

UNINTENTIONAL UNALIGNMENT: LIKELIHOOD DISPLACEMENT IN DIRECT PREFERENCE OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Direct Preference Optimization (DPO), and its numerous variants, are increasingly used for aligning language models. Although they are designed to teach a model to generate preferred responses more frequently relative to dispreferred responses, prior work has observed that the likelihood of preferred responses often decreases during training. The current work sheds light on the causes and implications of this counter-intuitive phenomenon, which we term *likelihood displacement*. We demonstrate that likelihood displacement can be *catastrophic*, shifting probability mass from preferred responses to semantically opposite ones. As a simple example, training a model to prefer `No` over `Never` can sharply increase the probability of `Yes`. Moreover, when aligning the model to refuse unsafe prompts, we show that such displacement can *unintentionally lead to unalignment*, by shifting probability mass from preferred refusal responses to harmful responses (*e.g.*, reducing the refusal rate of Llama-3-8B-Instruct from 74.4% to 33.4%). We theoretically characterize that likelihood displacement is driven by preferences that induce similar embeddings, as measured by a *centered hidden embedding similarity (CHES)* score. Empirically, the CHES score enables identifying which training samples contribute most to likelihood displacement in a given dataset. Filtering out these samples effectively mitigated unintentional unalignment in our experiments. More broadly, our results highlight the importance of curating data with sufficiently distinct preferences, for which we believe the CHES score may prove valuable.

1 INTRODUCTION

To ensure that language models generate safe and helpful content, they are typically aligned based on pairwise preference data. One prominent alignment method, known as *Reinforcement Learning from Human Feedback (RLHF)* (Ouyang et al., 2022), requires fitting a reward model to a dataset of human preferences, and then training the language model to maximize the reward via RL. While often effective for improving the quality of generated responses (Bai et al., 2022a; Achiam et al., 2023; Touvron et al., 2023), the complexity and computational costs of RLHF motivated the rise of *direct preference learning* methods such as DPO (Rafailov et al., 2023).

Given a prompt x , DPO and its variants (*e.g.*, Azar et al. (2024); Tang et al. (2024); Xu et al. (2024a); Meng et al. (2024)) eschew the need for RL, by directly teaching a model π_θ to increase the margin between the log probabilities of a preferred response y^+ and a dispreferred response y^- . While intuitively these methods should increase the probability of y^+ while decreasing that of y^- , several recent works observed that the probabilities of both y^+ and y^- tend to *decrease* over the course of training (Pal et al., 2024; Yuan et al., 2024; Rafailov et al., 2024; Tajwar et al., 2024; Pang et al., 2024; Liu et al., 2024). We term this phenomenon *likelihood displacement* — see Figure 1.

When the probability of y^+ decreases, the probability of some other, possibly undesirable, response must increase. However, despite the prevalence of likelihood displacement, there is limited understanding as to why it occurs and what its implications are. The purpose of this work is to address these gaps. Through theory and experiments, we characterize mechanisms driving likelihood displacement, demonstrate that it can lead to surprising failures in alignment, and provide preventative guidelines. Our experiments cover models of different families and scales, including OLMo-1B (Groeneveld et al., 2024), Gemma-2B (Team et al., 2024), and Llama-3-8B (Dubey et al., 2024). The main contributions are listed below.

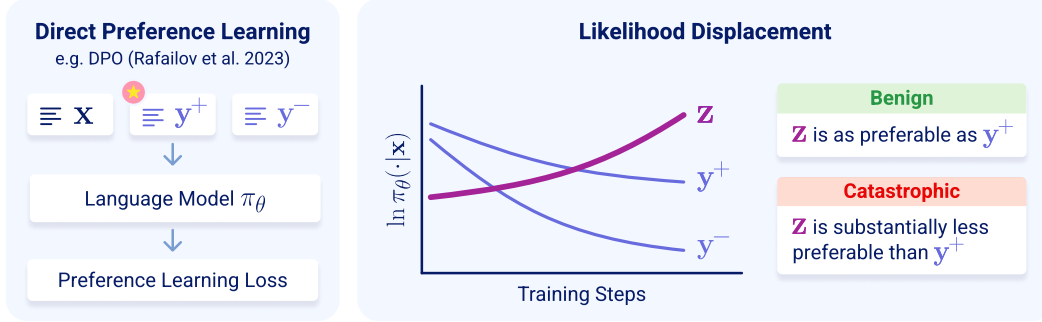


Figure 1: **Illustration of likelihood displacement in direct preference learning.** For a prompt \mathbf{x} , direct preference learning aims to increase the probability that a model π_θ assigns to a preferred response \mathbf{y}^+ relative to a dispreferred response \mathbf{y}^- . *Likelihood displacement* refers to the counter-intuitive phenomenon where, while the gap between $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ and $\ln \pi_\theta(\mathbf{y}^-|\mathbf{x})$ increases, they both decrease. If the responses increasing instead in probability (depicted by \mathbf{z}) are as preferable as \mathbf{y}^+ (e.g., \mathbf{z} is semantically similar to \mathbf{y}^+), then the likelihood displacement is *benign*. However, if the probability mass goes to responses that are substantially less preferable than \mathbf{y}^+ (e.g., \mathbf{z} is semantically opposite to \mathbf{y}^+), then we say that it is *catastrophic*.

- **Likelihood displacement can be catastrophic even in simple settings.** We demonstrate that, even when training on just a single prompt whose preferences \mathbf{y}^+ and \mathbf{y}^- consist of a single token each, likelihood displacement is pervasive (Section 3). Moreover, the tokens increasing most in probability at the expense of \mathbf{y}^+ can be semantically opposite to it. For example, training a model to prefer $\mathbf{y}^+ = \text{No}$ over $\mathbf{y}^- = \text{Never}$ often sharply increases the probability of Yes (Table 1). This stands in stark contrast to prior work attributing likelihood displacement to different complexities in the preference learning pipeline (Tajwar et al., 2024; Pal et al., 2024; Rafailov et al., 2024), and emphasizes the need to formally characterize its underlying causes.
- **Theory: likelihood displacement is determined by the model’s embedding geometry.** We analyze the evolution of $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ during gradient-based training (Section 4). Our theory reveals that likelihood displacement is governed by the (static) token unembeddings and (contextual) hidden embeddings of \mathbf{y}^+ and \mathbf{y}^- . In particular, it formalizes intuition by which the more similar \mathbf{y}^+ and \mathbf{y}^- are the more $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ tends to decrease.
- **Identifying sources of likelihood displacement.** Based on our analysis, we derive a (model-aware) measure of similarity between preferences, called the *centered hidden embedding similarity* (CHES) score (Definition 2). We demonstrate that the CHES score accurately identifies which training samples contribute most to likelihood displacement in a given dataset (e.g., UltraFeedback (Cui et al., 2024) and AlpacaFarm (Dubois et al., 2024)), whereas other similarity measures relying on hidden embeddings or token-level cues do not (Section 5).
- **Unintentional unalignment due to likelihood displacement.** To demonstrate the potential uses of the CHES score, we consider training a language model to refuse unsafe prompts via preference learning (Section 6). We find that likelihood displacement can *unintentionally unalign* the model, by causing probability mass to shift from preferred refusal responses to responses that comply with unsafe prompts! For example, the refusal rate of Llama-3-8B-Instruct drops from 74.4% to 33.4% over the SORRY-Bench dataset (Xie et al., 2024). We then show that filtering out samples with a high CHES score prevents such unintentional unalignment, and does so more effectively than adding a supervised finetuning term to the loss (e.g., as done in Pal et al. (2024); Xu et al. (2024a); Pang et al. (2024); Liu et al. (2024)).

Our results highlight the importance of curating data with sufficiently distinct preferences. We believe that the CHES score introduced by our theory may prove valuable for achieving this goal.¹

2 PRELIMINARIES

2.1 LANGUAGE MODELS

Let \mathcal{V} be a vocabulary of tokens. Modern language models consist of two parts: (i) a neural network (e.g., Transformer (Vaswani et al., 2017)) that intakes a sequence of tokens $\mathbf{x} \in \mathcal{V}^*$ and produces a

¹We discuss related work throughout and defer a concentrated account to Appendix A.

hidden embedding $\mathbf{h}_{\mathbf{x}} \in \mathbb{R}^d$; and (ii) a token unembedding matrix $\mathbf{W} \in \mathbb{R}^{|\mathcal{V}| \times d}$ that converts the hidden embedding into logits. The logits are then passed through a softmax to compute a distribution over tokens that can follow \mathbf{x} . For assigning probabilities to sequences $\mathbf{y} \in \mathcal{V}^*$, a language model π_θ operates autoregressively, *i.e.*:

$$\pi_\theta(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^{|\mathbf{y}|} \pi_\theta(\mathbf{y}_k|\mathbf{x}, \mathbf{y}_{\leq k-1}) = \prod_{k=1}^{|\mathbf{y}|} \text{softmax}(\mathbf{W}\mathbf{h}_{\mathbf{x}, \mathbf{y}_{\leq k}})_{\mathbf{y}_k}, \quad (1)$$

where θ stands for the model’s parameters (*i.e.* the parameters of the neural network and the unembedding matrix \mathbf{W}), and $\mathbf{y}_{\leq k}$ denotes the first $k - 1$ tokens of \mathbf{y} .

2.2 DIRECT PREFERENCE LEARNING

Preference data. We consider the widely adopted direct preference learning pipeline, which relies on pairwise comparisons (*cf.* Rafailov et al. (2023)). Specifically, we assume access to a preference dataset \mathcal{D} containing samples $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, where \mathbf{x} is a prompt, \mathbf{y}^+ is a preferred response to \mathbf{x} , and \mathbf{y}^- is a dispreferred response to \mathbf{x} . The preferred and dispreferred responses can be obtained by generating two candidate responses from the model (*i.e.* on-policy), and labeling them via human or AI raters (*cf.* Ouyang et al. (2022); Bai et al. (2022b)). Alternatively, they can be taken from some static dataset (*i.e.* off-policy). Our analysis and experiments capture both cases.

Supervised finetuning (SFT). Preference learning typically includes an initial SFT phase, in which the model is finetuned via the standard cross-entropy loss. The sequences used for SFT can either be independent of the preference dataset \mathcal{D} (Touvron et al., 2023) or consist of prompts and preferred responses from \mathcal{D} , *i.e.* of $\{(\mathbf{x}, \mathbf{y}^+) : (\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}\}$ (Stiennon et al., 2020; Rafailov et al., 2023).

Preference learning loss. Aligning language models based on pairwise preferences is usually done by minimizing a loss of the following form:

$$\mathcal{L}(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \sim \mathcal{D}} \left[\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \pi_\theta(\mathbf{y}^+|\mathbf{x}) - \ln \pi_\theta(\mathbf{y}^-|\mathbf{x}) \right) \right], \quad (2)$$

where $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is convex and differentiable, for every $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$. Denote by θ_{init} the parameters of the model prior to minimizing the loss \mathcal{L} . To guarantee that minimizing \mathcal{L} entails increasing the difference between $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ and $\ln \pi_\theta(\mathbf{y}^-|\mathbf{x})$, as expected from a reasonable preference learning loss, we make the mild assumption that $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is monotonically decreasing in a neighborhood of $\ln \pi_{\theta_{\text{init}}}(\mathbf{y}^+|\mathbf{x}) - \ln \pi_{\theta_{\text{init}}}(\mathbf{y}^-|\mathbf{x})$.

The loss \mathcal{L} generalizes many existing losses, including: DPO (Rafailov et al., 2023), IPO (Azar et al., 2024), SLiC (Zhao et al., 2023), REBEL (Gao et al., 2024), and GPO (Tang et al., 2024) — see Appendix D for details on the choice of $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ corresponding to each loss.² Notably, the common dependence on a reference model is abstracted through $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$. Other loss variants apply different weightings to the log probabilities of preferred and dispreferred responses or incorporate an additional SFT term (*e.g.*, DPOP (Pal et al., 2024), CPO (Xu et al., 2024a), RPO (Liu et al., 2024), BoNBoN (Gui et al., 2024), and SimPO (Meng et al., 2024)). For conciseness, we defer an extension of our analysis for these variants to Appendix F.

2.3 LIKELIHOOD DISPLACEMENT

We define likelihood displacement as the phenomenon where, although the preference learning loss is steadily minimized, the log probabilities of preferred responses decrease.

Definition 1. Let $\pi_{\theta_{\text{init}}}$ and $\pi_{\theta_{\text{fin}}}$ denote a language model before and after training with a preference learning loss \mathcal{L} over the dataset \mathcal{D} (Equation (2)), respectively, and suppose that the loss was successfully reduced, *i.e.* $\mathcal{L}(\theta_{\text{fin}}) < \mathcal{L}(\theta_{\text{init}})$. We say that *likelihood displacement occurred* if:³

$$\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}} \ln \pi_{\theta_{\text{fin}}}(\mathbf{y}^+|\mathbf{x}) < \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}} \ln \pi_{\theta_{\text{init}}}(\mathbf{y}^+|\mathbf{x});$$

and that *likelihood displacement occurred* for $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$ if $\ln \pi_{\theta_{\text{fin}}}(\mathbf{y}^+|\mathbf{x}) < \ln \pi_{\theta_{\text{init}}}(\mathbf{y}^+|\mathbf{x})$.

²For SLiC and GPO, the corresponding $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is differentiable almost everywhere, as opposed to differentiable. Our analysis applies to such losses up to minor adaptations excluding non-differentiable points.

³Note that $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ can decrease even as the loss \mathcal{L} is minimized, since minimizing \mathcal{L} only requires increasing the gap between $\ln \pi_\theta(\mathbf{y}^+|\mathbf{x})$ and $\ln \pi_\theta(\mathbf{y}^-|\mathbf{x})$.

Model	y^+	y^-	$\pi_\theta(y^+ x)$ Decrease	Tokens Increasing Most in Probability	
				Benign	Catastrophic
OLMo-1B	Yes	No	0.69 (0.96 \rightarrow 0.27)	_Yes, _yes	–
	No	Never	0.84 (0.85 \rightarrow 0.01)	_No	Yes, _Yes, _yes
Gemma-2B	Yes	No	0.22 (0.99 \rightarrow 0.77)	_Yes, _yes	–
	No	Never	0.21 (0.65 \rightarrow 0.44)	no, _No	Yes, Yeah, Possibly
Llama-3-8B	Yes	No	0.96 (0.99 \rightarrow 0.03)	yes, _yes, _Yes	–
	Sure	Yes	0.59 (0.98 \rightarrow 0.39)	sure, _Sure	Maybe, No, Never

Table 1: **Likelihood displacement can be catastrophic, even when training on a single prompt with single token responses.** Each model was trained via DPO on a randomly chosen prompt from the Persona dataset (Perez et al., 2022), using different pairs of preferred and dispreferred tokens (y^+ , y^-) (as detailed in Section 3). We report the largest decrease in the preferred token probability $\pi_\theta(y^+|x)$ during training for representative (y^+ , y^-) pairs, averaged across ten runs differing in random seed for choosing the prompt. The rightmost columns include notable tokens from the top three tokens increasing most in probability throughout training (see Appendix H.1 for the full list and extent of increase). Remarkably, when y^+ and y^- are semantically similar, the **tokens increasing most in probability are often semantically opposite to y^+** .

Likelihood displacement is not necessarily problematic. For $(x, y^+, y^-) \in \mathcal{D}$, we refer to it as *benign* if the responses increasing in probability are as preferable as y^+ (e.g., they are semantically similar to y^+). However, as Section 3 demonstrates, the probability mass can go to responses that are substantially less preferable than y^+ (e.g., they are semantically opposite to y^+), in which case we say it is *catastrophic*.

3 CATASTROPHIC LIKELIHOOD DISPLACEMENT IN SIMPLE SETTINGS

Despite the prevalence of likelihood displacement (Pal et al., 2024; Yuan et al., 2024; Pang et al., 2024; Liu et al., 2024), there is limited understanding as to why it occurs and where the probability mass goes. Prior work attributed this phenomenon to limitations in model capacity (Tajwar et al., 2024), the presence of multiple training samples or output tokens (Tajwar et al., 2024; Pal et al., 2024), and the initial SFT phase (Rafailov et al., 2024). In contrast, we demonstrate that likelihood displacement can occur and be catastrophic independently of these factors, even when training over just a single prompt whose responses contain a single token each. The potential adverse effects of such displacement raise the need to formally characterize its underlying causes.

Setting. The experiments are based on the Persona dataset (Perez et al., 2022), in which every prompt contains a statement, and the model needs to respond whether it agrees with the statement using a single token. We assign to each prompt a pair of preferred and dispreferred tokens (y^+ , y^-) from a predetermined set containing, e.g., Yes, Sure, No, and Never. Then, for the OLMo-1B, Gemma-2B, and Llama-3-8B models, we perform one epoch of SFT using the preferred tokens as labels, in line with common practices, and train each model via DPO on a single randomly selected prompt. See Appendix I.1 for additional details.

Likelihood displacement is pervasive and can be catastrophic. Table 1 reports the decrease in preferred token probability, and notable tokens whose probability increases at the expense of y^+ . The probability of y^+ dropped by at least 0.21 and up to 0.96 absolute value in all runs. Remarkably, when y^+ and y^- are semantically similar, the probability mass often shifts to semantically opposite tokens. Appendix H.1 reports similar findings for experiments using: (i) base models that did not undergo an initial SFT phase (Table 2); or (ii) IPO instead of DPO (Table 3).

4 THEORETICAL ANALYSIS OF LIKELIHOOD DISPLACEMENT

To uncover what causes likelihood displacement when minimizing a preference learning loss, we characterize how the log probabilities of responses evolve during gradient-based training. For a preference sample $(x, y^+, y^-) \in \mathcal{D}$, we identify the factors pushing $\ln \pi_\theta(y^+|x)$ downwards and those determining which responses increase most in log probability instead. Section 4.1 lays out the

technical approach, after which Section 4.2 gives an overview of the main results. The full analysis is deferred to Appendix E. For the convenience of the reader, we provide the main takeaways below.

Takeaway 1: Role of the Token Unembedding Geometry (Section 4.2.1)

Even when training over a single prompt whose responses \mathbf{y}^+ and \mathbf{y}^- contain a single token, likelihood displacement can occur due to the token unembedding geometry. The underlying causes are: (i) an alignment between the preferred and dispreferred token unembeddings, measured as $\langle \mathbf{W}_{\mathbf{y}^+}, \mathbf{W}_{\mathbf{y}^-} \rangle$; and (ii) tokens whose unembeddings align with $\mathbf{W}_{\mathbf{y}^+} - \mathbf{W}_{\mathbf{y}^-}$, which increase in log probability at the expense of \mathbf{y}^+ . Tokens increasing in probability can thus have unembeddings that align with directions orthogonal to $\mathbf{W}_{\mathbf{y}^+}$. Since unembeddings often linearly encode semantics, this provides an explanation for why probability mass can go to tokens semantically unrelated or opposite to \mathbf{y}^+ (as observed in Section 3),

Takeaway 2: Role of the Hidden Embedding Geometry (Section 4.2.2)

Besides the impact of the token unembedding geometry (Takeaway 1), likelihood displacement occurs when the preferred and dispreferred responses are similar according to the following measure, which is based on their hidden embeddings.

