

Sequence-level Large Language Model Training with Contrastive Preference Optimization

Anonymous ACL submission

Abstract

The next token prediction loss is the dominant self-supervised training objective for large language models and has achieved promising results across a variety of downstream tasks. However, upon closer investigation of this objective, we find that it lacks an understanding of sequence-level signals, leading to a mismatch between the training and inference processes. To bridge this gap, we introduce a contrastive preference optimization procedure that can inject sequence-level signals into the language model at any training stage without expensive human labeling. Notably, our experiments revealed that the proposed objective surpasses the next-token prediction in terms of GPT winning rate on both instruction-following and text generation. Specifically, using OpenLlama-3B, our method achieves a 13% improvement on an instruction-following task, and a 3% increase on a text generation task.

1 Introduction

Next token prediction¹ is now the prevalent way to the pretraining and supervised finetuning (SFT) of large language models (LLM). This loss function can be easily scaled up to train models with trillions of parameters and tokens, and it has also demonstrated the ability to generate coherent and contextually relevant text. Let P be the unknown target language distribution and Q be the distribution of our model at hand. The goal of next token prediction is to minimize the *forward-KL* divergence between P and Q ; during test time, we usually first generate a set of samples using the trained model, and evaluate the quality of these generations using a certain metric², for example

¹This term is subsequently used interchangeably with maximum likelihood estimation (MLE).

²Metric in the sense of a quality measure, rather than the mathematical concept.

the *reverse-KL*. This training process only supervise on predicting one token at a time given the full context from groundtruth. On the other hand during inference, the model needs to generate a whole sequence (for a given prompt) relying on its own prior predictions. This mismatch between the training and inference stage is also known as *exposure-bias* in the literature of RNN and sequence-to-sequence model (Bengio et al., 2015; Ranzato et al., 2015).

In other words, the next token prediction based training injects only *token-level* information into the model, but missing *sequence-level* signal. Of course, such discrepancies can be mitigated by the subsequent reinforcement learning with human feedback (RLHF) step (Ouyang et al., 2022) in LLM training. In RLHF, a reward signal is enforced on the generated sequence of the language model and guides the model generation to align with human preference. RLHF is computationally intensive and often faces instability issues. Therefore, many open-sourced LLMs do not incorporate this discipline. Direct preference optimization (DPO) (Rafailov et al., 2023) is a recently proposed alternative to RLHF, that enables sequence-level LLM training without the need for costly model generations. One drawback of both DPO and RLHF methods is that they require expensive human labeling to score the LLM training samples. RLHF requires human preference data to train the reward-model, and DPO-training requires a supervised pair of positive and negative completions for each given prompt. However, the majority of existing LLM training data does not consist of such human preference information. Therefore, in this work, we ask the following question:

Can we introduce sequence-level information in LLM training even in the absence of human-preference data?

We answer the question affirmatively with our proposed **CONTRASTIVE PREFERENCE**

OPTIMIZATION (CPO) method. CPO shares a similar principle to RLHF/DPO in the sense that they all parameterize (perhaps implicitly) the optimal model (with respect to a certain sequence-level signal) with an energy-based model (EBM). However, the goal of CPO is **not for alignment, but for generation quality**. Therefore unlike RLHF and DPO, the proposed CPO method does not require human preference information as the training signal. Another related method that optimizes the language quality is BRIO (Liu et al., 2022). Although unlike BRIO, the proposed CPO method does not rely on autoregressively sampled negative sequences from the model, and therefore is much more computational efficient and easier to scale up. The experiments in this paper demonstrate that CPO is able to improve the quality of text generation in terms of reward model scores and reverse-KL divergence.

2 Related work

LLMs trained with next token prediction loss (Radford et al., 2019; Chung et al., 2022; Sanh et al., 2021; Zhou et al., 2023) have demonstrated many fascinating capabilities, including the ability to perform zero-shot or few-shot tasks (Radford et al., 2019; Brown et al., 2020), the ability to improve the robustness of visual learning in multimodal models (Menon and Vondrick, 2022; Feng et al., 2023), and the ability to reason (Wei et al., 2022).

Several works have investigated the shortcoming of MLE and exposure bias. Arora et al. (2022) measured the error accumulation of language generation due to exposure bias. Schmidt (2019) connected exposure bias to generalization. Wang and Sennrich (2020) studied how exposure bias leads to hallucination in neural machine translation. To mitigate exposure bias, there exists a long line of work that have explored sequence level training methods. Bengio et al. (2015); Ranzato et al. (2015) proposed to train RNN with RL or RL-related algorithms rather than teacher-forcing. BRIO Liu et al. (2022) targeted the summarization task with the ROUGE signal. Pang and He (2020) trained the language models with an offline learning algorithm.

