

Preference Tuning with Human Feedback on Language, Speech, and Vision Tasks: A Survey

Genta Indra Winata

*AI Foundations, Capital One
New York, United States*

GENTA.WINATA@CAPITALONE.COM

Hanyang Zhao

*Department of IEOR, Columbia University
New York, United States*

HZ2684@COLUMBIA.EDU

Anirban Das

*AI Foundations, Capital One
New York, United States*

ANIRBAN.DAS3@CAPITALONE.COM

Wenpin Tang

David D. Yao
*Department of IEOR, Columbia University
New York, United States*

WT2319@COLUMBIA.EDU

DDY1@COLUMBIA.EDU

Shi-Xiong Zhang

Sambit Sahu
*AI Foundations, Capital One
New York, United States*

SHIXIONG.ZHANG@CAPITALONE.COM

SAMBIT.SAHU@CAPITALONE.COM

Abstract

Preference tuning is a crucial process for aligning deep generative models with human preferences. This survey offers a thorough overview of recent advancements in preference tuning and the integration of human feedback. The paper is organized into three main sections: 1) **introduction and preliminaries**: an introduction to reinforcement learning frameworks, preference tuning tasks, models, and datasets across various modalities: language, speech, and vision, as well as different policy approaches, 2) **in-depth exploration of each preference tuning approach**: a detailed analysis of the methods used in preference tuning, and 3) **applications, discussion, and future directions**: an exploration of the applications of preference tuning in downstream tasks, including evaluation methods for different modalities, and an outlook on future research directions. Our objective is to present the latest methodologies in preference tuning and model alignment, enhancing the understanding of this field for researchers and practitioners. We hope to encourage further engagement and innovation in this area. Additionally, we provide a GitHub link <https://github.com/hanyang1999/Preference-Tuning-with-Human-Feedback>.

1. Introduction

Learning from human feedback is a crucial step in aligning generative models with human preferences to generate output that closely resembles human speech and writing. Despite the powerful learning capabilities of generative models in self-supervised learning, these models frequently misinterpret instructions, leading to hallucinations in generation (Ji et al., 2023a; Yao et al., 2023a). Additionally, ensuring the safety of the generated content remains

a significant challenge for these models. Extensive research on preference tuning using human feedback has demonstrated that adversarial samples can be utilized to jailbreak systems (Rando & Tramèr, 2023; Wei et al., 2024). Ideally, generative models need to be controlled to ensure that their outputs are safe and do not cause harm. Models often exhibit unintended behaviors, such as fabricating facts (Chen & Shu, 2023; Sun et al., 2024), producing biased or toxic text (Hartvigsen et al., 2022), or failing to follow user instructions (Ji et al., 2023b; Tonmoy et al., 2024). Additionally, maintaining the privacy of data is crucial to ensure the safe operation of models and protect user privacy (Brown et al., 2022). In the text-to-image generation task, large-scale models often struggle to produce images that are well-aligned with text prompts (Feng et al., 2022), particularly in compositional image generation (Liu et al., 2022; Lee et al., 2023), object recognition (Qiao et al., 2024), and coherent generation (Liu et al., 2023). Similarly, in text-to-speech tasks, Zhang, Li, Li, Zhang, Wang, Zhou, and Qiu (2024), Chen, Hu, Wu, Wang, Chng, and Zhang (2024a) integrate subjective human evaluation into the training loop to better align synthetic speech with human preferences.

The application of preference tuning has been widely used in language tasks by training instruction-tuned large language models (LLMs), such as Llama (Touvron et al., 2023b; Dubey et al., 2024), Phi (Abdin et al., 2024), Mistral (Jiang et al., 2023), Nemotron (Parmar et al., 2024; Adler et al., 2024), Gemma (Team et al., 2024). Commercial models like GPT-4 (Achiam et al., 2023), Gemini (Team et al., 2023; Reid et al., 2024), Claude (Anthropic, 2024), Command-R, and Reka (Ormazabal et al., 2024) have leveraged human preference alignment to enhance their performance. Alignment of LLM improves task-specific skills, coherence, fluency, and helps avoid undesired outputs. Additionally, alignment research has benefited multilingual LLMs, such as Aya (Aryabumi et al., 2024; Üstün et al., 2024), BLOOMZ, and mT0 (Muennighoff et al., 2023), as well as regional LLMs like Cendol (Cahyawijaya et al., 2024) and SEALLM (Nguyen et al., 2023). Common approaches to achieving LLM alignment involve reinforcement learning techniques that guide language models to follow preferred samples by maximizing rewards. Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017) is the initial approach that is used to align models with human preference, which is further applied to the deep learning space that has been popularized by its successes in LLMs (Ouyang et al., 2022; Bai et al., 2022a) via PPO (Schulman et al., 2017), REINFORCE (Kool et al., 2019), Online Directed Preference Optimization (online DPO) (Guo et al., 2024a), and Supervised Fine-Tuning (SFT)-like approach (Dong et al., 2023). It typically involves three key aspects: human feedback collection, reward modeling, and *online* RL for policy optimization. Recent methods, however, allow for training the reward model alongside the policy model in an *offline* manner, as demonstrated by DPO (Rafailov et al., 2024), and SLIC-HF (Zhao et al., 2023). Moreover, preference tuning is also applied to vision-text tasks, and has been shown to improve the representation of both image and text using the alignment score of image and text embeddings (Ramesh et al., 2022; Saharia et al., 2022; Yu et al., 2022b) measured by pre-trained vision-text models, such as CLIP (Radford et al., 2021) and CoCa (Yu et al., 2022a). Wu, Sun, Zhu, Zhao, and Li (2023d) utilize LoRA (Hu et al., 2021) to align Stable Diffusion (Lee et al., 2023), a vision-text pre-trained model. The application in speech has not been much explored, and there is only a handful works in the literature. Zhang et al. (2024) focus on investigating alignment between codes and the text.

