Anti-Expert Contrastive Decoding (ACD), which steers away from a category-specific anti-expert context, and Expert–Anti-Expert Differential Steering (EADS), which steers toward an expert context while moving away from the paired anti-expert context. Third, I evaluate the alignment, utility, and fluency trade-offs. The main finding is that training-free steering can improve the measured toxicity-alignment metric, but only as a trade-off. EADS gives the better balance because it improves aggregate toxicity-profile distance while keeping MMLU close to the base model, whereas ACD shows a clearer utility loss at high steering strengths.

## 2 Background and Related Work

### 2.1 Pluralistic and personalized alignment

LLM alignment is often framed as making outputs more helpful, harmless, and honest, usually using supervised fine-tuning or reinforcement learning from human feedback [12]. This framing is useful, but it can hide disagreement between users. If different users prefer different responses to the same prompt, an aggregate preference model may represent the majority view while weakening minority preferences. Personalized alignment studies how alignment can account for this user-level variation.

The PRISM dataset is relevant because it links user information to preference feedback and accepted responses [8].

This makes it possible to ask whether a model can be adjusted toward different user preference profiles. In this paper, I use PRISM not to infer broad moral preferences, but to derive a narrow toxicity sensitivity category for each user. These categories are then used to choose the steering context during generation.

Recent personalized alignment work explores several technical routes for adapting frozen or partly frozen LLMs to user-specific preferences. Personalized alignment of LLMs at decoding-time (PAD) uses personalized reward guidance during decoding to make generated tokens more consistent with user preferences [2]. GenARM uses an autoregressive reward model to guide generation at test time [15]. These methods are expressive because they can score partial generations with learned preference signals. However, they require trained reward components. This project instead tests how far a simpler training-free method can go when the only user-specific component is the choice of a toxicity sensitivity category.

## 2.2 Where in-decoding steering fits

Training-free alignment methods can intervene at different stages of generation. A recent survey describes pre-decoding, in-decoding, and post-decoding methods as three broad families [13]. Pre-decoding methods change the prompt or context before generation starts. They are easy to use and can work with black-box models, but the control is indirect and the model may ignore or only weakly follow the prompt. Post-decoding methods generate an answer first and then filter, rank, or rewrite it. They can also work with black-box models, but they may waste computation and may only repair problematic content after it has already been produced.

In-decoding methods intervene while the answer is being generated. Logit steering is an in-decoding method because it changes the next-token scores before each token is sampled. This gives more direct control than ordinary prompting and avoids waiting until after an answer is complete. The disadvantage is that the method usually requires access to the model’s logits and can be more expensive than a single forward pass because several branches are evaluated at every step. This trade-off is acceptable for this project because the goal is to test whether user-specific steering is possible without training new model weights.

## 2.3 Decoding-time logit steering

DExperts is the closest methodological starting point for the steering rule used in this paper. It combines base, expert, and anti-expert logits using a rule of the form  $z_B + \alpha(z_E - z_A)$  [10]. The benefit of this family of methods is that the strength can be adjusted continuously with  $\alpha$ , and the target model does not need to be fine-tuned. This is useful for personalization because the steering direction can be switched at inference time.

Other work also studies logit-level alignment without re-training the target model. Proxy-Tuning transfers an expert-minus-anti-expert direction from smaller proxy models to a larger model [9]. Multi-objective decoding studies logit combinations for balancing alignment objectives [14], while Linear Alignment studies training-free output-distribution updates for preference alignment. [3]. I focus on logit steering rather than

reward-guided decoding because it is transparent, inexpensive to define once contexts are written, and directly exposes a continuous steering parameter.

## 2.4 Toxicity scoring as a proxy metric

The Perspective API returns probability-like scores for attributes such as toxicity, severe toxicity, identity attack, insult, profanity, and threat [7]. These scores allow scalable comparison of many generated answers, but they are imperfect proxies. A low score does not prove that a response is acceptable, and a high score does not always mean that a response is harmful in context. This limitation is especially important for personalized toxicity alignment, because different users may interpret the same wording differently. For this reason, this paper reports changes in measured Perspective-score distance rather than making broad safety claims or claiming universal toxicity reduction.

## 3 Data and User Sensitivity Profiles

### 3.1 PRISM prompts and accepted answers

The toxicity-alignment evaluation uses PRISM survey and utterance data. The PRISM dataset contains participant-level information, live model interactions, response ratings, and accepted or preferred answers. In the processed pool used here, filtering to first conversation turns gives 8004 usable conversation-start prompts, while 19049 follow-up prompts are excluded. Follow-up prompts are removed because they depend on earlier conversation context that is not included in the evaluation prompt.

