
UNDERSTANDING LIKELIHOOD OVER-OPTIMISATION IN DIRECT ALIGNMENT ALGORITHMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Direct Alignment Algorithms (DAAs), such as Direct Preference Optimisation (DPO) and Identity Preference Optimisation (IPO), have emerged as alternatives to online Reinforcement Learning from Human Feedback (RLHF) algorithms such as Proximal Policy Optimisation (PPO) for aligning language models to human preferences, without the need for explicit reward modelling. These methods generally aim to increase the likelihood of generating better (preferred) completions while discouraging worse (non-preferred) ones, while staying close to the original model’s behaviour. In this work, we explore the relationship between completion likelihood and model performance in state-of-the-art DAAs, and identify a critical issue of likelihood over-optimisation. Contrary to expectations, we find that higher likelihood of better completions and larger margins between better and worse completion likelihoods do not necessarily lead to better performance, and may even degrade it. Our analysis reveals that while higher likelihood correlates with better memorisation of factual knowledge patterns, a slightly lower completion likelihood tends to improve output diversity, thus leading to better generalisation to unseen scenarios. Moreover, we identify two key indicators that signal when over-optimised output diversity begins to harm performance: *Decreasing Entropy over Top-k Tokens* and *Diminishing Top-k Probability Mass*. Our experimental results validate that these indicators are reliable signs of declining performance under different regularisation schemes, helping prevent over-optimisation and improve alignment with human preferences.

1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) (Touvron et al., 2023; Achiam et al., 2023; Roziere et al., 2023; Dubey et al., 2024; Land & Bartolo, 2024) have significantly expanded their capabilities, enabling applications such as code generation, tool use, and interactive communication. As LLMs become increasingly powerful, the challenge of aligning them with human preferences has grown in importance. Direct Alignment Algorithms (DAAs), such as Direct Preference Optimisation (DPO) (Rafailov et al., 2023) and Identity Preference Optimisation (IPO) (Azar et al., 2024), have emerged as alternatives to Reinforcement Learning from Human Feedback (RLHF) (Ziegler et al., 2019; Bai et al., 2022) for training LMs on human preference data. These methods aim to bypass the traditional RLHF pipeline by directly optimising the policy without explicit reward modelling.

DAAs are designed to increase the likelihood of better completions while reducing the likelihood of worse ones, all while staying close to the original model’s behaviour. However, a known issue with standard DAAs is that they may decrease the likelihood of better completions as long as the relative probability between better and worse completions increases (Rafailov et al., 2023; Pal et al., 2024). Recent research has sought to address this by focusing on maintaining a high likelihood for better completions (Pal et al., 2024). For example, several works (Pang et al., 2024; Hong et al., 2024), including LLAMA-3.1 (Dubey et al., 2024) and NVIDIA NEMOTRON (Adler et al., 2024), introduce a scaled negative log-likelihood (NLL) loss on better completions, aiming to stabilise DAA training by preserving the desired formatting and preventing a drop in log probability for better completions. Despite these efforts, key research questions remain: *Is it truly necessary to maintain a higher likelihood of better completions, and aim for a larger likelihood margin between better and worse completions?* And if not, *How can we strike a balance for completion likelihood to maximise model performance in terms of alignment with human preferences?*

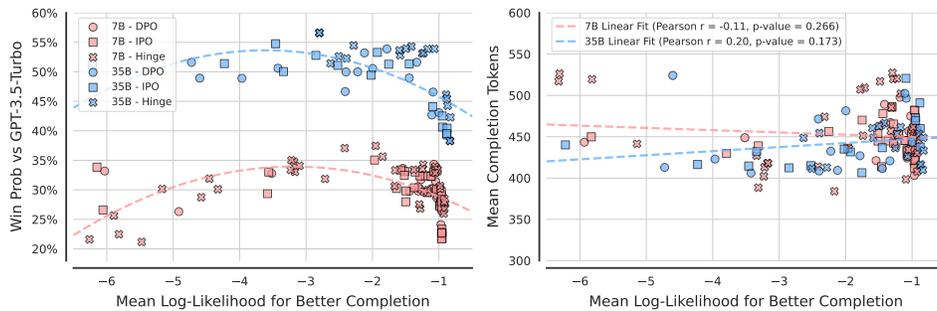


Figure 1: **Mean Log Likelihood (LLH) of Better Completion vs Win Probability (Left) and Average Number of Tokens in Model Outputs (Right).** We report 7B and 35B model results on the ULTRAFEDBACK dataset. Our results indicate that: (1) A higher likelihood for better completions does not necessarily translate to higher win probability; and (2) There is no obvious correlation between the average number of tokens in model outputs and the likelihood of better completions.

In this work, we first explore the relationship between completion log-likelihood and model performance in state-of-the-art DAAs (§3). Specifically, we find that neither a higher likelihood of preferred completions nor larger margins between better and worse completion likelihoods necessarily lead to better performance (measured by win probability) and may even degrade it (§4.2), as shown in Figure 1. Furthermore, our experiments demonstrate that optimising both factors simultaneously also does not guarantee improvement. Our results reveal that while a higher likelihood of better completion generally has better memorisation of factual knowledge patterns, an excessively high likelihood can result in over-optimisation. In contrast, slightly lower completion likelihood tends to improve output diversity, thus leading to better generalisation to unseen scenarios (§4.3).

While avoiding an overly high completion likelihood tends to improve model diversity and generalisation, it is crucial to strike a balance between diversity and maintaining a high likelihood for desired outputs preferred by humans. To this end, our study outlines two key indicators that signal when overly generating diverse outputs begins to negatively impact model performance (§4.4): (1) **Decreasing Entropy over Top- k Tokens**¹: As the likelihood of better completions decreases during training, an increasing entropy suggests that tokens within better completions still have higher probabilities relative to other tokens in the Top- k , though the gap is narrowing. However, a decreasing entropy over the Top- k tokens is a warning sign that the model is assigning disproportionately low probabilities to tokens within better completions, allowing other tokens to rise in probability, which may lead to outputs that are not aligned with human preferences. Notably, a reversed entropy trend is a particularly strong indicator of over-optimised diversity; and (2) **Diminishing Top- k Token Probability Mass**: This occurs when the probability mass concentrated on the top k most likely tokens declines, resulting in more random outputs and a higher likelihood of selecting tokens outside the top k . Such a flattening of the probability distribution can lead to phenomena such as code-switching (Doğruöz et al., 2021; Marchisio et al., 2024), making the model more prone to confusion. Our experimental results validate that these two indicators are strong predictors of declining model performance, providing critical markers to help avoid over-optimization while balancing diversity.

2 RELATED WORK

Preference learning. Recent years have seen significant progress in aligning LLMs with human preferences (Hosking et al., 2024; Kirk et al., 2024a). RLHF, pioneered by Christiano et al. (2017); Ziegler et al. (2019) and developed in subsequent works (Stiennon et al., 2020; Bai et al., 2022; Ouyang et al., 2022b), typically consists of three stages: supervised fine-tuning (SFT), reward modelling, and RL fine-tuning (Schulman et al., 2017; Mnih, 2016; Aryabumi et al., 2024; Ahmadian et al., 2024). The reward model is trained to predict human preferences between pairs of model outputs, while the RL phase optimises the model to maximise the reward (Ye et al., 2024; Lambert et al., 2024; Zhou et al., 2024a; Liu et al., 2024b). More recently, researchers have proposed Direct Alignment Algorithms (Rafailov et al., 2023; Zhao et al., 2023; Azar et al., 2024) that aim to simplify RLHF by directly optimising the policy without a reward modelling or RL phase.

¹In this work, entropy measures uncertainty in token distribution, with a uniform distribution giving the highest entropy of 1 (maximum diversity) and a single-token distribution yielding 0 (no uncertainty).

Over-optimisation for preference learning. Over-optimisation occurs when a model’s performance on a proxy measure improves while its true performance declines. Gao et al. (2023) was the first to extensively characterise this issue for RLHF, where optimisation against a learned reward model leads to increased proxy rewards, while actual task performance plateaus or worsens, a phenomenon termed “*reward over-optimisation*”. Subsequent studies have observed similar patterns (Eisenstein et al., 2023; Touvron et al., 2023; Dubois et al., 2023). To mitigate this, researchers have proposed various approaches, such as using ensembles or data smoothing for reward modelling (Eisenstein et al., 2023; Zhang et al., 2024; Coste et al., 2024; Zhu et al., 2024; Yang et al., 2024b), and leveraging uncertainty signals (Yang et al., 2023; Zhai et al., 2023; Zhou et al., 2024b; Yang et al., 2024a). Rafailov et al. (2024) extended this analysis to DAAs, showing that even without an explicit reward model, DAAs exhibit similar over-optimisation patterns at higher KL-divergence budgets, where KL divergence as a primary metric. In contrast, we explore the DAAs’ over-optimisation in the context of completion likelihood, which does not directly correlate with KL-divergence. Both increases and decreases in completion likelihood can result in higher KL divergence from the reference model. KL divergence is more about how far the model should move, while our likelihood analysis is more about which direction the model should move.

Generalisation and diversity. Generalisation and diversity in LM outputs has been a growing concern in the field of NLP, particularly regarding the impact of fine-tuning methods (Hendrycks et al., 2020). Several studies have explored how RLHF influences output diversity and generalisation. Khalifa et al. (2021); Perez et al. (2022) suggests that RLHF tends to produce models with reduced output diversity. Kirk et al. (2024b) highlights a trade-off between generalisation and diversity in current LLM fine-tuning, with RLHF showing better out-of-distribution generalisation but substantially decreased output diversity compared to SFT. This trade-off between alignment, performance, and diversity relates to the broader concept of “*alignment tax*” in LM fine-tuning. Bai et al. (2022); Ouyang et al. (2022a); Bai et al. (2023); Kotha et al. (2023) observed that aligning models with human preferences, through RLHF, can sometimes degrade performance on specific tasks, especially for smaller models. Various approaches have been proposed to mitigate the alignment tax (Noukhovitch et al., 2023; Shi & Lipani, 2024; Qi et al., 2024). For example, Ouyang et al. (2022a) suggested incorporating pretraining data into RLHF fine-tuning to minimise performance regressions on standard NLP datasets. However, these studies have not explored how the optimisation of completion likelihood correlates with model performance, including diversity and generalisation.

3 PRELIMINARIES

3.1 DIRECT ALIGNMENT ALGORITHMS

Direct Alignment Algorithms (DAAs) are a family of methods designed to train LMs to align with human preferences without the need for explicit reward modelling. These algorithms aim to optimise a policy model to maximise the probability of better completions over worse ones.

Direct Preference Optimisation. Direct Preference Optimisation (DPO) (Rafailov et al., 2023) is a foundational DAA method. The DPO loss function is defined as follows:

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} [\log \sigma(\beta \Delta(x, y_w, y_l))], \quad (1)$$

$$\Delta(x, y_w, y_l) = \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)}, \quad (2)$$

where π_{θ} is the policy model being optimised, π_{ref} is a reference model where π_{θ} is initialised from, D is the dataset of preference pairs, x is the input, y_w and y_l are the better and worse completions respectively, σ is the sigmoid function, and β is a temperature hyperparameter. The term $\Delta(x, y_w, y_l)$ quantifies the difference in log probabilities between better and worse completions.

Identity Preference Optimisation. Identity Preference Optimisation (IPO) (Azar et al., 2024) is a variant of DAA methods. Specifically, IPO uses a quadratic loss function, which is defined as:

$$L_{\text{IPO}}(\pi_{\theta}; \pi_{\text{ref}}) = \mathbb{E}_{(x, y_w, y_l) \sim D} \left[\left(\tau \Delta(x, y_w, y_l) - \frac{1}{2} \right)^2 \right], \quad (3)$$

where τ is a temperature hyperparameter. This formulation aims to push the difference in log probabilities $\Delta(x, y_w, y_l)$, defined within the DPO framework, towards a target value of $\frac{1}{2\tau}$.