Recently, RLHF (Stiennon et al., 2020; Ouyang et al., 2022) is developed. While the primary goal of RLHF for model alignment, it is effectively a sequence-level training technique. For the RLHF training, we usually need to gather a pair of continuations for each prefix, where one continuation

aligns with human preference and the other does not. This pair of sequences is used to train a reward model, which is later used to supervise the samples generated by the RL-trained model. The model is typically optimized by REINFORCE (Williams, 1992) or PPO (Schulman et al., 2017).

RLHF process is also closed related to energy-based models (EBM) (Korbak et al., 2022), and RLHF training can be reframed as a supervised learning algorithm coined as direct preference optimization (DPO) (Rafailov et al., 2023) under the assumption of the Bradley-Terry model (Bradley and Terry, 1952) or Plackett-Luce model (Plackett, 1975; Luce, 2012). The particular formulation of the EBM that minimizes the RLHF objective exactly mimics the formulation in Deng et al. (2020), with the reward function being the energy. However, Deng et al. (2020) directly formulate the EBM as a language model, which is computationally heavy for sampling and inference (due to the estimation of the partition function). This EBM form has also been studied in controlled text generation. Kumar et al. (2022) adopted the Langevin dynamics technique to directly sample from the EBM, with different energy functions that characterize toxicity, fluency, and diversity. These methods can all be viewed as sequence-level algorithms for different purposes.

3 Preliminary

Notation Consider a sentence of T tokens $x = \{x_1, \dots, x_T\} \in \mathcal{X}$, and let P be the unknown target language distribution, $\tilde{P}(x)$ be the empirical distribution of the training data (which is an approximation of P), and Q be the distribution of our model at hand. Since our paper is also closely related to RLHF, we will also use π to represent the distributions. In particular, we sometimes write π_θ for a distribution that is parameterized by θ , where θ is usually a subset of trainable parameters of the LLM; we write π_{ref} for a reference distribution that should be clear given the context. The next token prediction loss is minimizing the forward-KL between P and Q .

Forward-KL vs. reverse-KL The forward-KL is formally defined as the following:

$$\begin{aligned}
& \arg \min_Q D_{\text{KL}}(P||Q) \\
& \approx \arg \min_Q D_{\text{KL}}(\tilde{P}||Q) \\
& = \arg \min_Q -\frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} \log Q(\mathbf{x}).
\end{aligned}$$

Since we are only optimizing Q , minimizing the forward-KL is equivalent to the maximum likelihood estimation (MLE) $\max \log Q(X)$. Further decomposing $Q(\mathbf{x}) = \prod_i Q(\mathbf{x}_i | \mathbf{x}_1^{i-1})$, we get the next token prediction loss function

$$\arg \max_Q \sum_{\mathbf{x} \in \mathcal{X}} \sum_{x_t \in \mathcal{X}} \log Q(x_t | \mathbf{x}_1^{t-1}). \quad (1)$$

To actually measure the quality of the generated text, typically we will first generate several sequences and then evaluate the quality of these generated sequences. Here we look closely at the reverse-KL:

$$D_{\text{KL}}(Q||P) = \sum_{\mathbf{x} \in \mathcal{X}} Q(\mathbf{x}) \log \left(\frac{Q(\mathbf{x})}{P(\mathbf{x})} \right), \quad (2)$$

however, since $\mathbf{x} \sim Q$, and we do not have access to P , the reverse-KL cannot be computed exactly.

The equivalence of RLHF and EBM For the completeness of this paper, we include the result on the equivalence between RLHF and EBM. For the full proofs, we refer the reader to (Rafailov et al., 2023; Korbak et al., 2022).

The RLHF objective is the following:

$$\begin{aligned}
& \max_{\pi_\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_\theta(\mathbf{y}|\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] \\
& - \beta D_{\text{KL}}(\pi_\theta(\mathbf{y}|\mathbf{x}) || \pi_{\text{ref}}(\mathbf{y}|\mathbf{x})),
\end{aligned} \quad (3)$$

where $\mathbf{x} \sim \mathcal{D}$ is a given prefix, $\mathbf{y} \sim \pi_\theta(\mathbf{y}|\mathbf{x})$ is a sampled continuation from the model at training π_θ . Meanwhile we want to control the divergence between π_θ and π_{ref} , the latter is usually an already pretrained or finetuned LLM. Its optimum is achieved at the following EBM:

$$\pi^*(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \exp \left(\frac{1}{\beta} r(\mathbf{x}, \mathbf{y}) \right), \quad (4)$$

where $Z(\mathbf{x}) = \sum_{\mathbf{y}} \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \exp \left(\frac{1}{\beta} r(\mathbf{x}, \mathbf{y}) \right)$ is the partition function.

4 Our approach

While the RL penalty with KL control eq. (3) is widely adopted in RLHF, it can also be used directly to train LLMs: instead of a preference reward $r(\cdot)$, we can use any metric that measures general text qualities, including ROUGE, BLEU, MAUVE, etc. The benefit of eq. (3) over eq. (1) is that r guides the model over a whole sequence \mathbf{y} , rather than just a single token. This motivates our work to investigate the possibility of using such objective in pretraining and SFT stage of LLMs.