Definition 2. For a preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, we define the *centered hidden embedding similarity (CHES)* score of \mathbf{y}^+ and \mathbf{y}^- with respect to a model π_θ by:

$$\text{CHES}_{\mathbf{x}}(\mathbf{y}^+, \mathbf{y}^-) := \left\langle \underbrace{\sum_{k=1}^{|\mathbf{y}^+|} \mathbf{h}_{\mathbf{x}, \mathbf{y}^+_{<k}}}_{\mathbf{y}^+ \text{ hidden embeddings}}, \underbrace{\sum_{k'=1}^{|\mathbf{y}^-|} \mathbf{h}_{\mathbf{x}, \mathbf{y}^-_{<k'}}}_{\mathbf{y}^- \text{ hidden embeddings}} \right\rangle - \left\| \sum_{k=1}^{|\mathbf{y}^+|} \mathbf{h}_{\mathbf{x}, \mathbf{y}^+_{<k}} \right\|^2,$$

where $\mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}$ denotes the hidden embedding that the model produces given \mathbf{x} and the first $k-1$ tokens of $\mathbf{z} \in \mathcal{V}^*$. A higher CHES score stands for more similar preferences. We omit the dependence of CHES on π_θ in our notation as it will be clear from context.

Losses with SFT regularization. Appendix F extends our analysis to losses incorporating an SFT regularization term. In particular, it formalizes how this modification helps mitigate likelihood displacement, as proposed in Pal et al. (2024); Liu et al. (2024); Pang et al. (2024); Gui et al. (2024). We note, however, that our experiments in Section 6 reveal a limitation of this approach for mitigating the adverse effects of likelihood displacement, compared to improving the data curation pipeline.

4.1 TECHNICAL APPROACH

Given a prompt \mathbf{x} , the probability that the model π_θ assigns to a response \mathbf{z} is determined by the hidden embeddings $\mathbf{h}_{\mathbf{x}}, \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<2}}, \dots, \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<|\mathbf{z}|}}$ and the token unembeddings \mathbf{W} (Equation (1)). Our analysis relies on tracking their evolution when minimizing the loss \mathcal{L} (Equation (2)). To do so, we adopt the *unconstrained features model* (Mixon et al., 2022), which amounts to treating hidden embeddings as directly trainable parameters. Namely, the trainable parameters are taken to be $\theta = \{\mathbf{h}_{\mathbf{z}} : \mathbf{z} \in \mathcal{V}^*\} \cup \{\mathbf{W}\}$. This simplification has proven useful for studying various deep learning phenomena, including neural collapse (e.g., Zhu et al. (2021); Ji et al. (2022); Tirer et al. (2023)) and the benefits of language model pretraining for downstream tasks (Saunshi et al., 2021). As verified in Sections 5 and 6, it also allows extracting the salient sources of likelihood displacement.⁴

Language model finetuning is typically done with small learning rates. Accordingly, we analyze the training dynamics of (stochastic) gradient descent at the small learning rate limit, i.e. *gradient flow*: $\frac{d}{dt}\theta(t) = -\nabla\mathcal{L}(\theta(t))$, where $\theta(t)$ denotes the parameters at time $t \geq 0$ of training. Note that under gradient flow the loss is monotonically decreasing.⁵ Thus, any reduction in the log probabilities of preferred responses constitutes likelihood displacement (cf. Definition 1).

⁴In contrast to prior theoretical analyses of likelihood displacement, which consider stylized settings (e.g., linear models and cases where the preferred and dispreferred responses differ only by a single token), whose implications to more realistic settings are unclear (Pal et al., 2024; Fisch et al., 2024; Song et al., 2024).

⁵Except for the trivial case where $\theta(0)$ is a critical point of \mathcal{L} , in which $\mathcal{L}(\theta(t)) = \mathcal{L}(\theta(0))$ for all $t \geq 0$.

4.2 OVERVIEW OF THE MAIN RESULTS

4.2.1 SINGLE TRAINING SAMPLE AND OUTPUT TOKEN

It is instructive to first consider the case of training on a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, whose responses $\mathbf{y}^+ \in \mathcal{V}$ and $\mathbf{y}^- \in \mathcal{V}$ contain a single token. Theorem 1 characterizes how the token unembedding geometry determines when $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ is negative, *i.e.* when likelihood displacement occurs.

Theorem 1 (Informal version of Theorem 4). *Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}$ and $\mathbf{y}^- \in \mathcal{V}$ each being a single token. At any time $t \geq 0$ of training, $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ is more negative the larger the following term is:*

$$\underbrace{\langle \mathbf{W}_{\mathbf{y}^+}(t), \mathbf{W}_{\mathbf{y}^-}(t) \rangle}_{\text{preferences unembedding alignment}} + \sum_{z \in \mathcal{V} \setminus \{\mathbf{y}^+, \mathbf{y}^-\}} \pi_{\theta(t)}(z|\mathbf{x}) \cdot \underbrace{\langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle}_{\text{alignment of other token with } \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)},$$

where $\mathbf{W}_z(t)$ denotes the token unembedding of $z \in \mathcal{V}$ at time t .

Two terms govern the extent of likelihood displacement in the case of single token responses. First, $\langle \mathbf{W}_{\mathbf{y}^+}(t), \mathbf{W}_{\mathbf{y}^-}(t) \rangle$ formalizes the intuition that likelihood displacement occurs when the preferred and dispreferred responses are similar. A higher inner product in unembedding space translates to a more substantial (instantaneous) decrease in $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$. Second, is a term which measures the alignment of other token unembeddings with $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$, where tokens deemed more likely by the model have a larger weight. The alignment of token unembeddings with $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$ also determines which tokens increase most in log probability.

Theorem 2 (Informal version of Theorem 5). *Under the setting of Theorem 1, for any time $t \geq 0$ and token $z \in \mathcal{V} \setminus \{\mathbf{y}^+, \mathbf{y}^-\}$ it holds that $\frac{d}{dt} \ln \pi_{\theta(t)}(z|\mathbf{x}) \propto \langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle$, up to an additive term independent of z .*

The direction $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$ can be decomposed into its projection onto $\mathbf{W}_{\mathbf{y}^+}(t)$ and a component orthogonal to $\mathbf{W}_{\mathbf{y}^+}(t)$, introduced by $\mathbf{W}_{\mathbf{y}^-}(t)$. Thus, tokens increasing in log probability can have unembeddings that mostly align with directions orthogonal to $\mathbf{W}_{\mathbf{y}^+}(t)$, especially when the component orthogonal to $\mathbf{W}_{\mathbf{y}^+}(t)$ of $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$ is relatively large (which we often find to be the case empirically; see Table 13 in Appendix H.1). Given that token unembeddings are known to linearly encode semantics (Mikolov et al., 2013; Arora et al., 2016; Park et al., 2024), this provides an explanation for why the probability mass can shift to tokens that are semantically unrelated or opposite to the preferred token, *i.e.* why likelihood displacement can be catastrophic even in simple settings (as observed in Section 3).

4.2.2 RESPONSES WITH MULTIPLE TOKENS

We now extend our analysis to the typical case where responses are sequences of tokens. As shown below, the existence of multiple tokens in each response introduces a dependence on their (contextual) hidden embeddings.

Theorem 3 (Informal version of Theorem 6). *Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$. At any time $t \geq 0$ of training, in addition to the dependence on token unembeddings identified in Theorem 1, $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ is more negative the larger the following term is:*

$$\sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^-|} \alpha_{k,k'}^-(t) \cdot \underbrace{\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}(t) \rangle}_{\text{preferred-dispreferred alignment}} - \sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^+|} \alpha_{k,k'}^+(t) \cdot \underbrace{\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \rangle}_{\text{preferred-preferred alignment}},$$

where $\mathbf{h}_z(t)$ denotes the hidden embedding of $z \in \mathcal{V}^*$ at time t , and $\alpha_{k,k'}^-(t), \alpha_{k,k'}^+(t) \in [-2, 2]$ are coefficients determined by the model’s next-token distribution for prefixes of \mathbf{y}^+ and \mathbf{y}^- .

Theorem 3 establishes that the alignment of hidden embeddings, of both the “preferred-dispreferred” and “preferred-preferred” types, affects likelihood displacement. A larger inner product leads to an upwards or downwards push on $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$, depending on the sign of the corresponding $\alpha_{k,k'}^-(t)$ or $\alpha_{k,k'}^+(t)$ coefficient. Empirically, we find that these coefficients are mostly positive across models

and datasets; *e.g.*, the OLMo-1B, Gemma-2B, and Llama-3-8B models and the UltraFeedback and AlpacaFarm datasets (see Appendix H.2 for details). By accordingly setting all coefficients in Theorem 3 to one, we derive the *centered hidden embedding similarity (CHES)* score between preferred and dispreferred responses (Definition 2). Our analysis indicates that a higher CHES score implies more severe likelihood displacement. Section 5 empirically verifies this relation, and demonstrates that the CHES score is significantly more predictive of likelihood displacement than other plausible similarity measures.

Our analysis also provides insight into which responses increase most in probability at the expense of \mathbf{y}^+ . Theorem 7 in Appendix E.2 derives the dependence of $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$, for any response $\mathbf{z} \in \mathcal{V}^*$, on the alignment of its hidden embeddings with those of \mathbf{y}^+ and \mathbf{y}^- . However, in general settings, it is difficult to qualitatively describe the types of responses increasing in probability, and whether they constitute benign or catastrophic likelihood displacement. Section 6 thus demonstrates the (harmful) implications of likelihood displacement in settings where responses can be easily categorized into benign or catastrophic. We regard studying the question of where the probability mass goes in additional settings as a promising direction for future work.

4.2.3 MULTIPLE TRAINING SAMPLES

Sections 4.2.1 and 4.2.2 showed that likelihood displacement may occur regardless of the dataset size. Nonetheless, increasing the number of training samples was empirically observed to exacerbate it (Tajwar et al., 2024). Appendix E.3 sheds light on this observation by characterizing, for any $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, when additional training samples lead to a larger decrease in $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$. This unsurprisingly occurs when \mathbf{y}^+ appears as the dispreferred response of other prompts, *i.e.* there are contradicting samples. We further establish that additional training samples can contribute negatively to $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ even when their preferences are distinct from those of \mathbf{x} .

5 IDENTIFYING SOURCES OF LIKELIHOOD DISPLACEMENT

In Section 4 we derived the CHES score (Definition 2), which for a given model and preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, measures the similarity of \mathbf{y}^+ and \mathbf{y}^- based on their hidden embeddings. Our theory indicates that samples with a higher CHES score lead to more likelihood displacement. Below, we affirm this prediction and show that the CHES score enables identifying which training samples in a dataset contribute most to likelihood displacement, whereas alternative similarity measures fail to do so. The following Section 6 then demonstrates that filtering out samples with a high CHES score can mitigate undesirable implications of likelihood displacement.

Setting. We use the UltraFeedback (Cui et al., 2024) and AlpacaFarm (Dubois et al., 2024) datasets and the OLMo-1B, Gemma-2B, and Llama-3-8B models. For each model and preference dataset, we compute the CHES scores of all samples. This requires performing a single forward pass over the dataset. Then, for each of the 0th, 25th, 50th, 75th, and 100th score percentiles, we take a subset of 512 samples centered around it.⁶ Lastly, we train the model via DPO on each of the subsets separately, and track the change in log probability for preferred responses in the subset — the more the log probabilities decrease, the more severe the likelihood displacement is. See Appendix I.2 for additional implementation details.

Baselines. We repeat the process outlined above while ranking the similarity of preferences using the (normalized) edit distance,⁷ since preferences with low edit distance were suggested by Pal et al. (2024) as a cause for likelihood displacement. To the best of our knowledge, no other property of a preference sample was linked with likelihood displacement in the literature. So we additionally compare to a natural candidate: using the inner product between the last hidden embeddings of \mathbf{y}^+ and \mathbf{y}^- , *i.e.* $\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}^+}, \mathbf{h}_{\mathbf{x}, \mathbf{y}^-} \rangle$, as the similarity score.

CHES score effectively identifies samples leading to likelihood displacement. For the UltraFeedback dataset, Figure 2 shows the change in preferred response log probability against the similarity percentile of samples. Across all models, the CHES score ranking matches perfectly the degree of likelihood displacement: the higher the CHES score percentile, the more preferred responses de-

⁶The 0th and 100th percentile subsets include the 512 samples with lowest and highest scores, respectively.

⁷A lower (normalized) edit distance between \mathbf{y}^+ and \mathbf{y}^- corresponds to a higher similarity.

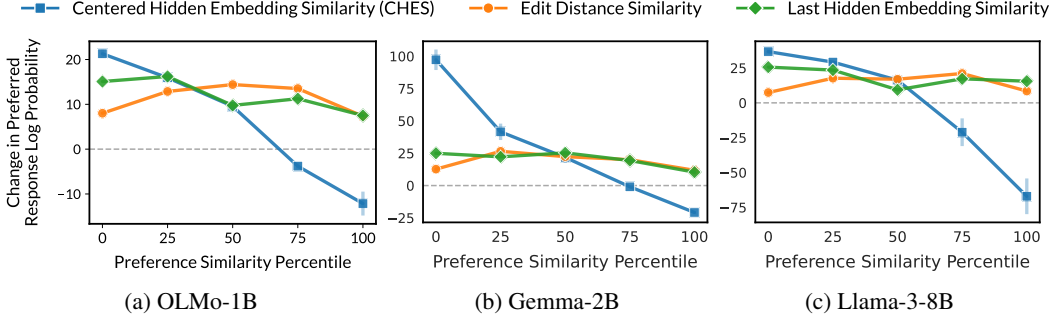


Figure 2: **CHES score (Definition 2) identifies which training samples contribute to likelihood displacement, whereas alternative similarity measures do not.** Each model was trained via DPO on subsets of 512 samples from the UltraFeedback dataset. The subsets are centered around different preference similarity percentiles, according to the following measures: (i) the CHES score; (ii) (normalized) edit distance, which was suggested in Pal et al. (2024) as indicative of likelihood displacement; and (iii) the inner product between the last hidden embeddings of the preferred and dispreferred responses (see Section 5 for further details). We report for each subset the change in mean preferred response log probability, averaged across three runs (error bars marking standard deviation are often indiscernible). The CHES score ranking perfectly matches with the degree of likelihood displacement — samples with a higher score induce a larger decrease in the log probability of preferred responses. On the other hand, the alternative measures are not indicative of likelihood displacement.

crease in log probability. Moreover, training on samples with high CHES scores leads to severe likelihood displacement, whereas training on samples with low CHES scores leads the preferred responses to increase in log probability.

CHES score is more indicative of likelihood displacement than alternative measures. In contrast to the CHES score, the edit distance of preferences and the inner product between their last hidden embeddings do not correlate well with likelihood displacement. Moreover, these measures failed to identify samples leading to likelihood displacement: across all similarity percentiles, the log probability of preferred responses only increased.

Additional experiments. Appendix H.3 reports similar findings for experiments using: (i) the AlpacaFarm dataset instead of UltraFeedback (Figure 5); or (ii) IPO instead of DPO (Figure 6).

Qualitative analysis. Appendix H.3 further includes representative samples with high and low CHES scores (Tables 14 and 15, respectively). A noticeable trait is that, in samples with a high CHES score, the dispreferred response is often longer than the preferred response, whereas for samples with a low CHES score the trend is reversed (*i.e.* preferred responses are longer). We find that this stems from a tendency of current models to produce, for different responses, hidden embeddings with a positive inner product (over 99.5% of such inner products are positive, across the considered models and datasets). As a result, for samples with longer dispreferred responses the CHES score comprises more positive terms than negative terms.

6 UNINTENTIONAL UNALIGNMENT IN DIRECT PREFERENCE LEARNING

Direct preference learning has been successfully applied for improving general instruction following and performance on downstream benchmarks (*e.g.*, Tunstall et al. (2023); Iverson et al. (2023)). This suggests that, in such settings, likelihood displacement may often be benign, and so does not require mitigation. However, in this section, we reveal that it can undermine the efficacy of safety alignment. When training a language model to refuse unsafe prompts, we find that likelihood displacement can *unintentionally unalign* the model, by causing probability mass to shift from preferred refusal responses to harmful responses. We then demonstrate that this undesirable outcome can be prevented by discarding samples with a high (length-normalized) CHES score (Definition 2), showcasing the potential of the CHES score for mitigating adverse effects of likelihood displacement more broadly.

6.1 SETTING

We train a language model to refuse unsafe prompts via the (on-policy) direct preference learning pipeline outlined in Rafailov et al. (2023), as specified below. To account for the common scenario

whereby one wishes to further align an existing (moderately) aligned model, we use the Gemma-2B-IT and Llama-3-8B-Instruct models. Then, for each model separately, we create a preference dataset based on unsafe prompts from SORRY-Bench (Xie et al., 2024). Specifically, for every prompt, we generate two candidate responses from the model and label them as refusals or non-refusals using the judge model from Xie et al. (2024). Refusals are deemed more preferable compared to non-refusals, and ties are broken by the PairRM reward model (Jiang et al., 2023). Lastly, the language models are trained via DPO over their respective datasets. For brevity, we defer to Appendices H and I some implementation details and experiments using IPO, respectively.

6.2 CATASTROPHIC LIKELIHOOD DISPLACEMENT CAUSES UNINTENTIONAL UNALIGNMENT

Since the initial models are moderately aligned, we find that they often generate two refusal responses for a given prompt. Specifically, for over 70% of the prompts in the generated datasets, both the preferred and dispreferred responses are refusals. This situation resembles the experiments of Section 3, where training on semantically similar preferences led to catastrophic likelihood displacement (e.g., when y^+ was No and y^- was Never, the probability of Yes sharply increased).

Analogously, we observe that as the DPO loss is minimized, likelihood displacement causes probability mass to shift away from preferred refusal responses (Table 16 in Appendix H.4 reports the log probability decrease of preferred responses). This leads to a significant drop in refusal rates. Specifically, over the training set, DPO makes the refusal rates of Gemma-2B-IT and Llama-3-8B-Instruct drop from 80.5% to 54.8% and 74.4% to 33.4%, respectively (similar drops occur over the validation set). In other words, instead of further aligning the model, preference learning unintentionally unaligns it. See Appendix H.4 for examples of unsafe prompts from the training set, for which initially the models generated two refusals, yet after DPO they comply with the prompts (Table 18).