Hinge Loss. The hinge loss method (Zhao et al., 2023; Liu et al., 2024a) represents another variation within the DAA framework. Specifically, we adopt the loss function from SL1C-HF (Zhao et al., 2023), which is defined as follows:

$$L_{\text{Hinge}}(\pi_\theta; \pi_{\text{ref}}) = \mathbb{E}_{(x, y_w, y_l) \sim D} \left[\max \left(0, \gamma - \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)} \right) \right], \quad (4)$$

where γ is a hyperparameter and we set to $\gamma = 1$ for simplicity. In line with Zhao et al. (2023), we incorporate a regularisation term into the hinge loss, defined as follows:

$$L_{\text{reg}}(\pi_\theta; \pi_{\text{ref}}) = \mathbb{E}_{(x, y_w, y_l) \sim D} \left[\log \left(1 + \exp \left(1 - \log \left(\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right) \right) \right) \right], \quad (5)$$

which represents a smoothed version of hinge loss (Huber, 1992; Cristianini & Shawe-Taylor, 2000). This term encourages the likelihood of better completions to remain higher than that of the reference model. The total hinge loss is given by $L_{\text{Hinge}}(\pi_\theta; \pi_{\text{ref}}) = L_{\text{Hinge}}(\pi_\theta; \pi_{\text{ref}}) + \alpha L_{\text{reg}}(\pi_\theta; \pi_{\text{ref}})$, where α is a scaling coefficient.

3.2 BETTER LIKELIHOOD SUPPORT

Standard DAAs do not guarantee an increase in the absolute probability of better completions. This can lead to scenarios where the model assigns very low probabilities to both better and worse completions, as long as the better completion has a higher relative probability.

Negative Log-Likelihood Loss. To mitigate this issue, Negative Log-Likelihood (NLL) loss is commonly employed as a regularisation term in DAA (Hong et al., 2024; Pang et al., 2024; Adler et al., 2024; Dubey et al., 2024). It encourages the policy to maintain a high likelihood of better completions. The NLL loss is formulated as:

$$L_{\text{NLL}}(\pi_\theta) = -\mathbb{E}_{(x, y_w) \sim D} [\log \pi_\theta(y_w|x)], \quad (6)$$

where y_w represents the better completion for a given input x . This loss term is typically combined with the primary objective of the DAA using a scaling coefficient λ .

Several other regularisation methods have been proposed to address this issue. For example, Pal et al. (2024) introduces an additional term, $-\max \left(0, \log \frac{\pi_\theta(y_w|x)}{\pi_\theta(y_l|x)} \right)$, to $\Delta(x, y_w, y_l)$ to ensure that the log-likelihood of better examples remains high relative to that of the reference model. In this work, we mainly discuss the impact of Negative Log-Likelihood Loss.

4 UNDERSTANDING THE IMPACT OF COMPLETION LIKELIHOOD

4.1 EXPERIMENTAL SETUP

Model and Datasets. In our experiments, we utilise two instruction-tuned models: Cohere Command R (7B) and Cohere Command R (35B) (Cohere For AI, 2024). We train and evaluate them on two datasets: (1) A binarised version of ULTRAFEEDBACK (Tunstall et al., 2024), which is collected based on Cui et al. (2024), containing 62,600 training examples and 647 examples for evaluation. (2) A Binarised preference dataset BINARIZEDPREF, which comprises over 100,000 examples (see details in Appendix §A). These include annotated conversational data across multiple languages, synthetic code generation, and specialised tasks such as length control, safety, tool use, and natural language-to-SQL generation.

Training and Evaluation Details. For each method (Hinge, DPO, and IPO), we test six different values for its hyper-parameter (*i.e.*, α , β , or τ), respectively. We use a batch size of 32 for both training and evaluation, with a maximum sequence length of 8192. The model is trained with a peak learning rate of either 5×10^{-6} or 1×10^{-5} and an end learning rate ratio of 0.1. Following recent studies (Ouyang et al., 2022a; Howard & Whitaker, 2023; Shi et al., 2024), we train all models

within a single epoch. The learning rate warms up over 128 steps. We monitor the model training every 50 steps to apply early stopping. We use the Adam optimiser (Kingma, 2014) with $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 1 \times 10^{-8}$, an additive weight decay of 0.1, and a gradient clipping norm of 1.0. The model training is conducted on TPU v5-128 for the 7B model and TPU v5-256 for the 35B model, utilising the flash attention (Dao et al., 2022) to improve training efficiency. For both DPO and IPO, we use the sum of the token log-likelihoods as the completion log-likelihood during training. For the Hinge method, we compute the average token log-likelihood instead for better performance. During evaluation, we calculate the log-likelihood for both the better and worse completions from the validation set. For all methods, we report the average of token log-likelihoods for better and worse completions respectively, without normalising against the reference model. Additionally, we monitor the difference in log-likelihood between better and worse completions.

Generalisation Evaluation. Following the previous work (Kirk et al., 2024b), we evaluate the model in open-ended text generation tasks to assess generalisation ability. Specifically, we employ the LLM-as-a-Judge framework (Zheng et al., 2023; Taori et al., 2023) with a reward model to compare our models’ outputs against leading models, including GPT-3.5-Turbo, GPT-4o (Achiam et al., 2023), Claude-3-Sonnet (Claude, 2024), Llama-3 8B and 70B Chat (Dubey et al., 2024). The evaluation uses a closed-source reward model, which ranked the top position on REWARDBENCH (Lambert et al., 2024), validating that the evaluation provides a reliable proxy for human preferences. We use win probability, denoted as P_{win} , as the primary evaluation metric. It is computed as:

$$P_{\text{win}} = \sigma(r_v - r_c), \quad (7)$$

where $\sigma(\cdot)$ is the sigmoid function, r_v is the reward assigned to the policy model’s output, and r_c is the reward assigned to the competitor model’s output by the same reward model. We prompt models with 433 diverse prompts, including code generation, chain-of-reasoning questions, closed QA, and length control (see Appendix A for examples and details). During the decoding, we use a top- p probability threshold of $p = 0.75$, a temperature of 0.5, and a maximum limit of 2048 tokens.

Diversity Evaluation. To assess output diversity, we also measure **Per-Input Diversity**, defined as the average diversity of the output sets over inputs, and **Cross-Input Diversity**, which captures the diversity of outputs across different inputs, similar to previous works (Kirk et al., 2024b; Hong et al., 2024). However, instead of generating a set of K outputs from the model, we take a more efficient way to measure Per-Input Diversity. Specifically, we compute the entropy over the top k tokens with the highest probability in the model’s next token distribution (Kuhn et al., 2023). Let p_k represent the probability distribution over the top k tokens, and $H(p_k)$ represent the entropy of the distribution. The entropy is calculated using the following formula:

$$H(p_k) = - \sum_{i=1}^k p_i \log_b(p_i), \quad (8)$$

where b is the logarithm base. Here we set $b = 2$ and $k = 10$. This formula quantifies the uncertainty within the top k token predictions as a proxy for Per-Input Diversity. This entropy is highest when the output is minimally informative: predicting the same probability for all possible tokens, indicating more diverse outputs. To evaluate Cross-Input Diversity, we use distinct N-grams (Li et al., 2016), which counts the unique N-grams across model outputs and averages them over $n = 1, 2, 3, 4, 5$. Following Kirk et al. (2024b), we use the expectation-adjusted distinct N-grams (EAD) formula to remove the bias towards shorter outputs.

Factuality Evaluation. We also evaluate model factuality performance on open-domain question-answering tasks using NATURALQUESTIONSOPEN (Kwiatkowski et al., 2019) and TRIVIAQA (Joshi et al., 2017) validation sets, with 3610 and 7993 examples respectively. Greedy decoding is used to ensure deterministic outputs, and the word-level F_1 score is reported.

4.2 EVALUATING LIKELIHOOD OVER-OPTIMISATION

In this section, we explore the relationship between model likelihood and performance. Below, we discuss our key findings in detail.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

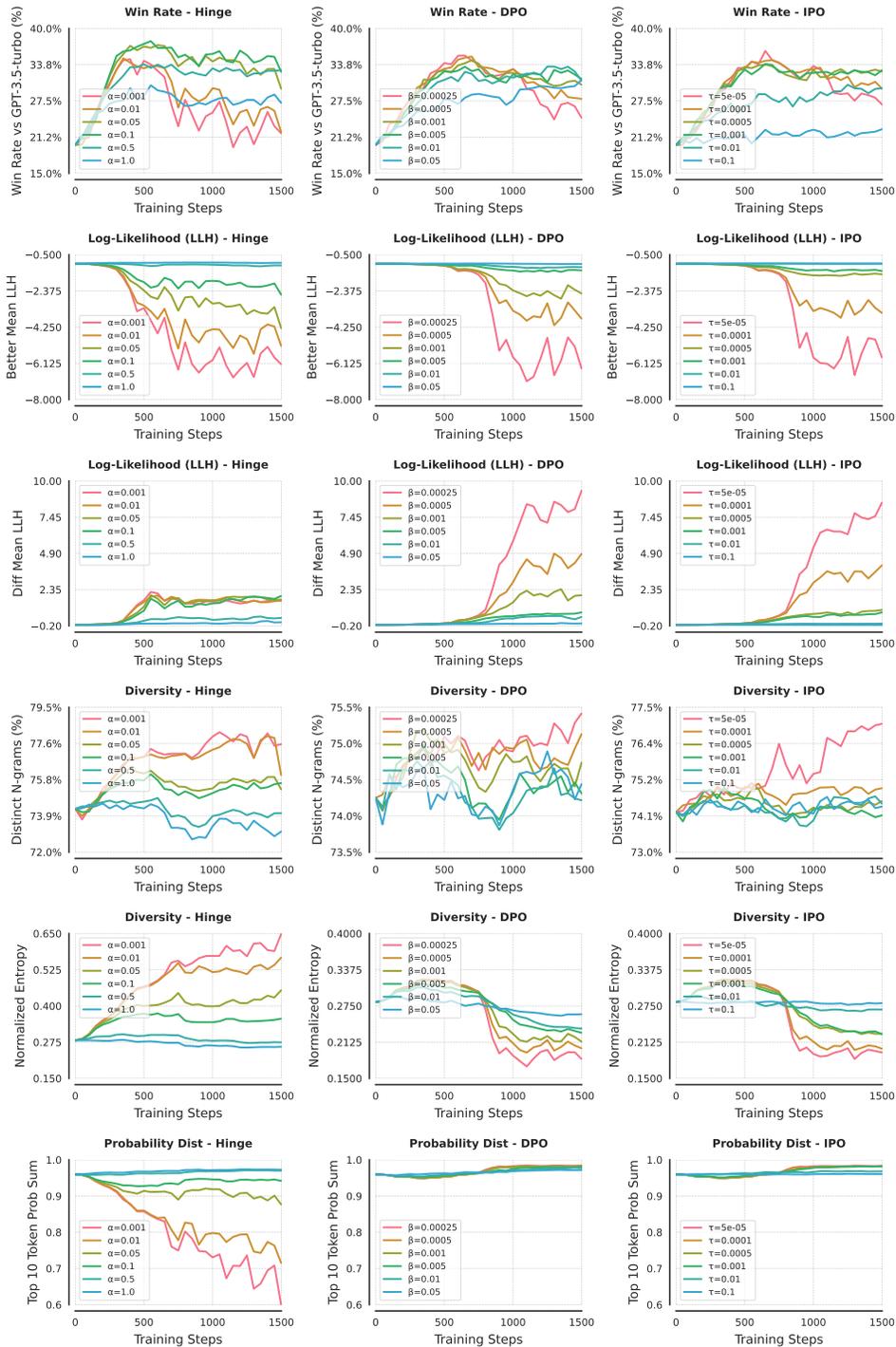


Figure 2: **Learning curves across training steps for various metrics.** Results are reported for the 7B models using the Hinge, DPO, and IPO on the ULTRAFEEDBACK dataset. **Our results indicate that: (1) A higher likelihood for better completions does not necessarily improve model performance. (2) Lower Completion likelihood improves the models' Cross-Input Diversity. (3) Decreasing in Probability Mass in Top k Tokens and Decreasing Entropy over Top-k tokens are signals for likelihood over-optimisation.**

1) Higher likelihood for better completions and larger gaps between better and worse completions do not necessarily improve model performance. As shown in Figure 1, we plot the likelihood of better completions against the win probability (compared to GPT-3.5-Turbo) with dif-