Interestingly, while most preference tuning methods rely on training-based optimization that modifies model weights, recent works are starting to explore training-free, inference-time techniques that align model outputs with human preferences without altering parameters. These include prompt and in-context alignment, feedback-driven refinement, latent vector-based modulation, and decoding-time adjustments—offering lightweight, flexible alternatives for dynamic preference adaptation (Huang et al., 2024; Lin et al., 2024; Liu et al., 2024; Li et al., 2025). While the majority of the paper focuses on categorizing and analyzing various preference optimization techniques, for the sake of comprehensiveness, we also examine how these methods are evaluated in practice. A preference optimization technique is only useful if it results in *measurable* improvements in alignment with human intent or task-specific goals. Therefore, we also include a survey of the current landscape of complementary evaluation strategies or metrics in speech, language, vision and reward modeling domain, that enable systematic comparison and benchmarking.

In this paper, we survey the recent advances of preference tuning with human feedback in different modalities. It provides not only a comprehensive introduction including preliminaries to get readers familiar with the topic, but also an in-depth review on the latest proposed approaches and in-depth discussions. To summarize, the paper comprises the following contributions:

- We provide a comprehensive overview of preference tuning for models on different modalities, such as language, speech, and vision tasks, and expand our survey to all existing preference tuning methods, including reinforcement learning (RL) approaches.
- We formulate and taxonomize a systematic framework and classification for preference tuning for deep generative models from the existing literature.
- We present various applications of preference tuning to improve generation aspects using human feedback. We also describe the automatic and human-based evaluations to measure the quality of generation in deep generative models.
- We discuss the opportunities and future directions for preference tuning.

Through this survey, we aim to present the recent methodologies on preference tuning and alignment for deep generative models, enabling researchers and practitioners to better understand this topic and further innovate.

Survey Paper Organization We structure the paper as follows: In Section 2 We introduce the formal definitions of the tasks and the notations used throughout. We then examine preference tuning pipelines, datasets, and the various generative models employed in preference optimization in Section 3. In Section 4, we explore online alignment methods in depth, and examine offline alignment approaches in Section 5. In Section 6 we discuss combined policies and sampling-agnostic alignment strategies. We then analyze training-free optimization techniques in Section 7 and compare with training-based optimization. We review evaluation methods for preference-optimized models in Section 8 and finally conclude with a discussion on emerging topics in preference optimization and potential future research directions in Section 9.

2. Preliminaries

This section outlines the preliminaries of preference tuning, including the formal definitions of the tasks and the notations used throughout this paper. Additionally, we provide a taxonomy for classifying preference tuning methods.

2.1 Tasks and Definition

In general, the entire preference tuning mechanism for generative models can be formulated as a RL problem described as follows.

2.1.1 RL FRAMEWORK CONCEPTS

Policy Model The policy model π_θ is a generative model that takes in an input prompt x and returns a sequence of output or probability distributions y . We define a generative model as a policy model π_θ where it is parameterized by θ with a policy model π . Given a prompt x , a generative model generates an output y as following:

$$\pi_\theta(y|x) = \prod_t \pi_\theta(y_t|x, y_{<t}), \quad (1)$$

where y_t is the t -th token in the response and $y_{<t}$ is tokens in the response before y_t . For example, for the text-based tasks, the input prompt is a text sequence x and the output is a probability distribution over text vocabulary of LLM y ; and for the vision-text-based tasks, such as text-to-image tasks, the input x is the text sequence, and y is the generated image.

Reward Model The reward model (RM) processes both the input x and the target y , passing them through the model to obtain a reward $r_\theta(y|x)$, which reflects the notion of preferability. This preferability score can also be interpreted as a relative score assigned to the target y given the input x . Less preferred outcomes receive a lower score compared to more preferred samples.

Action Space The action refers to all tokens corresponding to the vocabulary of generative models. For text tasks, the action space encompasses the entire vocabulary of the LLM. For vision tasks (similarly for speech tasks), the action space consists of real values representing the image, for example, the next hierarchy in diffusion generative models (if understanding diffusion models as Hierarchical Variational Autoencoders (Luo, 2022)).

Environment The distribution encompasses all possible input token sequences for generative models. In text-based tasks, these input token sequences correspond to text sequences, highly depending on the sampling methods for the inference. In vision tasks, they correspond to possible images.

2.1.2 PREFERENCE DATA

In the preference tuning pipeline, we utilize the supervised data \mathcal{D}_{sft} and the preference data $\mathcal{D}_{\text{pref}}$. We denote the supervised data $\mathcal{D}_{\text{sft}} = [(x^1, y^1), \dots, (x^M, y^M)]$ as a list of input and label pairs. Specifically for the text SFT data, x can be represented as prompts. The prompt $x^i = (I^i, F^i, Q^i)$ consists of the concatenation of an instruction I^i , few-shot samples F^i , and

a query Q^i . Then, we denote the preference data $\mathcal{D}_{\text{pref}} = [(x^1, y_w^1, y_l^1), \dots, (x^N, y_w^N, y_l^N)]$, a list of input x^i with preferred response y_w^i and dispreferred response y_l^i , and they are either sampled from the reference policy model π_{ref} or collected by human annotation. Generally, given the preference data, we can obtain a reward r associated to the response with the input.

2.1.3 TERMINOLOGY AND NOTATION

Table 1 lists the common notations used in this survey paper. The table serves as a quick reference guide for understanding the mathematical expressions and technical terms used throughout the paper.

Name	Notation	Description
Input Sequence	x	Input sequence that is passed to the model.
Output Sequence	y	Expected label or output of the model.
Dispreferred Response	y_l	Negative samples for reward model training.
Preferred Response	y_w	Positive samples for reward model training.
Optimal Policy Model	π^*	Optimal policy model.
Policy Model	π_θ	Generative model that takes the input prompt and returns a sequence of output or probability distribution.
Reference Policy Model	π_{ref}	Generative model that is used as a reference to ensure the policy model is not deviated significantly.
Preference Dataset	$\mathcal{D}_{\text{pref}}$	Dataset with a set of preferred and dispreferred responses to train a reward model.
SFT Dataset	\mathcal{D}_{sft}	Dataset with a set of input and label for supervised fine-tuning.
Loss Function	\mathcal{L}	Loss function.
Regularization Hyper-parameters	$\alpha, \beta_{\text{reg}}$	Regularization Hyper-parameters for preference tuning.
Reward	r	Reward score.
Target Reward Margin	γ	The margin separating the winning and losing responses.
Variance	β_i	Variance (or noise schedule) used in diffusion models.

Table 1: Table of Terminology and Notation.