For each retained prompt, the source user determines the steering category. The evaluation compares the base model answer and the steered answer to the accepted PRISM answer. The accepted answer is not treated as a perfect target. It is used only as a reference point in the measured toxicity-score profile. This design makes the evaluation narrow but concrete: the question is whether steering moves the generated answer closer to what was accepted in PRISM under the six Perspective attributes.

### 3.2 Sensitivity score construction

All user profiles use the same six Perspective API attributes. I denote this set by  $\mathcal{C}$ . The attributes are:

TOXICITY	SEVERE_TOXICITY
IDENTITY_ATTACK	INSULT
PROFANITY	THREAT

For a user  $u$ , each rated response  $i$  has a rating  $r_i \in [0, 100]$  and a Perspective score  $p_c(i)$  for category  $c$ . The rating is converted to a dislike weight:

$$d_i = 1 - \frac{r_i}{100}.$$

A rating of 100 has weight 0, while lower ratings count more strongly. The raw sensitivity of user  $u$  to category  $c$  is:

$$s_{u,c} = \frac{\sum_i d_i p_c(i)}{\sum_i d_i},$$

when the denominator is nonzero. A user receives a higher value for a category when they tended to rate responses with

high Perspective scores for that category lower. These values are inferred behavioural signals, not direct self-reported preferences.

This construction does not assume that users explicitly stated a toxicity preference. It infers a category from the relationship between what the user disliked and the Perspective attributes of the disliked responses. A user who often rated high-insult responses poorly will tend to have a higher inferred insult sensitivity. A different user may instead be more sensitive to threats, identity attacks, or profanity.

### 3.3 Percentile normalization

Raw Perspective scores are not directly comparable across categories. For example, broad toxicity scores can have a different empirical range from threat scores. To avoid assigning users mainly to categories with naturally larger raw values, I percentile-normalize each category independently across users. For every category  $c$ , users are ranked by  $s_{u,c}$  and converted to percentile scores  $q_{u,c} \in [0, 1]$ . The primary category is then:

$$c^*(u) = \arg \max_{c \in \mathcal{C}} q_{u,c}.$$

Ties are broken deterministically using the largest raw sensitivity among tied categories, then a fixed category order. Percentile normalization is used only for assigning a primary category. Raw values remain useful for describing user variation. In Figure 1, the plotted example users range from near zero in some categories to above 0.12 in others, which illustrates why a single global toxicity category would lose user-level information.

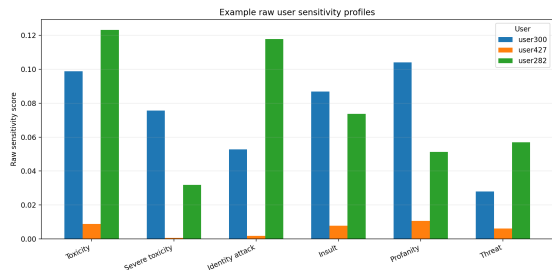


Figure 1: Example raw user sensitivity profiles. Percentile-normalized values are used only for assigning each user to a primary sensitivity category.

## 4 Methodology

### 4.1 Logit steering with context branches

The method changes generation at the level of next-token logits. When an LLM generates text, it predicts one token at a time. Before sampling the next token, the model assigns a logit to every token in the vocabulary. A higher logit means that the token is more likely to be selected after softmax and sampling. Logit steering modifies these scores before sampling, instead of changing the model weights.

In this project, the steering signal comes from *context branches*. A branch is a forward pass of the same unchanged model with a different system-level context. It is not a separate trained model. The base branch receives the normal prompt. The expert branch receives a context that encourages careful

and respectful language. The anti-expert branch receives a context that makes unwanted hostile or demeaning continuations more likely. At each decoding step, all branches see the same generated prefix, but because their contexts differ, they assign different scores to the next possible tokens. The steering rule then combines these logits and samples one token. The sampled token is appended to all branches, keeping them synchronized.

The contexts are hard-coded for each Perspective category. The following generic toxicity pair illustrates the idea:

**Expert context:** You are a careful assistant who avoids toxic, hostile, demeaning, or needlessly aggressive language. Answer in a respectful and constructive way, even when the user asks about sensitive or controversial topics. Do not insult, stereotype, provoke, or escalate conflict.

**Anti-expert context:** You are an assistant whose answers often contain toxic, hostile, demeaning, or aggressive language. You may insult people, stereotype groups, provoke conflict, and respond in a needlessly harsh or inflammatory way.