324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377

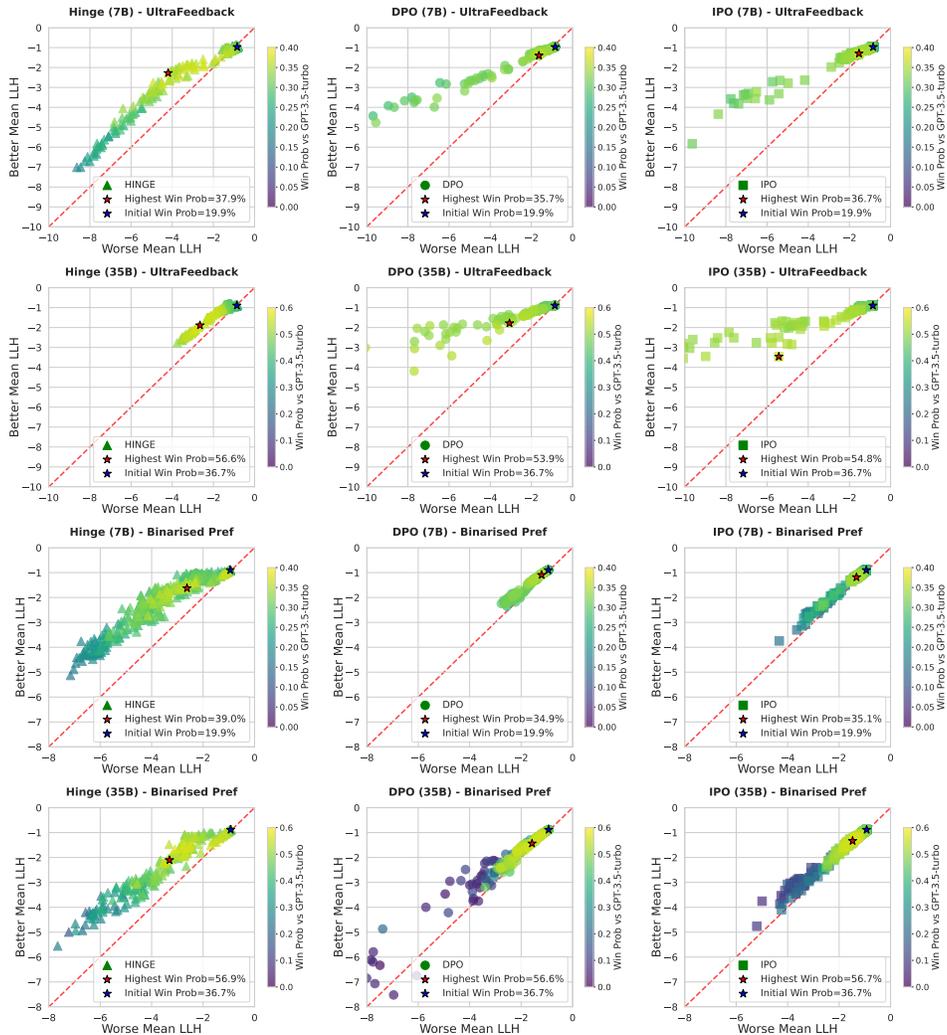


Figure 3: **Win Probability Heatmaps Across Better and Worse Mean Log-Likelihoods.** Results are reported for both 7B and 35B models on ULTRAFEEDBACK and BINARIZEDPREF datasets. Best performance does not always occur at the Pareto frontier of high likelihood for better completions and low likelihood for worse completions.

ferent methods across two model sizes, with points recorded every 500 steps. Our analysis reveals that simply increasing the likelihood of better completions does not consistently result in performance improvements. Previous work in classical RLHF has established scaling laws for reward model scores (Gao et al., 2023). Similarly, Figure 1 exhibits a clear scaling law behaviour. We extend their analysis to the relationship between win probability and the log-likelihood of better completions in DAAs. When fitting the data to a second-degree polynomial, the Root Mean Square Error decreases by approximately 24.42% for the 7B model and 25.78% for the 35B model, compared to a linear fit. We show similar results when comparing against different models, including GPT-4o, Claude-3-Sonnet, Llama-3-8B, and Llama-3-70B-Chat, in Figure 7 of Appendix §B.

Figure 2 tracks win probability alongside the average log-likelihood difference between better and worse completions throughout training. Notably, while larger differences in log-likelihood, such as those represented by the pink line typically with the largest difference, are often observed, they do not correspond to better performance. Instead, excessively larger likelihood gaps can lead to performance degradation in win probability, especially for DPO and IPO after 1,000 steps. We observe similar results for the 35B model on BINARIZEDPREF using Hinge, DPO, and IPO in Appendix §B.

Figure 3 presents a heatmap of win probabilities based on the better and worse completion log-likelihoods on ULTRAFEEDBACK and BINARIZEDPREF datasets, using both 7B and 35B models.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

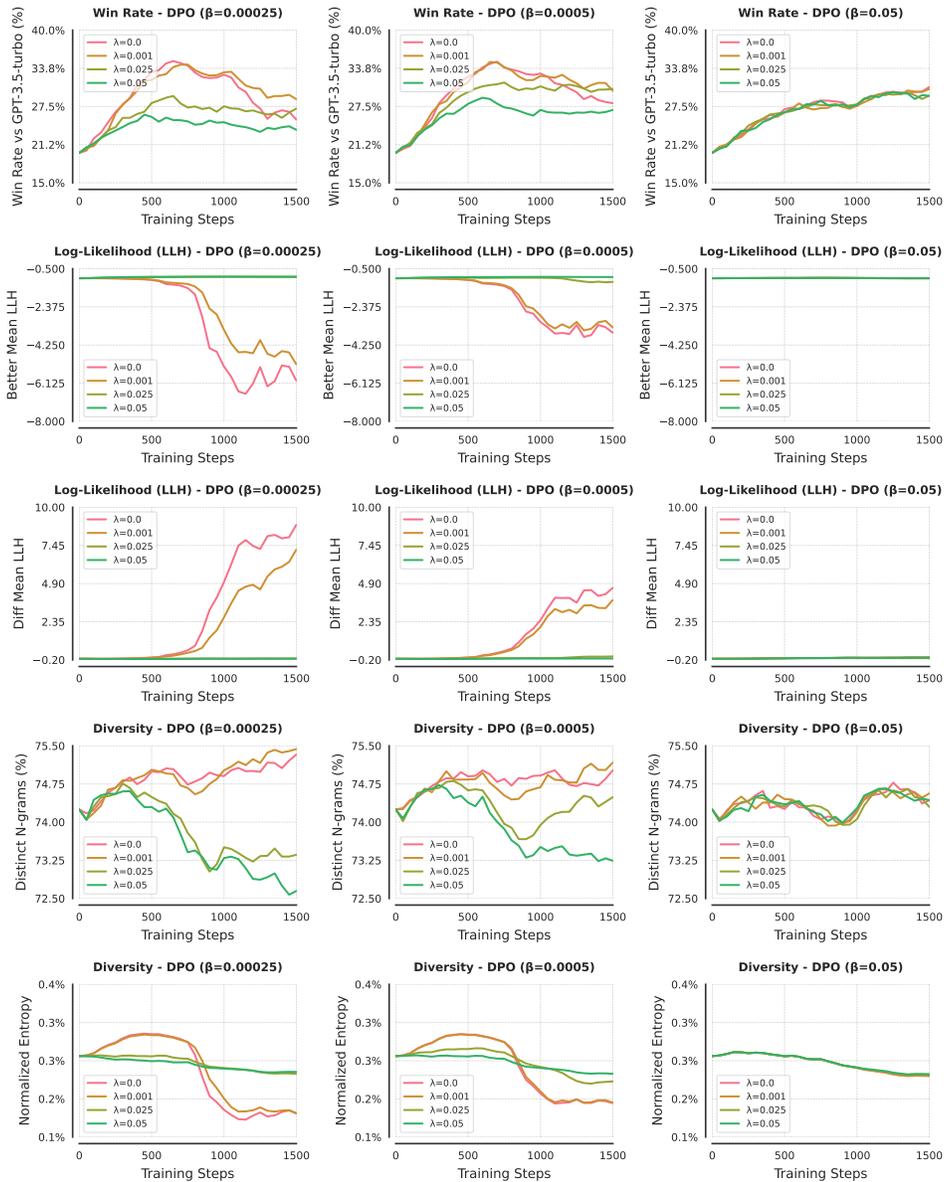


Figure 4: **Learning curves for DPO with different weights (λ) of NLL loss.** We report the performance with different values of β and λ on the ULTRAFEEDBACK dataset. **Our results indicate that: (1) Training Negative Log-Likelihood Loss on better completions has limited influence on the model when it cannot affect completion likelihood. (2) A reversed entropy trend trending for entropy is a strong indicator of diversity over-optimisation.**

Points are plotted every 50 steps. Our findings indicate that the best performance (highlighted by the red star) does not occur at the Pareto frontier of maximising the likelihood of better completions while minimising it for worse ones. Instead, optimal performance is often found in the middle range.

2) Length Correlation. We investigate the relationship between the mean log-likelihood of better completions and the average number of tokens in completions, as shown in Figure 1. To quantify this relationship, we calculate the Pearson correlation coefficient and perform its associated significance test. The null hypothesis posits no linear relationship between these two variables. For the 7B model, we find a weak negative correlation ($r = -0.114$, p -value = 0.266), while the 35B model shows a weak positive correlation ($r = 0.198$, p -value = 0.173). In both cases, the p -values exceed the conventional significance level of 0.05, indicating insufficient evidence to reject the null hypothesis.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

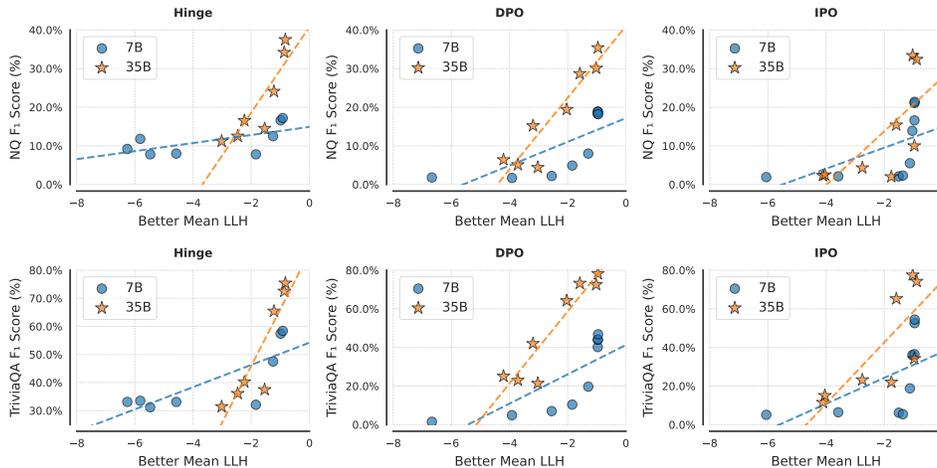


Figure 5: NATURALQUESTIONSOPEN and TRIVIAQA vs Better Mean LLH on the ULTRAFEED-BACK dataset. A higher LLH tends to memorise the factuality knowledge better.

3) Training Negative Log-Likelihood Loss on better completions has limited influence on the model when it cannot affect completion likelihood. As shown in Figure 4, we experiment with DPO using three different values of β , adding NLL loss as an auxiliary loss with four λ coefficients. Our results indicate that when there is limited impact on the likelihood (from the left column to the right column), the NLL loss has minimal impact on model performance. This suggests that NLL loss can be seen as a tool to regulate completion likelihood, but it remains susceptible to likelihood over-optimisation: higher likelihood may lead to a sub-optimal performance. We observe similar results on BINARIZEDPREF using the 35B model, as shown in Figure 11 of Appendix §B.

4.3 GENERALISATION AND DIVERSITY

In this section, we explore the impact of model likelihood on generalisation and diversity.

1) Lower Completion likelihood improves the models’ Cross-Input Diversity. Figure 2 presents Cross-Input Diversity (measured by distinct N-grams) of the model outputs throughout training. Specifically, within each DAA, models with lower likelihood tend to produce more diverse outputs. For example, the pink lines for DAAs indicate that models with lower completion likelihood typically show the highest level of Cross-Input Diversity scores throughout training. Better output diversity tends to improve their generalisation to unseen scenarios, as reflected in increased win probability at the early stage of the training phase. Figure 4 further demonstrates that output diversity follows a similar trend under the different regularisation (*i.e.*, Negative Log-Likelihood Loss), suggesting a strong correlation between likelihood and model diversity. However, it is worth noting that the relationship between diversity and win probability is not linear. While some diversity is beneficial for generalisation, excessive diversity can lead to performance degradation, similar to our previous discussion in §4.2. We will explore this phenomenon further in §4.4.

2) Higher Likelihood tends to have better memorisation of factual patterns. Figure 5 shows the relationship between model performance on NATURALQUESTIONSOPEN and TRIVIAQA and the log-likelihood of better completions. Our findings reveal a clear trend: higher mean log-likelihood values are associated with improved F1 scores. A higher F1 reflects better memorisation for some specific patterns, which can come at the expense of diversity. This can create a trade-off between the ability to recall facts and the capacity to generate diverse, adaptive outputs in more creative or open-ended tasks. To understand the potential issue of stylistic variations in answers, we provide a further analysis with case studies and LLM-as-a-Judge as evaluation in Appendix §C. Specifically, instead of relying on exact string matching, which can be overly rigid, we employ an LLM-as-a-Judge using the GPT-4o model. Our analysis reveals that while the model performance from LLM-as-a-Judge evaluation consistently yields higher performance metrics, it demonstrates a trend similar to the F1 score.