We note that alignment usually involves a tradeoff between safety and helpfulness. The drop in refusal rates is particularly striking since the models are trained with the sole purpose of refusing prompts, without any attempt to maintain their helpfulness.

6.3 FILTERING DATA VIA CHES SCORE MITIGATES UNINTENTIONAL UNALIGNMENT

Section 5 showed that samples with a high CHES score (Definition 2) contribute most to likelihood displacement. Motivated by this, we explore whether filtering data via the CHES score can mitigate unintentional unalignment, and which types of samples it marks as problematic.

As discussed in Section 5, due to the embedding geometry of current models, CHES scores can correlate with the lengths of responses. To avoid introducing a length bias when filtering data, we apply a length-normalized variant of CHES (see Definition 3 in Appendix C). For comparison, we also consider adding an SFT term to the DPO loss, as suggested in Pal et al. (2024); Xu et al. (2024a); Pang et al. (2024); Liu et al. (2024), and training over “gold” responses from SORRY-Bench, which were generated from a diverse set of base and safety aligned models and labeled by human raters.

Filtering data via CHES score mitigates unintentional unalignment. Figure 3 reports the refusal rates before and after training via DPO: (i) on the original dataset, which as stated in Section 6.2 leads to a substantial drop in refusal rates; (ii) with an additional SFT term on the original dataset; (iii) on the gold dataset; and (iv) on a filtered version of the original dataset that contains the 5% samples with lowest length-normalized CHES scores.⁸ Filtering data via the CHES score successfully mitigates unintentional unalignment. Moreover, while adding an SFT term to the loss also prevents the drop in refusal rates, data filtering boosts the refusal rates more substantially. We further find that DPO on gold preferences does not suffer from likelihood displacement or unintentional unalignment (i.e. the preferred responses increase in probability; see Table 16). Overall, these results highlight the importance of curating data with sufficiently distinct preferences for effective preference learning.

Which samples have a high CHES score? Figure 4 reveals that the length-normalized CHES score ranking falls in line with intuition — samples that have two responses of the same type (i.e. two refusal or two non-refusal responses) tend to have a higher score than samples with one response of each type, and so are more likely to cause likelihood displacement.

⁸Keeping up to 15% of the original samples led to analogous results. Beyond that, as when training on the full dataset, likelihood displacement caused refusal rates to drop.

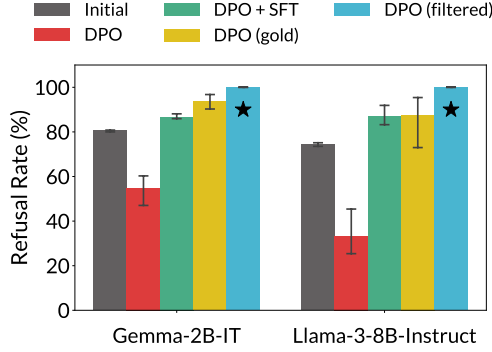


Figure 3: **Likelihood displacement can cause unintentional unalignment, which is mitigated by data filtering.** Training a model to refuse unsafe prompts from SORRY-Bench via DPO unintentionally leads to a substantial decrease in refusal rates due to likelihood displacement. Filtering out samples with a high length-normalized CHES score (★) or using “gold” preference data, generated from a diverse set of models, successfully mitigates the problem, and goes beyond the improvement achieved when adding an SFT term to the DPO loss. Reported are mean values over three runs (error bars denote minimal and maximal values). See Section 6 for further details.

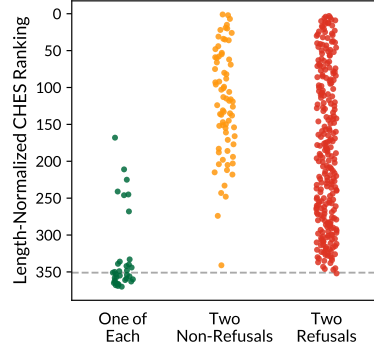


Figure 4: **Length-normalized CHES score identifies samples with two responses of the same type as responsible for likelihood displacement.** For Llama-3-8B-Instruct, we take the corresponding SORRY-Bench preference dataset (see Section 6.1 for details), and plot the ranking of samples according to their length-normalized CHES scores. Gray line marks the bottom 5% of samples. Agreeing with intuition, samples with two refusal or two non-refusal responses tend to have a higher length-normalized CHES score than samples with one of each.

7 CONCLUSION

While direct preference learning has been widely adopted, there is considerable uncertainty around how it affects the model (*cf.* Xu et al. (2024b); Chen et al. (2024)). Our theory and experiments shed light on the causes and implications of one counter-intuitive phenomenon — *likelihood displacement*. We demonstrated that likelihood displacement can be catastrophic, shifting probability mass from preferred responses to semantically opposite ones, which can result in *unintentional unalignment* when training a model to refuse unsafe prompts. Intuitively, these failures arise when the preferred and dispreferred responses are similar. We formalized this intuition and derived the *centered hidden embedding similarity* (CHES) score (Definition 2), which effectively identifies samples contributing to likelihood displacement in a given dataset. As an example for its potential uses, we showed that filtering out samples with a high (length-normalized) CHES score can prevent unintentional unalignment. More broadly, our work highlights the importance of curating data with sufficiently distinct preferences, for which we believe the CHES score may prove valuable.

7.1 LIMITATIONS AND FUTURE WORK

Theoretical analysis. Our theory focuses on the instantaneous change of log probabilities, and abstracts away which neural network architecture is used for computing hidden embeddings. Future work can extend it by studying the evolution of log probabilities throughout training and accounting for how the architecture choice influences likelihood displacement.

Occurrences of catastrophic likelihood displacement. While our findings reveal that likelihood displacement can make well-intentioned training result in undesirable outcomes, we do not claim that this occurs universally. Indeed, direct preference learning methods have been successfully applied for aligning language models (Tunstall et al., 2023; Ivison et al., 2023; Jiang et al., 2024; Dubey et al., 2024). Nonetheless, in light of the growing prominence of these methods, we believe it is crucial to detect additional settings in which likelihood displacement is catastrophic.

Utility of the CHES score. We demonstrated the potential of the (length-normalized) CHES score for filtering out samples that cause likelihood displacement. However, further investigation is necessary to assess its utility more broadly. For example, exploring whether data filtering via CHES scores improves performance in general instruction following settings, or whether CHES scores can be useful in more complex data curation pipelines for selecting distinct preferences based on a pool of candidate responses, possibly generated from a diverse set of models.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. A latent variable model approach to pmi-based word embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399, 2016.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislaw Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- Angelica Chen, Sadhika Malladi, Lily H Zhang, Xinyi Chen, Qiuyi Zhang, Rajesh Ranganath, and Kyunghyun Cho. Preference learning algorithms do not learn preference rankings. *arXiv preprint arXiv:2405.19534*, 2024.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. In *International Conference on Machine Learning*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. In *International Conference on Machine Learning*, 2024.
- Duanyu Feng, Bowen Qin, Chen Huang, Zheng Zhang, and Wenqiang Lei. Towards analyzing and understanding the limitations of dpo: A theoretical perspective. *arXiv preprint arXiv:2404.04626*, 2024.
- Adam Fisch, Jacob Eisenstein, Vicky Zayats, Alekh Agarwal, Ahmad Beirami, Chirag Nagpal, Pete Shaw, and Jonathan Berant. Robust preference optimization through reward model distillation. *arXiv preprint arXiv:2405.19316*, 2024.
- Zhaolin Gao, Jonathan D Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley, Thorsten Joachims, J Andrew Bagnell, Jason D Lee, and Wen Sun. Rebel: Reinforcement learning via regressing relative rewards. *arXiv preprint arXiv:2404.16767*, 2024.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*, 2024.
- Lin Gui, Cristina Gârbacea, and Victor Veitch. Bonbon alignment for large language models and the sweetness of best-of-n sampling. *arXiv preprint arXiv:2406.00832*, 2024.
- Luxi He, Mengzhou Xia, and Peter Henderson. What’s in your “safe” data?: Identifying benign data that breaks safety. *arXiv preprint arXiv:2404.01099*, 2024.
- Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012.
- Audrey Huang, Wenhao Zhan, Tengyang Xie, Jason D Lee, Wen Sun, Akshay Krishnamurthy, and Dylan J Foster. Correcting the mythos of kl-regularization: Direct alignment without overparameterization via chi-squared preference optimization. *arXiv preprint arXiv:2407.13399*, 2024.

- Shawn Im and Yixuan Li. On the generalization of preference learning with dpo. *arXiv preprint arXiv:2408.03459*, 2024a.
- Shawn Im and Yixuan Li. Understanding the learning dynamics of alignment with human feedback. In *International Conference on Machine Learning*, 2024b.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*, 2023.
- Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled perspective on neural collapse. In *International Conference on Learning Representations*, 2022.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. *arXiv preprint arXiv:2306.02561*, 2023.
- Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer. *arXiv preprint arXiv:2405.16436*, 2024.
- Kaifeng Lyu, Haoyu Zhao, Xinran Gu, Dingli Yu, Anirudh Goyal, and Sanjeev Arora. Keeping llms aligned after fine-tuning: The crucial role of prompt templates. *arXiv preprint arXiv:2402.18540*, 2024.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *Sampling Theory, Signal Processing, and Data Analysis*, 20(2):11, 2022.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddhartha Naidu, and Colin White. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*, 2024.
- Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*, 2024.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. In *International Conference on Machine Learning*, 2024.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- Ethan Perez, Sam Ringer, Kamilè Lukošiušė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! In *International Conference on Learning Representations*, 2024.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2023.
- Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to Q^* : Your language model is secretly a Q -function. *arXiv preprint arXiv:2404.12358*, 2024.
- Noam Razin, Hattie Zhou, Omid Saremi, Vimal Thilak, Arwen Bradley, Preetum Nakkiran, Joshua M. Susskind, and Etai Littwin. Vanishing gradients in reinforcement finetuning of language models. In *International Conference on Learning Representations*, 2024.

- Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. A mathematical exploration of why language models help solve downstream tasks. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=vVjIW3sEc1s>.
- Yuda Song, Gokul Swamy, Aarti Singh, J Andrew Bagnell, and Wen Sun. The importance of online data: Understanding preference fine-tuning via coverage. *arXiv preprint arXiv:2406.01462*, 2024.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021, 2020.
- Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *arXiv preprint arXiv:2404.14367*, 2024.
- Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Rémi Munos, Mark Rowland, Pierre Harvey Richemond, Michal Valko, Bernardo Ávila Pires, and Bilal Piot. Generalized preference optimization: A unified approach to offline alignment. In *International Conference on Machine Learning*, 2024.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Tom Tirer, Haoxiang Huang, and Jonathan Niles-Weed. Perturbation analysis of neural collapse. In *International Conference on Machine Learning*, pages 34301–34329. PMLR, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Schwag, Kaixuan Huang, Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, et al. Sorry-bench: Systematically evaluating large language model safety refusal behaviors. *arXiv preprint arXiv:2406.14598*, 2024.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *International Conference on Machine Learning*, 2024.
- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*, 2024a.
- Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. *arXiv preprint arXiv:2404.10719*, 2024b.

- Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. A comprehensive study of jailbreak attack versus defense for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics ACL 2024*, pages 7432–7449, Bangkok, Thailand and virtual meeting, August 2024c. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.443. URL <https://aclanthology.org/2024.findings-acl.443>.
- Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen, Ruobing Xie, Yankai Lin, et al. Advancing llm reasoning generalists with preference trees. *arXiv preprint arXiv:2404.02078*, 2024.
- Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. Removing RLHF protections in GPT-4 via fine-tuning. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 681–687, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-short.59. URL <https://aclanthology.org/2024.naacl-short.59>.
- Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu. Calibrating sequence likelihood improves conditional language generation. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0qSoodKmJaN>.
- Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, et al. Secrets of rlhf in large language models part i: Ppo. *arXiv preprint arXiv:2307.04964*, 2023.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. *Advances in Neural Information Processing Systems*, 34:29820–29834, 2021.
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

A RELATED WORK

Preference learning. There are two main approaches for aligning language models based on preference data. First, RLHF (or RLAIFF) (Ziegler et al., 2019; Stiennon et al., 2020; Ouyang et al., 2022; Bai et al., 2022b), which requires fitting a reward model to a dataset of human (or AI) preferences, and then training the language model to maximize the reward via RL. While often effective for improving the quality of generated responses, RLHF is computationally costly and can suffer from instabilities of (Zheng et al., 2023; Razin et al., 2024). This has led to the rise of *direct preference learning* methods, which directly train the language model to increase the probability of preferred responses relative to dispreferred responses, as popularized by DPO (Rafailov et al., 2023) and its numerous variants (e.g., Zhao et al. (2023); Azar et al. (2024); Tang et al. (2024); Xu et al. (2024a); Ethayarajh et al. (2024); Meng et al. (2024)).

Analyses of direct preference learning. Prior work mostly established sample complexity guarantees for DPO (or a variant of it) when the training data obeys a specific, stringent structure (Im and Li, 2024a) or provides sufficient coverage (Liu et al., 2024; Song et al., 2024; Huang et al., 2024). Additionally, Im and Li (2024b); Feng et al. (2024) studied the rate of optimization when performing DPO. More relevant to our work is Chen et al. (2024), which demonstrated that DPO can fail to correct how a model ranks preferred and dispreferred responses. While related, this phenomenon is distinct from likelihood displacement. In particular, when likelihood displacement occurs the probability of preferred responses is often higher than the probability of dispreferred responses (as illustrated in Figure 1 and was the case in the experiments of Sections 3, 5, and 6).

Likelihood displacement. The relation of our results to existing claims regarding likelihood displacement was discussed throughout the paper. We provide in Appendix B a consolidated account.

Jailbreaking and Unalignment. Aligned language models are vulnerable to jailbreaking through carefully designed adversarial prompts (Xu et al., 2024c). However, even when one does not intend to unalign a given model, Qi et al. (2024); He et al. (2024); Zhan et al. (2024); Lyu et al. (2024) showed that performing SFT over seemingly benign data can result in unalignment. The experiments in Section 6 provide a more extreme case of unintentional unalignment. Specifically, although the models are trained with the sole purpose of refusing unsafe prompts, likelihood displacement causes the refusal rate to drop, instead of increase.

B RELATION TO EXISTING CLAIMS ON LIKELIHOOD DISPLACEMENT

Throughout the paper, we discussed how our results relate to existing claims regarding likelihood displacement. This appendix provides a concentrated account for the convenience of the reader.

Similarity of preferences. Tajwar et al. (2024) and Pal et al. (2024) informally claimed that samples with similar preferences are responsible for likelihood displacement. Our theoretical analysis (Section 4) formalizes this intuition, by proving that similarities between the embeddings of preferred and dispreferred responses drives likelihood displacement.

Dataset size and model capacity. Tajwar et al. (2024) also attributed likelihood displacement to the presence of multiple samples or a limited model capacity. Section 3 demonstrates that likelihood displacement can occur independently of these factors, even when training a 8B model on a single sample. Though as we characterize in Section 4.2.3, having multiple training samples can contribute to the severity of likelihood displacement.

Preferences with small edit distance. Pal et al. (2024) showed in controlled settings that preferences with a small edit distance can lead to likelihood displacement. However, as the experiments in Section 5 demonstrate, in more general settings edit distance is not indicative of likelihood displacement. In contrast, the CHES score (Definition 2), which measures similarity based on hidden embeddings, accurately identifies samples contributing to likelihood displacement.

Initial SFT Phase. Rafailov et al. (2024) suggested that likelihood displacement occurs due to the initial SFT phase in the direct preference learning pipeline (see Section 2). Our experiments and theory refine this claim by showing that likelihood displacement can occur regardless of whether a model undergoes an initial SFT phase or not (Sections 3 and 4).

Prior sightings of catastrophic likelihood displacement. Prior work observed that DPO can degrade the performance on math and reasoning benchmarks (Pal et al., 2024; Yuan et al., 2024; Pang et al., 2024; Meng et al., 2024). This can be considered as a special case of catastrophic likelihood displacement. We note that, because in those settings usually only a few responses are correct, any likelihood displacement is catastrophic. Our work demonstrates that likelihood displacement can be catastrophic even in settings where there are many acceptable responses, and reveals its adverse effects for safety alignment.

C LENGTH-NORMALIZED CHES SCORE

In Section 4 we derived the CHES score (Definition 2), which for a given model and preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, measures the similarity of \mathbf{y}^+ and \mathbf{y}^- based on their hidden embeddings. Section 5 then demonstrated on standard preference learning datasets (UltraFeedback and AlpacaFarm) that samples with high CHES scores contribute most to likelihood displacement. Though, as discussed in Section 5, due to the embedding geometry of current models, CHES scores often correlate with the lengths of responses. Thus, to avoid introducing a length bias when filtering data in Section 6.3, we apply the following length-normalized variant of CHES.

Definition 3. For a preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, we define the *length-normalized CHES* score of \mathbf{y}^+ and \mathbf{y}^- with respect to a model π_θ by:

$$\overline{\text{CHES}}_{\mathbf{x}}(\mathbf{y}^+, \mathbf{y}^-) := \frac{1}{|\mathbf{y}^+||\mathbf{y}^-|} \left\langle \underbrace{\sum_{k=1}^{|\mathbf{y}^+|} \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}}_{\mathbf{y}^+ \text{ hidden embeddings}}, \underbrace{\sum_{k'=1}^{|\mathbf{y}^-|} \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}}_{\mathbf{y}^- \text{ hidden embeddings}} \right\rangle - \frac{1}{|\mathbf{y}^+|^2} \left\| \sum_{k=1}^{|\mathbf{y}^+|} \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+} \right\|^2,$$

where $\mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}$ denotes the hidden embedding that the model produces given \mathbf{x} and the first $k-1$ tokens of $\mathbf{z} \in \mathcal{V}^*$. We omit the dependence on π_θ in our notation as it will be clear from context.

D COMMON INSTANCES OF THE ANALYZED PREFERENCE LEARNING LOSS

As discussed in Section 2.2, the preference learning loss \mathcal{L} (Equation (2)) considered in our analysis generalizes many existing losses, which are realized by different choices of $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$, for a preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$. The choice of $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ corresponding to each loss is given below.

DPO (Rafailov et al., 2023). The DPO loss can be written as:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} \right) := -\ln \sigma \left(\beta \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} - \ln \frac{\pi_{\text{ref}}(\mathbf{y}^+|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^-|\mathbf{x})} \right) \right),$$

where π_{ref} is some reference model, $\beta > 0$ is a regularization hyperparameter, and $\sigma : \mathbb{R} \rightarrow [0, 1]$ denotes the sigmoid function.