2.2 Taxonomy

We define the following categories for all of the preference tuning approaches as shown in Table 2. Figure 1 shows the five categories we study in this survey paper and described in the following:

Sampling Likewise in the literature of RL, we categorize the methods based on how we sample the data and use them to train or obtain the reward: *offline* and *online* human alignments. The categorization is related to how we compute the reward and use it in the policy models. In online human alignment setting, the agent that collects a batch of examples by interacting with the environment and uses them to update the policy. The reward of the examples can be collected by the reward model or samples generated by the policy model. While for the offline human alignment setting, the data are collected from offline human demonstrations. In online reinforcement learning, we say a method is

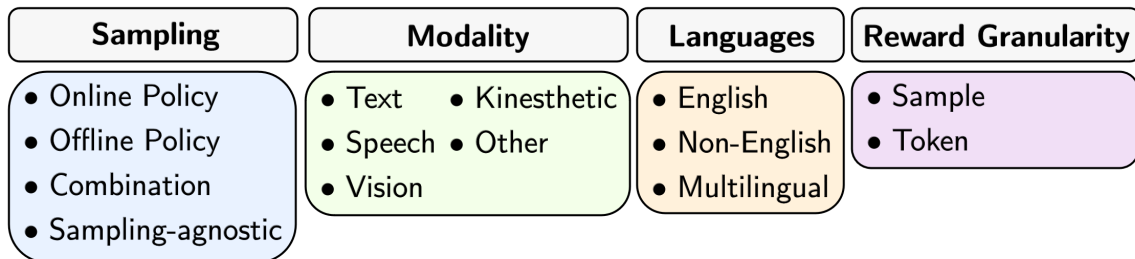


Figure 1: Taxonomy of the Preference Tuning methods.

on-policy if the agent collects data using the same policy that it is currently improving. Conversely, a method is *off-policy* if the data is gathered by a different behavior policy than the one being optimized.

Modality We study the use of preference tuning on various modality, such as text, speech, vision, kinesthetic and others if we are not able to classify them. In the latest advancement of NLP, the idea of RL has been further explored to language and speech tasks, even in multi-modal tasks, such as vision-text. Thus, it is essential to categorize the papers by the extend of the study in terms of the modality, such as text, speech, vision, and vision-text.

Language We explore the preference tuning application on different languages. In this case, we categorize the method by English, non-English, and multilingual.

Reward Granularity In the preference tuning, the reward can be computed in different granularity levels. The granularity levels can be expanded into two: sample- and token-level. The token-level for each modality may differ, for example, in text tasks, we can use subwords from vocabulary as tokens. And, in vision tasks, patches of image are tokens.

3. Preference Tuning

In this section, we cover the general framework to train preference-tuned generative models. As shown in Table 3, the preference tuning training framework typically begins with the supervised fine-tuning (SFT) stage, during which the generative model is trained to excel at next-token prediction or use an instruction-tuned model as the base initialized model. The SFT focuses on improving the model capability to generate tokens as it guides the model on how an generative model should response to a prompt input. Once the model is able to properly generate fluent text sequences, the model is further aligned by further policy optimization via RL. The alignment is useful to guide the model to answer with a appropriate manner based on the preference objective. This step is a necessary training stage to make sure the model generation aligned to human preference, thus, the model will act more human-like. Notably, the human alignment stage can also be jointly trained alongside SFT. As a disclaimer, rigorously speaking, preference learning is only one direction among different directions and approaches for alignment (though the most dominant and popular one for now), but in this paper, since we will only focus on the preference line methods, we will use these two terminologies interchangeably.

Method	Modality					Languages			Reward Granularity	
	Text	Speech	Vision	Kinesthetic	Other	EN	Non-EN	Multi.	Sample	Token
Online Methods										
RLHF (Christiano et al., 2017)										
PPO (Schulman et al., 2017)										
AI Feedback (Bai et al., 2022b)	✓	×	×	×	×	✓	×	×	✓	×
P3O (Wu et al., 2023c)	✓	×	×	×	×	✓	×	×	✓	×
MaxMin-RLHF (Chakraborty et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
Multi-Ling RLHF (Dang et al., 2024)	✓	×	×	×	×	✓	✓	✓	×	×
RLHF-PPO (Ouyang et al., 2022)	✓	×	×	×	×	✓	×	×	✓	×
RLHF Workflow (Dong et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
REINFORCE (Williams, 1992)										
ReMax (Li et al., 2023c)	✓	×	×	×	×	✓	×	×	✓	×
RLOO (Ahmadian et al., 2024a)	✓	×	×	×	×	✓	×	×	✓	×
GRPO (Shao et al., 2024)	✓	×	×	×	×	✓	✓	✓	✓	×
Online DPO										
Iterative DPO (Xu et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
OAIF (Guo et al., 2024a)	✓	×	×	×	×	✓	×	×	✓	×
OPTune (Chen et al., 2024d)	✓	×	×	×	×	✓	×	×	✓	×
Self-Rewarding (Yuan et al., 2024b)	✓	×	×	×	×	✓	×	×	✓	×
Nash-Learning										
NLHF (Mimos et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
SPPO (Wu et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
SFT-like										
RAFT (Dong et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
ReST (Gulcehre et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
RRHF (Yuan et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
SuperHF (Mukobi et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
Multi-Modal Models										
Diffusion (Schulman et al., 2017)										
AlignProp (Prabhudesai et al., 2023)	✓	×	✓	×	×	✓	×	×	✓	×
DDPO (Black et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
DPOK (Fan et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
DRaFT (Clark et al., 2023)	✓	×	✓	×	×	✓	×	×	✓	×
PRDF (Deng et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
ReFL (Xu et al., 2024b)	✓	×	✓	×	×	✓	×	×	✓	×
VLLM (Liu et al., 2024b)										
DLPO (Chen et al., 2024c)	✓	×	✓	×	×	✓	×	×	✓	×
HIVE (Zhang et al., 2024b)	✓	×	✓	×	×	✓	×	×	✓	×
LLaVA-rlhf (Sun et al., 2023)	✓	×	✓	×	×	✓	×	×	✓	×
RLHF-V (Yu et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
Rich Feedback (Liang et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
Offline Methods										
BPPO (Zhuang et al., 2023)	×	×	×	✓	×	×	×	×	✓	×
Multi-Modal Models										
Diffusion-DPO (Wallace et al., 2024)	✓	×	✓	×	×	✓	×	×	✓	×
POVID (Zhou et al., 2024b)	✓	×	✓	×	×	✓	×	×	✓	×
Offline DPO (Rafailov et al., 2024)										
ALLO (Chen et al., 2024)	✓	×	×	×	×	✓	×	×	×	✓
CPO (Guo et al., 2024b)	✓	×	×	×	×	✓	✓	✓	✓	×
GPO (Tang et al., 2024b)	✓	×	×	×	×	✓	×	×	✓	×
IPO (Azar et al., 2024)	×	×	×	×	✓	×	×	×	✓	×
KTO (Ethayarajh et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
ODPO (Amini et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
ORPO (Hong et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
PRO (Song et al., 2024b)	✓	×	×	×	×	✓	×	×	✓	×
R-DPO (Park et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
rDPO (Chowdhury et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
sDPO (Kim et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
VPO (Chen et al., 2024b)	✓	×	×	×	×	✓	×	×	✓	×
Mallows-DPO (Chen et al., 2024b)	✓	×	×	×	×	✓	×	×	✓	×
RainbowPO (Zhao et al., 2024a)	✓	×	×	×	×	✓	×	×	✓	×
SimPO (Meng et al., 2024)	✓	×	×	×	×	✓	×	×	✓	×
(Li et al., 2024b)	✓	×	×	×	×	✓	✓	✓	✓	×
SLiC-HF (Zhao et al., 2023)	✓	×	×	×	×	✓	×	×	✓	×
Combination										
P3O (Fakoor et al., 2020)	×	×	✓	×	×	×	×	×	✓	×
RTO (DPO + PPO) (Zhong et al., 2024)	✓	×	×	×	×	✓	×	×	×	✓
Sampling-Agnostic										
ExPO (Zheng et al., 2024a)	✓	×	×	×	×	✓	×	×	✓	×