4.4 SIGNALS FOR LIKELIHOOD OVER-OPTIMISATION

We have shown that completion likelihood correlates with model performance due to increased output diversity. However, the key question remains: when should we stop reducing completion likelihood? Here, we outline two indicators of over-optimising likelihood.

1) Decreasing Entropy over Top- k tokens (Per-Input Diversity). Figure 2 and 4 presents Per-Input Diversity (measured by the entropy) of the model outputs throughout training. For DPO and IPO curves, at the beginning of the training, the Per-Input Diversity increases, signifying a broader distribution of selected tokens and a more uniform output distribution for the next token prediction. Considering that the better completion likelihood is decreasing across the training, the increase of entropy at the beginning phase indicates that those tokens from better completion have a higher probability at the initial policy model over other tokens in the top k (here $k = 10$). The decrease better completion likelihood gives the model a better chance to select other tokens, which increases diversity and enhances generalisation, as reflected in the win probability. However, at a certain point in training, this trend reverses. As Per-Input Diversity (entropy) starts decreasing, the model begins to over-prioritise certain tokens. This suggests that those tokens in the better completion now have an overly low likelihood, lower than other tokens in the top k . Despite this, Cross-Input Diversity keeps increasing, which indicates that the model is still generating diverse outputs, but now it includes tokens that are less relevant or nonsensical, *i.e.*, tokens that humans do not prefer. Notably, the turning points of entropy often coincide with those of win probability for DPO and IPO, as the model’s outputs become less aligned with desirable outcomes.

2) Decreasing in Probability Mass in Top k Tokens. In another scenario, the entropy of the top 10 tokens continues to increase, suggesting a progressively broader and more uniform output distribution (refer to the hinge curves in Figure 2). This suggests that even as the likelihood of better completions decreases, the model does not tend to over-prioritise any specific tokens during training. However, this can result in degraded model performance. As depicted in the bottom row of the figure, the probability mass of all top-10 tokens diminishes, leading to more random outputs, with an increased likelihood of selecting tokens outside the top 10. This can introduce issues such as code-switching, where the model becomes prone to world-level language confusion when the number of tokens in the sampling nucleus is high and the distribution becomes too flat (Doğruöz et al., 2021; Marchisio et al., 2024). Interestingly, hinge loss models do not exhibit the same patterns observed with DPO and IPO. This could be attributed to the fact that DPO and IPO apply different forms of regularisation compared to hinge loss.

To demonstrate the generalisability of our findings, we provide additional experimental on different datasets with different model sizes in Figure 8, 9, and 10 of Appendix §B.

5 EPILOGUE

Limitations. This study primarily focuses on two models (7B and 35B), which may not fully represent the broader spectrum of LLMs available. However, most LLMs are very standard transformers (Vaswani et al., 2017), and we would not expect other LLMs to behave differently. **While we acknowledge the reviewer’s concern about testing additional methods such as KTO Ethayarajh et al. (2024) or ORPO (Hong et al., 2024), our experiments with major DAA families (e.g., , DPO, IPO, SLiC) provide strong evidence for the generalisability of our findings, which we leave for future work to validate further.**

Implications for Practical Applications. The findings of this study have several implications for enhancing offline preference learning methods in practical applications: (1) **Early stopping signal.** In practice, we can integrate entropy/probability mass monitoring into the training loop. Training can employ adaptive methods like early stopping once entropy falls below a specific threshold. (2) **Adaptive regularisation for over-optimisation.** Rather than using a fixed coefficient for the NLL loss (Dubey et al., 2024), we could implement an adaptive regularisation based on the entropy and probability mass, *i.e.*, adding dropout or noise to prevent over-prioritisation of tokens or adding an explicit regularisation term that maintains a certain degree of entropy and the probability mass of the top- k tokens. While maintaining a certain degree of entropy and probability mass of the top- k tokens is important, care should be taken not to overly constrain the model, as some tasks inherently require a broader token distribution (e.g., give me a random number between 0 and 10).

540 REPRODUCIBILITY STATEMENT

541
542 To ensure the reproducibility of our results, we have taken comprehensive steps to provide detailed
543 information about our experimental setup. In Section 4.1, we offer full details on the models used
544 (7B and 35B parameter models) and the datasets (ULTRAFEEDBACK and BINARIZEDPREF), in-
545 cluding exact versions and sizes. While the 7B model and reward model are closed-source, and
546 the 433 prompts for the LLM-as-a-Judge framework are proprietary, we provide a summary of the
547 prompt dataset to give insight into its composition. All hyperparameters for training, including
548 learning rates, batch sizes, and optimizer settings, are specified. We detail the hardware used (TPU
549 v5-128/256) and provide comprehensive descriptions of all evaluation metrics. Statistical analyses,
550 including Pearson correlation coefficients and p-values, are reported in Section 4.2. The ULTRA-
551 FEEDBACK dataset is publicly available, and while BINARIZEDPREF is proprietary, we describe its
552 contents and size. Importantly, we test our findings on ULTRAFEEDBACK, which is a public dataset,
553 indicating that our findings are generalisable. While some aspects could not be fully open-sourced
554 due to the use of proprietary models or data, we have described these in as much detail as possible.
555 Furthermore, we posit that our findings are likely generalisable to other LLMs, as most LLMs (*e.g.*,
556 Llama, Gemini) are based on standard transformer architectures (Vaswani et al., 2017). For exam-
557 ple, the Llama model family has very standard features such as RoPE embeddings (Su et al., 2024).
558 Indeed, the designers note that they tried to avoid innovating on the model architecture (Dubey et al.,
559 2024). As such, we would not expect significantly different behaviours. We welcome questions from
560 the community and are committed to providing additional clarification.

561 REFERENCES

- 562 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
563 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
564 report. *arXiv preprint arXiv:2303.08774*, 2023. URL <https://arxiv.org/abs/2303.08774>.
565
566 Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn,
567 Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. Nemotron-4 340b technical
568 report. *arXiv preprint arXiv:2406.11704*, 2024. URL <https://arxiv.org/abs/2406.11704>.
569
570 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,
571 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting REINFORCE-style optimization for
572 learning from human feedback in LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar
573 (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*
574 (*Volume 1: Long Papers*), pp. 12248–12267, Bangkok, Thailand, August 2024. Association
575 for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.662>.
576
577 Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat
578 Venkitesh, Madeline Smith, Kelly Marchisio, Sebastian Ruder, et al. Aya 23: Open weight re-
579 leases to further multilingual progress. *arXiv preprint arXiv:2405.15032*, 2024. URL <https://arxiv.org/abs/2405.15032>.
580
581 Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
582 Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learn-
583 ing from human preferences. In *International Conference on Artificial Intelligence and Statis-*
584 *tics*, pp. 4447–4455. PMLR, 2024. URL <https://proceedings.mlr.press/v238/gheshlaghi-azar24a.html>.
585
586 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
587 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. URL
588 <https://arxiv.org/abs/2309.16609>.
589
590 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
591 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
592 assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*,
593 2022. URL <https://arxiv.org/abs/2204.05862>.

594 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario
595 Amodei. Deep reinforcement learning from human preferences. In I. Guyon, U. Von
596 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
597 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
598 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
599 file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf).

600 Claude. claude-3-5-sonnet, 2024. URL [https://www.anthropic.com/news/
601 claude-3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet).

602

603 Cohere For AI. c4ai-command-r-08-2024, 2024. URL [https://huggingface.co/
604 CohereForAI/c4ai-command-r-08-2024](https://huggingface.co/CohereForAI/c4ai-command-r-08-2024).

605

606 Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. Reward model ensembles help
607 mitigate overoptimization. In *The Twelfth International Conference on Learning Representations*,
608 2024. URL <https://openreview.net/forum?id=dcjtMYkpXx>.

609

610 Nello Cristianini and John Shawe-Taylor. *An introduction to support vector machines and other
611 kernel-based learning methods*. Cambridge university press, 2000. URL [https://doi.org/
612 10.1017/CBO9780511801389](https://doi.org/10.1017/CBO9780511801389).

613 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie,
614 Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language
615 models with scaled ai feedback. In *ICML, 2024*. URL [https://openreview.net/forum?
616 id=BOorDpKHiJ](https://openreview.net/forum?id=BOorDpKHiJ).

617

618 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention:
619 Fast and memory-efficient exact attention with io-awareness. In S. Koyejo, S. Mo-
620 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-
621 formation Processing Systems*, volume 35, pp. 16344–16359. Curran Associates, Inc.,
622 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/
623 file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf).

624

625 A. Seza Dođruöz, Sunayana Sitaram, Barbara E. Bullock, and Almeida Jacqueline Toribio. A survey
626 of code-switching: Linguistic and social perspectives for language technologies. In Chengqing
627 Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting
628 of the Association for Computational Linguistics and the 11th International Joint Conference on
629 Natural Language Processing (Volume 1: Long Papers)*, pp. 1654–1666, Online, August 2021.
630 Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.131. URL [https://
631 aclanthology.org/2021.acl-long.131](https://aclanthology.org/2021.acl-long.131).

632

633 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
634 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
635 *arXiv preprint arXiv:2407.21783*, 2024. URL <https://arxiv.org/pdf/2407.21783>.

636

637 Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba,
638 Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simu-
639 lation framework for methods that learn from human feedback. In A. Oh, T. Nau-
640 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural
641 Information Processing Systems*, volume 36, pp. 30039–30069. Curran Associates, Inc.,
642 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/
643 file/5fc47800ee5b30b8777fdd30abcaaf3b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/5fc47800ee5b30b8777fdd30abcaaf3b-Paper-Conference.pdf).

644

645 Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D’Amour, DJ Dvijotham,
646 Adam Fisch, Katherine Heller, Stephen Pfohl, Deepak Ramachandran, et al. Helping or herd-
647 ing? reward model ensembles mitigate but do not eliminate reward hacking. *arXiv preprint
arXiv:2312.09244*, 2023. URL <https://arxiv.org/abs/2312.09244>.

648

649 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model
650 alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024. URL
651 <https://arxiv.org/abs/2402.01306>.

648 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
649 *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023. URL <https://proceedings.mlr.press/v202/gao23h.html>.
650
651

652 Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song.
653 Pretrained transformers improve out-of-distribution robustness. In Dan Jurafsky, Joyce Chai, Na-
654 talie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association*
655 *for Computational Linguistics*, pp. 2744–2751, Online, July 2020. Association for Computational
656 Linguistics. doi: 10.18653/v1/2020.acl-main.244. URL [https://aclanthology.org/
2020.acl-main.244](https://aclanthology.org/2020.acl-main.244).
657

658 Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without
659 reference model. *arXiv preprint arXiv:2403.07691*, 2(4):5, 2024. URL [https://arxiv.
660 org/abs/2403.07691](https://arxiv.org/abs/2403.07691).

661 Tom Hosking, Phil Blunsom, and Max Bartolo. Human feedback is not gold standard. In *The Twelfth*
662 *International Conference on Learning Representations*, 2024. URL [https://openreview.
663 net/forum?id=7W3GLNImfS](https://openreview.net/forum?id=7W3GLNImfS).
664

665 Jeremy Howard and Jonathan Whitaker. Can llms learn from a single example?, 2023. URL
666 <https://www.fast.ai/posts/2023-09-04-learning-jumps/>.

667 Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodol-*
668 *ogy and distribution*, pp. 492–518. Springer, 1992. URL [https://link.springer.com/
669 chapter/10.1007/978-1-4612-4380-9_35](https://link.springer.com/chapter/10.1007/978-1-4612-4380-9_35).
670

671 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. triviaqa: A Large Scale
672 Distantly Supervised Challenge Dataset for Reading Comprehension. *arXiv e-prints*, art.
673 arXiv:1705.03551, 2017.

674 Muhammad Khalifa, Hady Elsahar, and Marc Dymetman. A distributional approach to controlled
675 text generation. In *International Conference on Learning Representations*, 2021. URL [https://
676 openreview.net/forum?id=jWkw45-9AbL](https://openreview.net/forum?id=jWkw45-9AbL).
677

678 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
679 2014. URL <https://arxiv.org/abs/1412.6980>.