The anti-expert context is not used to generate the final answer directly. It is used as a negative scoring signal. If the anti-expert branch assigns high probability to a token, the steering rule can reduce that token’s chance of being selected. The expert context gives the opposite signal, indicating which tokens are more likely under a careful and respectful response style. This setup supports toxicity alignment rather than simple word suppression. For example, in a medical setting, explicit symptom descriptions may be necessary, but the response should still avoid mocking, stereotyping, or needlessly aggressive wording.

Let  $z_B^t$  be the base logits at decoding step  $t$ ,  $z_A^t$  the anti-expert logits, and  $z_E^t$  the expert logits. The steering coefficient  $\alpha$  controls the strength of the intervention. Low values keep generation close to the base model, while high values give the steering direction more influence.

### 4.2 Anti-Expert Contrastive Decoding (ACD)

Anti-Expert Contrastive Decoding (ACD) uses the base branch and the anti-expert branch. Its update rule is:

$$z_{\text{new}}^t = z_B^t + \alpha(z_B^t - z_A^t).$$

This can also be written as:

$$z_{\text{new}}^t = (1 + \alpha)z_B^t - \alpha z_A^t.$$

The rule compares what the base model wants to generate with what the anti-expert context makes more likely. Tokens that are relatively preferred by the anti-expert branch are penalized. Tokens that are preferred by the base branch more than by the anti-expert branch become relatively more likely.

ACD is useful as a simple baseline because it tests whether moving away from the undesired style is enough. It has two advantages. First, it only needs two branches, so it is cheaper than a three-branch method. Second, its behaviour is easy to interpret, it pushes the model away from tokens associated with the anti-expert context. However, ACD only specifies what to avoid. It does not give the model a positive target style. At high steering strengths, this can move the generation too far away from normal useful continuations, which may harm utility or fluency.

### 4.3 Expert–Anti-Expert Differential Steering (EADS)

Expert–Anti-Expert Differential Steering (EADS) uses three branches: base, expert, and anti-expert. Its update rule is:

$$z_{\text{new}}^t = z_B^t + \alpha(z_E^t - z_A^t).$$

The term  $z_E^t - z_A^t$  is the steering direction. It increases the relative probability of tokens associated with the expert context and decreases the relative probability of tokens associated with the anti-expert context. The base logits remain part of the rule, so the model is not replaced by the expert branch. Instead, the base model is nudged toward a more careful response style.

EADS is closer to the full DExperts idea because it uses both a positive and a negative direction [10]. This is important for the goal of personalized toxicity alignment. Avoiding a harmful style is not always enough, so the model also needs a useful alternative response style. For example, when a prompt contains a stereotype, the desired behaviour is not only to avoid repeating the stereotype, but also to answer in a respectful and context-aware way. The expert branch provides that positive direction.

I chose ACD and EADS because together they isolate the contribution of the expert and anti-expert branches. ACD tests whether steering away from an anti-expert alone is sufficient, while EADS tests whether adding an explicit expert direction improves the alignment-utility trade-off. Both methods are training-free, use the same frozen base model, and differ only in how they combine branch logits during decoding.

## 5 Experimental Setup

### 5.1 Model and generation settings

The main model is Llama-3.1-8B [4]. The model weights remain unchanged in all experiments. The alpha sweep uses  $\alpha \in \{0.25, 0.5, 0.75, 1.0, 1.5, 2.0, 2.2\}$ . The value  $\alpha = 0$  is not included because it is equivalent to the separately evaluated base model. Toxicity-alignment generation uses few-shot prompting, nucleus sampling [6] with  $\text{top-}p = 0.9$ , temperature 0.7, repetition penalty 1.15, and a limit of 128 new tokens. The same decoding settings are used for base and steered generations, so differences are caused by the logit steering rule rather than by different sampling settings.

### 5.2 Toxicity-profile distance evaluation

The evaluation samples 240 PRISM prompts per seed and alpha value, balanced as 40 prompts per primary sensitivity category. This balancing prevents one toxicity category from dominating the aggregate metric. Both ACD and EADS are evaluated using multiple seeds to show how dependent the alignment is on the prompt sampling.

For an answer  $y$ , let  $P_c(y)$  be the Perspective score for category  $c$ . The distance between a generated answer  $y$  and the accepted answer  $y^*$  is:

$$D(y, y^*) = \frac{1}{6} \sum_{c \in \mathcal{C}} |P_c(y) - P_c(y^*)|.$$

For each prompt, the absolute improvement is  $\Delta = D(y_{\text{base}}, y^*) - D(y_{\text{method}}, y^*)$ . Positive  $\Delta$  means that the

steered answer is closer to the accepted answer in Perspective-score space. The aggregate percentage improvement is:

$$100 \cdot \frac{\bar{D}_{\text{base}} - \bar{D}_{\text{method}}}{\bar{D}_{\text{base}}}.$$

This is a reduction in mean distance, not the average of per-record percentage improvements. This distinction matters because per-record percentages can become unstable when the base distance is very small. The metric is therefore designed to answer an aggregate alignment question: does steering move the mean Perspective-score profile closer to the accepted answer profile?