Table 2: Preference Tuning methods. The categorization based on the methods under study and it does not limit the extension of the method to other domains or modalities.

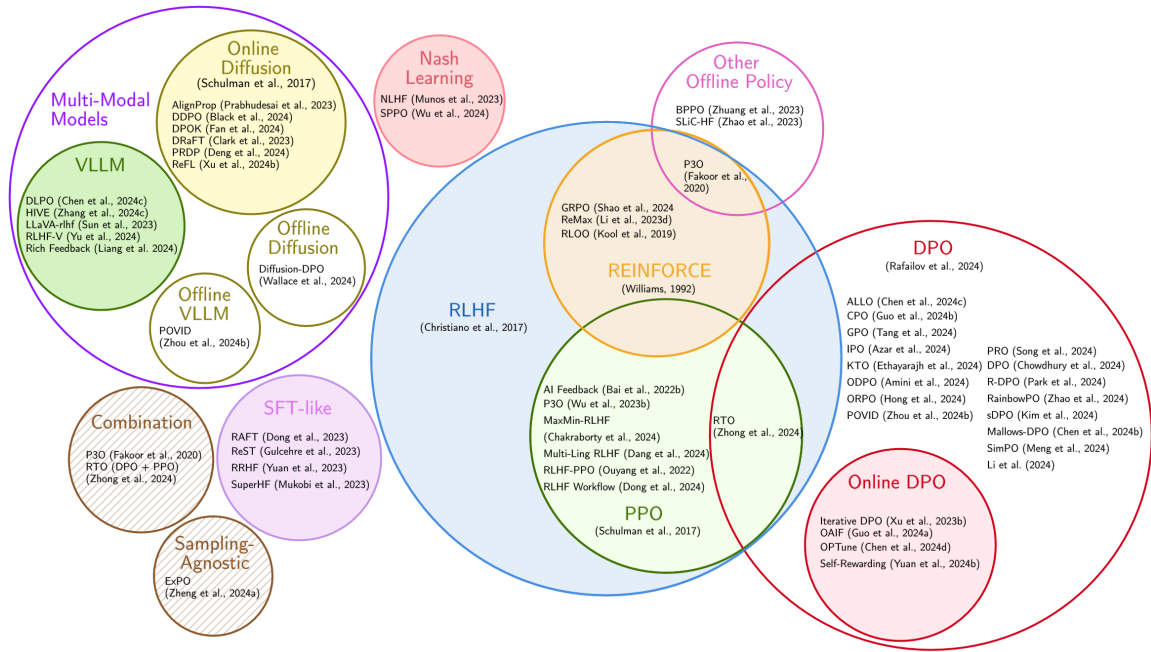


Figure 2: Preference Tuning methods. The circles with shaded areas represent off-policy methods, while the unshaded circles denote on-policy methods. The overlapping area signifies methods that incorporate both on-policy and off-policy approaches. The policy-agnostic circle indicates methods that are applicable to either on-policy or off-policy scenarios. The combination circle represents methods that integrate both online and off-policy strategies.

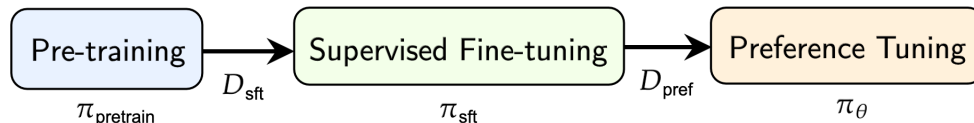


Figure 3: Training stages.

3.1 Training Phases

The training phases for preference tuning are described as follows.