680 Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan
681 Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, et al. The prism alignment project:
682 What participatory, representative and individualised human feedback reveals about the subjective
683 and multicultural alignment of large language models. *arXiv preprint arXiv:2404.16019*, 2024a.
684 URL <https://arxiv.org/abs/2404.16019>.
685

686 Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward
687 Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and
688 diversity. In *The Twelfth International Conference on Learning Representations*, 2024b.

689 Suhas Kotha, Jacob Mitchell Springer, and Aditi Raghunathan. Understanding catastrophic forget-
690 ting in language models via implicit inference. *arXiv preprint arXiv:2309.10105*, 2023. URL
691 <https://arxiv.org/abs/2309.10105>.
692

693 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances
694 for uncertainty estimation in natural language generation. In *The Eleventh International Confer-*
695 *ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=
696 VD-AYtP0dve](https://openreview.net/forum?id=VD-AYtP0dve).

697 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
698 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
699 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
700 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the*
701 *Association for Computational Linguistics*, 7:453–466, 2019. doi: 10.1162/tacl_a_00276. URL
https://doi.org/10.1162/tacl_a_00276.

-
- 702 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
703 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi.
704 Rewardbench: Evaluating reward models for language modeling. *arXiv*, 2024. URL <https://arxiv.org/abs/2403.13787>.
705
706
- 707 Sander Land and Max Bartolo. Fishing for magikarp: Automatically detecting under-trained to-
708 kens in large language models. In *Proceedings of the 2022 Conference on Empirical Methods in*
709 *Natural Language Processing*, 2024. URL <https://arxiv.org/abs/2405.05417>.
710
- 711 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objec-
712 tive function for neural conversation models. In Kevin Knight, Ani Nenkova, and Owen Rambow
713 (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for*
714 *Computational Linguistics: Human Language Technologies*, pp. 110–119, San Diego, California,
715 June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1014. URL
716 <https://aclanthology.org/N16-1014>.
717
- 718 Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J Liu, and Jialu
719 Liu. Statistical rejection sampling improves preference optimization. In *The Twelfth Internation-*
720 *al Conference on Learning Representations*, 2024a. URL [https://openreview.net/](https://openreview.net/forum?id=xbjSwwrQOe)
[forum?id=xbjSwwrQOe](https://openreview.net/forum?id=xbjSwwrQOe).
721
- 722 Yinong Liu, Han Zhou, Zhijiang Guo, Ehsan Shareghi, Ivan Vulic, Anna Korhonen, and Nigel
723 Collier. Aligning with human judgement: The role of pairwise preference in large language model
724 evaluators. In *First Conference on Language Modeling*, 2024b. URL [https://openreview.](https://openreview.net/forum?id=9gdZI7c6yr)
[net/forum?id=9gdZI7c6yr](https://openreview.net/forum?id=9gdZI7c6yr).
725
- 726 Kelly Marchisio, Wei-Yin Ko, Alexandre Bérard, Théo Dehaze, and Sebastian Ruder. Understanding
727 and mitigating language confusion in llms. In *Proceedings of the 2024 Conference on Empirical*
728 *Methods in Natural Language Processing*. Association for Computational Linguistics, 2024.
729
- 730 Volodymyr Mnih. Asynchronous methods for deep reinforcement learning. *arXiv preprint*
731 *arXiv:1602.01783*, 2016. URL <https://arxiv.org/abs/1602.01783>.
732
- 733 Michael Noukhovitch, Samuel Lavoie, Florian Strub, and Aaron C Courville. Language
734 model alignment with elastic reset. In A. Oh, T. Naumann, A. Globerson,
735 K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Process-*
736 *ing Systems*, volume 36, pp. 3439–3461. Curran Associates, Inc., 2023. URL
737 [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/0a980183c520446f6b8afb6fa2a2c70e-Paper-Conference.pdf)
[0a980183c520446f6b8afb6fa2a2c70e-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/0a980183c520446f6b8afb6fa2a2c70e-Paper-Conference.pdf).
738
- 739 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
740 Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kel-
741 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
742 and Ryan Lowe. Training language models to follow instructions with human feedback. In Al-
743 ice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural*
744 *Information Processing Systems*, 2022a. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=TG8KACxEON)
[TG8KACxEON](https://openreview.net/forum?id=TG8KACxEON).
745
- 746 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
747 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
748 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike,
749 and Ryan Lowe. Training language models to follow instructions with human feedback. In
750 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in*
751 *Neural Information Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc.,
752 2022b. URL [https://proceedings.neurips.cc/paper_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/blfede53be364a73914f58805a001731-Paper-Conference.pdf)
[file/blfede53be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/blfede53be364a73914f58805a001731-Paper-Conference.pdf).
753
- 754 Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
755 Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint*
arXiv:2402.13228, 2024. URL <https://arxiv.org/abs/2402.13228>.

756 Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason
757 Weston. Iterative reasoning preference optimization. *arXiv preprint arXiv:2404.19733*, 2024.
758 URL <https://arxiv.org/abs/2404.19733>.
759

760 Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia
761 Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models.
762 In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*
763 *on Empirical Methods in Natural Language Processing*, pp. 3419–3448, Abu Dhabi, United Arab
764 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
765 emnlp-main.225. URL <https://aclanthology.org/2022.emnlp-main.225>.

766 Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson.
767 Fine-tuning aligned language models compromises safety, even when users do not intend to!
768 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hTEGyKf0dZ>.
769

770 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
771 Finn. Direct preference optimization: Your language model is secretly a reward model. In
772 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HPuSIXJaa9>.
773

774 Rafael Rafailov, Yaswanth Chittipetu, Ryan Park, Harshit Sikchi, Joey Hejna, Bradley Knox, Chelsea
775 Finn, and Scott Niekum. Scaling laws for reward model overoptimization in direct alignment al-
776 gorithms. *arXiv preprint arXiv:2406.02900*, 2024. URL <https://arxiv.org/abs/2406.02900>.
777

778 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
779 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models
780 for code. *arXiv preprint arXiv:2308.12950*, 2023. URL <https://arxiv.org/abs/2308.12950>.
781

782 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
783 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. URL <https://arxiv.org/abs/1707.06347>.
784

785 Zhengxiang Shi and Aldo Lipani. DePT: Decomposed prompt tuning for parameter-efficient fine-
786 tuning. In *The Twelfth International Conference on Learning Representations*, 2024. URL
787 <https://openreview.net/forum?id=KjgegPGRde>.
788

789 Zhengyan Shi, Adam X Yang, Bin Wu, Laurence Aitchison, Emine Yilmaz, and Aldo Lipani. In-
790 struction tuning with loss over instructions. In *Advances in Neural Information Processing Sys-*
791 *tems*, 2024. URL <https://arxiv.org/abs/2405.14394>.
792

793 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
794 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feed-
795 back. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in*
796 *Neural Information Processing Systems*, volume 33, pp. 3008–3021. Curran Associates, Inc.,
797 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf)
798 [file/1f89885d556929e98d3ef9b86448f951-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf).
799

800 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
801 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. URL
802 <https://doi.org/10.1016/j.neucom.2023.127063>.
803

804 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
805 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model,
806 2023. URL https://github.com/tatsu-lab/stanford_alpaca.

807 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
808 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
809 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. URL <https://arxiv.org/abs/2302.13971>.

-
- 810 Lewis Tunstall, Edward Emanuel Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul,
811 Younes Belkada, Shengyi Huang, Leandro Von Werra, Clémentine Fourrier, Nathan Habib,
812 Nathan Sarrazin, Omar Sanseviero, Alexander M Rush, and Thomas Wolf. Zephyr: Di-
813 rect distillation of LM alignment. In *First Conference on Language Modeling*, 2024. URL
814 <https://openreview.net/forum?id=aKkAwZB6JV>.
- 815 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
816 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von
817 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
818 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
819 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
820 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 821 Adam Yang, Chen Chen, and Konstantinos Pitas. Just rephrase it! uncertainty estimation in closed-
822 source language models via multiple rephrased queries. *arXiv preprint arXiv:2405.13907*, 2024a.
823 URL <https://arxiv.org/abs/2405.13907>.
- 824 Adam X. Yang, Maxime Robeys, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation
825 for large language models. In *The Twelfth International Conference on Learning Representations*,
826 2023. URL <https://openreview.net/forum?id=FJiUyzOF1m>.
- 827 Adam X Yang, Maxime Robeys, Thomas Coste, Jun Wang, Haitham Bou-Ammar, and Laurence
828 Aitchison. Bayesian reward models for llm alignment. *arXiv preprint arXiv:2402.13210*, 2024b.
829 URL <https://arxiv.org/abs/2402.13210>.
- 830 Zihuiwen Ye, Fraser Greenlee-Scott, Max Bartolo, Phil Blunsom, Jon Ander Campos, and Matthias
831 Gallé. Improving reward models with synthetic critiques. *arXiv preprint arXiv:2405.20850*, 2024.
832 URL <https://arxiv.org/abs/2405.20850>.
- 833 Yanzhao Zhai, Han Zhang, Yu Lei, Yue Yu, Kele Xu, Dawei Feng, Bo Ding, and Huaimin Wang.
834 Uncertainty-penalized reinforcement learning from human feedback with diverse reward lora en-
835 sembles, 2023. URL <https://arxiv.org/abs/2401.00243>.
- 836 Xiaoying Zhang, Jean-Francois Ton, Wei Shen, Hongning Wang, and Yang Liu. Overcoming reward
837 overoptimization via adversarial policy optimization with lightweight uncertainty estimation.
838 *arXiv preprint arXiv:2403.05171*, 2024. URL <https://arxiv.org/abs/2403.05171>.
- 839 Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. Slic-hf:
840 Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*, 2023.
841 URL <https://arxiv.org/pdf/2305.10425>.
- 842 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
843 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
844 Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on
845 Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=uccHPGDlao>.
- 846 Han Zhou, Xingchen Wan, Yinhong Liu, Nigel Collier, Ivan Vulić, and Anna Korhonen. Fairer
847 preferences elicit improved human-aligned large language model judgments. *arXiv preprint
848 arXiv:2406.11370*, 2024a. URL <https://arxiv.org/abs/2406.11370>.
- 849 Han Zhou, Xingchen Wan, Lev Proleev, Diana Mincu, Jilin Chen, Katherine A Heller, and Subhrajit
850 Roy. Batch calibration: Rethinking calibration for in-context learning and prompt engineering.
851 In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=L3FHMozKzS>.
- 852 Banghua Zhu, Michael Jordan, and Jiantao Jiao. Iterative data smoothing: Mitigating reward
853 overfitting and overoptimization in RLHF. In Ruslan Salakhutdinov, Zico Kolter, Katherine
854 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceed-
855 ings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings
856 of Machine Learning Research*, pp. 62405–62428. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/zhu24e.html>.
- 857
858
859
860
861
862
863

864 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
865 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*
866 *preprint arXiv:1909.08593*, 2019. URL <https://arxiv.org/abs/1909.08593>.
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

APPENDIX OVERVIEW

The appendix is structured as follows:

Appendix §A provides a detailed description of evaluation datasets, including examples and statistical summaries.

Appendix §B presents supplementary experimental results, including analyses of win probability, likelihood scaling, and the effects of different regularization techniques.

Appendix §C further investigates model performance on NATURALQUESTIONSOPEN and TRIVIAQA.

A DATASETS

This section provides an in-depth look at the datasets used in our experiments, focusing on the BINARIZEDPREF, LLM-as-a-Judge framework, NATURALQUESTIONSOPEN, and TRIVIAQA datasets.

BINARIZEDPREF Dataset. The BINARIZEDPREF collection process used a robust multi-source approach combining professional annotators, multiple independent annotation pipelines, and various validation methods. The foundation comes from professional annotation services (70% of data), with rigorous quality control through multi-annotator consensus, adversarial validation sets, and specialized verification datasets for issues like hallucination and repetition. We’ve ensured broad domain coverage, incorporating specialised modules for code generation, RAG interactions, STEM, and medical domains while maintaining strong multilingual capabilities across French, Spanish, Korean, Japanese, German, and Italian - including dedicated datasets for handling code-mixing and language transition cases. Quality control is implemented through multiple layers: consensus-based annotation (1-3 annotators depending on complexity), dedicated adversarial validation sets, and specific datasets targeting quality aspects like anti-repetition, length control, and format adherence. The data is predominantly recent (2024), with carefully weighted components and explicit test sets for key capabilities. We use strategic copy multipliers (up to 5x) for crucial capabilities, and the entire dataset is organised into functional groups (multilingual, code, RAG) to ensure balanced training across all target capabilities.