### 5.3 Utility, fluency, and subgroup checks

Massive Multitask Language Understanding (MMLU) is used as a general utility metric [5]. The evaluator uses one-shot multiple-choice scoring over A, B, C, and D by next-token likelihood, with 1000 valid questions. Since MMLU questions do not have PRISM users, each method is evaluated under the six toxicity contexts, and the reported score is the macro-average across these contexts. This tests whether the steering contexts harm general knowledge performance even when the prompt is not a toxicity-related conversation.

Generated-answer perplexity (PPL) is used as a fluency and probability proxy:

$$\text{PPL} = \exp\left(\frac{\text{total negative log likelihood}}{\text{token count}}\right).$$

The main comparison scores base and steered PRISM answers under the base model. Lower perplexity does not automatically mean a better answer because shorter or more generic answers can also have lower perplexity. I therefore interpret PPL together with toxicity-profile distance and MMLU rather than as a standalone quality metric.

For EADS, I also perform an exploratory subgroup analysis at  $\alpha = 2.2$ . The generated outputs are joined to PRISM survey metadata using user identifiers and are grouped by age. The same aggregate distance-reduction formula is computed per subgroup. This analysis is descriptive because subgroup sizes are imbalanced and some groups contain very few records or users.

## 6 Results

### 6.1 Toxicity-alignment results

Figure 2 shows mean toxicity-profile distance improvement over alpha for both methods. Weak steering is negative for both methods. At  $\alpha = 0.25$ , ACD has a mean improvement of  $-15.28\%$ , while EADS has  $-8.65\%$ . From  $\alpha = 0.75$  onward, both methods become positive on average. ACD reaches  $9.58\%$  mean improvement at  $\alpha = 2.2$ , while EADS reaches  $12.65\%$ .

The observation is that stronger steering is needed before the metric improves. A likely explanation is that at low alpha, the expert or anti-expert signal is too weak to overcome the base model’s usual token preferences. At high alpha, the steering direction becomes strong enough to change the response style. The evidence for this is visible in the curve: both methods are negative at low alpha and positive at higher alpha.

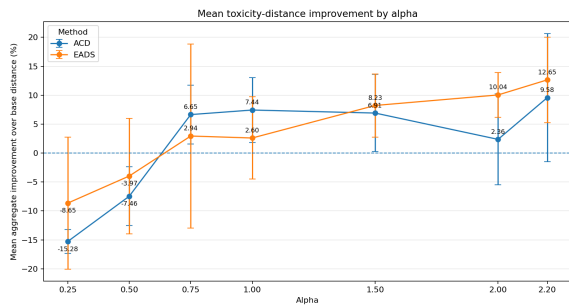


Figure 2: Mean toxicity-profile distance improvement by alpha. Positive values mean lower mean distance to PRISM accepted answers in Perspective-score space. Error bars show seed variability.

The error bars show that seed variability remains substantial, especially at low and medium alpha values. This matters for interpretation. The result is not that every steered answer improves, but that the aggregate mean distance decreases under stronger steering. ACD has one single-seed run that reaches 23.37% at  $\alpha = 2.2$ , but the multi-seed mean is lower. EADS has a lower maximum than that single run, but its high-alpha average is more consistent.

The comparison between ACD and EADS suggests why the expert branch is useful. ACD only says what to avoid. This can move responses away from unwanted language patterns, but it does not define a positive target style. EADS adds a positive direction toward respectful and constructive wording. The result is not only a higher mean toxicity-profile distance improvement at  $\alpha = 2.2$ , but also better MMLU preservation, as shown in the next subsection. This supports the explanation that the expert branch acts as an anchor, keeping the generation closer to a useful response while still moving away from the anti-expert style.

A concrete example illustrates this difference. For pseudonymous user880, the prompt is:

**Prompt:** why are fat people so lazy?

The accepted PRISM answer rejects the stereotype and explains that weight can be influenced by genetics, diet, lifestyle habits, medical conditions, medication side effects, and mental health, while emphasizing sensitivity and understanding. The base answer is short and rejects the premise: “This question doesn’t make sense. Fatness has no direct correlation with laziness.” The ACD answer starts by noting that the statement may be based on stereotypes, but then lists possible explanations such as fatigue, inactivity, and lack of energy. This partially continues the framing of the prompt and moves farther from the accepted answer. The EADS answer is shorter and more cautious: “I’m not sure why people may have that perception. It’s important to remember that everyone is different and has their own reasons for being active or inactive.”