Following Rafailov et al. (2023), we assume that the preference over two sequences \mathbf{y}_w and \mathbf{y}_l given \mathbf{x} is parameterized by the Bradley-Terry model:

$$P(\mathbf{y}_w \succ \mathbf{y}_l) = \frac{e^{r(\mathbf{x}, \mathbf{y}_w)}}{e^{r(\mathbf{x}, \mathbf{y}_l)} + e^{r(\mathbf{x}, \mathbf{y}_w)}}. \quad (5)$$

The optimal policy π^* takes the aforementioned EBM form eq. (4). and this EBM reparametrization establish the equivalence between the original RLHF object eq. (3) and the following supervised objective:

$$\begin{aligned}
& \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = \\
& \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(\mathbf{y}_w | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_w | \mathbf{x})} \right. \right. \\
& \left. \left. - \beta \log \frac{\pi_\theta(\mathbf{y}_l | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_l | \mathbf{x})} \right) \right],
\end{aligned} \quad (5)$$

where $\sigma(\cdot)$ is the Sigmoid function.

We can also generalize the formulation to the Plackett-Luce model, where we have a linear ordering $\tau(\cdot)$ among K sequences:

$$\begin{aligned}
& \mathcal{L}_{\text{DPO}}(\pi_\theta, \pi_{\text{ref}}) = \\
& \mathbb{E}_{\substack{\tau, \mathbf{x} \sim \mathcal{D} \\ \mathbf{y}_1, \dots, \mathbf{y}_K}} \left[\log \prod_{k=1}^K \frac{\exp \left(\beta \log \frac{\pi_\theta(\mathbf{y}_{\tau(k)} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{\tau(k)} | \mathbf{x})} \right)}{\sum_{j=k}^K \exp \left(\beta \log \frac{\pi_\theta(\mathbf{y}_{\tau(j)} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{\tau(j)} | \mathbf{x})} \right)} \right].
\end{aligned} \quad (6)$$

Here, $\tau(1), \dots, \tau(K)$ induce a ranking among K sequences. To ease the notation, from now on, we always assume that $\mathbf{y}_1 \sim \mathcal{D}$ is the natural text appeared in the training data.

Investigating the DPO objective, we notice two caveats for its use in the pretraining and SFT stages:

1. We need human labelers to gather $\mathbf{y}_l \sim \mathcal{D}$.
2. There may not be a natural ranking among negative sequences $\mathbf{y}_2, \dots, \mathbf{y}_K$ in terms of text quality. To tackle the first point, we sample $\mathbf{y}_l \sim \mathcal{A}$ where

239 \mathcal{A} is some noise distribution where it is cheap to
 240 sample; to tackle the second point, we provide a
 241 variant objective that models a “best-of- K ” event:
 242 \mathbf{y}_1 is the best among K sequences, rather than a
 243 linear ordering event $\mathbf{y}_{\tau(1)} \succ \mathbf{y}_{\tau(2)} \succ \dots \succ \mathbf{y}_{\tau(K)}$.
 244 These modifications lead to our proposed CPO ob-
 245 jective:

$$246 \mathcal{L}_{\text{CPO}}(\pi_\theta, \pi_{\text{ref}}) = \mathbb{E}_{\substack{(\mathbf{x}, \mathbf{y}_1) \sim \mathcal{D} \\ \mathbf{y}_2, \dots, \mathbf{y}_K \sim \mathcal{A}}} \left[\log \frac{\exp\left(\beta \log \frac{\pi_\theta(\mathbf{y}_1|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1|\mathbf{x})}\right)}{\sum_{j=1}^K \exp\left(\beta \log \frac{\pi_\theta(\mathbf{y}_j|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_j|\mathbf{x})}\right)} \right]. \quad (7)$$

247 If ranking information is desired, we have the
 248 following CPO objective with ranking:

$$249 \mathcal{L}_{\text{CPO}}(\pi_\theta, \pi_{\text{ref}}) = \mathbb{E}_{\substack{\tau, (\mathbf{x}, \mathbf{y}_1) \sim \mathcal{D} \\ \mathbf{y}_2, \dots, \mathbf{y}_K \sim \mathcal{A}}} \left[\log \frac{\prod_{k=1}^K \exp\left(\beta \log \frac{\pi_\theta(\mathbf{y}_{\tau(k)}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{\tau(k)}|\mathbf{x})}\right)}{\sum_{j=k}^K \exp\left(\beta \log \frac{\pi_\theta(\mathbf{y}_{\tau(j)}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{\tau(j)}|\mathbf{x})}\right)} \right]. \quad (8)$$

250 We will later discuss some possible choices of
 251 ranking signals, and show that the ranking can in-
 252 deed further improve the text generation quality.