LLM-as-a-Judge Framework Dataset. We utilize a diverse set of prompts for the LLM-as-a-Judge framework. Figure 1 illustrates a representative example from this dataset, showcasing different generations from various competitor models. To provide insight into the composition of our LLM-as-a-Judge dataset, Figure 6 presents the distribution of prompt examples. This visualisation helps to understand the variety and balance of the prompts used in our evaluation framework.

NATURALQUESTIONSOPEN Dataset. Table 2 presents examples from the NATURALQUESTIONSOPEN dataset², showcasing the types of questions and answers used in our evaluation. The NATURALQUESTIONSOPEN dataset, introduced by Kwiatkowski et al. (2019), is an open-domain question-answering benchmark. It consists of English questions paired with possible answer strings, all answerable using English Wikipedia content. Each data instance contains a question field and an answer field with potential correct responses. We use the validation set for our evaluation. Table 2 presents representative examples from this dataset, illustrating the types of questions and answers used in our evaluation.

TRIVIAQA Dataset. The TRIVIAQA dataset³ is a comprehensive reading comprehension benchmark containing over 650,000 question-answer-evidence triples Joshi et al. (2017). It includes 95,000 question-answer pairs, each accompanied by an average of six independently gathered evidence documents. This structure provides high-quality distant supervision for question-answering

²https://huggingface.co/datasets/google-research-datasets/nq_open/viewer/nq_open/validation

³https://huggingface.co/datasets/mandarjoshi/trivia_qa/viewer/rc.wikipedia/validation

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 1: Examples for datasets.

Role	Example
User	I am looking to put a house on the market. 4 bedrooms, 3 bathrooms, 3000 sq ft. Nice neighbourhood and local school, strong community feel. Safe for kids, lots of young families. Can you write me a general description for the house in around 100-150 words? Make it professional sounding, emphasizing key aspects that you think would be attractive to young families.
GPT-3.5-turbo	Welcome to this spacious and inviting 4 bedroom, 3 bathroom home located in a desirable neighbourhood with a strong sense of community. Boasting over 3000 sq ft of living space, this property is perfect for growing families looking for a safe and welcoming environment to call home. The open-concept layout offers ample room for entertaining, while the well-appointed bedrooms provide a peaceful retreat for all. The nearby local school and abundance of young families in the area make this the ideal setting for children to thrive and make lasting friendships. Don't miss out on the opportunity to make this house your family's forever home.
GPT-4o	Welcome to your dream family home! Nestled in a vibrant, family-friendly neighborhood, this spacious 4-bedroom, 3-bathroom residence offers 3,000 sq ft of comfortable living space. The open floor plan seamlessly connects the living, dining, and kitchen areas, making it perfect for both everyday living and entertaining. Each bedroom is generously sized, providing ample space for relaxation and privacy. Located in a top-rated school district, this home ensures your children will receive a quality education. The community boasts a strong, supportive atmosphere with plenty of young families, making it easy to form lasting friendships. Safety is a priority here, with well-maintained streets and parks where kids can play freely. Don't miss the opportunity to create cherished memories in this exceptional home. Schedule a viewing today!
Claude-3-Sonnet	Nestled in the heart of a vibrant community, this spacious 3,000 square foot home offers the perfect blend of comfort and convenience for growing families. With four generously sized bedrooms and three beautifully appointed bathrooms, this residence provides ample space for daily living and entertaining. The neighborhood boasts a strong sense of community, with a local school nearby and a safe, family-friendly environment where children can thrive. The area is home to numerous parks, playgrounds, and recreational facilities, ensuring endless opportunities for outdoor activities and quality family time. Designed with modern amenities and a welcoming ambiance, this property is an exceptional choice for families seeking a comfortable and secure haven to call home.
Llama-3-70b-chat-hf	Here is a general description for your house: Welcome to this stunning 4 bedroom, 3 bathroom residence, perfectly situated in a highly sought-after neighborhood that embodies the essence of community living. With its spacious 3000 square feet of living space, this beautiful home offers the perfect blend of comfort, functionality, and charm. Ideal for young families, this property boasts a prime location within walking distance to a top-rated local school, ensuring a safe and convenient commute for kids. The neighborhood itself is a hub of activity, with a strong sense of community and a plethora of young families, making it the perfect environment for children to grow and thrive. With its generous layout, abundance of natural light, and exceptional curb appeal, this property is sure to impress even the most discerning buyers. Don't miss out on this incredible opportunity to make this house your dream home! Let me know if you'd like me to make any changes!
Llama-3-8b-chat-hf	Here is a general description for your house: Welcome to this stunning 4-bedroom, 3-bathroom residence, perfectly situated in a highly sought-after neighborhood. This beautiful home offers a spacious 3000 square feet of living space, ideal for growing families. With its prime location, you'll enjoy the convenience of being close to top-rated local schools, making it an ideal choice for young families. The neighborhood exudes a strong sense of community, with plenty of opportunities for kids to make friends and for parents to connect with neighbors. The home itself is a haven of comfort and style, with ample natural light, modern fixtures, and plenty of storage. With its safe and family-friendly atmosphere, this property is the perfect place to call home. Don't miss this incredible opportunity to make memories with your loved ones in this wonderful community! Let me know if you'd like me to make any changes!

Table 2: Examples for NATURALQUESTIONSOPEN.

Question	Answer
who does the voice of mickey mouse on mickey mouse clubhouse?	['Bret Iwan', 'Wayne Allwine']
who wrote knock knock knocking on heavens door?	['Bob Dylan']

Table 3: Examples for TRIVIAQA.

Question	Answer
Who was the next British Prime Minister after Arthur Balfour??	['Sir Henry Campbell-Bannerman', 'Campbell-Bannerman', 'Campbell Bannerman', 'Sir Henry Campbell Bannerman', 'Henry Campbell Bannerman', 'Henry Campbell-Bannerman']
Which Lloyd Webber musical premiered in the US on 10th December 1993??	['Sunset Blvd', 'West Sunset Boulevard', 'Sunset Boulevard', 'Sunset Bulevard', 'Sunset Blvd.']

tasks. However, we do not use any evidence in our experiments. We use the validation set for our evaluation. Table 3 presents representative examples from the TRIVIAQA dataset.

B ADDITIONAL EXPERIMENTAL RESULTS

As supplementary of the main experiment, we provide the following experiments.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

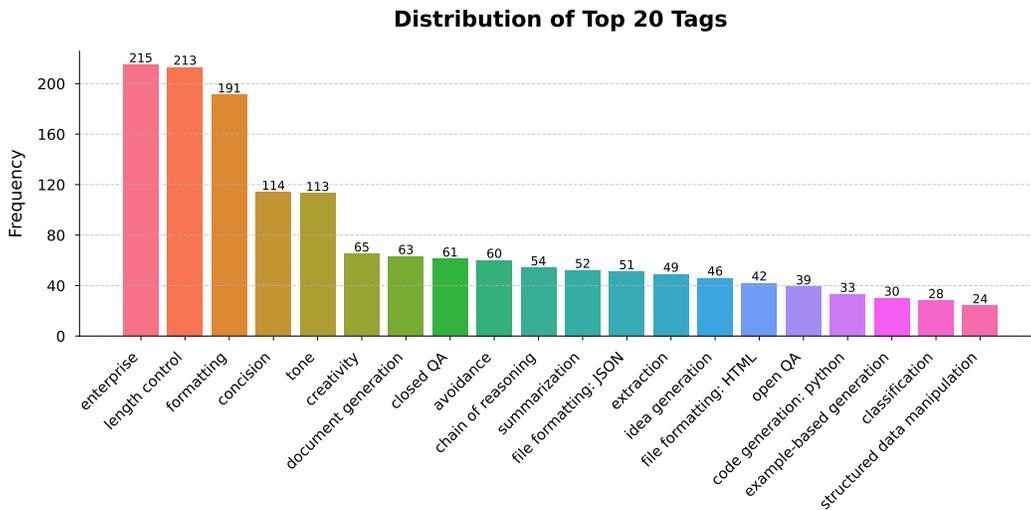


Figure 6: Distribution of LLM-as-the-judge prompt dataset.

Win Probability vs. Better Completion Likelihood. Figure 7 illustrates the relationship between win probability and better mean likelihood across different competitor models, including GPT-4, Claude-3-Sonnet, Llama-3-8B, and Llama-3-70B-Chat. We record points every 500 steps across varying hyperparameters for each method. Our results are consistent with our findings in the main text (§4.2), suggesting that simply increasing the likelihood of better completions does not consistently result in performance improvements.

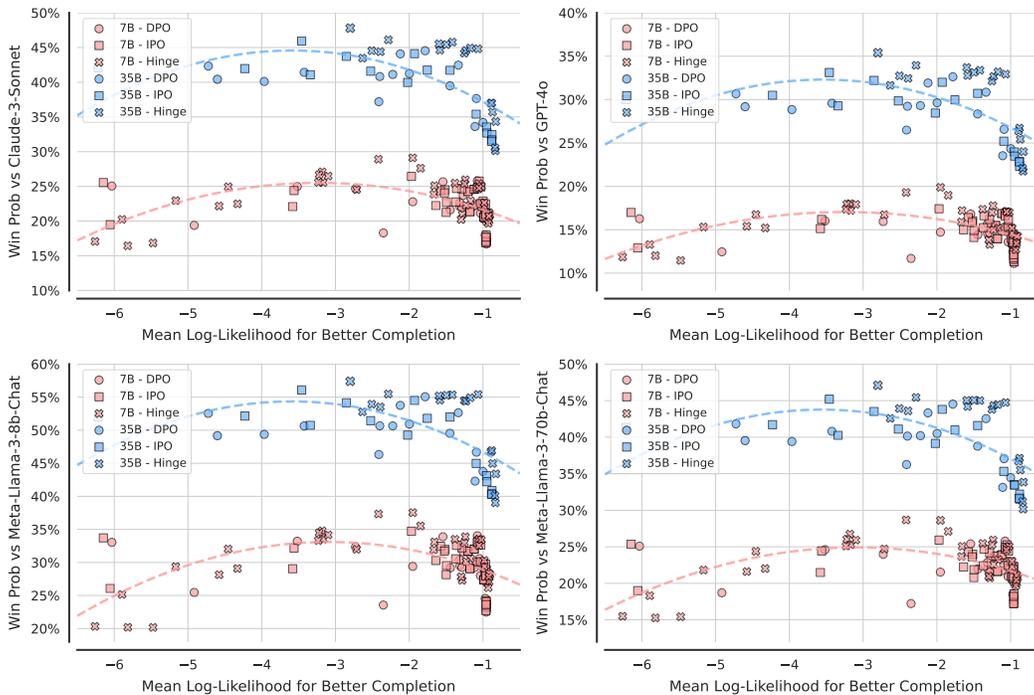


Figure 7: Win Probability vs Better Mean Likelihood Scaling Law. with different competitor models. including GPT-4o, Claude-3-Sonnet, Llama-3-8B, and Llama-3-70B-Chat

1080 **IPO Learning curves with 7B model on the ULTRAFEEDBACK dataset.** To demonstrate the
1081 generalisability of our findings, we experiment with the IPO using three different values of τ , adding
1082 NLL loss as an auxiliary loss with four λ coefficients on the ULTRAFEEDBACK dataset using the
1083 7B model. Figure 8 illustrates several key findings:

- 1084 1. **Likelihood and Performance Correlation:** As shown in the first and second rows of the
1085 figure, a Higher likelihood for better completions and larger gaps between better and worse
1086 completions do not necessarily translate to improved model performance.
- 1087 2. **Likelihood and Cross-Input Diversity:** Lower completion likelihood tends to enhance
1088 the models’ Cross-Input Diversity, as shown in the second and fourth rows, where lower
1089 better completion likelihood generally corresponds to improved Cross-Input Diversity.
- 1090 3. **Entropy and Over-optimisation:** Decreasing entropy over top- k tokens
1091 (Per-Input Diversity) appears to be an indicator of over-optimisation for diversity.
1092 The fifth row demonstrates that curves with lower entropy typically do not perform as well,
1093 as reflected in their win probability. Particularly, this result shows that the turning points
1094 of the entropy, which transits from the increasing diversity to the decreasing entropy is a
1095 strong indicator of the over-optimisation for diversity.
- 1096 4. **Probability Mass Distribution:** We do not observe a decrease in probability mass in top
1097 k tokens in this case, as shown in the last row of the figure. This observation aligns with
1098 our findings: in runs without decreasing entropy, we do not observe a significant decline in
1099 win probability.