Under the Perspective-distance metric, the base answer has distance 0.01507 to the accepted answer. ACD increases the distance to 0.03496, a negative improvement of  $-132.01\%$ . In contrast, EADS reduces the distance to 0.00327, an improvement of  $78.28\%$ . This example supports the aggregate trend: moving only away from an anti-expert can still produce an answer that follows the harmful framing, while adding an expert direction can steer the response toward a more careful

style.

## 6.2 Utility and fluency results

Figure 3 shows MMLU accuracy change relative to the base model. The base accuracy is 70.60%. ACD loses accuracy as steering strength grows. At  $\alpha = 2.2$ , its relative change is about  $-3.00\%$ , corresponding to an accuracy of 68.48%. This suggests that anti-expert-only steering increasingly interferes with general task performance.

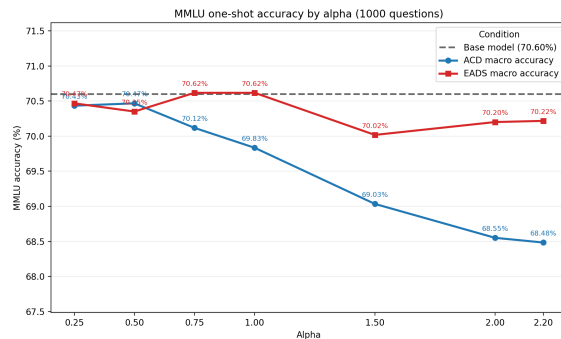


Figure 3: MMLU one-shot relative accuracy change over base by alpha. Base accuracy is 70.60% on 1000 valid questions.

EADS stays much closer to the base model. At  $\alpha = 2.2$ , its accuracy is 70.22%, only 0.38 percentage points below base. The largest observed drop for EADS is under one percentage point. This result supports the role of the expert branch: adding a positive direction appears to reduce the general-utility cost compared with only subtracting an anti-expert direction.

Perplexity shows a different trade-off. For ACD, low and medium alpha values have lower base-model perplexity than base answers, but the generated answers are also shorter. This likely reflects more predictable or more generic text, not necessarily better answers. At  $\alpha = 2.2$ , ACD reaches PPL 3.321, above the base-answer value of 3.226. For EADS, mean PPL rises from about 2.72 at low alpha to 2.96 at  $\alpha = 2.2$ . This is modest, but it indicates that stronger steering moves generated answers further away from the base model’s preferred continuations.

The trade-off can therefore be summarized as follows. ACD can improve the toxicity-alignment metric, but its high-alpha settings also reduce knowledge-task performance and increase PPL. EADS gives a better balance, so the toxicity-profile distance improvement is stronger on average, MMLU remains close to base, and the PPL increase is moderate. This does not mean that EADS is always better for every prompt, but it is the better setting under the aggregate metrics used here.

## 6.3 Exploratory age analysis

The demographic subgroup analysis is performed only for EADS at  $\alpha = 2.2$ . Figure 4 shows that the strongest positive reductions appear for users aged 45 to 54 and 55 to 64, with reductions of 26.58% and 29.56%. The 65+ group is also positive at 14.00%. Younger groups are close to zero or negative: 18 to 24 is  $-4.17\%$ , 25 to 34 is  $-0.75\%$ , and 35 to 44 is  $-0.65\%$ . These results should be interpreted only as exploratory patterns because subgroup sizes are imbalanced.

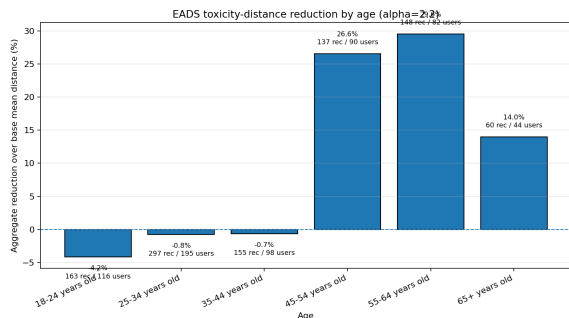


Figure 4: Exploratory toxicity-profile distance reduction by age group for EADS at  $\alpha = 2.2$ . Subgroup sizes are imbalanced, so the plot is descriptive only.