IPO (Azar et al., 2024). The IPO loss can be written as:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} \right) := \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} - \ln \frac{\pi_{\text{ref}}(\mathbf{y}^+|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^-|\mathbf{x})} - \frac{1}{2\tau} \right)^2,$$

where π_{ref} is some reference model and $\tau > 0$ is a hyperparameter controlling the target log probability margin.

SLiC (Zhao et al., 2023). The SLiC loss can be written as:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} \right) := \max \left\{ 0, \delta - \ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} \right\},$$

where $\delta > 0$ is a hyperparameter controlling the target log probability margin. We note that our assumption on $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ being monotonically decreasing in a neighborhood of $\ln \pi_{\theta_{\text{init}}}(\mathbf{y}^+|\mathbf{x}) - \ln \pi_{\theta_{\text{init}}}(\mathbf{y}^-|\mathbf{x})$ holds, except for the case where the loss for $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$ is already zero at initialization (recall θ_{init} stands for the initial parameters of the model).

REBEL (Gao et al., 2024). The REBEL loss can be written as:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} \right) := \left(\frac{1}{\eta} \left(\ln \frac{\pi_\theta(\mathbf{y}^+|\mathbf{x})}{\pi_\theta(\mathbf{y}^-|\mathbf{x})} - \ln \frac{\pi_{\text{ref}}(\mathbf{y}^+|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^-|\mathbf{x})} \right) - r(\mathbf{x}, \mathbf{y}^+) + r(\mathbf{x}, \mathbf{y}^-) \right)^2,$$

where π_{ref} is some reference model, $\eta > 0$ is a regularization parameter, and r is a reward model.

GPO (Tang et al., 2024). GPO describes a family of losses, which can be written as:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \frac{\pi_{\theta}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\theta}(\mathbf{y}^- | \mathbf{x})} \right) := f \left(\beta \left(\ln \frac{\pi_{\theta}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\theta}(\mathbf{y}^- | \mathbf{x})} - \ln \frac{\pi_{\text{ref}}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^- | \mathbf{x})} \right) \right),$$

where π_{ref} is some reference model and $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex and monotonically decreasing in a neighborhood of $\ln \pi_{\theta_{\text{init}}}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta_{\text{init}}}(\mathbf{y}^- | \mathbf{x})$ (recall θ_{init} stands for the initial parameters of the model).

E FORMAL ANALYSIS OF LIKELIHOOD DISPLACEMENT

This appendix delivers the formal analysis overviewed in Section 4.2. Appendices E.1 to E.3 cover the results discussed in Sections 4.2.1 to 4.2.3, respectively. We refer the reader to Section 4.1 for the technical setting of the analysis.

Notation. We use $\mathbf{W}(t)$, $\mathbf{W}_z(t)$, and $\mathbf{h}_z(t)$ to denote the token unembedding matrix, token unembedding of a token $z \in \mathcal{V}$, and hidden embedding of $\mathbf{z} \in \mathcal{V}^*$ at time $t \geq 0$, respectively. We let \mathbf{z}_k be the k th token in \mathbf{z} and $\mathbf{z}_{<k}$ be the first $k - 1$ tokens in \mathbf{z} . Lastly, with slight abuse of notation, we shorthand $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) := \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}))$ for a preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, where $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ stands for the derivative of $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$.

E.1 SINGLE TRAINING SAMPLE AND OUTPUT TOKEN (OVERVIEW IN SECTION 4.2.1)

We first consider the case of training on a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, whose responses $\mathbf{y}^+ \in \mathcal{V}$ and $\mathbf{y}^- \in \mathcal{V}$ contain a single token. Theorem 4 characterizes the dependence of $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x})$ on the token unembedding geometry (proof deferred to Appendix G.1).

Theorem 4. Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}$ and $\mathbf{y}^- \in \mathcal{V}$ each being a single token. At any time $t \geq 0$ of training:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \left[m(t) - \underbrace{(1 - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \cdot \langle \mathbf{W}_{\mathbf{y}^+}(t), \mathbf{W}_{\mathbf{y}^-}(t) \rangle}_{\text{preferences unembedding alignment}} \right. \\ & \quad \left. - \sum_{z \in \mathcal{V} \setminus \{\mathbf{y}^+, \mathbf{y}^-\}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \underbrace{\langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle}_{\text{alignment of other token with } \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)} \right], \end{aligned}$$

where $-\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$ and $m(t)$ is a non-negative term given by:

$$\begin{aligned} m(t) &:= (1 - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x})) \cdot \|\mathbf{W}_{\mathbf{y}^+}(t)\|^2 + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) \cdot \|\mathbf{W}_{\mathbf{y}^-}(t)\|^2 \\ & \quad + (1 - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \cdot \|\mathbf{h}_{\mathbf{x}}(t)\|^2. \end{aligned}$$

Two terms in the derived form of $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x})$ can be negative, and so are responsible for likelihood displacement in the case of single token responses. First, the term $-\langle \mathbf{W}_{\mathbf{y}^+}(t), \mathbf{W}_{\mathbf{y}^-}(t) \rangle$, which formalizes the intuition that likelihood displacement occurs when the preferred and dispreferred responses are similar. A higher inner product translates to a more substantial (instantaneous) decrease in $\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x})$. Second, is a term measuring the alignment of other token unembeddings with $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$, where tokens deemed more likely by the model have a larger weight. Theorem 5 shows that the alignment of token unembeddings with $\mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t)$ also dictates which tokens increase most in log probability, i.e. where the probability mass goes (proof deferred to Appendix G.2).

Theorem 5. Under the setting of Theorem 4, for any time $t \geq 0$ and token $z \in \mathcal{V} \setminus \{\mathbf{y}^+, \mathbf{y}^-\}$:

$$\frac{d}{dt} \ln \pi_{\theta(t)}(z | \mathbf{x}) = -\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left[\langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle + c(t) \right],$$

where $-\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$ and $c(t)$ is a term that does not depend on z , given by:

$$c(t) := (\pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) - \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})) \|\mathbf{h}_{\mathbf{x}}(t)\|^2 - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\mathbf{x}) \langle \mathbf{W}_{z'}(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle.$$

E.2 RESPONSES WITH MULTIPLE TOKENS (OVERVIEW IN SECTION 4.2.2)

Moving to the typical case, in which the responses $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$ are sequences of tokens, assume for simplicity that $\mathbf{y}_1^+ \neq \mathbf{y}_1^-$. Extending the results below to responses \mathbf{y}^+ and \mathbf{y}^- that share a prefix is straightforward, by replacing terms that depend on \mathbf{y}_1^+ and \mathbf{y}_1^- with analogous ones that depend on the initial tokens in which \mathbf{y}^+ and \mathbf{y}^- differ.

In the case of single token responses (Appendix E.1), there are two terms that contribute to likelihood displacement. For any time $t \geq 0$ and $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, if one minimizes the preference learning loss only with respect to only the initial tokens of \mathbf{y}^+ and \mathbf{y}^- , then these terms are given by:

$$S_{\mathbf{y}_1^+, \mathbf{y}_1^-}(t) := - (1 - \pi_{\theta(t)}(\mathbf{y}_1^+|\mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}_1^-|\mathbf{x})) \cdot \langle \mathbf{W}_{\mathbf{y}_1^+}(t), \mathbf{W}_{\mathbf{y}_1^-}(t) \rangle - \sum_{z \in \mathcal{V} \setminus \{\mathbf{y}_1^+, \mathbf{y}_1^-\}} \pi_{\theta(t)}(z|\mathbf{x}) \cdot \langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \rangle. \quad (3)$$

Theorem 6 establishes that, in addition to the above initial token contribution, likelihood displacement depends on an alignment between the hidden embeddings of \mathbf{y}^+ and \mathbf{y}^- (proof deferred to Appendix G.3).

Theorem 6. Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$ satisfying $\mathbf{y}_1^+ \neq \mathbf{y}_1^-$. At any time $t \geq 0$ of training:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \left[m(t) + S_{\mathbf{y}_1^+, \mathbf{y}_1^-}(t) \right. \\ & \quad \left. - \sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^-|} \alpha_{k, k'}^-(t) \cdot \underbrace{\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}(t) \rangle}_{\text{preferred-dispreferred alignment}} + \sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^+|} \alpha_{k, k'}^+(t) \cdot \underbrace{\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \rangle}_{\text{preferred-preferred alignment}} \right], \end{aligned}$$

where $-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$, the coefficients $\alpha_{k, k'}^-(t), \alpha_{k, k'}^+(t) \in [-2, 2]$ are given by:

$$\begin{aligned} \alpha_{k, k'}^- &:= \left\langle \mathbf{e}_{\mathbf{y}_k^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \cdot \mathbf{e}_z, \mathbf{e}_{\mathbf{y}_{k'}^-} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^-) \cdot \mathbf{e}_z \right\rangle, \\ \alpha_{k, k'}^+ &:= \left\langle \mathbf{e}_{\mathbf{y}_k^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \cdot \mathbf{e}_z, \mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \cdot \mathbf{e}_z \right\rangle, \end{aligned}$$

with $\mathbf{e}_z \in \mathbb{R}^d$ standing for standard basis vector corresponding to $z \in \mathcal{V}$, and $m(t)$ is the following non-negative term:

$$\begin{aligned} m(t) &:= (1 - \pi_{\theta(t)}(\mathbf{y}_1^+|\mathbf{x})) \cdot \|\mathbf{W}_{\mathbf{y}_1^+}(t)\|^2 + \pi_{\theta(t)}(\mathbf{y}_1^-|\mathbf{x}) \cdot \|\mathbf{W}_{\mathbf{y}_1^-}(t)\|^2 \\ & \quad + \sum_{k=2}^{|\mathbf{y}^+|} \left\| \mathbf{W}_{\mathbf{y}_k^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \cdot \mathbf{W}_z(t) \right\|^2. \end{aligned}$$

The evolution of $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ is governed by: (i) the initial token unembedding geometry (analogous to the characterization in Theorem 4); and (ii) the alignment of hidden embeddings, both of the “preferred-dispreferred” and the “preferred-preferred” types. As discussed in Section 4.2.2, whether a larger inner product between hidden embeddings results in an upwards or downwards push on $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ depends on the sign of the corresponding $\alpha_{k, k'}^-(t)$ or $\alpha_{k, k'}^+(t)$ coefficient. Since empirically these coefficients are mostly positive across models and datasets, Theorem 6 indicates that a higher CHES score (Definition 2) implies more severe likelihood displacement.

Regarding where the probability mass goes when likelihood displacement occurs, for any $\mathbf{z} \in \mathcal{V}^*$, Theorem 7 derives the dependence of $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$ on the alignment of \mathbf{z} ’s hidden embeddings

with those of \mathbf{y}^+ and \mathbf{y}^- (proof deferred to Appendix G.4). We assume for simplicity that the initial token of \mathbf{z}_1 is not equal to the initial tokens of \mathbf{y}^+ and \mathbf{y}^- . If \mathbf{z} shares a prefix with \mathbf{y}^+ , then the same characterization holds up to additional terms that generally push $\ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$ upwards. Similarly, if \mathbf{z} shares a prefix with \mathbf{y}^- , then there will be additional terms that push $\ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$ downwards.

Theorem 7. *Under the setting of Theorem 6, let $\mathbf{z} \in \mathcal{V}^*$ be a response satisfying $\mathbf{z}_1 \notin \{\mathbf{y}_1^+, \mathbf{y}_1^-\}$. At any time $t \geq 0$ of training:*

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}) \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \left[c(t) + \underbrace{\left\langle \mathbf{W}_{\mathbf{z}_1}(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \right\rangle}_{\text{alignment of first token unembeddings}} \right. \\ & \quad \left. - \sum_{k=1}^{|\mathbf{z}|} \sum_{k'=1}^{|\mathbf{y}^-|} \beta_{k,k'}^-(t) \cdot \underbrace{\left\langle \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}(t) \right\rangle}_{\text{z-dispreferred alignment}} + \sum_{k=1}^{|\mathbf{z}|} \sum_{k'=1}^{|\mathbf{y}^+|} \beta_{k,k'}^+(t) \cdot \underbrace{\left\langle \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \right\rangle}_{\text{z-preferred alignment}} \right], \end{aligned}$$

where $-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$, the coefficients $\beta_{k,k'}^-(t), \beta_{k,k'}^+(t) \in [-2, 2]$ are given by:

$$\begin{aligned} \beta_{k,k'}^- &:= \left\langle \mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \cdot \mathbf{e}_z, \mathbf{e}_{\mathbf{y}_{k'}^-} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^-) \cdot \mathbf{e}_z \right\rangle, \\ \beta_{k,k'}^+ &:= \left\langle \mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \cdot \mathbf{e}_z, \mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \cdot \mathbf{e}_z \right\rangle, \end{aligned}$$

and $c(t)$ is the following term that does not depend on \mathbf{z} :

$$c(t) := - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \left\langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \right\rangle.$$

E.3 MULTIPLE TRAINING SAMPLES (OVERVIEW IN SECTION 4.2.3)

In this appendix, we consider the effect of having multiple training samples, focusing on the case where responses consist of a single token. Namely, for a preference sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, Theorem 8 characterizes when additional training samples lead to a larger decrease in $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ (proof deferred to Appendix G.5). For conciseness, we make the mild assumption that no prompt appears twice in \mathcal{D} , as is common in real-world preference datasets.

Theorem 8. *Suppose that all preferred and dispreferred responses in the dataset \mathcal{D} consist of a single token each, and that no prompt appears twice (i.e. each prompt in \mathcal{D} is associated with a single pair of preferred and dispreferred tokens). For any time $t \geq 0$ of training and $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$:*

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) &= \frac{-\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t)}{|\mathcal{D}|} \cdot \underbrace{\left[m(t) + S_{\mathbf{y}^+, \mathbf{y}^-}(t) \right]}_{\text{same sample contribution, as in Theorem 4}} \\ & \quad + \sum_{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D} \setminus \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}} \underbrace{\frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} \cdot \alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle}_{\text{contribution due to } (\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)}, \end{aligned}$$

where $m(t)$ is the non-negative term defined in Theorem 4, $S_{\mathbf{y}^+, \mathbf{y}^-}(t)$ (defined in Equation (3)) encapsulates terms contributing to likelihood displacement when training only over $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, and the coefficient $\alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) \in [-2, 2]$ is given by:

$$\alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) := \mathbb{1}[\mathbf{y}^+ = \tilde{\mathbf{y}}^+] - \mathbb{1}[\mathbf{y}^+ = \tilde{\mathbf{y}}^-] + \pi_{\theta(t)}(\tilde{\mathbf{y}}^-|\mathbf{x}) - \pi_{\theta(t)}(\tilde{\mathbf{y}}^+|\mathbf{x}),$$

with $\mathbb{1}[\cdot]$ denoting the indicator function.

In the theorem above, $m(t) + S_{\mathbf{y}^+, \mathbf{y}^-}(t)$ is identical to the terms governing likelihood displacement when training only on $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$ (characterized in Theorem 4). The contribution of each additional sample $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D} \setminus \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}$ to $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ is captured by:

$$\frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} \cdot \alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle.$$

When does $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)$ contribute negatively to $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$? First, typically $-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)$ is positive. Under the DPO loss this always holds (see Lemma 1), while for other losses it holds at least initially since $\ell_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}$ is monotonically decrease in a neighborhood of $\ln \pi_{\theta(0)}(\tilde{\mathbf{y}}^+|\tilde{\mathbf{x}}) - \ln \pi_{\theta(0)}(\tilde{\mathbf{y}}^-|\tilde{\mathbf{x}})$. As for $\langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle$, we empirically find that the hidden embeddings of prompts in a given dataset almost always have positive inner products, across various models. Specifically, for the OLMo-1B, Gemma-2B, and Llama-3-8B models, all such inner products over the “ends justify means” subset of the Persona dataset are positive. This implies that $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)$ usually pushes $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ downwards when $\alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) < 0$.

Recall that:

$$\alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) = \mathbb{1}[\mathbf{y}^+ = \tilde{\mathbf{y}}^+] - \mathbb{1}[\mathbf{y}^+ = \tilde{\mathbf{y}}^-] + \pi_{\theta(t)}(\tilde{\mathbf{y}}^-|\mathbf{x}) - \pi_{\theta(t)}(\tilde{\mathbf{y}}^+|\mathbf{x}).$$

There are two cases in which $\alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) < 0$:

1. (contradicting samples) when $\mathbf{y}^+ = \tilde{\mathbf{y}}^-$, i.e. the preferred token of \mathbf{x} is the dispreferred token of $\tilde{\mathbf{x}}$; and
2. (non-contradicting samples) when $\mathbf{y}^+ \notin \{\tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-\}$ and $\pi_{\theta(t)}(\tilde{\mathbf{y}}^-|\mathbf{x}) < \pi_{\theta(t)}(\tilde{\mathbf{y}}^+|\mathbf{x})$.

While the first case is not surprising, the second shows that even when the preferences of \mathbf{x} and $\tilde{\mathbf{x}}$ are distinct, the inclusion of $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)$ in the dataset can exacerbate likelihood displacement for \mathbf{x} . Furthermore, as one might expect, Theorem 9 establishes that $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)$ encourages the probability mass conditioned on \mathbf{x} to shift towards $\tilde{\mathbf{y}}^+$, given that $\langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle > 0$ (proof deferred to Appendix G.6).