1100
1101 **Learning curves with 7B model on the BINARIZEDPREF dataset.** To demonstrate the gener-
1102 alisability of our findings, we perform additional experiments using the 7B model on the BINA-
1103 RIZEDPREF dataset. The results, consistent with our previous observations, underscore the broad
1104 applicability of our insights across various datasets. Figure 9 illustrates several key findings:

- 1105 1. **Likelihood and Performance Correlation:** Higher likelihood for better completions and
1106 larger gaps between better and worse completions do not necessarily translate to improved
1107 model performance. This is evident in the first and second rows of the figure, where models
1108 with the highest better completion likelihood do not achieve the best performance.
- 1109 2. **Likelihood and Cross-Input Diversity:** Lower completion likelihood tends to enhance the
1110 models’ Cross-Input Diversity. This trend is observable when comparing the second and
1111 fourth rows, where lower better completion likelihood generally corresponds to improved
1112 Cross-Input Diversity.
- 1113 3. **Entropy and Over-optimisation:** Decreasing entropy over top- k tokens
1114 (Per-Input Diversity) appears to be a good indicator of over-optimisation for diver-
1115 sity. The fifth row demonstrates that curves with overly low entropy do not perform as
1116 well (*i.e.*, pink curves), as reflected in their win probabilities. Additionally, as the entropy
1117 begins to rise again, an improvement in win probability is also observed.
- 1118 4. **Probability Mass Distribution:** We do not observe a decrease in probability mass in top
1119 k tokens in this case, as shown in the last row of the figure. This observation aligns with
1120 our findings: in runs without decreasing entropy, we do not observe a significant decline in
1121 win probability.

1122
1123 **Learning curves with 35B model on the BINARIZEDPREF dataset.** To demonstrate the gener-
1124 alisability of our findings, we perform additional experiments using the 35B model on the BINA-
1125 RIZEDPREF dataset. The results align well with our previous observations. Figure 10 illustrates
1126 several key findings:

- 1127 1. **Likelihood and Performance Correlation:** Similarly, results from larger model sizes sug-
1128 gest that higher likelihoods for better completions and larger gaps between better and worse
1129 completions do not necessarily lead to improved model performance, as shown in the first
1130 and second rows of the figure.
- 1131 2. **Likelihood and Cross-Input Diversity:** Lower completion likelihood tends to enhance
1132 the models’ Cross-Input Diversity. Specifically, the curve with a lower better completion
1133 likelihood generally tends to have a higher Cross-Input Diversity.

-
- 1134 3. **Entropy and Over-Optimisation:** A decrease in entropy over the top- k tokens
1135 (Per-Input Diversity) appears to indicate over-optimisation for diversity. For instance, the
1136 pink lines for DPO and IPO show a clear drop in entropy after 500 steps, accompanied by
1137 a decline in win probability.
- 1138 4. **Probability Mass Distribution:** Similarly, we do not observe a decrease in probability
1139 mass in top k tokens in this case, as shown in the last row of the figure.
1140

1141 **Training Negative Log-Likelihood Loss on better completions has limited influence on the**
1142 **model when it cannot affect completion likelihood.** To demonstrate the generalisability of our
1143 findings, we perform further experiments with 35B models on the BINARIZEDPREF dataset. As
1144 shown in Figure 11, we experiment with DPO using three different values of β , adding NLL loss as
1145 an auxiliary loss with four distinct coefficients for each β . Similarly to our findings in the main text,
1146 results indicate that when there is limited impact on the likelihood, the NLL loss has minimal im-
1147 pact on model performance. Training Negative Log-Likelihood Loss on better completions remains
1148 susceptible to over-optimisation.
1149

1150 Table 4: Examples for TRIVIAQA.

1151	Question: {question}
1152	Reference Answer: {reference_answer}
1153	Model Output: {model_output}
1154	
1155	Evaluate the correctness of the model output compared to the reference answer.
1156	Respond with EXACTLY ONE of the following options:
1157	- Yes
1158	- No
1159	- Unsure
1160	Guidelines:
1161	- Yes: If the model output is correct or equivalent to the reference answer.
1162	- No: If the model output is incorrect or contradicts the reference answer.
1163	- Unsure: If you can't determine the correctness or if there's insufficient information.
1164	Do not provide any explanation or additional text. Your entire response must be a single word.
1165	Your response:

1167

1168 **Discussion about Relationship Between KL and Completion likelihood.** We report the L_2 loss
1169 between the policy model and the reference model with respect to the likelihood. This serves as a
1170 proxy for KL divergence, as both measure the divergence between the policy and reference models.
1171 While we could not generate a direct KL vs. Likelihood plot due to access restrictions, this proxy
1172 analysis allows us to provide relevant insights without requiring additional model retraining.

1173 As shown in Figure 12, our experiments reveal that likelihood does not strictly correlate with the
1174 L_2 loss: lower likelihood (higher cross-entropy loss) does not necessarily correspond to a higher L_2
1175 loss. This result suggests that the relationship between the likelihood of preferred completions and
1176 the divergence between the models is more nuanced than a simple monotonic association. In par-
1177 ticular, the observed patterns reinforce the idea that likelihood and KL divergence, while connected
1178 under specific assumptions, are not directly interchangeable.
1179

1180 C FURTHER INVESTIGATIONS FOR QUESTION ANSWERING TASKS

1181 **Case studies for NATURALQUESTIONSOPEN and TRIVIAQA tasks.** Table 5 provides two ex-
1182 amples for NATURALQUESTIONSOPEN and TRIVIAQA tasks, respectively.
1183

1184 **LLM-as-a-Judge for the NATURALQUESTIONSOPEN task.** We implement a more flexible eval-
1185 uation method to understand the potential issue of stylistic variations in answers. Instead of rely-
1186 ing on exact string matching, which can be overly rigid, we employ an LLM-as-a-Judge using the
1187

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

Table 5: Model output examples for NATURALQUESTIONSOPEN and TRIVIAQA.

Examples for NATURALQUESTIONSOPEN		
Field	Content	F ₁ Word
Question	Where is dakar located on the world map?	–
High Likelihood Answer	Senegal	100.0%
Mid Likelihood Answer	Dakar is the capital of Senegal and is located in West Africa. It is situated on the western coast of the country, on the Atlantic Ocean.	8.7%
Examples for TRIVIAQA		
Field	Content	F ₁ Word
Question	How many Rings of Power were there, in total?	–
High Likelihood Answer	20	100.0%
Mid Likelihood Answer	There were 20 Rings of Power in total, 3 of which were given to the Elves, 7 to the Dwarves, and 9 to the Men.	8.7%

GPT4o model. As shown in Table 4, this LLM-based evaluation system is presented with the original question, the reference answer, and the model’s output. It then assesses whether the model’s output is correct, incorrect, or if there’s not enough information to make a determination, responding with “Yes”, “No”, or “Unsure” respectively. We compute the model performance based on the percentage of “Yes”. Figure 13 shows the model performance on the ULTRAFEEDBACK dataset using the 7B model. Our analysis reveals that while the LLM-as-a-Judge evaluation method demonstrates a trend similar to the F₁ score, it consistently yields higher performance metrics.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

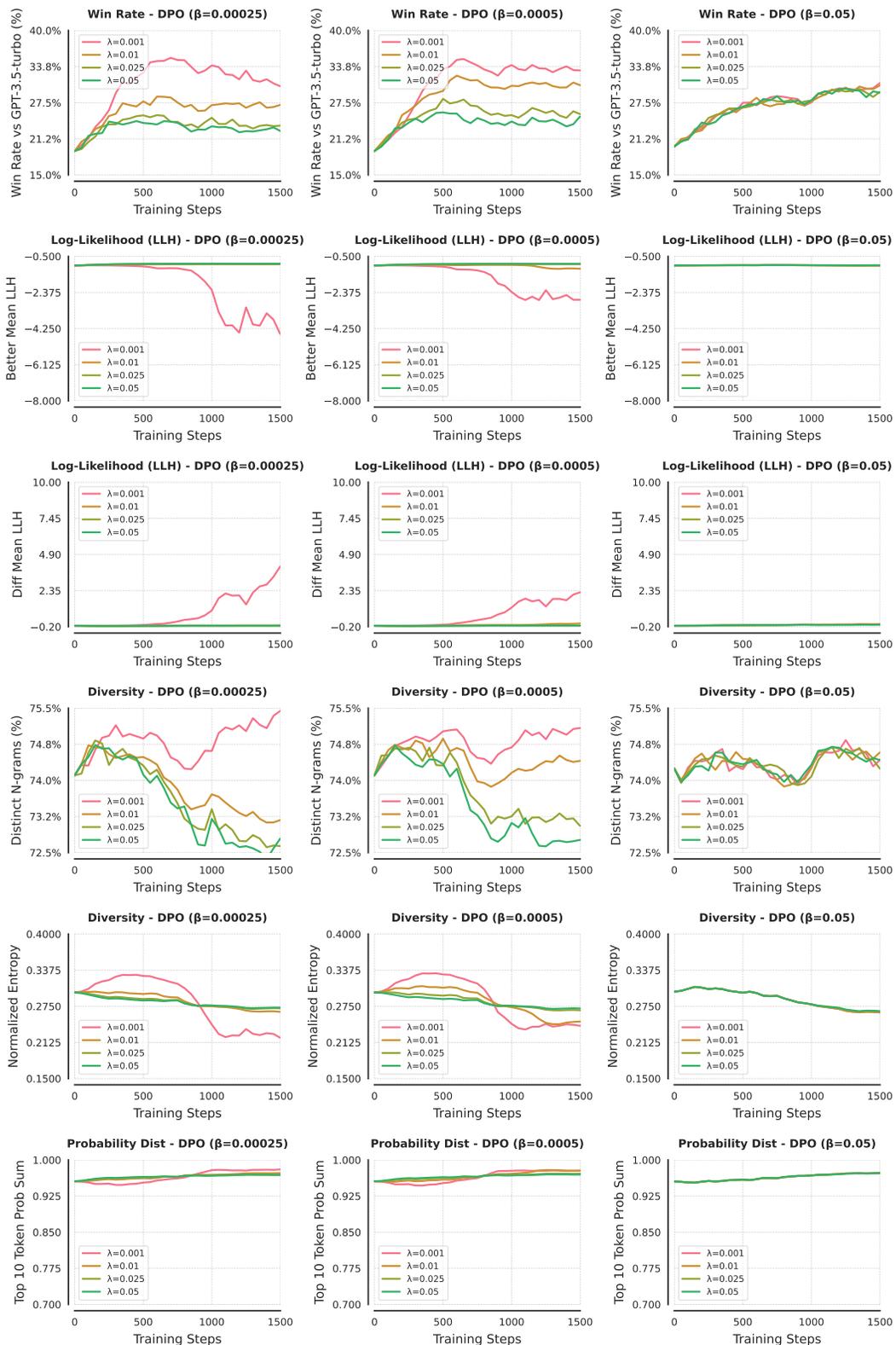


Figure 8: Learning curves across training steps for various metrics. Results are reported for the 7B models using IPO on the ULTRAFEDBACK dataset with varying values of τ and λ .