253 The crucial aspect of CPO is how to generate the
 254 negative sequences $\mathbf{y}_2, \dots, \mathbf{y}_k \sim \mathcal{A}$. For RLHF, the
 255 negative sequences are simply the ones that humans
 256 dislike. For the qualities of text generation, we
 257 implicitly model the sequence-level signal $r(\mathbf{x}, \mathbf{y})$
 258 such that $r(\mathbf{x}, \mathbf{y}_k) < r(\mathbf{x}, \mathbf{y}_1), \forall k \in \{2, \dots, K\}$.
 259 In other word, we use the reward $r(\cdot)$ that prefers
 260 the ground truth to any other sequences. Import-
 261 tantly, the actual signal r is not parameterized ex-
 262 plicitly, instead it is represented by the log density
 263 ratio $\log \frac{\pi_\theta}{\pi_{\text{ref}}}$.

264 4.1 Connection to noise contrastive 265 estimation

266 Noise contrastive estimation (NCE) (Gutmann and
 267 Hyvärinen, 2010) is a novel estimation technique
 268 introduced to tackle the computational infeasibility
 269 of traditional likelihood-based methods in large-
 270 scale machine learning models, particularly those
 271 involving high-dimensional data. NCE diverges
 272 from typical maximum likelihood estimation by
 273 transforming the problem into a classification task,

274 which is deeply connected to both DPO and CPO.
 275 In NCE, the model is trained to distinguish between
 276 real data and noise/synthetic data. Beyond binary
 277 classification, RankingNCE³ also train the model
 278 to rank the real data higher than all the noise sam-
 279 ples (Ma and Collins, 2018).

280 There are two important distinctions between
 281 CPO and NCE. First, instead of asking the model
 282 to distinguish between real data and noise (at which
 283 any reasonable language model should already be
 284 good), we ask the model to distinguish *better* than
 285 a reference model does, hence making the model
 286 better at recognizing natural text. Second, we in-
 287 corporate a denser signal by incorporating the simi-
 288 larity among embeddings of different samples. The
 289 experiments in this paper demonstrate that such a
 290 dense training signal consistently improves the text
 291 generation quality.

292 4.2 Synthetic negative samples

293 In this work, we propose four ways to generate
 294 synthetic negative samples. The first is to autore-
 295 regressively generate continuations to the training
 296 prefixes from a model trained with the next token
 297 prediction loss. We fix the synthetic data genera-
 298 tion strategy to be top- k sampling with $k = 50$.
 299 The advantage of this strategy to the forthcoming
 300 strategies is that the generated continuations are
 301 of higher quality and lead to better downstream
 302 performance, while the disadvantage is that sam-
 303 pling is slow. We denote these negative samples
 304 as **autoregressive negatives (AN)**. One can speed
 305 up the sampling process via speculative sampling
 306 (Chen et al., 2023) or using a smaller or distilled
 307 model, this direction is orthogonal to our approach
 308 and can be directly incorporated into our frame-
 309 work.

310 The second way is to directly use the continua-
 311 tions to other (possibly unrelated) prefixes within
 312 the same mini-batch as the negative samples. More
 313 specifically, given a batch of prefixes and continua-
 314 tions $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^b$, the negative samples to the prefix
 315 \mathbf{x}_i are composed of $\{\mathbf{y}_j\}_{j \neq i}$. Although these neg-
 316 ative samples are not difficult to distinguish, they
 317 are very simple to create and can be easily scaled
 318 up. We denote these as **batch negatives (BN)**.

319 The third way is to perform a token-level per-
 320 turbation. Given a sequence $\mathbf{y} = \{y_1, \dots, y_T\}$,
 321 we randomly select c percent of the positions

³Despite the name, it means the model is ranking the real data highest among all data, rather than learning a total ordering.

$\{t_1, \dots, t_j\} \subseteq [T]$, and substitute each y_{t_i} independently based on $\pi_\theta(y_{t_i}|y_1, \dots, y_{t_i})$. We call these **meanfield negatives (MN)**. The name is based on the fact that we use a fully separable distribution to approximate the AR distribution, but note that this is not the optimal mean-field approximation. Nevertheless, computing this particular meanfield approximation does not take an additional cost, compared to estimating the best meanfield approximation. This method does not generate semantically meaningful sentences, but it does generate hard negative samples, as the model tends to give them high probabilities.

Lastly, for each ground truth continuation, we can truncate the continuation at a random position and append an extra EOS token to the end. We denote this as **truncation negatives (TN)**.

4.3 Possible ranking signals

As mentioned before, our reward implicitly prefers the ground truth over other sequences, and we do not explicitly model the reward parametrically. The upside of the implicit representation of the reward is that it bypasses the shortcuts (e.g. Krishna et al., 2021) that are known to other explicit metrics⁴. However, since we do not have access to a concrete score for the text quality, when presented more than one negative samples, we do not have a direct ranking among them.