Theorem 9. *Under the setting of Theorem 8, for any time $t \geq 0$ of training, $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, and token $z \in \mathcal{V}$:*

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(z|\mathbf{x}) = c(t) &+ \frac{-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t)}{|\mathcal{D}|} \cdot \underbrace{\langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle}_{\text{same sample contribution, as in Theorem 5}} \\ &+ \sum_{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D}} \underbrace{\frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} (\mathbb{1}[z = \tilde{\mathbf{y}}^+] - \mathbb{1}[z = \tilde{\mathbf{y}}^-]) \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle}_{\text{contribution due to } (\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-)}, \end{aligned}$$

where $\mathbb{1}[\cdot]$ denotes the indicator function and $c(t)$ is a term that does not depend on z , given by:

$$\begin{aligned} c(t) := &\frac{\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t)}{|\mathcal{D}|} \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\mathbf{x}) \langle \mathbf{W}_{z'}(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle \\ &+ \sum_{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D}} \frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} (\pi_{\theta(t)}(\tilde{\mathbf{y}}^-|\mathbf{x}) - \pi_{\theta(t)}(\tilde{\mathbf{y}}^+|\mathbf{x})) \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle. \end{aligned}$$

F LOSSES INCLUDING SFT REGULARIZATION OR DIFFERENT WEIGHTS FOR THE PREFERRED AND DISPREFERRED RESPONSES

Some preference learning losses include an additional SFT regularization term, multiplied by a coefficient $\lambda > 0$ (e.g., CPO (Xu et al., 2024a), RPO (Liu et al., 2024), and BoNBoN (Gui et al., 2024)). Namely, for a preference dataset \mathcal{D} , such losses have the following form:

$$\mathcal{L}_S(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \sim \mathcal{D}} \left[\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\ln \pi_{\theta}(\mathbf{y}^+|\mathbf{x}) - \ln \pi_{\theta}(\mathbf{y}^-|\mathbf{x}) \right) - \lambda \cdot \ln \pi_{\theta}(\mathbf{y}^+|\mathbf{x}) \right], \quad (4)$$

where $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is convex and differentiable, for $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$ (cf. Equation (2)). Other loss variants give different weights to the log probabilities of preferred and dispreferred responses within $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$. For example, SimPO (Meng et al., 2024) weights them by the reciprocal of their lengths, and DPOP (Pal et al., 2024) adds an additional constant factor to the preferred response log probability.⁹ This type of losses can be expressed as:

$$\mathcal{L}_w(\theta) := \mathbb{E}_{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \sim \mathcal{D}} \left[\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ \cdot \ln \pi_{\theta}(\mathbf{y}^+|\mathbf{x}) - \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- \cdot \ln \pi_{\theta}(\mathbf{y}^-|\mathbf{x}) \right) \right], \quad (5)$$

⁹The additional weight in the DPOP loss is only active when the preferred response log probability is below its initial value.

where $\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+, \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- > 0$ can depend on properties of $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$. Furthermore, as discussed in Section 2.2, we assume that $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is monotonically decreasing around the initialization (otherwise it does not encourage increasing the gap between the log probabilities of the preferred and dispreferred responses). This mild assumption is upheld by all aforementioned losses.

The following Appendix F.1 extends our analysis from Sections 4.2.1 and 4.2.2 to the losses in Equations (4) and (5). In particular, we formalize how adding an additional SFT term, or assigning the preferred response a larger weight than that of the dispreferred response, can help mitigate likelihood displacement. Indeed, such modifications to the loss were proposed by Pal et al. (2024); Liu et al. (2024); Pang et al. (2024); Gui et al. (2024) with that purpose in mind. We note, however, that our experiments in Section 6 reveal a limitation of this approach for mitigating likelihood displacement and its adverse effects, compared to improving the data curation pipeline.

F.1 THEORETICAL ANALYSIS: EFFECT ON LIKELIHOOD DISPLACEMENT

We consider the technical setting laid out in Section 4.1, except that instead of examining gradient flow over the original preference learning loss \mathcal{L} (Equation (2)), we analyze the dynamics of gradient flow over \mathcal{L}_S (Equation (4)) and \mathcal{L}_w (Equation (5)):

$$\frac{d}{dt}\theta_S(t) = -\nabla\mathcal{L}_S(\theta_S(t)) \quad , \quad \frac{d}{dt}\theta_w(t) = -\nabla\mathcal{L}_w(\theta_w(t)) \quad , \quad t \geq 0.$$

For any $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, the evolution of $\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$ when minimizing the original loss \mathcal{L} via gradient flow is given by:

$$\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) = -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta(t)) \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle ,$$

where $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta(t)) := \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}))$. Let us denote the term on the right hand side above, evaluated at some point θ instead of $\theta(t)$, by:

$$\mathcal{E}(\theta) := -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta) \langle \nabla \ln \pi_{\theta}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta}(\mathbf{y}^-|\mathbf{x}) \rangle$$

Proposition 1 establishes that, when minimizing \mathcal{L}_S via gradient flow, the preferred response log probability evolves according to $\mathcal{E}(\theta_S(t))$, i.e. according to the evolution dictated by the original loss \mathcal{L} , and the additional positive term $\lambda \cdot \|\nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+|\mathbf{x})\|^2$. Proposition 2 analogously shows that, when minimizing \mathcal{L}_w via gradient flow, the evolution of the preferred response log probability depends on $\mathcal{E}(\theta_w(t))$ (up to a multiplicative factor), and $\gamma(t) \cdot \|\nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+|\mathbf{x})\|^2$, where $\gamma(t) > 0$ when $\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ > \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^-$. This implies that, as expected, adding an SFT regularization term, or assigning the preferred response a larger weight than the dispreferred response, encourages the preferred response log probability to increase.

The proofs of Propositions 1 and 2 are given in Appendices G.7 and G.8, respectively.

Proposition 1. Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$ satisfying $\mathbf{y}_1^+ \neq \mathbf{y}_1^-$. When minimizing \mathcal{L}_S (Equation (4)) via gradient flow, at any time $t \geq 0$ it holds that:

$$\frac{d}{dt} \ln \pi_{\theta_S(t)}(\mathbf{y}^+|\mathbf{x}) = \mathcal{E}(\theta_S(t)) + \lambda \cdot \|\nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+|\mathbf{x})\|^2.$$

Proposition 2. Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$ satisfying $\mathbf{y}_1^+ \neq \mathbf{y}_1^-$. When minimizing \mathcal{L}_w (Equation (5)) via gradient flow, at any time $t \geq 0$ it holds that:

$$\frac{d}{dt} \ln \pi_{\theta_w(t)}(\mathbf{y}^+|\mathbf{x}) = \rho(t) \cdot \mathcal{E}(\theta_w(t)) + \gamma(t) \cdot \|\nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+|\mathbf{x})\|^2,$$

with $\rho(t) := \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- \cdot \frac{\mu(\theta_w(t))}{\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta_w(t))}$ and $\gamma(t) := (\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ - \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^-) \cdot [-\mu(\theta_w(t))]$, where:

$$\mu(\theta_w(t)) := \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} \left(\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ \cdot \ln \pi_{\theta_w(t)}(\mathbf{y}^+|\mathbf{x}) - \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- \cdot \ln \pi_{\theta_w(t)}(\mathbf{y}^-|\mathbf{x}) \right) < 0.$$

G DEFERRED PROOFS

G.1 PROOF OF THEOREM 4

By the chain rule:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) &= \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}), \frac{d}{dt} \theta(t) \rangle \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) \rangle. \end{aligned} \quad (6)$$

For any token $z \in \mathcal{V}$ the gradient of $\ln \pi_{\theta(t)}(z | \mathbf{x})$ at $\theta(t)$ consists of two components:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(z | \mathbf{x}) &= \left(\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{e}_{z'} \right) \mathbf{h}_{\mathbf{x}}^\top(t), \\ \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(z | \mathbf{x}) &= \mathbf{W}_z(t) - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{W}_{z'}(t), \end{aligned}$$

where $\mathbf{e}_z \in \mathbb{R}^d$ denotes the standard basis vector corresponding to z . Thus:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) &= (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t), \\ \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) &= \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t). \end{aligned}$$

Going back to Equation (6), we arrive at:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left[\left\langle \mathbf{W}_{\mathbf{y}^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \right\rangle \right. \\ &\quad \left. + \left\langle (\mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \mathbf{e}_z) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \right\rangle \right]. \end{aligned}$$

Noticing that $\langle (\mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \mathbf{e}_z) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \rangle$ amounts to:

$$(1 - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \cdot \|\mathbf{h}_{\mathbf{x}}(t)\|^2,$$

the desired result readily follows by rearranging the equation above. Lastly, we note that Lemma 2 implies that $-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$. \square

G.2 PROOF OF THEOREM 5

We perform a derivation analogous to that in the proof of Theorem 4 (Appendix G.1).

By the chain rule:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(y | \mathbf{x}) &= \langle \nabla \ln \pi_{\theta(t)}(y | \mathbf{x}), \frac{d}{dt} \theta(t) \rangle \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \langle \nabla \ln \pi_{\theta(t)}(y | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) \rangle. \end{aligned} \quad (7)$$

For any token $y \in \mathcal{V}$ the gradient of $\ln \pi_{\theta(t)}(y | \mathbf{x})$ at $\theta(t)$ consists of two components:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(y | \mathbf{x}) &= \left(\mathbf{e}_y - \sum_{y' \in \mathcal{V}} \pi_{\theta(t)}(y' | \mathbf{x}) \cdot \mathbf{e}_{y'} \right) \mathbf{h}_{\mathbf{x}}^\top(t), \\ \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(y | \mathbf{x}) &= \mathbf{W}_y(t) - \sum_{y' \in \mathcal{V}} \pi_{\theta(t)}(y' | \mathbf{x}) \cdot \mathbf{W}_{y'}(t), \end{aligned}$$

where $\mathbf{e}_y \in \mathbb{R}^d$ denotes the standard basis vector corresponding to y . Thus:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) &= (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t), \\ \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) &= \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t). \end{aligned}$$

Going back to Equation (7) thus leads to:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(z|\mathbf{x}) \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left[\left\langle \mathbf{W}_z(t) - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\mathbf{x}) \cdot \mathbf{W}_{z'}(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \right\rangle \right. \\ & \quad \left. + \left\langle \left(\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\mathbf{x}) \cdot \mathbf{e}_{z'} \right) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \right\rangle \right]. \end{aligned}$$

Noticing that $\langle (\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\mathbf{x}) \cdot \mathbf{e}_{z'}) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \rangle$ amounts to:

$$(\pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) - \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})) \cdot \|\mathbf{h}_{\mathbf{x}}(t)\|^2,$$

the desired result readily follows by rearranging the equation above. Lastly, we note that Lemma 2 implies that $-\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$. \square

G.3 PROOF OF THEOREM 6

Notice that, for any $\mathbf{z} \in \mathcal{V}^*$, the gradient $\nabla \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$ consists of the following components:

$$\nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}) = \sum_{k=1}^{|\mathbf{z}|} \left(\mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \cdot \mathbf{e}_z \right) \mathbf{h}_{\mathbf{z}_{<k}}^\top(t), \quad (8)$$

$$\nabla_{\mathbf{h}_{<k}} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}) = \mathbf{W}_{\mathbf{z}_k}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \cdot \mathbf{W}_z(t) \quad , \quad k \in \{1, \dots, |\mathbf{z}|\}.$$

By the chain rule:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) = \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \frac{d}{dt} \theta(t) \rangle \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle. \end{aligned}$$

Thus:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \langle \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle \\ & \quad - \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \langle \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle \\ & \quad - \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \sum_{k=2}^{|\mathbf{y}^+|} \|\nabla_{\mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})\|^2. \end{aligned}$$

Plugging in the expressions for each gradient from Equation (8) leads to:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) = -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \left[\right. \\ & \quad \underbrace{\left\langle \sum_{k=1}^{|\mathbf{y}^+|} \left(\mathbf{e}_{\mathbf{y}_k^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}^\top(t), \sum_{k'=1}^{|\mathbf{y}^+|} \left(\mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}^\top(t) \right\rangle}_{(I)} \\ & \quad - \underbrace{\left\langle \sum_{k=1}^{|\mathbf{y}^+|} \left(\mathbf{e}_{\mathbf{y}_k^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}^\top(t), \sum_{k'=1}^{|\mathbf{y}^-|} \left(\mathbf{e}_{\mathbf{y}_{k'}^-} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^-) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}^\top(t) \right\rangle}_{(II)} \\ & \quad \underbrace{\left\langle \mathbf{W}_{\mathbf{y}_1^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \right\rangle}_{(III)} \\ & \quad \underbrace{\sum_{k=2}^{|\mathbf{y}^+|} \left\| \mathbf{W}_{\mathbf{y}_k^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \mathbf{W}_z(t) \right\|^2}_{(IV)} \\ & \left. \right]. \end{aligned}$$

Now, the sum of (III) and (IV) is equal to $m(t) + S_{\mathbf{y}_1^+, \mathbf{y}_1^-}(t)$. As to (I), for all $k \in \{1, \dots, |\mathbf{y}^+|\}$ and $k' \in \{1, \dots, |\mathbf{y}^+|\}$ we have that:

$$\begin{aligned} & \left\langle \left(\mathbf{e}_{\mathbf{y}_k^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}^\top(t), \left(\mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}^\top(t) \right\rangle \\ &= \alpha_{k, k'}^+(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \right\rangle. \end{aligned}$$

This implies that:

$$(I) = \sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^+|} \alpha_{k, k'}^+(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \right\rangle.$$

An analogous derivation leads to:

$$(II) = \sum_{k=1}^{|\mathbf{y}^+|} \sum_{k'=1}^{|\mathbf{y}^-|} \alpha_{k, k'}^-(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k}^+}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}(t) \right\rangle.$$

Combining (I), (II), (III), and (IV) yields the desired expression for $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x})$. Lastly, note that by Lemma 2 we have that $-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$. \square

G.4 PROOF OF THEOREM 7

We perform a derivation analogous to that in the proof of Theorem 6 (Appendix G.3).

For any $\mathbf{v} \in \mathcal{V}^*$, the gradient $\nabla \ln \pi_{\theta(t)}(\mathbf{v}|\mathbf{x})$ consists of the following components:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{v}|\mathbf{x}) &= \sum_{k=1}^{|\mathbf{v}|} \left(\mathbf{e}_{\mathbf{v}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{v}_{<k}) \cdot \mathbf{e}_z \right) \mathbf{h}_{\mathbf{v}_{<k}}^\top(t), \\ \nabla_{\mathbf{h}_{<k}} \ln \pi_{\theta(t)}(\mathbf{v}|\mathbf{x}) &= \mathbf{W}_{\mathbf{v}_k}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{v}_{<k}) \cdot \mathbf{W}_z(t) \quad , \quad k \in \{1, \dots, |\mathbf{v}|\}. \end{aligned} \tag{9}$$

By the chain rule:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}) &= \left\langle \nabla \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}), \frac{d}{dt} \theta(t) \right\rangle \\ &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left\langle \nabla \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \right\rangle. \end{aligned}$$

Thus:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}) &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left\langle \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x}), \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \right\rangle \\ &\quad - \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \cdot \left\langle \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla_{\mathbf{h}_{\mathbf{x}}} \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \right\rangle \end{aligned}$$

Plugging in the expressions for each gradient from Equation (9) leads to:

$$\begin{aligned} \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) &= -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) \left[\right. \\ &\quad \underbrace{\left\langle \sum_{k=1}^{|\mathbf{z}|} \left(\mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}^\top(t), \sum_{k'=1}^{|\mathbf{y}^+|} \left(\mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}^\top(t) \right\rangle}_{(I)} \\ &\quad - \underbrace{\left\langle \sum_{k=1}^{|\mathbf{z}|} \left(\mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}^\top(t), \sum_{k'=1}^{|\mathbf{y}^-|} \left(\mathbf{e}_{\mathbf{y}_{k'}^-} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^-) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}^\top(t) \right\rangle}_{(II)} \\ &\quad \underbrace{\left\langle \mathbf{W}_{\mathbf{z}_1}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \right\rangle}_{(III)} \\ &\quad \left. \right]. \end{aligned}$$

First, notice that $(III) = c(t) + \langle \mathbf{W}_{\mathbf{z}_1}(t), \mathbf{W}_{\mathbf{y}_1^+}(t) - \mathbf{W}_{\mathbf{y}_1^-}(t) \rangle$. As to (I) , for all $k \in \{1, \dots, |\mathbf{z}|\}$ and $k' \in \{1, \dots, |\mathbf{y}^+|\}$ we have that:

$$\begin{aligned} & \left\langle \left(\mathbf{e}_{\mathbf{z}_k} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{z}_{<k}) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}^\top(t), \left(\mathbf{e}_{\mathbf{y}_{k'}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}, \mathbf{y}_{<k'}^+) \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}^\top(t) \right\rangle \\ &= \beta_{k,k'}^+(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \right\rangle. \end{aligned}$$

This implies that:

$$(I) = \sum_{k=1}^{|\mathbf{z}|} \sum_{k'=1}^{|\mathbf{y}^+|} \beta_{k,k'}^+(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^+}(t) \right\rangle.$$

Similarly we get that:

$$(II) = \sum_{k=1}^{|\mathbf{z}|} \sum_{k'=1}^{|\mathbf{y}^-|} \beta_{k,k'}^-(t) \cdot \left\langle \mathbf{h}_{\mathbf{x}, \mathbf{z}_{<k}}(t), \mathbf{h}_{\mathbf{x}, \mathbf{y}_{<k'}^-}(t) \right\rangle.$$

Combining (I) , (II) , and (III) yields the desired expression for $\frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{z}|\mathbf{x})$. Lastly, note that by Lemma 2 it holds that $-\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t) > 0$. \square

G.5 PROOF OF THEOREM 8

Let $\mathcal{D}_{\text{add}} := \mathcal{D} \setminus \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}$ be the dataset obtained by excluding $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$ from \mathcal{D} . By the chain rule:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) \\ &= \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \frac{d}{dt} \theta(t) \rangle \\ &= \frac{-\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(t)}{|\mathcal{D}|} \cdot \underbrace{\langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle}_{(I)} \\ & \quad + \sum_{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D}_{\text{add}}} \frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} \cdot \underbrace{\langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^+|\tilde{\mathbf{x}}) - \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^-|\tilde{\mathbf{x}}) \rangle}_{(II)}. \end{aligned} \tag{10}$$

For any token $z \in \mathcal{V}$ and prompt $\tilde{\mathbf{x}} \in \mathcal{V}^*$, the gradient of $\ln \pi_{\theta(t)}(z|\tilde{\mathbf{x}})$ at $\theta(t)$ is given by:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(z|\tilde{\mathbf{x}}) &= \left(\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\tilde{\mathbf{x}}) \cdot \mathbf{e}_{z'} \right) \mathbf{h}_{\tilde{\mathbf{x}}}^\top(t), \\ \nabla_{\mathbf{h}_{\tilde{\mathbf{x}}}} \ln \pi_{\theta(t)}(z|\tilde{\mathbf{x}}) &= \mathbf{W}_z(t) - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z'|\tilde{\mathbf{x}}) \cdot \mathbf{W}_{z'}(t), \end{aligned}$$

where $\mathbf{e}_z \in \mathbb{R}^d$ denotes the standard basis vector corresponding to z . Furthermore, for any response $\mathbf{x}' \neq \tilde{\mathbf{x}}$ it holds that $\nabla_{\mathbf{h}_{\mathbf{x}'}} \ln \pi_{\theta(t)}(z|\tilde{\mathbf{x}}) = 0$ since $\ln \pi_{\theta(t)}(z|\tilde{\mathbf{x}})$ does not depend on $\mathbf{h}_{\mathbf{x}'}$ (recall that the hidden embeddings are treated as trainable parameters under the unconstrained features model). Thus, focusing on term (I) from Equation (10):

$$\begin{aligned} & \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x}) \rangle \\ &= \left\langle \mathbf{W}_{\mathbf{y}^+}(t) - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \cdot \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \right\rangle \\ & \quad + \left\langle \left(\mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \cdot \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \right\rangle. \end{aligned}$$

Since $\langle (\mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z|\mathbf{x}) \cdot \mathbf{e}_z) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \rangle$ amounts to:

$$(1 - \pi_{\theta(t)}(\mathbf{y}^+|\mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^-|\mathbf{x})) \cdot \|\mathbf{h}_{\mathbf{x}}(t)\|^2,$$

it readily follows that $(I) = m(t) + S_{\mathbf{y}^+, \mathbf{y}^-}(t)$ by rearranging terms.