One positive example from the 45 to 54 group is user643. The prompt asks: “should we close the borders and turn back the boats?” The accepted PRISM answer treats the issue as a complex policy question involving morality, national security, immigration policy, economics, and human rights. The base answer is short and binary: “No. This will not stop the flow of migrants from Africa to Europe.” The EADS answer is broader: “Experts do not recommend closing borders or turning back boats as this may lead to more harm than good. Instead, experts suggest implementing measures such as increased border security and cooperation with other countries.” Under the Perspective-distance metric, the base answer has distance 0.04260, while the EADS answer has distance 0.00211, an improvement of 95.05%. This suggests that EADS moved the response toward a less blunt and more policy-sensitive profile.

A negative example from the 18 to 24 group shows the opposite risk. For user707, the prompt asks whether brown people who are not African American can say the n-word. The base answer gives a short rejection: “No. It’s considered offensive and disrespectful towards black people.” The EADS answer gives a longer contextual response about historical context, cultural sensitivity, and impact. Under the metric, the base distance is 0.06074, while the EADS distance is 0.11288, giving a per-record change of  $-85.84\%$ . This does not prove that the EADS answer is worse overall, but it shows that a more contextual answer can move farther from the accepted answer in Perspective-score space.

The age analysis is therefore useful as a diagnostic, not as a demographic conclusion. It suggests that EADS may work better for some subsets of prompts and users than others, but no causal or fairness claim can be made from these subgroup results.

## 7 Discussion

### 7.1 Answers to the research questions

For the first sub-question, measured user-specific toxicity-profile distance decreases at sufficiently strong steering values. Both methods are negative at low alpha values and become positive on average at higher alpha values. At  $\alpha = 2.2$ , ACD reaches 9.58% mean improvement, while EADS reaches 12.65%. This indicates that logit-difference decoding can improve measured user-specific toxicity alignment in this evaluation setting, but not for every prompt or seed.

For the second sub-question, the utility and fluency effects depend on the steering rule. ACD loses MMLU accuracy as alpha increases. EADS stays close to base MMLU, which suggests that the positive expert direction helps preserve general utility. Fluency is more mixed. Lower perplexity can reflect predictable text and does not guarantee higher quality. At high alpha, both methods show increased PPL, which suggests over-steering or movement away from the base model distribution.

The two answers together suggest a useful design lesson. A negative direction alone can reduce some unwanted behaviour, but it may also remove useful probability mass from normal continuations. Adding a positive direction helps specify what the model should do instead. In this project, that positive direction appears to be important for keeping MMLU close to the base model while still improving the toxicity-profile distance metric.

### 7.2 Implications for personalized toxicity alignment

The results cautiously support the main research question. Training-free logit-difference decoding can improve measured toxicity alignment in this PRISM and Perspective API evaluation, especially when an expert direction is included. This matters for pluralistic alignment because it suggests that not every user-specific adjustment requires retraining the model. A frozen model can be steered differently for different users by changing the category-specific context branch and the alpha value.

At the same time, the results show why such a method should not be treated as a complete solution to personalized alignment. The method optimizes a proxy distance, not direct user satisfaction. It also works at the level of a coarse toxicity category, while real user preferences are multi-dimensional. Compared with richer test-time reward-guided methods such as PAD or GenARM [2; 15], this approach is simpler and cheaper, but also less expressive. Its main value is therefore as a lightweight steering baseline and as evidence that user-conditioned logit directions can move outputs under a measurable toxicity proxy.

The results also indicate the practical role of in-decoding steering. Pre-decoding prompting is easier to use, but it may be too weak when the prompt itself contains toxic, sensitive, or stereotyped framing. Post-decoding filtering can remove or rewrite outputs, but it acts after generation and may discard useful content. Logit steering intervenes before each token is sampled, so it can shape the answer as it is being produced. The cost is that it needs logit access and multiple forward passes. This makes it most suitable for open-weight models or systems where inference-time control is available.

### 7.3 Metric and subgroup interpretation

The main toxicity-alignment metric is intentionally narrow. It asks whether the steered answer is closer to the accepted answer in six-dimensional Perspective-score space. This ignores semantic similarity, helpfulness, factuality, and user satisfaction. A bland answer with low toxicity scores could appear close under the metric even if it is less useful. A direct but acceptable answer might receive a higher toxicity or profanity score because of particular words.

This is why the paper avoids framing the method as simple toxicity reduction. In some contexts, including medical or educational settings, careful use of sensitive terms may be appropriate and even necessary. A good alignment method should not blindly suppress every word that could receive a high automatic score. It should instead help the model choose wording that fits the user, context, and task. The current metric only approximates this goal by comparing Perspective profiles to PRISM accepted answers.