3.1.1 SUPERVISED FINE-TUNING (SFT)

On the preference tuning, a generative model with trainable weights θ normally starts by SFT via maximum likelihood (MLE) using teacher forcing and cross-entropy loss. The training is done using the supervised fine-tuning dataset \mathcal{D}_{sft} . The objective is to maximize the log probability of a set of human demonstrations. The generative model is trained to generate the label by predicting the next token y_{t+1} given the input x , current and previous

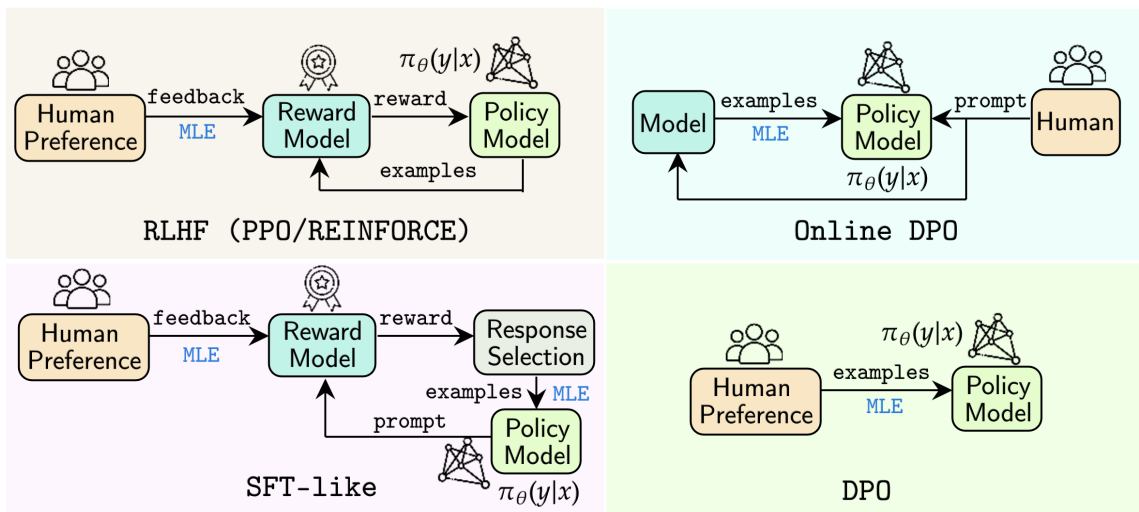


Figure 4: Preference Tuning methods for online algorithms, such as RLHF, Online DPO, and SFT-like, and offline methods, such as DPO.

Reward Model	Sizes	Model Base	Datasets
Single Objective			
BTRM Qwen2	7B Δ	Qwen2	UNK
Eurus-RM (Yuan et al., 2024a)	7B Δ	Mistral	UltraInteract, UltraFeedback, UltraSafety
FsfairX-LLama3-v0.1 (Dong et al., 2023)	8B Δ	Llama3	UNK
GRM-llama3-8B-sftreg (Yang et al., 2024a)	8B Δ	Llama3	Preference 700K
GRM-llama3-8B-distill (Yang et al., 2024a)	8B Δ	Llama3	Preference 700K
InternLM2 (Cai et al., 2024)	1.8B Δ , 7B Δ , 20B Δ	UNK	UNK
SteerLM-Llama3 (Wang et al., 2024)	70B Δ	Llama3	HelpSteer2
Nemotron-4-340B-Reward (Adler et al., 2024)	340B Δ	Nemotron4	HelpSteer2
Pair-preference-model-LLamA3-8B (Dong et al., 2024)	8B Δ	LLama3	RLHFFlow Pair Preference
Starling-RM-34B	34B Δ	Yi-34B-Chat	Nectar
UltraRM (Cui et al., 2023)	13B Δ	Llama2	UltraFeedback
Multi-Objective			
ArmoRM-Llama3-8B-v0.1 (Wang et al., 2024)	8B Δ	Llama3	HelpSteer, UltraFeedback, BeaverTails-30k CodeUltraFeedback, Prometheus, Argilla-Capybara Argilla-OpenOrca, Argilla-Math-Preference
Multi-Model			
MetaMetrics-RM (Winata et al., 2024)	Multiple	Multiple	Skywork Preference Data and AllenAI Preference Data

Table 3: Reward Models.

label tokens $y_{t:<t}$. During the SFT, we utilize an attention mask applying to the entire context x and $y_{t:<t}$, and avoid applying attention to future tokens. The trained model denoted π_θ^{sft} and it is often to be used to initialize reward model and policy model π_θ .

3.1.2 REWARD MODELING

The reward model $r_\phi(x, y)$ can be trained either separately (offline) or jointly trained with the policy model π_θ (online). Table 3 shows the list of reward models.

Single Objective Reward Model Bradley-Terry Reward Model (Bradley & Terry, 1952) is a pairwise comparison between two samples. It estimates the probability that

the pairwise comparison $i \succ j$, which indicates a strong preference of i over j , is true as:

$$P(i \succ j) = \frac{\exp s_i}{\exp s_i + \exp s_j}, \quad (2)$$

where s_i and s_j are latent variables representing sample i and sample j , respectively. Thus, given the preference dataset $\mathcal{D}_{\text{pref}} = \{x^i, y_w^i, y_l^i\}_{i=1}^N$, we could obtain an estimation of the reward model $r_\phi(x, y)$ by minimizing the negative log-likelihood loss:

$$\mathcal{L}(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_{\text{pref}}} \log P(y_w \succ y_l | x) \quad (3)$$

$$= -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_{\text{pref}}} \log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l)), \quad (4)$$

which σ denotes the logistic function, i.e., $\sigma(x) := (1 + e^{-x})^{-1}$.

Multi-Objective Reward Model Absolute-Rating Multi-Objective Reward Model (ArmoRM) (Wang et al., 2024) is a two-stage approach that first trains a multi-objective RM and then learns a gating layer that scalarizes reward objectives in a mixture-of-experts way. Each example consists of an input x and output y with k -dimensional rating vector, where each dimension corresponds to a reward objective. A concatenation of input and output $x \oplus y$ is passed through the model f_θ with a linear regression layer w , which outputs a k -dimensional rating prediction. The model is trained with regression loss:

$$\min_{\theta, w} \mathbb{E}_{x, y, r \in \mathcal{D}} \|w^\top f_\theta(x \oplus y) - r\|_2^2. \quad (5)$$

Then, it learns a mixture-of-experts gating function, g_ϕ , which is implemented as a shallow MLP. This MLP takes the representation of the input x and outputs a k -dimensional vector, which is then processed by a softmax function. During the training of the gating layer, the backbone and the regression layer are kept frozen. Only the gating layer is trained using the Bradley-Terry loss, augmented with an additional scaling variable.