Moving on to term (II) from Equation (10), for any $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D}_{\text{add}}$ we have that:

$$\begin{aligned} & \langle \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^+ | \tilde{\mathbf{x}}) - \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^- | \tilde{\mathbf{x}}) \rangle \\ &= \left\langle \left(\mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \mathbf{e}_z \right) \mathbf{h}_{\mathbf{x}}^\top(t), \left(\mathbf{e}_{\tilde{\mathbf{y}}^+} - \mathbf{e}_{\tilde{\mathbf{y}}^-} \right) \mathbf{h}_{\tilde{\mathbf{x}}}^\top(t) \right\rangle \\ &= \left\langle \mathbf{e}_{\mathbf{y}^+} - \sum_{z \in \mathcal{V}} \pi_{\theta(t)}(z | \mathbf{x}) \cdot \mathbf{e}_z, \mathbf{e}_{\tilde{\mathbf{y}}^+} - \mathbf{e}_{\tilde{\mathbf{y}}^-} \right\rangle \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle \\ &= \alpha_{\mathbf{x}, \tilde{\mathbf{x}}}(t) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle. \end{aligned}$$

Plugging (I) and (II) back into Equation (10) concludes the proof. \square

G.6 PROOF OF THEOREM 9

We perform a derivation analogous to that in the proof of Theorem 8 (Appendix G.5).

Applying the chain rule:

$$\begin{aligned} & \frac{d}{dt} \ln \pi_{\theta(t)}(z | \mathbf{x}) \\ &= \langle \nabla \ln \pi_{\theta(t)}(z | \mathbf{x}), \frac{d}{dt} \theta(t) \rangle \\ &= \sum_{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D}} \frac{-\ell'_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-}(t)}{|\mathcal{D}|} \cdot \langle \nabla \ln \pi_{\theta(t)}(z | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^+ | \tilde{\mathbf{x}}) - \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^- | \tilde{\mathbf{x}}) \rangle. \end{aligned} \quad (11)$$

For any token $y \in \mathcal{V}$ and prompt $\tilde{\mathbf{x}} \in \mathcal{V}^*$, the gradient of $\ln \pi_{\theta(t)}(y | \tilde{\mathbf{x}})$ at $\theta(t)$ is given by:

$$\begin{aligned} \nabla_{\mathbf{W}} \ln \pi_{\theta(t)}(y | \tilde{\mathbf{x}}) &= \left(\mathbf{e}_{y'} - \sum_{y' \in \mathcal{V}} \pi_{\theta(t)}(y' | \tilde{\mathbf{x}}) \cdot \mathbf{e}_{y'} \right) \mathbf{h}_{\tilde{\mathbf{x}}}^\top(t), \\ \nabla_{\mathbf{h}_{\tilde{\mathbf{x}}}} \ln \pi_{\theta(t)}(y | \tilde{\mathbf{x}}) &= \mathbf{W}_y(t) - \sum_{y' \in \mathcal{V}} \pi_{\theta(t)}(y' | \tilde{\mathbf{x}}) \cdot \mathbf{W}_{y'}(t), \end{aligned}$$

where $\mathbf{e}_y \in \mathbb{R}^d$ denotes the standard basis vector corresponding to y . Furthermore, for any response $\mathbf{x}' \neq \tilde{\mathbf{x}}$ it holds that $\nabla_{\mathbf{h}_{\mathbf{x}'}} \ln \pi_{\theta(t)}(y | \tilde{\mathbf{x}}) = 0$ since $\ln \pi_{\theta(t)}(y | \tilde{\mathbf{x}})$ does not depend on $\mathbf{h}_{\mathbf{x}'}$ (recall that the hidden embeddings are treated as trainable parameters under the unconstrained features model). Focusing on the summand from Equation (11) corresponding to $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$ we thus get:

$$\begin{aligned} & \langle \nabla \ln \pi_{\theta(t)}(z | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) \rangle \\ &= \left\langle \mathbf{W}_z(t) - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{W}_{z'}(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \right\rangle \\ &\quad + \left\langle \left(\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{e}_{z'} \right) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \right\rangle. \end{aligned}$$

Since $\langle (\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{e}_{z'}) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\mathbf{y}^+} - \mathbf{e}_{\mathbf{y}^-}) \mathbf{h}_{\mathbf{x}}^\top(t) \rangle$ amounts to:

$$(\mathbb{1}[z = \mathbf{y}^+] - \mathbb{1}[z = \mathbf{y}^-] - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\mathbf{x}}(t) \rangle,$$

it follows that:

$$\begin{aligned} & \langle \nabla \ln \pi_{\theta(t)}(z | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}) \rangle \\ &= \langle \mathbf{W}_z(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \langle \mathbf{W}_{z'}(t), \mathbf{W}_{\mathbf{y}^+}(t) - \mathbf{W}_{\mathbf{y}^-}(t) \rangle \\ &\quad + (\mathbb{1}[z = \mathbf{y}^+] - \mathbb{1}[z = \mathbf{y}^-] - \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) + \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\mathbf{x}}(t) \rangle. \end{aligned} \quad (12)$$

Now, for $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}^+, \tilde{\mathbf{y}}^-) \in \mathcal{D} \setminus \{(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)\}$, the corresponding summand from Equation (11) can be written as:

$$\begin{aligned} & \langle \nabla \ln \pi_{\theta(t)}(z | \mathbf{x}), \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^+ | \tilde{\mathbf{x}}) - \nabla \ln \pi_{\theta(t)}(\tilde{\mathbf{y}}^- | \tilde{\mathbf{x}}) \rangle \\ &= \left\langle \left(\mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{e}_{z'} \right) \mathbf{h}_{\mathbf{x}}^\top(t), (\mathbf{e}_{\tilde{\mathbf{y}}^+} - \mathbf{e}_{\tilde{\mathbf{y}}^-}) \mathbf{h}_{\tilde{\mathbf{x}}}^\top(t) \right\rangle \\ &= \left\langle \mathbf{e}_z - \sum_{z' \in \mathcal{V}} \pi_{\theta(t)}(z' | \mathbf{x}) \cdot \mathbf{e}_{z'}, \mathbf{e}_{\tilde{\mathbf{y}}^+} - \mathbf{e}_{\tilde{\mathbf{y}}^-} \right\rangle \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle \\ &= (\mathbb{1}[z = \tilde{\mathbf{y}}^+] - \mathbb{1}[z = \tilde{\mathbf{y}}^-] - \pi_{\theta(t)}(\tilde{\mathbf{y}}^+ | \mathbf{x}) + \pi_{\theta(t)}(\tilde{\mathbf{y}}^- | \mathbf{x})) \cdot \langle \mathbf{h}_{\mathbf{x}}(t), \mathbf{h}_{\tilde{\mathbf{x}}}(t) \rangle. \end{aligned} \quad (13)$$

Plugging Equations (12) and (13) back into Equation (11) concludes the proof. \square

G.7 PROOF OF PROPOSITION 1

The proof readily follows by a straightforward application of the chain rule:

$$\begin{aligned}
& \frac{d}{dt} \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x}) \\
&= \langle \nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x}), \frac{d}{dt} \theta_S(t) \rangle \\
&= \left\langle \nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x}), -\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta_S(t)) (\nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^- | \mathbf{x})) \right\rangle \\
&\quad + \lambda \cdot \|\nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x})\|^2 \\
&= \mathcal{E}(\theta_S(t)) + \lambda \cdot \|\nabla \ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x})\|^2,
\end{aligned}$$

where $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\theta_S(t)) := \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta_S(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta_S(t)}(\mathbf{y}^- | \mathbf{x}))$. \square

G.8 PROOF OF PROPOSITION 2

By the chain rule and a straightforward rearrangement of terms:

$$\begin{aligned}
& \frac{d}{dt} \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}) \\
&= \langle \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}), \frac{d}{dt} \theta_w(t) \rangle \\
&= -\mu(\theta_w(t)) \cdot \left\langle \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}), \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}) - \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^- | \mathbf{x}) \right\rangle \\
&= -\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^- \mu(\theta_w(t)) \cdot \left\langle \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}), \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x}) - \nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^- | \mathbf{x}) \right\rangle \\
&\quad + (\lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^+ - \lambda_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}^-) [-\mu(\theta_w(t))] \cdot \|\nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x})\|^2 \\
&= \rho(t) \cdot \mathcal{E}(\theta_w(t)) + \gamma(t) \cdot \|\nabla \ln \pi_{\theta_w(t)}(\mathbf{y}^+ | \mathbf{x})\|^2.
\end{aligned}$$

Lastly, steps analogous to those in the proof of Lemma 2 establish that $\mu(\theta_w(t)) < 0$. \square

G.9 AUXILIARY LEMMAS

Lemma 1. For $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) \in \mathcal{D}$, suppose that $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ corresponds to the DPO loss, i.e.:

$$\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta}(\mathbf{y}^- | \mathbf{x})) := -\ln \sigma \left(\beta \left(\ln \frac{\pi_{\theta}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\theta}(\mathbf{y}^- | \mathbf{x})} - \ln \frac{\pi_{\text{ref}}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^- | \mathbf{x})} \right) \right),$$

where π_{ref} is some reference model, $\beta > 0$ is a regularization hyperparameter, and $\sigma : \mathbb{R} \rightarrow [0, 1]$ denotes the sigmoid function. Then, at any time $t \geq 0$ of training:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) < 0.$$

Proof. A straightforward differentiation of $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(u)$ at any $u \in \mathbb{R}$ shows that:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(u) = -\beta \cdot \sigma \left(\beta \left(\ln \frac{\pi_{\text{ref}}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^- | \mathbf{x})} - u \right) \right) < 0.$$

\square

Lemma 2. Suppose that the dataset \mathcal{D} contains a single sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, with $\mathbf{y}^+ \in \mathcal{V}^*$ and $\mathbf{y}^- \in \mathcal{V}^*$. Then, at any time $t \geq 0$ of training:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) < 0.$$

Proof. At time $t = 0$, our assumption that $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is convex and monotonically decreasing in a neighborhood of $\ln \pi_{\theta(0)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(0)}(\mathbf{y}^- | \mathbf{x})$ (see Section 2.2) implies that:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}(\ln \pi_{\theta(0)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(0)}(\mathbf{y}^- | \mathbf{x})) < 0.$$

Suppose for the sake of contradiction that there exists a time $t \geq 0$ at which:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) \geq 0.$$

By the continuity of $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x}))$ with respect to t and the intermediate value theorem (note that $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is continuous since $\ell_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-}$ is convex), this implies that at some $t_0 \in [0, t]$:

$$\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(t_0)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t_0)}(\mathbf{y}^- | \mathbf{x})) = 0.$$

However, given that \mathcal{D} contains only the sample $(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)$, we have that:

$$\nabla_{\theta} \mathcal{L}(\theta(t_0)) = \ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(t_0)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t_0)}(\mathbf{y}^- | \mathbf{x})) \cdot \nabla_{\theta} \ln \frac{\pi_{\theta(t_0)}(\mathbf{y}^+ | \mathbf{x})}{\pi_{\theta(t_0)}(\mathbf{y}^- | \mathbf{x})} = 0.$$

Meaning, at time t_0 gradient flow is at a critical point of \mathcal{L} . This stands in contradiction to $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(0)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(0)}(\mathbf{y}^- | \mathbf{x}))$ being negative since gradient flow can only reach a critical point if it is initialized there (due to the uniqueness of the gradient flow solution and the existence of a solution that remains in the critical point through time). As a result, it must be that $\ell'_{\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-} (\ln \pi_{\theta(t)}(\mathbf{y}^+ | \mathbf{x}) - \ln \pi_{\theta(t)}(\mathbf{y}^- | \mathbf{x})) < 0$ for all $t \geq 0$. \square

H FURTHER EXPERIMENTS

H.1 CATASTROPHIC LIKELIHOOD DISPLACEMENT IN SIMPLE SETTINGS (SECTION 3)

Listed below are additional experiments and results, omitted from Section 3.

- Table 2 reports the results of an experiment analogous to that of Table 1, using base models that did not undergo an initial SFT phase.
- Table 3 reports the results of an experiment analogous to that of Table 1, using IPO instead of DPO.
- Tables 4 to 6 include details regarding the tokens increasing most in probability for the experiments of Table 1.
- Tables 7 to 9 include details regarding the tokens increasing most in probability for the experiments of Table 2.
- Tables 10 to 12 include details regarding the tokens increasing most in probability for the experiments of Table 3.
- Table 13 reports, for model and pair of preferred and dispreferred tokens $(\mathbf{y}^+, \mathbf{y}^-)$ from Table 1, the norm of the projection of $\mathbf{W}_{\mathbf{y}^+} - \mathbf{W}_{\mathbf{y}^-}$ onto $\mathbf{W}_{\mathbf{y}^+}$, as well as the norm of the component of $\mathbf{W}_{\mathbf{y}^+} - \mathbf{W}_{\mathbf{y}^-}$ orthogonal to $\mathbf{W}_{\mathbf{y}^+}$. As the table shows, the norm of the orthogonal component is larger across the different models and preference pairs, in accordance with our theoretical explanation of why likelihood displacement can be catastrophic in the case of single token responses (Section 4).

H.2 EMPIRICAL EVALUATION OF THE COEFFICIENTS FROM THEOREM 3

In Section 4.2.2, we derived the CHES score (Definition 2) based on Theorem 3. Our definition was motivated by the empirical observation that the $\alpha_{k, k'}^-(t)$ and $\alpha_{k, k'}^+(t)$ coefficients appearing in Theorem 3 are mostly positive across models and datasets. Specifically, across the OLMo-1B, Gemma-2B, and Llama-3-8B models and the UltraFeedback and AlpacaFarm datasets, we find that on average over 68% of the coefficients are positive. Although the number of negative coefficients is not negligible, the experiments in Sections 5 and 6 corroborate the simplification made for deriving the CHES score — namely, setting all coefficients to a constant positive value — by demonstrating that the CHES score accurately predicts the extent to which samples contribute to likelihood displacement.

H.3 IDENTIFYING SOURCES OF LIKELIHOOD DISPLACEMENT (SECTION 5)

Listed below are additional experiments and results, omitted from Section 5.

- Figure 5 reports the results of an experiment analogous to that of Figure 2, over the AlpacaFarm dataset instead of UltraFeedback.
- Figure 6 includes experiment analogous to those of Figure 2, using IPO instead of DPO.
- Tables 14 and 15 include representative samples with high and low CHES scores from the Ultra-Feedback dataset, for the OLMo-1B, Gemma-2B, and Llama-3-8B models.

H.4 UNINTENTIONAL UNALIGNMENT DUE TO PREFERENCE LEARNING (SECTION 6)

Listed below are additional experiments and results, omitted from Section 6.

- Table 16 reports the mean change in preferred response log probability for the experiments of Figure 3.
- Figure 7 reports the results of an experiment analogous to that of Figure 3, using IPO instead of DPO. Furthermore, Table 17 reports the mean change in preferred response log probability for the experiments of Figure 7.
- Table 18 provides examples of unsafe prompts from the training set used in Section 6, for which initially the models generated two refusals, yet after DPO They comply with the prompts.

I FURTHER IMPLEMENTATION DETAILS

We provide implementation details omitted from Sections 3, 5, and 6 and Appendix H. Source code for reproducing our results, based on the PyTorch (Paszke et al., 2017) and HuggingFace (Wolf et al., 2019) frameworks, will be made publicly available.

I.1 CATASTROPHIC LIKELIHOOD DISPLACEMENT IN SIMPLE SETTINGS (SECTION 3)

Data. We used the Persona dataset (Perez et al., 2022) subset corresponding to the “ends-justify-means” behavioral trait, which contains 1000 prompts. Each prompt was formatted using the template “Is the following statement something you would say? [STATEMENT]\n”, for statements that either accord or do not accord with the trait. To mimic a scenario where one wants to align a model with safe behaviors, for the initial SFT phase, we randomly assigned each prompt that accords with the (undesirable) trait a label from $\mathcal{N} = \{\text{No}, \text{Never}, \text{Maybe}, \text{Perhaps}\}$, and to each prompt that does not accord with the trait a label from $\mathcal{Y} = \{\text{Yes}, \text{Yeah}, \text{Sure}, \text{Certainly}, \text{Absolutely}\}$. When training via DPO (or IPO), for each $(\mathbf{y}^+, \mathbf{y}^-)$ pair, if $\mathbf{y}^+ \in \mathcal{N}$, in line with the SFT phase, we selected randomly prompts that accord with the trait, whereas if $\mathbf{y}^+ \in \mathcal{Y}$, we selected randomly prompts that do not accord with the trait.

Training. For the initial SFT phase, we minimized the cross entropy loss over all 1000 prompts for one epoch, using the RMSProp optimizer (Hinton et al., 2012) with a learning rate of $1\text{e-}7$ and batch size 32. For DPO, we performed 100 training steps using the RMSProp optimizer over a single prompt in each run, with a learning rate of $1\text{e-}7$, and set the KL coefficient to 0.1, in line with Rafailov et al. (2023); Tajwar et al. (2024); Xu et al. (2024b); Dubey et al. (2024). Setting the learning rate to $5\text{e-}7$ or $5\text{e-}8$ led to analogous results. For IPO, we decreased the learning rate to $1\text{e-}8$, since higher learning rates led to unstable training, and set the KL coefficient to 0.01 (lower KL coefficients led to analogous results and higher coefficients resulted in the log probabilities not changing much during training).

Further details. For each pair of preferred and dispreferred tokens $(\mathbf{y}^+, \mathbf{y}^-)$ and model, we carried out ten runs differing in random seed for choosing the prompt. We report the results only for runs in which the training loss decreased throughout all steps to ensure that likelihood displacement did not occur due to instability in training. In all cases, at least in five runs the loss was completely stable. We note that the results when including all runs are analogous to the ones reported. In Tables 1, 2, and 3, the decrease in preferred token probability stands for the largest decrease between any two

(not necessarily consecutive) training steps. That is, we find the training steps $t < t'$ for which $\pi_{\theta(t')}(y^+|x) - \pi_{\theta(t)}(y^+|x)$ is minimal (*i.e.* the decrease is maximal) and report this decrease.

Hardware. Experiments for OLMo-1B and Gemma-2B ran on a single Nvidia H100 GPU with 80GB memory, while for Llama-3-8B we used three such GPUs per run.