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

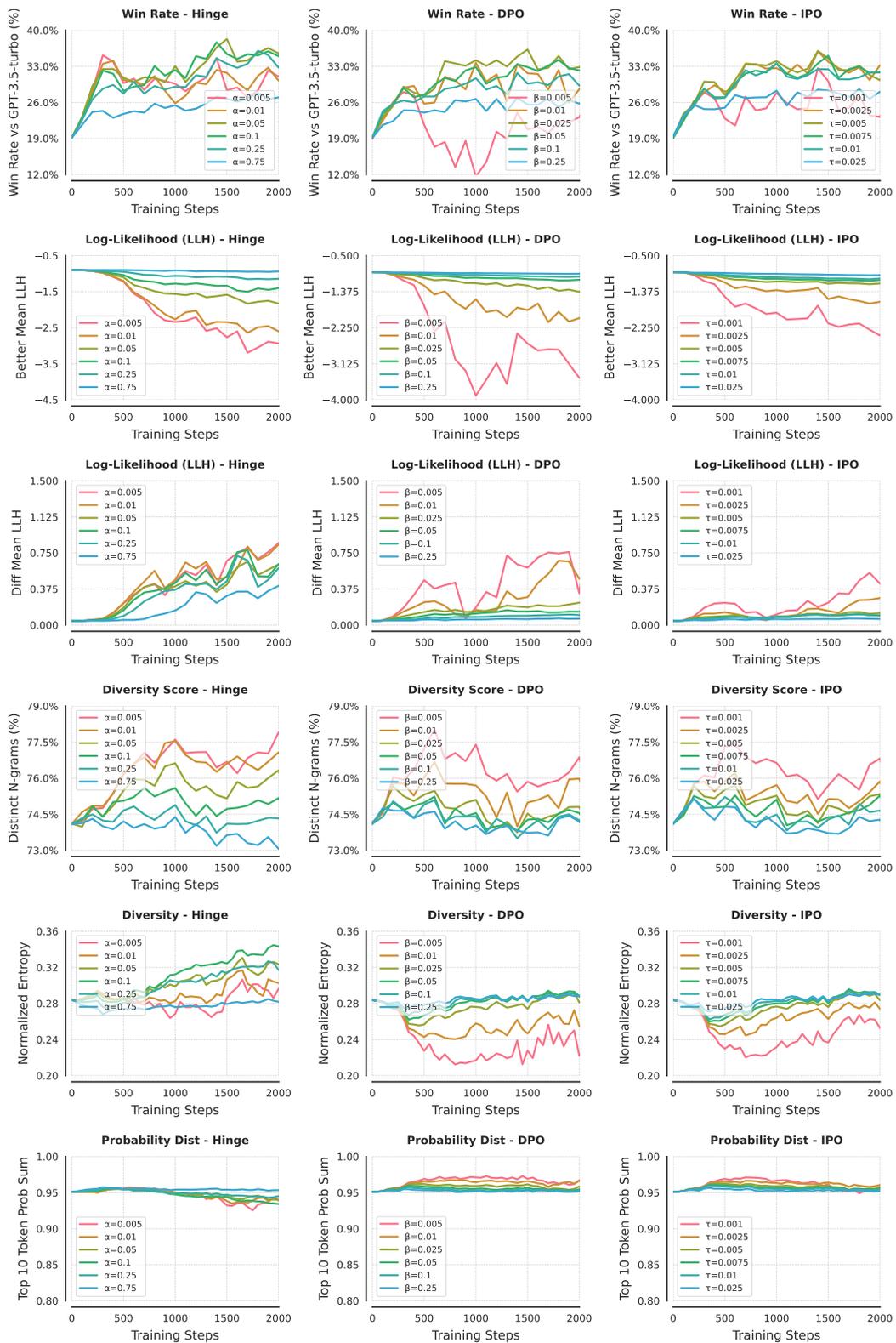


Figure 9: Learning curves across training steps for various metrics. Results are reported for the 7B models using the Hinge, DPO, and IPO on the BinarizedPref dataset.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

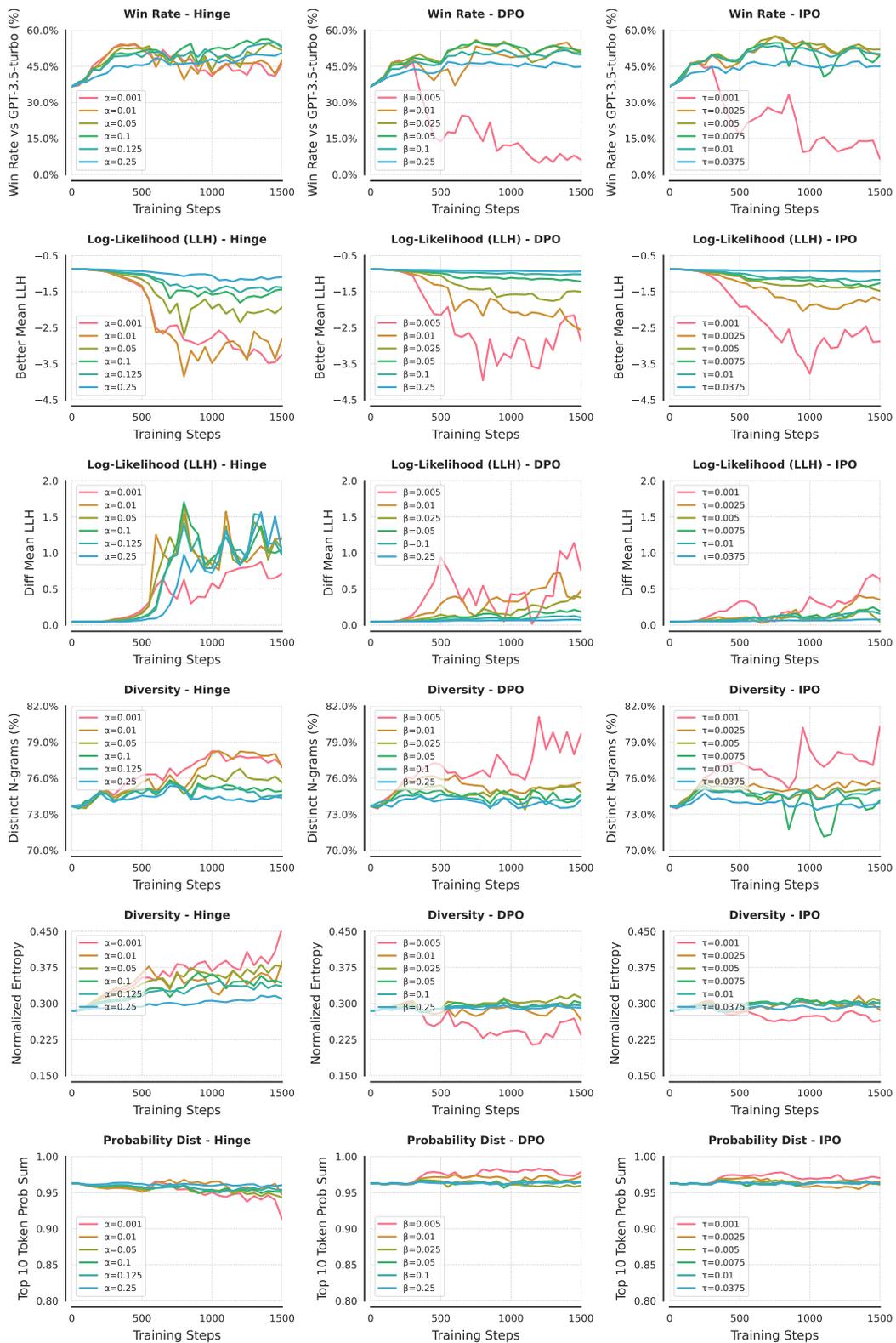


Figure 10: Learning curves across training steps for various metrics. Results are reported for the 35B models using the Hinge, DPO, and IPO on the BINARIZEDPREF dataset.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

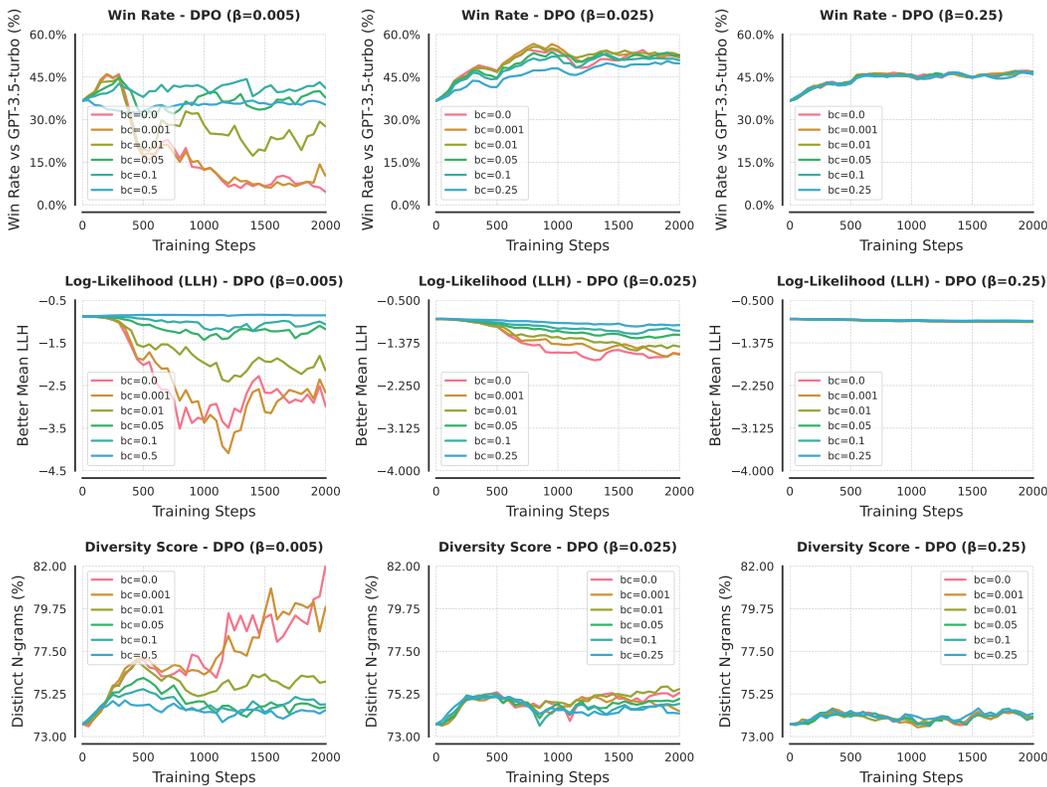


Figure 11: Control Likelihood via training on better completion on the BINARIZEDPREF dataset, using the 35B model. When different runs have similar likelihoods, the win probability and diversity of their model outputs tend to follow the same trend throughout training.

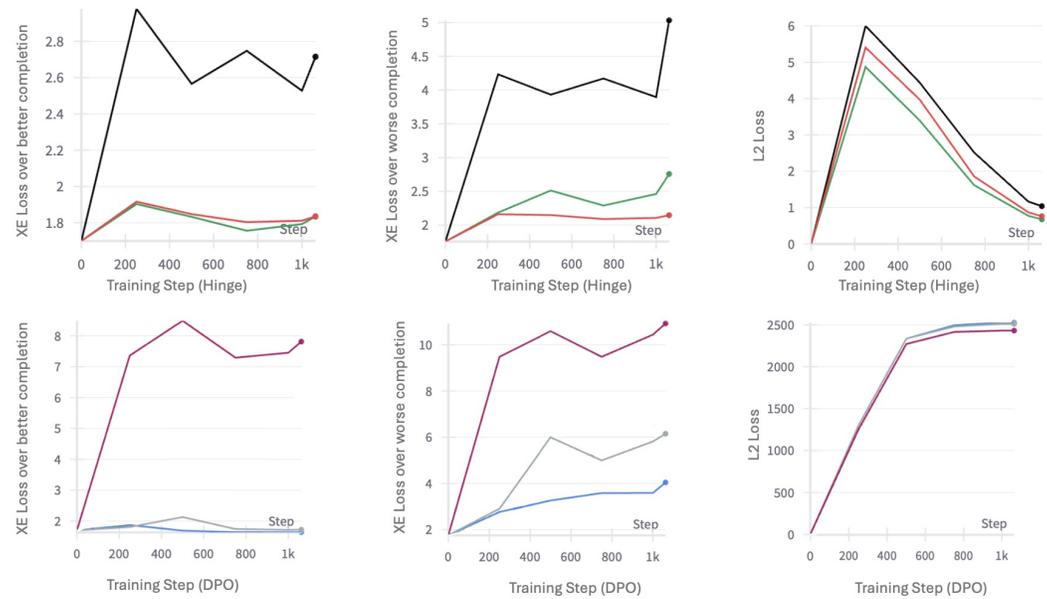


Figure 12: Our results indicate that completion likelihood does not strictly correlate with the L_2 loss: lower likelihood (higher cross-entropy loss) does not necessarily correspond to a higher L_2 loss.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

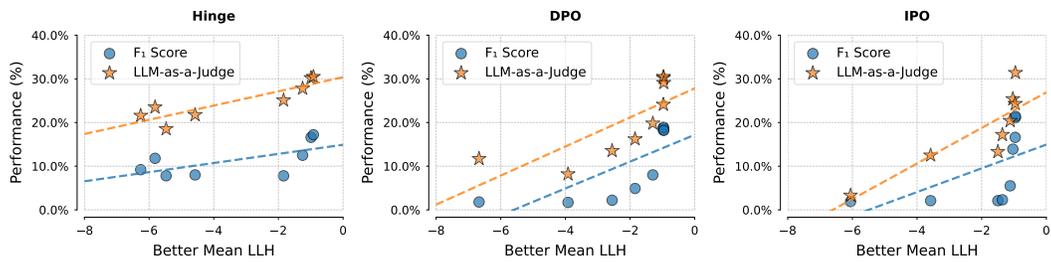


Figure 13: NATURALQUESTIONSOPEN vs Better Mean LLH on the ULTRAFEEDBACK dataset using the 7B model. The F_1 score and LLM-as-a-Judge results are reported.