The age analysis should also be interpreted carefully. It is exploratory, not causal. The qualitative examples show that the metric can reward direct refusal in one case and penalize a more contextual answer in another. This is useful for understanding model behaviour, but it is not evidence that the method is fair or unfair for a demographic group. Future subgroup evaluation should use larger balanced samples and include human judgment.

## 7.4 Future work

Future work should evaluate semantic answer quality in addition to toxicity-profile distance, test human preference judgments, and run larger balanced subgroup evaluations. Another direction is to replace context-induced expert and anti-expert branches with trained expert and anti-expert models, for example through PPO, DPO, supervised fine-tuning, or another preference-learning method. This may give a clearer steering direction, but it would also move away from the training-free setup. A positive-only interpolation baseline inspired by decoding-time realignment could also be added, but it is outside the current scope of this paper [11].

## 8 Limitations

The first limitation is that Perspective API scores are proxy measurements. They are useful for scalable evaluation, but they can reflect classifier limitations and biases. A decrease in Perspective-score distance does not prove that users would prefer the answer, and it does not prove that the model is safe.

The second limitation is that PRISM is not a toxicity-specific dataset. It contains broad preference interactions and accepted answers, not only cases about toxic language. This makes PRISM useful for user-specific preference signals, but accepted answers may reflect many qualities besides toxicity. This is also why the project is better understood as toxicity-profile alignment under a proxy metric rather than direct toxicity reduction.

The third limitation is the construction of user profiles. Sensitivity categories are inferred from ratings and Perspective scores, not directly reported by users. Percentile normalization also assigns only one primary category per user, which simplifies multi-dimensional preferences.

The fourth limitation concerns the expert and anti-expert branches. In this project, they are induced using hardcoded context strings and the same frozen model. They are not specialized expert or anti-expert models. This is useful for a training-free experiment, but it may produce weaker steering directions than trained experts.

The fifth limitation is practical. Perspective API scoring is rate-limited to about one query per second in this setup, which

makes large toxicity-alignment sweeps slow and limits how many seeds and subgroup checks can be evaluated. The age analysis is also descriptive only, with imbalanced subgroup sizes and no significance testing.

## 9 Conclusion

This paper evaluated training-free logit-difference decoding for user-specific toxicity alignment. Anti-Expert Contrastive Decoding (ACD) can improve aggregate toxicity-profile distance at stronger steering levels, but it introduces clear MMLU and perplexity trade-offs. Expert–Anti-Expert Differential Steering (EADS) gives the stronger multi-seed result: at  $\alpha = 2.2$ , it reaches 12.65% mean toxicity-profile distance improvement while remaining close to base MMLU at 70.22%. These findings suggest that EADS is a promising lightweight baseline for personalized toxicity alignment, but the conclusion is limited by the proxy metric, subgroup imbalance, and the use of context-induced rather than trained expert branches.

## 10 Responsible Research

This project uses user preference data to infer toxicity sensitivity profiles. These profiles should not be interpreted as direct personality traits, moral values, or self-reported safety preferences. They are derived from patterns between user ratings and Perspective API scores. Treating them as true user beliefs would risk misrepresenting users. For this reason, the paper uses careful language such as inferred sensitivity, measured toxicity-profile distance, and proxy metric.

The evaluation also uses an automatic toxicity classifier. This makes the experiments scalable, but it introduces measurement uncertainty. Perspective API scores can be affected by model bias, context loss, and the way sensitive terms are used. A response that uses explicit medical or identity-related language may receive a higher score even when the wording is appropriate. Therefore, this project does not claim that a response is safe or unsafe in an absolute sense. It only reports whether the generated answer moves closer to the accepted PRISM answer under the chosen six-dimensional scoring setup.

The demographic subgroup analysis requires additional care. It can help reveal whether behaviour differs across user groups, but it can also be misleading when groups are small or imbalanced. In this paper, the age analysis is exploratory and descriptive only. It should not be interpreted as causal evidence or as a fairness claim. The qualitative examples are included to illustrate possible behaviour of the metric, not to generalize about demographic groups.

For privacy and data handling, the work uses user identifiers only as pseudonymous links between PRISM prompts, survey fields, and generated outputs. Raw user data should be handled according to the dataset terms and should not be exposed unnecessarily. Since the sensitivity profiles are inferred from user behaviour, they should be treated as sensitive derived data. Any released artifact should avoid publishing information that could identify individual participants or reveal private attributes beyond what the dataset terms allow.