Prior works on sequence-level training (Liu et al., 2022; Bengio et al., 2015) have suggested a variety set of signals, including BLEU, ROGUE, and BertScore. These signals are usually specific for certain downstream tasks like translation or summarization. In modern era of LLMs, they have been shown to not align with human evaluations anymore (Goyal et al., 2022). Since our goal is to improve text generation or instruction following, the cosine similarity between embeddings is a more intuitive signal to measure the distance between sequences. The usage of embedding for generation quality measurement is also suggested in the MAUVE metric (Pillutla et al., 2021).

When presented K sequences and ranking is desired, the sequences are ranked based on its cosine similarity with the ground truth. Let e_1, \dots, e_K be the embeddings of given sequences and without loss of generality assume e_1 is the ground truth,

⁴Note that here we are not claiming CPO does not exist any shortcuts. The statement here simply means other existing metrics have known shortcuts.

we define $\tau(e_i) < \tau(e_j)$ if $\frac{\langle e_i, e_1 \rangle}{\|e_i\| \|e_1\|} > \frac{\langle e_j, e_1 \rangle}{\|e_j\| \|e_1\|}$, with the lower ranking index indicating the better sample. Using the objective eq. (8), this process gives us denser signals during training, and can lead to better downstream performance.

Another good candidate for ranking signal is the reward model score. In fact, since the downstream performance is judged by a reward model, this will probably yield the best test performance as well. However, one has to train and host an extra reward model, creating extra memory and computation overhead. Therefore, we did not include such signal during training in this work.

4.4 Approximate reverse-KL

In the subsequence experiment, we show how CPO improves reverse-KL. As we discuss previously, an unavoidable issue of calculating the reverse-KL is we do not have access to the probability of the generated sequences under the true language distribution. However, if we agree that the capability of imitating true language scales with the model size, then we can approximate the true language distribution P with a more capable model \hat{P} , hence approximating the reverse-KL divergence. Since many of our tasks are conditional by nature, for example, the instruction following task is to generate a response, condition on the input instruction, we further consider the expected reverse-KL divergence:

$$\begin{aligned} & \mathbb{E}_x \left[D_{\text{KL}} \left(Q(\cdot|x) \parallel \hat{P}(\cdot|x) \right) \right] \\ & \approx \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} Q(y|x) \log \left(\frac{Q(y|x)}{\hat{P}(y|x)} \right), \end{aligned} \quad (9)$$

where \mathcal{X} is the set of inputs (e.g. instructions) in the test set, and \mathcal{Y} is the set of generated continuations (e.g. responses). During our evaluation, we also notice that a more capable Q tends to generate sequences y with lower probability $Q(y|x)$, compare to a less capable Q . This phenomenon is indeed expected, since a more capable model should be able to generate more diverse continuations. To overcome the numerical instability with a vanishing Q , we also use the following surrogate:

$$\frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \frac{-\log \hat{P}(y|x)}{|y|}, \quad (10)$$

where $|y|$ is the length of y .

This is the log conditional probability normalized by length, and its usage has been justified in Cho et al. (2014); Liu et al. (2022); Fu et al. (2023). Particularly, Fu et al. (2023) has discussed the use of the normalized conditional probability of a capable evaluator.

5 Experiment

Throughout this section, we use **BN** for models trained with batch negatives, **MN** for models trained with meanfield negatives, **TN** for models trained with truncation negatives, **MixN** for a mixed negative sampling strategies which the details should be found in its context, and **AN** for autoregressive negatives. We use **ANR** for models trained with autoregressive negatives and ranking signals, similarly we can define **MixNR**, etc.

Task and model. We consider two tasks in this paper. The first is an instruction-following task, trained and evaluated on the Dolly dataset (Conover et al., 2023). This dataset is composed of 15011 total instruction and response pairs. We train with 7505 sequences and test with the rest 7506. We use the pretrained GPT2-XL (Radford et al., 2019) and OpenLlama-3B (Touvron et al., 2023; Geng and Liu, 2023) as the base model. The second task is an open-ended text generation task on the Wikidump data (Foundation). We train the OpenLlama-3B model to predict the rest 85% tokens given the leading 15% tokens.

Training details. Throughout the experiment, we fix the learning rate to be $1e - 5$, we use the AdamW optimizer with weight decay of 0.05. We keep the batch size to be 64. Unless otherwise specified, for the baseline model, we train the GPT2-XL and OpenLlama-3B with the next token prediction loss for 2000 steps, and denote them as the MLE models. Using these models as the reference model π_{ref} , we continue to train with the CPO objective either with or without ranking signals, with $\beta = 5$, for 1000 steps. For both models, each batch during training contains 11 negative samples in total. For MixN and MixNR, we also use a negative sample size of 11, consisting 3 BN, 5 MN, and 3 TN.