Multi-Model Reward Model MetaMetrics (Winata et al., 2024) is a method to combine multiple existing reward models into a more powerful reward model by calibrating them using the preference data. The method is a systematic way to identify reward models that can be used complementary without blindly use the models. There are two methods introduced to calibrate the models using Bayesian optimization and boosting method. Thus, the approach is highly efficient and they are aspect-agnostic, thus allowing flexibility to use them in any preference data.

3.1.3 PREFERENCE ALIGNMENT

In the domain of preference alignment, methods can be primarily categorized into two types based on how data is utilized during training: online and offline preference alignment.

Online Preference Alignment Online preference alignment is a method in which data for preference tuning is continuously sampled throughout the training process. This approach effectively addresses the misalignment between the Supervised Fine-Tuning (SFT) objective and the ultimate goal of generating high-quality outputs as determined by human evaluation. While SFT significantly enhances performance, online preference alignment

ensures that the model remains aligned with human preferences. Stiennon, Ouyang, Wu, Ziegler, Lowe, Voss, Radford, Amodei, and Christiano (2020), Ouyang et al. (2022) propose reinforcement learning from human feedback (RLHF) to further align language models with human intent. RLHF pipeline starts with the stage of modeling the rewards from human preferences, known as reward modeling stage, by maximizing the likelihood of preferences under the ground truth assumption. After obtaining the RM, RLHF further trains the Language Model policy via Reinforcement Learning to maximize the score given by the RM. Proximal Policy Optimization (PPO) was commonly chosen as the RL algorithm to update the policy because of its great sample efficiency. We provide a detailed explanation in Section 4.

Offline Preference Alignment Offline preference alignment utilizes pre-existing data for preference tuning prior to the start of training. This approach is generally more efficient than online methods, as it removes the need for continuous data sampling and reward modeling. A prominent example is Directed Preference Optimization (DPO) (Rafailov et al., 2024). A detailed explanation is provided in Section 5.

3.1.4 JOINT TRAINING

Recent works also proposed that two stages of SFT and RLHF can be simplified as one stage with a weighted combination of the two loss functions and even lead to better performance. The key takeaway is to treat the preferred answer in the Human Alignment/RLHF stage as the SFT target, e.g., SLiC-HF (Zhao et al., 2023).

3.2 Datasets

The dataset sources for SFT and preference tuning can be collected from various sources, such as human and LLMs feedback. Table 4 shows the list of SFT and alignment text data labeled by the data source either they are collected by human or synthetically generated by LLM.

3.2.1 SFT DATASETS

The SFT data is useful for training LM on high-quality input-output demonstration pairs. This is usually conducted for the foundation model as initialization. The SFT data can be in the form of prompts with various format.

LLM-Generated Datasets Taori et al. (2023) propose Alpaca, a dataset with demonstrations generated using OpenAI’s GPT-3 text-davinci-003 model. The instruction data can be used to conduct instruction tuning for LLMs and allow them to follow instruction better. A version of Alpaca dataset with Chain-of-Thought (CoT) (Wei et al., 2022) and it is introduced to further improve the LLM’s reasoning ability. Multi-turn datasets generated using LLMs are also created, such as ChatAlpaca, UltraChat (Ding et al., 2023), and WildChat (Zhao et al., 2024c).

Human-Generated and Human-Annotated Datasets Using human-generated and human-annotated data are essential in training high-quality models. Zhou et al. (2024a) have shown quality is more important than quantity, as shown as using LIMA datasets

Dataset	# Samples or (# Tokens) or [Byte Size]	Usecase		Data Source		Annotation
		SFT	Alignment	Human	LLM	Human
Alpaca (Taori et al., 2023)	52k	✓	×	✓	✓	×
Alpaca-CoT [△]	127.5M [†]	✓	×	✓	✓	✓
Aya Dataset (Singh et al., 2024)	202k	✓	×	✓	×	✓
ChatAlpaca [△]	20k [‡]	✓	×	✓	✓	×
BeaverTails (Ji et al., 2024)	30k, 330k	✓	✓	✓	✓	✓
Code-Alpaca [△]	20k	✓	×	✓	✓	×
CodeUltraFeedback [△]	10k	✓	✓	✓	✓	✓
Dolly (Conover et al., 2023)	15k	✓	×	✓	×	✓
FLAN collection (Longpre et al., 2023)	UNK [‡]	✓	×	✓	×	✓
HC3 (Guo et al., 2023)	24.3k	✓	✓	✓	✓	✓
HelpSteer2 (Wang et al., 2024)	21k	✓	✓	✓	✓	✓
HH-RLHF (Bai et al., 2022a)	170k	×	✓	✓	×	✓
InstructionWild v2 (Ni et al., 2023)	110k	✓	×	✓	×	✓
LIMA (Zhou et al., 2024a)	1.3k	✓	×	✓	×	✓
Magpie (Air) (Xu et al., 2024b)	300k, 3M	✓	✓	✓	✓	✓
Magpie (Pro) (Xu et al., 2024b)	300k, 1M	✓	✓	✓	✓	✓
M2Lingual (Maheshwary et al., 2024)	174k	✓	×	✓	✓	×
Natural Questions (Kwiatkowski et al., 2019)	323k	✓	×	✓	×	✓
Oasst1 (Köpf et al., 2024)	88.8k	✓	✓	✓	×	✓
Okapi (Lai et al., 2023)	4.3M*	✓	✓	✓	✓	×
P3 (Sanh et al., 2021)	122M	✓	×	✓	×	✓
Preference 700K [△]	700K	×	✓	UNK	UNK	UNK
Prometheus2 (Kim et al., 2024)	200k	✓	✓	✓	✓	×
Prosocial-Dialog (Kim et al., 2022)	165.4k	✓	✓	✓	×	✓
RLHFFlow Pair Preference [△]	700k	×	✓	✓	✓	✓
Self-instruct (Wang et al., 2023b)	197k	✓	×	✓	×	✓
ShareGPT	Multiple Versions	✓	✓	✓	✓	✓
StackExchange [△]	10.8M	✓	✓	✓	×	✓
Super-Natural Instructions (Wang et al., 2022)	5M	✓	✓	✓	✓	✓
UltraChat (Ding et al., 2023)	1.5M	✓	×	×	✓	✓
UltraFeedback (Cui et al., 2023)	64k	✓	✓	✓	✓	✓
WildChat (Zhao et al., 2024c)	652k	✓	✓	✓	✓	✓
WizardLM (Xu et al., 2023)	250k	✓	✓	✓	✓	✓
xP3 (Muennighoff et al., 2023)	78.8M	✓	×	✓	×	✓

Table 4: SFT and alignment text datasets. [†]The dataset is updated over the time and the number placed on the table is from the latest dataset released by the authors. [‡]The exact size is unknown and some the datasets are no longer accessible. *The estimated number of translated and English instructions.