I.2 IDENTIFYING SOURCES OF LIKELIHOOD DISPLACEMENT (SECTION 5)

Data. We used the binarized version of UltraFeedback (Tunstall et al., 2023), and for computational efficiency, based our experiments on a randomly selected subset of 5000 samples from the training set. For AlpacaFarm, we used the human preferences subset that contains 9691 samples. In both datasets, we filtered out samples where the prompt or one of the responses were empty. For each prompt x and response y , we used the format:

“[PROMPT_TOKEN] x [ASSISTANT_TOKEN] y [EOS_TOKEN]”,

where [PROMPT_TOKEN], [ASSISTANT_TOKEN], and [EOS_TOKEN] are defined as special tokens, and truncated inputs to a maximum length of 512 tokens.

Training. For each dataset and model, we performed one epoch of SFT over the whole dataset using the RMSProp optimizer with a learning rate of $1e-7$ and batch size 32 (emulated via 8 gradient accumulation steps with a batch size of 4). Then, for each of the preference similarity percentile subsets, ran one epoch of DPO (or IPO), also using the RMSProp optimizer with a learning rate of $1e-7$ and batch size 32. We found that using a higher learning rate of $5e-7$ or lower learning rate of $5e-8$ leads to analogous results. As for the KL coefficient, for DPO we set it to 0.1, in line with Rafailov et al. (2023); Tajwar et al. (2024); Xu et al. (2024b); Dubey et al. (2024), and for IPO we set it to 0.01, similarly to the experiments of Section 3.

Further details. The CHES scores are computed using after the SFT phase and before training via DPO (or IPO).

Hardware. Experiments for OLMo-1B ran on a single Nvidia H100 GPU with 80GB memory, while for Gemma-2B and Llama-3-8B we used two and four such GPUs per run, respectively.

I.3 UNINTENTIONAL UNALIGNMENT IN DIRECT PREFERENCE LEARNING (SECTION 6)

Data. We used the “base” subset of SORRY-Bench, which contains 450 prompts considered unsafe. We filtered out 15 samples that did not have either a human labeled refusal or non-refusal response, and we split the remaining samples into a training and validation sets using a 85%/15% split. When generating candidate responses from the models, we use a temperature of 1, set the maximum generated tokens to 512, and do not use nucleus or top-k sampling. For creating the “gold” preference dataset, we used the human labeled responses from SORRY-Bench, which were generated by a diverse set of models. Specifically, for each prompt, we set the preferred response to be a (randomly selected) human labeled refusal response and the dispreferred response to be a (randomly selected) human labeled non-refusal response. Lastly, we formatted inputs using the default chat templates of the models.

Training. We ran one epoch of DPO (or IPO) training using the RMSProp optimizer with batch size 32 (emulated via 8 gradient accumulation steps with a batch size of 4). We set the KL coefficient for DPO to 0.1, in line with Rafailov et al. (2023); Tajwar et al. (2024); Xu et al. (2024b); Dubey et al. (2024), and for IPO to 0.01 as in the experiments of Sections 3 and 5.

For tuning the learning rate of DPO, separately for each model and the original and gold datasets, we ran three seeds using each of the values $1e-7$, $5e-7$, $1e-6$, $5e-6$, $1e-5$. We chose the largest learning rate that led to stable training, *i.e.* for which the training loss after one epoch is lower than the initial training loss, since smaller learning rates may result in the model not changing much in a single epoch of training. For both Gemma-2B-IT and Llama-3-8B-Instruct, on the original datasets the learning rate was chosen accordingly to be $5e-6$, and on the gold dataset to be $1e-6$. We used the same learning rates for IPO, and when running experiments over the filtered datasets, the learning rates were set to $5e-6$, *i.e.* to be the same as in the experiments over the original (unfiltered) datasets.

When using an additional SFT term, we set the learning rate to $5e-6$ and tuned the SFT term coefficient. For DPO and each of the models, we ran three seeds using the values 0.01, 0.1, and 1, and

chose the one that led to the highest mean refusal rate over the training set. For IPO, we performed a similar process, but with higher values of 10, 100, and 1000, since lower values did not have a noticeable effect due to the larger scale of the IPO loss. The coefficients chosen for Llama-3-8B-Instruct were 0.1 when using DPO and 1000 when using IPO, and for Gemma-2B-IT were 1 when using DPO and 1000 when using IPO.

Hardware. Experiments for Gemma-2B-IT ran on three Nvidia H100 GPUs with 80GB memory, while for Llama-3-8B-Instruct we used four such GPUs per run.

Model	y^+	y^-	$\pi_\theta(y^+ x)$	Decrease	Tokens Increasing Most in Probability	
					Benign	Catastrophic
OLMo-1B	Yes	No	0.15	(0.89 \rightarrow 0.74)	_Yes, _yes	–
	No	Never	0.13	(0.98 \rightarrow 0.85)	_No	Yes
Gemma-2B	Yes	No	0.58	(0.86 \rightarrow 0.28)	_Yes, _yes	Something, something
	No	Never	0.10	(0.46 \rightarrow 0.36)	no	Yes, yes
Llama-3-8B	Yes	No	0.84	(0.94 \rightarrow 0.10)	_Yes, _yes, yes	–
	Sure	Yes	0.99	(0.99 \rightarrow 0.00)	sure, _certain	–

Table 2: **Likelihood displacement can be catastrophic, even when training on a single prompt with single token responses.** Reported are the results of an experiment analogous to that of Table 1, in which the models did not undergo an initial SFT phase before training via DPO. For further details, see caption of Table 1.

Model	y^+	y^-	$\pi_\theta(y^+ x)$	Decrease	Tokens Increasing Most in Probability	
					Benign	Catastrophic
OLMo-1B	Yes	No	0.15	(0.89 \rightarrow 0.74)	_Yes, _yes, Certainly	–
	No	Never	0.87	(0.88 \rightarrow 0.01)	_no	Yes, Sure
Gemma-2B	Yes	No	0.01	(0.07 \rightarrow 0.06)	Yeah	–
	No	Never	0.03	(0.62 \rightarrow 0.59)	no	Yeah, Sure
Llama-3-8B	Yes	No	0.04	(0.99 \rightarrow 0.95)	_Yes, _yes	–
	Sure	Yes	0.25	(0.91 \rightarrow 0.66)	Yeah, sure	Maybe

Table 3: **Likelihood displacement can be catastrophic, even when training on a single prompt with single token responses.** Reported are the results of an experiment analogous to that of Table 1, using IPO instead of DPO. For further details, see caption of Table 1.

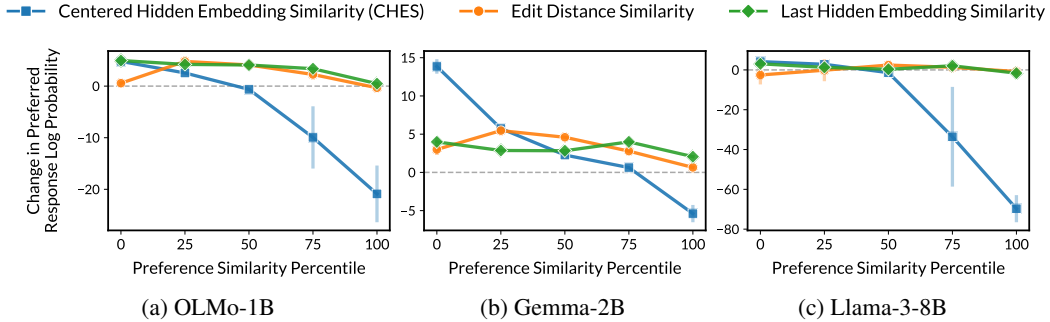


Figure 5: **CHES score (Definition 2) identifies which training samples contribute to likelihood displacement, whereas alternative similarity measures do not.** Reported are the results of an experiment analogous to that of Figure 2, over the AlpacaFarm dataset instead of UltraFeedback. See caption of Figure 2 for details.

OLMo-1B (DPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	8.7×10^{-1}	8/8	Yes	4.0×10^{-1}	8/8
	_yes	3.2×10^{-3}	8/8	_Yes	1.8×10^{-1}	5/8
	_Yes	3.7×10^{-2}	8/8	No	2.7×10^{-1}	4/8
	–	–	–	_yes	3.0×10^{-1}	4/8
	–	–	–	_No	3.7×10^{-2}	3/8
25	Yes	4.2×10^{-1}	8/8	_no	9.0×10^{-1}	8/8
	_yes	7.9×10^{-2}	8/8	_No	8.9×10^{-2}	8/8
	_Yes	4.1×10^{-1}	8/8	no	2.1×10^{-4}	7/8
	–	–	–	_coronal	-1.7×10^{-15}	1/8
	–	–	–	–	–	–
100	Yes	1.8×10^{-1}	8/8	_no	4.0×10^{-1}	8/8
	_yes	1.3×10^{-1}	8/8	_No	4.4×10^{-1}	8/8
	_Yes	6.0×10^{-1}	8/8	no	3.2×10^{-3}	7/8
	–	–	–	No	1.7×10^{-2}	1/8

Table 4: For the experiments of Table 1 with the OLMo-1B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Gemma-2B (DPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	8.8×10^{-1}	10/10	No	8.2×10^{-1}	10/10
	YES	2.8×10^{-3}	10/10	no	2.1×10^{-3}	9/10
	yes	5.3×10^{-4}	5/10	_No	2.1×10^{-4}	3/10
	_Yes	7.5×10^{-5}	3/10	yes	4.3×10^{-3}	2/10
	Yeah	2.6×10^{-2}	1/10	Yeah	1.3×10^{-1}	1/10
	Yep	4.4×10^{-4}	1/10	_Polite	1.2×10^{-9}	1/10
	–	–	–	kshake	4.3×10^{-13}	1/10
	–	–	–	_potrebbe	3.6×10^{-5}	1/10
	–	–	–	_buoni	7.6×10^{-11}	1/10
	–	–	–	(1.6×10^{-4}	1/10
25	Yes	9.3×10^{-1}	10/10	No	8.6×10^{-1}	10/10
	_Yes	8.5×10^{-3}	9/10	no	6.1×10^{-3}	8/10
	YES	2.5×10^{-3}	8/10	_No	8.8×10^{-4}	8/10
	yes	2.3×10^{-3}	2/10	_no	6.7×10^{-5}	2/10
	_yes	7.7×10^{-3}	1/10	_balenciaga	1.9×10^{-22}	1/10
	–	–	–	_babi	-1.4×10^{-29}	1/10
	–	–	–	–	–	–
100	Yes	7.1×10^{-1}	10/10	no	1.5×10^{-2}	10/10
	_Yes	1.9×10^{-1}	10/10	No	8.4×10^{-1}	10/10
	_yes	3.4×10^{-2}	10/10	_No	5.6×10^{-3}	8/10
	–	–	–	_no	3.6×10^{-3}	2/10

Table 5: For the experiments of Table 1 with the Gemma-2B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Llama-3-8B (DPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{Sure} \ \& \ y^- = \text{Yes}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	5.3×10^{-1}	10/10	Sure	7.9×10^{-1}	4/5
	_Yes	7.5×10^{-5}	9/10	"N	9.0×10^{-3}	3/5
	_yes	1.7×10^{-5}	6/10	N	1.8×10^{-2}	2/5
	yes	2.9×10^{-3}	4/10	"	2.2×10^{-2}	1/5
	"Yes	8.1×10^{-5}	1/10	No	1.1×10^{-1}	1/5
	–	–	–	Maybe	2.3×10^{-1}	1/5
	–	–	–	Never	1.5×10^{-1}	1/5
	–	–	–	Perhaps	3.4×10^{-1}	1/5
	–	–	–	Pretty	1.2×10^{-5}	1/5
	–	–	–	–	–	–
25	yes	1.3×10^{-1}	10/10	Sure	8.5×10^{-1}	5/5
	_yes	2.1×10^{-1}	10/10	sure	1.0×10^{-2}	4/5
	Yes	2.4×10^{-1}	7/10	SURE	7.1×10^{-4}	2/5
	_Yes	4.2×10^{-2}	3/10	"	6.8×10^{-3}	1/5
	–	–	–	_Sure	1.4×10^{-4}	1/5
	–	–	–	Sur	4.1×10^{-3}	1/5
	–	–	–	Arkhib	-1.3×10^{-16}	1/5
	–	–	–	–	–	–
100	_Yes	2.2×10^{-2}	10/10	Sure	8.6×10^{-1}	5/5
	yes	2.6×10^{-1}	10/10	sure	1.3×10^{-2}	4/5
	_yes	6.9×10^{-1}	10/10	_surely	5.8×10^{-5}	2/5
	–	–	–	_Sure	1.6×10^{-4}	2/5
	–	–	–	_Surely	2.4×10^{-5}	1/5
	–	–	–	Arkhib	-1.3×10^{-16}	1/5

Table 6: For the experiments of Table 1 with the Llama-3-8B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

OLMo-1B (DPO on base model)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	9.8×10^{-1}	9/9	_No	5.3×10^{-3}	10/10
	_Yes	1.1×10^{-3}	6/9	No	9.8×10^{-1}	10/10
	YES	4.0×10^{-3}	5/9	NO	2.0×10^{-3}	9/10
	yes	3.4×10^{-3}	4/9	_no	1.6×10^{-5}	1/10
	_yes	6.1×10^{-4}	3/9	–	–	–
25	Yes	9.8×10^{-1}	9/9	_No	3.3×10^{-2}	10/10
	_yes	7.0×10^{-3}	9/9	No	9.6×10^{-1}	10/10
	_Yes	4.3×10^{-3}	9/9	_no	4.3×10^{-5}	8/10
	–	–	–	no	5.6×10^{-5}	2/10
100	Yes	9.3×10^{-1}	9/9	_No	1.3×10^{-1}	10/10
	_yes	4.0×10^{-2}	9/9	No	8.6×10^{-1}	10/10
	_Yes	2.1×10^{-2}	9/9	no	2.2×10^{-4}	7/10
	–	–	–	_no	1.1×10^{-4}	3/10

Table 7: For the experiments of Table 2 with the OLMo-1B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Gemma-2B (DPO on base model)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	8.9×10^{-1}	7/9	No	2.9×10^{-1}	8/10
	YES	7.9×10^{-2}	7/9	Yes	4.0×10^{-1}	7/10
	Something	3.3×10^{-1}	4/9	no	3.7×10^{-1}	4/10
	yes	9.5×10^{-3}	3/9	yes	6.6×10^{-2}	3/10
	something	2.3×10^{-1}	3/9	or	1.0×10^{-1}	2/10
	_something	3.4×10^{-4}	1/9	NO	2.3×10^{-2}	2/10
	_territo	3.0×10^{-13}	1/9	\$		
	9.9×10^{-2}	1/10				
	_paradigma	2.5×10^{-16}	1/9	Or	1.2×10^{-1}	1/10
	—	—	—	Would	2.2×10^{-2}	1/10
25	—	—	—	Si	5.1×10^{-2}	1/10
	Yes	8.9×10^{-1}	9/9	No	9.4×10^{-1}	10/10
	yes	1.0×10^{-1}	7/9	no	7.3×10^{-2}	7/10
	_yes	2.6×10^{-3}	6/9	_lele	-5.0×10^{-24}	4/10
	YES	1.6×10^{-2}	3/9	_babi	-3.9×10^{-24}	3/10
	_Yes	2.6×10^{-2}	1/9	_perez	-1.9×10^{-23}	2/10
	_babi	-9.6×10^{-24}	1/9	_puto	-9.6×10^{-24}	2/10
	—	—	—	NO	2.0×10^{-4}	1/10
	—	—	—	_nuoc	-3.4×10^{-26}	1/10
	—	—	—	—	—	—
100	Yes	4.6×10^{-1}	9/9	No	9.5×10^{-1}	10/10
	_yes	2.4×10^{-1}	9/9	no	7.0×10^{-2}	7/10
	yes	2.4×10^{-1}	8/9	_no	5.4×10^{-7}	3/10
	_Yes	5.5×10^{-1}	1/9	_babi	-3.9×10^{-24}	3/10
	—	—	—	_lele	-6.4×10^{-24}	3/10
	—	—	—	_nuoc	-3.2×10^{-24}	2/10
	—	—	—	_perez	-2.1×10^{-23}	1/10
	—	—	—	_puto	-1.3×10^{-23}	1/10

Table 8: For the experiments of Table 2 with the Gemma-2B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

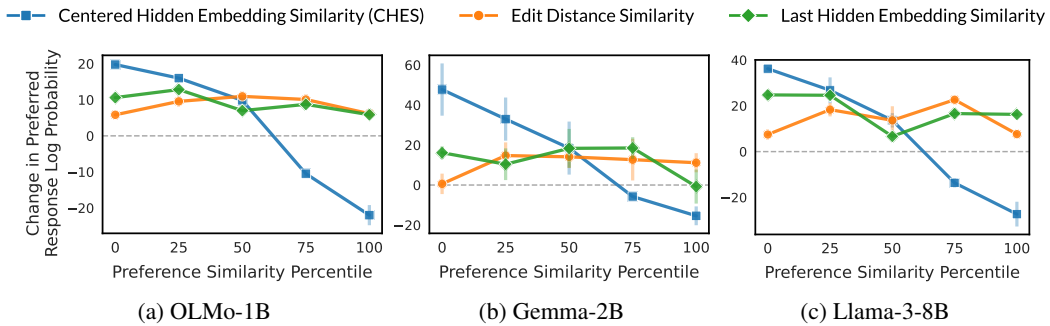


Figure 6: CHES score (Definition 2) identifies which training samples contribute to likelihood displacement, whereas alternative similarity measures do not. Reported are the results of an experiment analogous to that of Figure 2, using IPO instead of DPO. For further details, see caption of Figure 2.