Evaluation. As discussed in Goyal et al. (2022), almost all automated evaluation metrics have been shown to not align with human evaluations in modern era of LLMs, hence we decide to use GPT-3.5 (Brown et al., 2020) as the evaluation tool. For

each test instruction, we ask the models to generate continuations with various generation configurations, and query the reward model whether it prefers the generated continuations or the ground truth. A winning rate is then computed over all test instructions. As pointed out in Wang et al. (2023), GPT models are prone to position bias. When evaluating by asking GPT which of the two inputs it prefers, one can easily manipulate the result by exchanging the input positions. To counter this bias, for each test instruction, we ask both the CPO model and baseline model to generate continuations, and we compare each of them to the ground truth to calculate the winning rate. Now since both models’ generations suffer from the same position bias, we can meaningfully compare the difference between their winning rates against the ground truth.

The query template is the following: “For the following query to a chatbot, which response is more helpful?
Query: {}
Response A: {}
Response B: {}
State only "A" or "B" to indicate which response is more helpful.
More helpful:” For efficiency, we generate and evaluate 1000 samples in the test set.

In addition to winning rate, we also evaluate the model performance based on reverse-KL and normalized log conditional probability, as described in section 4.4.

Weight-space ensemble. Previous works (Liu et al., 2022) have also suggested to combine the newly proposed loss function with the MLE training objective $\mathcal{L}_{\text{MLE}} + \alpha \mathcal{L}_{\text{CPO}}$, the downside of combining loss functions in this way is that for a different choice of α one will have to retrain the model. To investigate the importance of loss combination, we instead take a similar approach to WISE-FT (Wortsman et al., 2022) and perform a weight-space ensemble. In particular, denote θ_{CPO} and θ_{MLE} the model parameters trained solely with CPO or MLE respectively, we generate with the interpolated weights $\theta = \alpha \theta_{\text{MLE}} + (1 - \alpha) \theta_{\text{CPO}}$.

5.1 Instruction-following task

On the Dolly instruction-following task, our proposed CPO method with various negative sampling strategies consistently outperforms the MLE baseline models. Using greedy sampling with GPT2-XL, the CPO model has a clear margin over the MLE model, and CPO+ANR has a 3.5% higher

Table 1: Study of the effect of different negative samples.

BNR	MNR	TNR
0.599	0.567	0.601

Table 2: The winning rate of GPT2-XL against the ground truth, samples generated by greedy decoding, evaluated by GPT-3.5.

	MLE	ANR	MixNR			
α	-	-	0	0.5	0.7	0.9
WinRate	0.471	0.506	0.476	0.479	0.487	0.485

Table 3: The winning rate of OpenLlama-3B trained with either MLE or CPO+MixNR against the ground truth, evaluated by GPT-3.5. The samples are generated by various strategies, we only present MLE and MixNR models here.

Model \ Config	$k = 50, p = 1$	$k = 50, p = 0.7$	beam= 2	beam= 4
	MixNR	0.591	0.611	0.607
MLE	0.497	0.517	0.532	0.514

winning rate, see table 2. Keep in mind that the CPO process only incur very little computation overhead during the actual training. Even if we generate the negative samples autoregressively, this cost is only offline and is one-time.

The improvement on OpenLlama-3B is more significant: CPO+ANR has a 13.8% higher winning rate than the MLE baseline, and CPO+MixNR has a 9.8% higher winning rate in table 4. We also observe that WISE-FT has a positive impact on the model. Heuristically, for OpenLlama-3B, a smaller α is preferred (more emphasis on the CPO weights) (table 4), and the reverse holds for GPT2-XL (table 2). We hypothesize that the choice of α should depend on model parameters: if the model is more capable, then it can benefit more from CPO. Here we show the existence of a good α , and we leave further exploration to future research.

Generation configuration. In addition to greedy decoding, we also experiment with different choice of sampling strategies. We test with various settings of top-k top-p sampling, as well as different length of beam search. In all settings, CPO has consistently demonstrated a superior performance to MLE table 3.

Effect of different negative samples. We perform a study on the effects of different negative sampling strategies, the results are presented in table 1. We first train the OpenLlama-3B model with MLE loss for 1000 steps, then we continually train with CPO for 200 steps. For every ground truth sequences, we use 4 negative sequences. In this setting, we always use the ranking information to

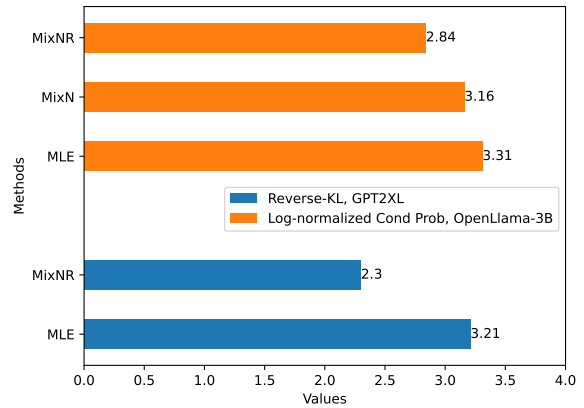


Figure 1: Reverse-KL and negative log-normalized conditional probability of different negative sampling strategies. Smaller numbers indicate better quality. CPO consistently outperforms MLE. CPO+ranking information also outperforms CPO without ranking information.

train CPO. We observe that the effects of BNR and TNR on the reward model preference is similar, and they perform slightly better than MNR.