Dataset	# Samples or (# Tokens) or [Byte Size]	Usecase		Data Source		Annotation
		SFT	Alignment	Human	LLM	Human
ImageRewardDB (Xu et al., 2024b)	137k+	✓	✓	✓	×	✓
Pick-a-pic (Kirstain et al., 2023)	500k+	×	✓	✓	✓	✓
RichHF-18K (Liang et al., 2024)	18k	×	✓	✓	×	✓

Table 5: SFT and alignment vision datasets. [†]The dataset is updated over the time and the number placed on the table is from the latest dataset released by the authors. [‡]The exact size is unknown and some the datasets are no longer accessible. *The estimated number of translated and English instructions.

that models trained only consist of 1,000 carefully human curated prompts and responses, without any reinforcement learning or human preference modeling can outperform models with much larger instruction-tuned datasets.

Dataset Collection FLAN collection (Longpre et al., 2023) is introduced to train a collection of tasks on top of T5 and PaLM models (Raffel et al., 2020). For training multilingual LMs, Cendol Collection (Cahyawijaya et al., 2024), ROOTS (Laurencon et al., 2022), and xP3 (Muennighoff et al., 2023) are used in SFT. Other potential datasets are crowd-sourcing datasets, although they are designed for SFT, but they can be useful resources for SFT, such as NusaCrowd (Cahyawijaya et al., 2023) and SEACrowd (Lovenia et al., 2024).

3.2.2 HUMAN PREFERENCE ALIGNMENT DATASETS

The human alignment data can be in the form of pair-wise or ranking format. We can have a set of preferred and dispreferred data $\mathcal{D}_{\text{pref}}$ for each input sample. For pairwise dataset, we collect pairs of preferred response y_w and dispreferred response y_l . In case of multiple responses, we can gather responses y_0, y_1, y_2, \dots and ask humans to pick the best y_i from each. These datasets have been used to train reward models.

Conversational Datasets Several existing conversational datasets are instrumental in evaluating the quality of dialogue system or chatbot responses. Notable examples include HelpSteer2 (Wang et al., 2024) and UltraFeedback (Cui et al., 2023). HelpSteer2 provides alignment scores across five different aspects—helpfulness, correctness, coherence, complexity, and verbosity—collected from human evaluators. UltraFeedback offers alignment scores for four aspects: instruction-following, truthfulness, honesty, and helpfulness. Additionally, HH-RLHF (Bai et al., 2022a) introduces datasets labeled with scores for helpfulness and harmlessness.

Code Datasets CodeUltraFeedback comprises 10,000 coding instructions, each annotated with four responses generated by a diverse pool of 14 LLMs (Weyssow et al., 2024). These responses are ranked based on five distinct coding preferences: instruction-following, complexity, style, readability, and another instance of instruction-following. The rankings are determined using GPT-3.5 as a judge, providing both numerical scores and detailed textual feedback.

3.3 Pre-trained Generative Models

We categorize pre-trained generative models into three main types: LMs, VLMs, and SLMs. Additionally, we classify these models based on their accessibility: **(1) Open Source:** The model and data are open and accessible, **(2) Open-Weight:** Only the model is accessible and some or all data are inaccessible, **(3) Close-weight and Close-source:** The model is a black-box and may only be accessible by API or service, and **(4) Close Access:** The model is inaccessible. We also categorize these models based on the datasets used for pre-training, specifically noting whether they are trained with Supervised Fine-Tuning (SFT) datasets or Human Preference Tuning datasets.

3.3.1 LANGUAGE MODELS (LMs)

Table 6 shows the list of LMs categorized by the model accessibility and annotated with the model sizes, languages, model base, and fine-tuning methods applied to the model.