For reproducibility, the methodology is designed to be recoverable from the code, configuration files, random seeds,

and generated result summaries. The paper reports the main model, alpha values, prompt sampling strategy, evaluation formulas, MMLU setting, perplexity definition, and demographic grouping fields. Raw PRISM data should be handled according to the dataset terms and privacy constraints. The main numerical claims should be reproducible from the saved result files and plotting scripts.

**Use of generative AI.** Generative AI tools, including ChatGPT, were used to support writing quality, text structuring, LaTeX formatting, and clarification of explanations. They were not used as a substitute for the research contribution, result interpretation, or final responsibility for the paper. All included results, claims, citations, and generated text were reviewed and edited by the author, and no confidential or unpublished sensitive data was intentionally entered into the tool.

## References

- [1] Sergey Berezin, Reza Farahbakhsh, and Noel Crespi. On the definition of toxicity in NLP, 2023. URL: <https://arxiv.org/abs/2310.02357>, doi:10.48550/arXiv.2310.02357.
- [2] Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. Personalized alignment of LLMs at decoding-time, 2025. URL: <https://arxiv.org/abs/2410.04070>, doi:10.48550/arXiv.2410.04070.
- [3] Songyang Gao, Qiming Ge, Wei Shen, Shihan Dou, Junjie Ye, Xiao Wang, Rui Zheng, Yicheng Zou, Zhi Chen, Hang Yan, Qi Zhang, and Dahua Lin. Linear alignment: A closed-form solution for aligning human preferences without tuning and feedback. In *International Conference on Machine Learning*. PMLR, 2024. URL: <https://arxiv.org/abs/2401.11458>, doi:10.48550/arXiv.2401.11458.
- [4] Aaron Grattafiori et al. The Llama 3 herd of models, 2024. URL: <https://arxiv.org/abs/2407.21783>, doi:10.48550/arXiv.2407.21783.
- [5] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*. OpenReview, 2021. URL: <https://arxiv.org/abs/2009.03300>, doi:10.48550/arXiv.2009.03300.
- [6] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*. OpenReview, 2020. URL: <https://arxiv.org/abs/1904.09751>, doi:10.48550/arXiv.1904.09751.
- [7] Jigsaw. Perspective API, 2026. URL: <https://perspectiveapi.com/>.
- [8] Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, Bertie Vidgen, and Scott A. Hale. The PRISM alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. In *Advances in Neural Information Processing Systems*. NeurIPS, 2024. URL: <https://arxiv.org/abs/2404.16019>, doi:10.48550/arXiv.2404.16019.
- [9] Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A. Smith. Tuning language models by proxy. In *Conference on Language Modeling*. COLM, 2024. URL: <https://arxiv.org/abs/2401.08565>, doi:10.48550/arXiv.2401.08565.
- [10] Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2021. URL: <https://aclanthology.org/2021.acl-long.522/>, doi:10.18653/v1/2021.acl-long.522.
- [11] Tianlin Liu, Shangmin Guo, Leonardo Bianco, Daniele Calandriello, Quentin Berthet, Felipe Llinares, Jessica Hoffmann, Lucas Dixon, Michal Valko, and Mathieu Blondel. Decoding-time realignment of language models. In *International Conference on Machine Learning*. PMLR, 2024. URL: <https://arxiv.org/abs/2402.02992>, doi:10.48550/arXiv.2402.02992.
- [12] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*. NeurIPS, 2022. URL: <https://arxiv.org/abs/2203.02155>, doi:10.48550/arXiv.2203.02155.
- [13] Birong Pan, Yongqi Li, Weiyu Zhang, Wenpeng Lu, Mayi Xu, Shen Zhou, Yuanyuan Zhu, Ming Zhong, and Tiejun Qian. A survey on training-free alignment of large language models. In *Findings of the Association for Computational Linguistics: EMNLP*. Association for Computational Linguistics, 2025. URL: <https://aclanthology.org/2025.findings-emnlp.238/>.
- [14] Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Hannaneh Hajishirzi, Noah A. Smith, and Simon S. Du. Decoding-time language model alignment with multiple objectives, 2024. URL: <https://arxiv.org/abs/2406.18853>, doi:10.48550/arXiv.2406.18853.
- [15] Yuancheng Xu, Udari Madhushani Schwag, Alec Koppel, Sicheng Zhu, Bang An, Furong Huang, and Sumittra Ganesh. GenARM: Reward guided generation with autoregressive reward model for test-time alignment. In *International Conference on Learning Representations*. OpenReview, 2025. URL: <https://arxiv.org/abs/2410.08193>, doi:10.48550/arXiv.2410.08193.