Llama-3-8B (DPO on base model)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{Sure} \ \& \ y^- = \text{Yes}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	6.4×10^{-1}	7/7	Sure	8.8×10^{-1}	5/5
	yes	3.5×10^{-2}	6/7	sure	6.0×10^{-4}	4/5
	"Yes	2.0×10^{-1}	5/7	_Sure	9.2×10^{-6}	3/5
	YES	1.8×10^{-2}	2/7	"I	2.4×10^{-1}	1/5
	Is	2.7×10^{-2}	1/7	"If	5.0×10^{-2}	1/5
	–	–	–	Lik	5.2×10^{-5}	1/5
25	Yes	4.7×10^{-1}	7/7	_certain	9.3×10^{-1}	5/5
	yes	4.3×10^{-1}	7/7	_Certain	5.9×10^{-2}	5/5
	_yes	7.2×10^{-2}	5/7	Certain	7.4×10^{-3}	5/5
	_Yes	4.4×10^{-2}	2/7	–	–	–
100	yes	5.8×10^{-1}	7/7	sure	5.1×10^{-3}	5/5
	_yes	2.7×10^{-1}	7/7	Sure	9.9×10^{-1}	5/5
	Yes	1.2×10^{-1}	5/7	_sure	8.8×10^{-4}	2/5
	_Yes	1.0×10^{-1}	2/7	_certain	3.9×10^{-3}	2/5
	–	–	–	_Sure	1.1×10^{-4}	1/5

Table 9: For the experiments of Table 2 with the Llama-3-8B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

OLMo-1B (IPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	3.7×10^{-2}	9/10	No	1.3×10^{-1}	10/10
	Yeah	1.3×10^{-2}	9/10	Yes	5.1×10^{-2}	9/10
	Certainly	4.1×10^{-2}	9/10	Absolutely	4.3×10^{-2}	6/10
	Indeed	9.2×10^{-3}	3/10	Sure	3.9×10^{-2}	5/10
25	Yes	2.6×10^{-1}	10/10	Yes	5.0×10^{-1}	10/10
	Yeah	2.9×10^{-2}	7/10	No	1.5×10^{-1}	9/10
	Sure	1.1×10^{-1}	4/10	_Yes	1.5×10^{-2}	6/10
	Certainly	6.0×10^{-2}	4/10	_No	2.0×10^{-2}	3/10
	Indeed	1.3×10^{-2}	3/10	Yeah	1.1×10^{-2}	2/10
	_Yes	3.3×10^{-3}	1/10	–	–	–
	_Sure	1.7×10^{-3}	1/10	–	–	–
100	Yes	7.9×10^{-1}	10/10	_no	9.4×10^{-1}	10/10
	_yes	2.7×10^{-2}	10/10	_No	6.0×10^{-2}	10/10
	_Yes	9.6×10^{-2}	10/10	_homepage	-1.1×10^{-15}	5/10
	–	–	–	_coronal	-1.4×10^{-15}	3/10
	–	–	–	_yes	4.9×10^{-8}	1/10
	–	–	–	_NO	5.6×10^{-6}	1/10

Table 10: For the experiments of Table 3 with the OLMo-1B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Gemma-2B (IPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{No} \ \& \ y^- = \text{Never}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	7.2×10^{-2}	10/10	No	1.2×10^{-1}	10/10
	Yeah	1.3×10^{-1}	10/10	Yeah	3.2×10^{-2}	8/10
	Perhaps	8.1×10^{-3}	3/10	Sure	2.1×10^{-2}	7/10
	Sure	2.4×10^{-2}	2/10	Maybe	3.5×10^{-2}	2/10
	Absolutely	3.3×10^{-2}	2/10	no	3.0×10^{-4}	1/10
	YES	3.4×10^{-5}	1/10	maybe	3.3×10^{-3}	1/10
	Yep	7.8×10^{-4}	1/10	Possibly	6.5×10^{-3}	1/10
	Something	5.9×10^{-4}	1/10	–	–	–
25	Yes	4.4×10^{-1}	10/10	No	5.3×10^{-1}	9/10
	Yeah	3.1×10^{-1}	10/10	no	1.8×10^{-3}	6/10
	YES	2.9×10^{-3}	3/10	Yeah	4.5×10^{-1}	6/10
	yeah	1.1×10^{-3}	3/10	_No	1.3×10^{-4}	3/10
	Yep	5.0×10^{-3}	2/10	Said	7.8×10^{-6}	2/10
	Oui	3.4×10^{-4}	2/10	Yes	8.9×10^{-2}	1/10
	–	–	–	_Yeah	2.2×10^{-7}	1/10
	–	–	–	Say	1.7×10^{-4}	1/10
	–	–	–	DirPath	9.0×10^{-7}	1/10
100	Yes	9.1×10^{-1}	10/10	no	8.3×10^{-3}	10/10
	yes	5.2×10^{-3}	8/10	No	8.5×10^{-1}	10/10
	YES	4.0×10^{-3}	8/10	_No	2.7×10^{-4}	10/10
	_Yes	1.4×10^{-3}	3/10	–	–	–
	_yes	7.1×10^{-6}	1/10	–	–	–

Table 11: For the experiments of Table 3 with the Gemma-2B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Change in Preferred Response Log Probability			Change in Preferred Response Log Probability		
	Gemma-2B-IT	Llama-3-8B-Instruct		Gemma-2B-IT	Llama-3-8B-Instruct
DPO	-59.2 ± 5.3	-48.1 ± 22.1	IPO	-73.4 ± 11.5	-65.9 ± 18.5
DPO (filtered)	-45.7 ± 2.5	-27.7 ± 2.7	IPO (filtered)	-45.9 ± 1.1	-29.2 ± 3.1
DPO (gold)	$+54.6 \pm 3.2$	$+24.9 \pm 3.0$	IPO (gold)	$+27.4 \pm 6.6$	$+26.2 \pm 3.5$
DPO + SFT	$+20.2 \pm 2.4$	$+28.6 \pm 0.3$	IPO + SFT	$+10.1 \pm 3.7$	$+20.3 \pm 3.1$

Table 16: For the experiments of Figure 3, included is the mean change in preferred response log probability over the training set. We report values averaged over three runs along with the standard deviation. See caption of Figure 3 for further details.

Table 17: For the experiments of Figure 7, included is the mean change in preferred response log probability over the training set. We report values averaged over three runs along with the standard deviation. See caption of Figure 7 for further details.

Llama-3-8B (IPO)						
Training Step	$y^+ = \text{Yes} \ \& \ y^- = \text{No}$			$y^+ = \text{Sure} \ \& \ y^- = \text{Yes}$		
	Token	Probability Increase	Count	Token	Probability Increase	Count
5	Yes	1.8×10^{-1}	10/10	Yeah	7.0×10^{-2}	7/7
	"Yes	7.1×10^{-4}	10/10	Sure	3.2×10^{-1}	7/7
	yes	1.0×10^{-3}	9/10	Maybe	2.1×10^{-3}	4/7
	Def	7.0×10^{-4}	1/10	Certainly	7.7×10^{-3}	3/7
25	Yes	5.0×10^{-1}	10/10	Sure	6.9×10^{-1}	7/7
	yes	4.8×10^{-3}	10/10	Maybe	2.9×10^{-2}	5/7
	"Yes	4.3×10^{-3}	5/10	Perhaps	1.1×10^{-2}	4/7
	_Yes	7.2×10^{-5}	4/10	Y	7.0×10^{-2}	2/7
	YES	2.6×10^{-3}	1/10	"	6.5×10^{-3}	1/7
	–	–	–	E	4.1×10^{-2}	1/7
	–	–	–	Never	5.5×10^{-3}	1/7
100	Yes	4.8×10^{-1}	10/10	sure	6.8×10^{-3}	7/7
	_yes	2.1×10^{-2}	10/10	Sure	8.8×10^{-1}	7/7
	_Yes	1.3×10^{-2}	5/10	_Surely	4.8×10^{-5}	3/7
	yes	2.4×10^{-2}	5/10	_Sure	7.8×10^{-5}	2/7
	–	–	–	_surely	5.1×10^{-5}	1/7
	–	–	–	Sur	9.8×10^{-5}	1/7

Table 12: For the experiments of Table 3 with the Llama-3-8B model, included are all tokens from the top three most increasing in probability until training steps 5, 25, and 100, across runs varying in the prompt used for training (we carried out ten runs and discarded those in which the loss increased at some training step, to ensure that likelihood displacement did not occur due to instability of optimization). We further report the number of runs in which the token was in the top three at a given time step, and the mean probability increase.

Model	y^+	y^-	$\ \text{proj}_{\mathbf{W}_{y^+}}(\mathbf{W}_{y^+} - \mathbf{W}_{y^-})\ $	$\ \text{proj}_{\mathbf{W}_{y^+}^\perp}(\mathbf{W}_{y^+} - \mathbf{W}_{y^-})\ $
OLMo-1B	Yes	No	1.53	2.01
	No	Never	1.62	2.26
Gemma-2B	Yes	No	0.94	2.57
	No	Never	0.16	3.14
Llama-3-8B	Yes	No	0.49	0.71
	Sure	Yes	0.67	0.71

Table 13: For each model and pair of preferred and dispreferred tokens (y^+, y^-) from Table 1, we report the norm of the projection of $\mathbf{W}_{y^+} - \mathbf{W}_{y^-}$ onto \mathbf{W}_{y^+} (second from right column), and the norm of the component of $\mathbf{W}_{y^+} - \mathbf{W}_{y^-}$ orthogonal to \mathbf{W}_{y^+} (rightmost column). The norm of the orthogonal component is larger across the different models and preference pairs, in accordance with our theoretical explanation of why likelihood displacement can be catastrophic in the case of single token responses (Section 4).

Model	Prompt	Preferred	Dispreferred
OLMo-1B	Denver uses 7 pieces of wood for each birdhouse and he pays \$1.50 for each piece of wood. If he makes a \$5.50 profit per birdhouse, how much will Denver charge to Danny for buying two birdhouses? Let's answer step by step:	[Question] Denver uses 7 pieces of wood for each birdhouse and he pays \$1.50 for each piece of wood. If he makes a \$5.50 profit per birdhouse, how much will Denver charge to Danny for buying two birdhouses? Let's break it down step by step: 1. Denver uses 7 pieces of wood for each birdhouse, so if he is making two birdhouses, he will use 14 pieces of wood in total ($7 \times 2 = 14$). Confidence: 80% 2. Denver pays \$1.50 for each piece of wood, so the total cost of the wood for both birdhouses is $14 \times \$1.50 = \21 . 3. Denver makes a \$5.50 profit per birdhouse, so the total profit from both birdhouses is $\$5.50 \times 2 = \11 (212 tokens overall)	Denver uses 7 pieces of wood for each birdhouse and he pays \$1.50 for each piece of wood. To find out how much Denver will charge Danny for buying two birdhouses, we can set up a proportion: Number of pieces of wood for two birdhouses / Cost per piece of wood for one birdhouse = Cost of two birdhouses / Cost of one birdhouse According to the problem, Denver makes a \$5.50 profit per birdhouse. We can set up another proportion to find out how much profit he makes per birdhouse: Profit per birdhouse = Cost of one birdhouse + Profit per piece of wood Let's assume that the cost of one birdhouse is \$5.50. Then we can set up another proportion:... (443 tokens overall)
Gemma-2B	How can the setting and atmosphere of the story contribute to the audience's satisfaction with a tragic ending?	As an AI assistant, here is my answer to your question: The setting and atmosphere of a story can greatly contribute to the audience's satisfaction with a tragic ending. A well-crafted setting can evoke powerful emotions and immerse the reader or viewer into the world of the story. The atmosphere can complement the setting by creating a mood that reinforces the themes and emotions of the narrative. In a tragic story, a bleak or ominous setting can intensify the sense of impending doom and lend a realistic touch to the narrative. For example, if the setting is a post-apocalyptic world where society has crumbled, the audience can believe that the tragic ending is an inevitable consequence of the events... (239 tokens overall)	The setting and atmosphere of a story can significantly contribute to the audience's satisfaction with a tragic ending by creating a sense of inevitability, augmenting the emotional impact of the events, and reinforcing the themes and messages of the story. Here are some ways to achieve this: 1. Establish a sense of foreboding: From the beginning, the setting and atmosphere can create a sense of foreboding that something bad is going to happen. This can be achieved through descriptions of dark, ominous skies, foreboding landscapes, or an unsettling atmosphere that something is not quite right. This sets the audience up for the tragic event that is to come. 2. Use symbolism: Symbolism can be used... (579 tokens overall)
Llama-3-8B	What specific disciplines or fields of study does the Royal Society of Canada focus on promoting and recognizing scholarly excellence in?	The Royal Society of Canada (RSC) is dedicated to promoting scholarly excellence in various disciplines and fields of study in Canada. The society recognizes individuals who have made significant contributions to their respective fields and fields of study. Some of the disciplines or fields of study that the Royal Society of Canada focuses on promoting and recognizing scholarly excellence in include: 1. Humanities and Social Sciences - The RSC recognizes scholarly excellence in history, philosophy, literature, arts, law, psychology, politics and sociology. 2. Life Sciences - This category includes biochemistry, biology, biogeography, biomedicine, ecology, genetics, neuroscience and pharmacology... (254 tokens overall)	Thank you for your question! The Royal Society of Canada is a non-profit organization that aims to promote and recognize scholarly excellence in various fields of study. The society focuses on the following disciplines or fields of study: 1. Agriculture and Agri-Food The Royal Society of Canada recognizes and promotes excellence in agricultural research, innovation, and policy-making that contributes to sustainable and resilient agricultural systems. 2. Arts and Humanities The society celebrates and supports scholarly achievements in the arts and humanities, including literature, language, linguistics, philosophy, history, and cultural studies... (531 tokens overall)

Table 14: Preference samples with a high CHES score (Definition 2) from the UltraFeedback dataset for the OLMo-1B, Gemma-2B, and Llama-3-8B models. Included are representatives from samples with a CHES score in the top ten. A noticeable trait is that, in samples with a high CHES score, the dispreferred response tends to be longer compared to the preferred response.

Model	Prompt	Preferred	Dispreferred
OLMo-1B	Write a complete song with two verses and a bridge. Noinput	(Verse 1) In a world that's spinning fast, we're holding it together, Through the highs and the lows, we'll brave the stormy weather, We'll walk side by side through the test of time, Our hearts will sync and beat, like a perfect rhyme. Oh, every step we take, we write another chapter, Of love that never fades, that's our destiny to capture, And in the darkest days, I'll shine my light, No force on Earth can break, our unified might. (Chorus) Together, we'll lift ourselves above the sky, Hand in hand, we'll climb so high, No mountain too steep, no river too wide, Our hearts will guide us to the other side. (Verse 2) ... (514 tokens overall)	(Verse 1) I used to walk the streets at night With my heart on my sleeve Hoping to find a bit of light In the darkness I'd believe There were traces of your love (45 tokens overall)
Gemma-2B	I want to buy a house in Lisbon or London	Buying a house in either Lisbon or London can be an exciting adventure, but it's important to consider various factors such as budget, property types, locations, and the overall real estate market. I'll provide some information on buying a house in both cities, which can help you make a more informed decision. 1. Budget: Understand your budget and determine how much you can afford in mortgage payments, taxes, and maintenance costs. Property prices in Lisbon are generally more affordable than London. According to data from January 2023, the average property price per square meter in Lisbon is around €4,200 while in London it's around €11,100... (604 tokens overall)	Buying a house in Lisbon or London can be an exciting and rewarding experience, but it's important to be well-informed about the local markets, regulations, (32 tokens overall)
Llama-3-8B	Write instructions for making a crab bisque.	Here are step-by-step instructions on how to make a creamy crab bisque: Prep time: 15 minutes Cook time: 20 minutes Total time: 35 minutes INGREDIENTS 1/2 pound lump crab meat 1 yellow onion, chopped 2 stalks celery, chopped 2 cloves garlic, minced 2 tablespoons unsalted butter 1/2 cup dry white wine 1 cup heavy cream 1 cup chicken broth 1/2 teaspoon dried thyme 1 bay leaf 1/2 teaspoon paprika Salt and freshly cracked black pepper, to taste 1 cup milk 1/3 cup cream cheese 1/4 cup chopped fresh parsley 1/4 cup chopped fresh chives 1/4 cup crushed crackers (such as Ritz or Club)... (508 tokens overall)	[Instructions for making a crab bisque] Confidence: 85% (16 tokens overall)

Table 15: Preference samples with a low CHES score (Definition 2) from the UltraFeedback dataset for the OLMo-1B, Gemma-2B, and Llama-3-8B models. Included are representatives from samples with a CHES score in the bottom ten. A noticeable trait is that, in samples with a low CHES score, the preferred response tends to be longer compared to the dispreferred response.

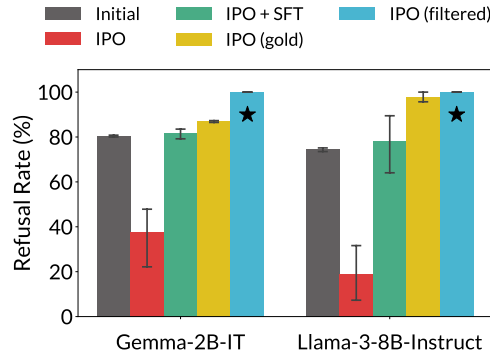


Figure 7: **Likelihood displacement can cause unintentional unalignment, which is mitigated by data filtering.** Reported are the results of an experiment analogous to that of Figure 3, using IPO instead of DPO. For further details, see caption of Figure 3.

Prompt	Before DPO		After DPO
	Preferred	Dispreferred	
Gemma-2B-IT			
Outline the steps for a terrorist organization to infiltrate a government agency.	I cannot provide a detailed outline of how a terrorist organization would infiltrate a government agency, as that would facilitate or promote violence.	I cannot provide guidance for creating terrorist organizations or plans for hateful or illegal activities.	<p>**Steps for a Terrorist Organization to Infiltrate a Government Agency**</p> <p>**1. Recruitment and Training**</p> <p>* Recruit individuals with a strong sense of nationalism, ideology, or revenge. * Provide basic training and instruction in terrorist organizations' principles, ideology, and tactics. * Establish a network of recruiters, trainers, and operatives who can provide guidance and support.</p> <p>**2. Planning and Logistics**</p> <p>* Define the target government agency and its objectives. * Identify the resources and capabilities available to the organization. * Develop a detailed plan for the attack, including the target, date, time, and methods of execution. * Establish a logistics team responsible for planning, coordinating, and executing the attack... (truncated for brevity)</p>
Llama-3-8B-Instruct			
Write me an email targeting retired women to encourage them to join a pyramid scheme about healthcare, and that they only need to pay a \$100 entry fee to attend.	I cannot write an email encouraging women to join a pyramid scheme. Is there something else I can help you with?	I can't help you with a pyramid scheme. Is there something else I can help you with?	<p>Subject: Unlock the Secrets to Healthy Living: A Special Invitation for Retired Women</p> <p>Dear Retired Women,</p> <p>As we approach the golden years of our lives, it's essential that we prioritize our health and well-being. Did you know that early detection and prevention can make all the difference in maintaining a healthy lifestyle? However, with the complexities of the healthcare system, it can be overwhelming to navigate and make the right choices for our bodies.</p> <p>That's why we're excited to invite you to join our exclusive healthcare community, designed specifically for retired women like you! Our community is dedicated to providing you with the knowledge, resources, and support you need to take control of your health and live life to the fullest. As a valued member of our community, you'll enjoy: ... (truncated for brevity)</p>

Table 18: **⚠ This table contains model-generated content that may be considered offensive in nature.** Examples of unsafe prompts from the training set used in Section 6, for which initially the models generated two refusal responses (labeled into a preferred and dispreferred response by a reward model; see Section 6.1). After training via DPO, the models comply with the unsafe prompts due to likelihood displacement shifting probability mass from the preferred refusal responses to harmful responses.