Reverse-KL. The reverse-KL (eq. (9)) and negative log-normalized conditional probability (eq. (10)) metrics are reported in fig. 1. Smaller numbers indicate higher quality for both metrics. MixNR and MixN both demonstrate improvements over the MLE model. Since these metrics measure how likely the generated texts are under the (approximated) language distribution, these results serve as complementary explanation to CPO’s improvement in reward model scores.

Table 4: The winning rate of OpenLlama-3B trained with CPO and MLE against the ground truth data in Dolly, sampled by greedy decoding, evaluated by GPT-3.5. MLE₁, ANR and AN are trained for 200 steps, the rest models are trained for 1000 steps. The best CPO model outperforms the MLE baseline by 13.8% winning rate.

	MLE ₁	ANR	AN	MLE ₂	MIXNR						MIXN
α	-	-	-	-	0	0.1	0.3	0.5	0.7	0.9	-
WinRate	0.505	0.643	0.56	0.522	0.608	0.620	0.614	0.610	0.601	0.550	0.576

Instruction: Do you prefer a Tesla model S or a Porsche Taycan

MLE response: S

CPO response: ? While both are excellent cars, if I had to choose between the two, I would go with the Tesla model S. The Tesla model S has a larger battery size (92 kWh vs 79 kWh), faster charging capabilities (230 kW vs 180 kW), and has a better range (418 miles vs 318 miles). The Tesla model S also feels more powerful and is easier to drive.

Table 5: OpenLlama-3B’s winning rate against the ground truth continuation on Wikidump. The model is trained with either MLE or CPO+BNR. Weight ensemble is adopted. The best CPO model outperforms the MLE baseline by 3% winning rate.

	MLE	BNR			
α	-	0	0.5	0.7	0.9
WinRate	0.508	0.455	0.505	0.5	0.538

5.2 Open-ended text generation task

We further test OpenLlama-3B’s ability on an open-ended text generation task with CPO. Using the Wikidump data (Foundation), for each test sample, we take its first 15% tokens as the prefix and train the model with CPO on the rest 85%. For negative sampling, we use four BNR examples. The results table 5 indicates that with a correctly picked weight interpolation coefficient α , CPO can greatly improve the model’s winning rate against the MLE baseline by 3%. The results also have a different pattern compared to the instruction-following task: the optimal choice of α shows a reverse trend. With the Dolly dataset we observe a small optimal α , but on the Wiki dataset we see a large optimal α .

5.3 What type of generations do CPO tend to create?

Investigating the generations of CPO vs those of MLE, we notice that CPO model tends to create more detailed continuations/responses to given prefixes/instructions, partly explaining why these generations are preferred by GPT reward. As the sam-

ple demonstrates, the CPO response appears to be more helpful with more details.

6 Conclusion

In this paper, we propose an auxiliary CPO loss function for SFT, it can be used with or without ranking signals depending on the quality of the negative samples. We investigated several ways to generate the negative samples, each with its own pros and cons. Experimentally, we show that both GPT2-XL and OpenLlama-3B models benefit from training with our proposed CPO objectives. On Dolly instruction-following task, OpenLlama-3B+CPO has a winning rate 13.8% higher than MLE; GPT2-XL has a winning rate 3.5% higher. On Wikipedia text generation task, OpenLlama-3B+CPO has a winning rate 3% higher than the MLE baseline model. It is interesting to explore other ways to efficiently generate high-quality negative data beyond the autoregressive fashion. One possible direction is to consider Langevin dynamic sampling, which samples all tokens in parallel.

References

- Kushal Arora, Layla El Asri, Hareesh Bahuleyan, and Jackie Chi Kit Cheung. 2022. Why exposure bias matters: An imitation learning perspective of error accumulation in language generation. *arXiv preprint arXiv:2204.01171*.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28.

608	Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. <i>Biometrika</i> , 39(3/4):324–345.	660	Tomasz Korbak, Ethan Perez, and Christopher L Buckley. 2022. RL with kl penalties is better viewed as bayesian inference. <i>arXiv preprint arXiv:2205.11275</i> .	661
609		662		663
610		664	Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering. <i>arXiv preprint arXiv:2103.06332</i> .	665
611		666		667
612	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	668	Sachin Kumar, Biswajit Paria, and Yulia Tsvetkov. 2022. Gradient-based constrained sampling from language models. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2251–2277.	669
613		670		671
614		672	Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. Brio: Bringing order to abstractive summarization. <i>arXiv preprint arXiv:2203.16804</i> .	673
615		674		675
616		676	R Duncan Luce. 2012. <i>Individual choice behavior: A theoretical analysis</i> . Courier Corporation.	677
618	Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. Accelerating large language model decoding with speculative sampling. <i>arXiv preprint arXiv:2302.01318</i> .	678	Zhuang Ma and Michael Collins. 2018. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. <i>arXiv preprint arXiv:1809.01812</i> .	679
619		680		681
620		682	Sachit Menon and Carl Vondrick. 2022. Visual classification via description from large language models. <i>arXiv preprint arXiv:2210.07183</i> .	683
621		684	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	685
622		686		687
623	Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. <i>arXiv preprint arXiv:1409.1259</i> .	688		689
624		690	Richard Yuanzhe Pang and He He. 2020. Text generation by learning from demonstrations. <i>arXiv preprint arXiv:2009.07839</i> .	691
625		692		693
626		694	Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. <i>Advances in Neural Information Processing Systems</i> , 34:4816–4828.	695
627	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	696		697
628		698		699
629		700	Robin L Plackett. 1975. The analysis of permutations. <i>Journal of the Royal Statistical Society Series C: Applied Statistics</i> , 24(2):193–202.	701
630		702	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	703
631		704		705
632	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm .	706	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>arXiv preprint arXiv:2305.18290</i> .	707
633		708		709
634		710		711
635		712	Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. <i>arXiv preprint arXiv:1511.06732</i> .	713
636		714		715
637	Yuntian Deng, Anton Bakhtin, Myle Ott, Arthur Szlam, and Marc’Aurelio Ranzato. 2020. Residual energy-based models for text generation. <i>arXiv preprint arXiv:2004.11714</i> .			
638				
639				
640				
641	Zhili Feng, Anna Bair, and J Zico Kolter. 2023. Leveraging multiple descriptive features for robust few-shot image learning. <i>arXiv preprint arXiv:2307.04317</i> .			
642				
643				
644				
645	Wikimedia Foundation. Wikimedia downloads .			
646	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. <i>arXiv preprint arXiv:2302.04166</i> .			
647				
648				
649	Xinyang Geng and Hao Liu. 2023. Openllama: An open reproduction of llama .			
650				
651	Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. <i>arXiv preprint arXiv:2209.12356</i> .			
652				
653				
654	Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In <i>Proceedings of the thirteenth international conference on artificial intelligence and statistics</i> , pages 297–304. JMLR Workshop and Conference Proceedings.			
655				
656				
657				
658				
659				

715 Victor Sanh, Albert Webson, Colin Raffel, Stephen H
716 Bach, Lintang Sutawika, Zaid Alyafeai, Antoine
717 Chaffin, Arnaud Stiegler, Teven Le Scao, Arun
718 Raja, et al. 2021. Multitask prompted training en-
719 ables zero-shot task generalization. *arXiv preprint*
720 *arXiv:2110.08207*.

721 Florian Schmidt. 2019. Generalization in generation:
722 A closer look at exposure bias. *arXiv preprint*
723 *arXiv:1910.00292*.

724 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec
725 Radford, and Oleg Klimov. 2017. Proximal
726 policy optimization algorithms. *arXiv preprint*
727 *arXiv:1707.06347*.

728 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel
729 Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
730 Dario Amodei, and Paul F Christiano. 2020. Learn-
731 ing to summarize with human feedback. *Advances*
732 *in Neural Information Processing Systems*, 33:3008–
733 3021.

734 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
735 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
736 Baptiste Rozière, Naman Goyal, Eric Hambro,
737 Faisal Azhar, et al. 2023. Llama: Open and effi-
738 cient foundation language models. *arXiv preprint*
739 *arXiv:2302.13971*.

740 Chaojun Wang and Rico Sennrich. 2020. On exposure
741 bias, hallucination and domain shift in neural ma-
742 chine translation. *arXiv preprint arXiv:2005.03642*.

743 Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai
744 Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang
745 Sui. 2023. Large language models are not fair evalu-
746 ators. *arXiv preprint arXiv:2305.17926*.

747 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
748 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
749 et al. 2022. Chain-of-thought prompting elicits
750 reasoning in large language models. *Advances in*
751 *Neural Information Processing Systems*, 35:24824–
752 24837.

753 Ronald J Williams. 1992. Simple statistical gradient-
754 following algorithms for connectionist reinforce-
755 ment learning. *Machine learning*, 8:229–256.

756 Mitchell Wortsman, Gabriel Ilharco, Jong Wook
757 Kim, Mike Li, Simon Kornblith, Rebecca Roelofs,
758 Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali
759 Farhadi, Hongseok Namkoong, et al. 2022. Robust
760 fine-tuning of zero-shot models. In *Proceedings of*
761 *the IEEE/CVF Conference on Computer Vision and*
762 *Pattern Recognition*, pages 7959–7971.

763 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao
764 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
765 Lili Yu, et al. 2023. Lima: Less is more for align-
766 ment. *arXiv preprint arXiv:2305.11206*.