Model	Sizes	SFT/Pref. Tuning Langs. [†]	Model Base	SFT	Pref. Tuning
Open-source LM					
Aya-23 (Aryabumi et al., 2024)	8B, 35B	Multi. (23)	Dec-Only; Command R	✓	×
Aya-101 (Üstün et al., 2024)	13B	Multi. (101)	Enc-Dec; mT5	✓	×
Bactrian-X (Li et al., 2023a)	7B	Multi. (52)	Dec-Only; Llama1	✓	×
BART (Lewis et al., 2020)	139M, 406M	English	Enc-Dec	×	×
BLOOM (Le Scao et al., 2023)	560M, 1.1B, 1.7B, 3B, 7.1B, 176B	Multi. (46) + Code (13)	Dec-Only	✓	×
BLOOMZ (Muennighoff et al., 2023)	560M, 1.1B, 1.7B, 3B, 7.1B, 176B	Multi. (108) + Code (13)	Dec-Only; BLOOM	✓	×
Cendol (Cahyawijaya et al., 2024)	7B, 13B	Multi. (10)	Dec-Only; Llama2	✓	×
FLAN-T5 (Longpre et al., 2023)	300M, 580M, 1.2B, 3.7B, 13B	Multi. (10)	Enc-Dec; mT5	✓	×
Llama1 (Touvron et al., 2023a)	80M, 250M, 780M, 3B, 11B	English	Enc-Dec; T5	✓	×
M2M-100 (Fan et al., 2021)	6.7B, 13B, 32.5B, 65.2B	English	Dec-Only	×	×
mBART (Liu, 2020)	418M, 1.2B, 12B	Multi. (100)	Enc-Dec	✓	×
Megatron-LM (Shoeybi et al., 2019)	406M	Multi. (25), Multi. (50)	Enc-Dec	✓	×
MPT (Instruct/Chat) [△]	1.2B, 2.5B, 4.2B, 8.3B	English	Dec-Only; GPT-2	×	×
mT0 (Muennighoff et al., 2023)	7B, 30B	English	Dec-Only	✓	✓
OLMo (Groeneveld et al., 2024)	560M, 1B7, 3B, 7B1	Multi. (108) + Code (13)	Enc-Dec; mT5;	✓	×
OPT (Zhang et al., 2022)	1B, 7B	English + Code	Dec-Only	×	×
Phi1 (Gunasekar et al., 2023)	125M, 350M, 1.3B, 2.7B, 6.7B, 13B, 30B, 66B, 175B	English	Dec-Only; Megatron-LM	×	×
Phi1.5 (Li et al., 2023a)	1.3B	English	Dec-Only	×	×
Pythia (Biderman et al., 2023)	70M, 160M, 410M, 1B, 1.4B, 2.8B, 6.9B, 12B	English	Decoder-Only; GPT-NeoX	×	×
SantaCoder (Allal et al., 2023)	1.1B	Code (3)	Dec-Only	✓	×
StarCoder (Li et al., 2023d)	15.5B	Code (80+)	Dec-Only	✓	×
T0 (Sanh et al., 2021)	3B, 11B	English	Enc-Dec; T5	✓	×
T5 (Raffel et al., 2020)	80M, 250M, 780M, 3B, 11B	English	Enc-Dec	✓	×
T5v1.1 (Raffel et al., 2020; Shazeer, 2020)	80M, 250M, 780M, 3B, 11B	English	Enc-Dec	×	×
WizardCoder (Luo et al., 2023)	7B, 13B, 15B, 33B	Code	Dec-Only	✓	×
Open-weight LM					
Alpaca [△]	7B	English	Dec-Only; Llama1	✓	×
C4AI Command-R (incl. Plus) [△]	35B, 104B	Multi. (13)	Dec-Only	✓	×
DBRX [△]	132B	Multi. (UNK) + Code	MoE	✓	✓
DeepSeek-V2 [△]	16B, 236B	Multi. (UNK) + Code	MoE	✓	✓
Falcon (Almazrouei et al., 2023)	7B, 40B, 180B	Multi. (2) + Code	Dec-Only	✓	×
Falcon2 [△]	11B	Multi. (11) + Code	Dec-Only	✓	×
Gemma (Team et al., 2024)	2B, 7B	Multi. (UNK) + Code	Dec-Only	✓	✓
Gemma2 [△]	9B, 27B	Multi. (UNK) + Code	Dec-Only	✓	✓
Llama2 (Touvron et al., 2023b)	7B, 13B, 70B	Multi. (UNK) + Code	Dec-Only	✓	✓
Llama3, Llama3.1 (Dubey et al., 2024) [△]	8B, 70B	Multi. (UNK) + Code	Dec-Only	✓	✓
LlaMAX (Lu et al., 2024)	7B, 8B	Multi. (102)	Dec-Only; Llama2, Llama3	✓	×
Mistral (Jiang et al., 2023)	7B	Multi. (UNK) + Code	Dec-Only	✓	✓
Mixtral-MoE (Jiang et al., 2024)	8×7B, 8×22B	Multi. (UNK) + Code	MoE; Mistral	✓	✓
Nemotron-4 (15B) (Parmar et al., 2024)	15B	Multi. (53) + Code (43)	Dec-Only	×	×
Nemotron-4 (340B) (Adler et al., 2024)	340B	Multi. (53) + Code (43)	Dec-Only; Nemotron-4 (15B)	✓	✓
NLLB (Costa-jussà et al., 2022)	600M, 1.3B, 3.3B, 54.5B (MoE)	Multi. (200+)	Enc-Dec; M2M-100, MoE	✓	×
Phi3 (Abdin et al., 2024)	3.8B, 7B, 14B	Multi. (UNK) + Code	Dec-Only	✓	✓
Qwen (Bai et al., 2023)	1.8B, 7B, 14B, 72B	Multi. (100) + Code	Dec-Only	✓	✓
Snowflake Arctic [△]	128 × 3.66B	Multi. (UNK) + Code	MoE	✓	✓
StableLM 2 (1.6B) (Bellagente et al., 2024)	1.6B	Multi. (7) + Code	Dec-Only; Vicuna	✓	✓
StableVicuna [△]	13B	English	Dec-Only; Vicuna	✓	✓
Vicuna (Chiang et al., 2023)	7B, 13B	English	Dec-Only; Llama1, Llama2	✓	×
Close-weight and Close-source LM					
Bard (Manyika & Hsiao, 2023)	UNK	UNK	UNK	✓	✓
Chinchilla (Hoffmann et al., 2022)	70B	English + Code	Dec-Only	×	×
Claude 3.5 Sonnet (Anthropic, 2024)	UNK	UNK	UNK	✓	✓
Command R (Plus) [△]	UNK	UNK	UNK	✓	✓
Gemini 1.0 (Team et al., 2023)	UNK	UNK	Dec-Only	✓	✓
Gemini 1.5 (Reid et al., 2024)	UNK	UNK	MoE; Gemini 1.0	✓	✓
Gopher (Rae et al., 2021)	280B	English + Code	Dec-Only	×	×
GPT-3 (Brown et al., 2020)	125M, ..., 175B	Multi. (UNK)	Dec-Only; GPT-2	×	×
GPT-3.5 (Instruct GPT) (Ouyang et al., 2022)	1.3B	UNK	Enc-Dec; GPT-3	✓	✓
GPT-4 (Achiam et al., 2023)	UNK	Multi. (UNK)	UNK	✓	✓
Reka (Ormazabal et al., 2024)	7B (Edge), 21B (Flash), UNK (Core)	Multi. (110)	Enc-Dec	✓	✓
Close-access LM					
AlexaTM (Soltan et al., 2022)	20B	Multi. (12)	Enc-Dec; BART	×	×
BloombergGPT (Wu et al., 2023b)	50.6B	English	Dec-Only; BLOOM	×	×
FLAN-PaLM (Longpre et al., 2023)	8B, 62B, 540B	Multi. (124+) + Code (24+)	UNK	✓	×
PaLM (Chowdhery et al., 2023)	8B, 62B, 540B	Multi. (124) + Code (24)	Dec-Only	×	×
PaLM2 (Anil et al., 2023)	400M, ..., 15B	Multi. (124+) + Code (24+)	UNK	✓	×

Table 6: Pre-trained Generative Language Models. [†]The languages do not include the languages seen by the base model.

3.3.2 SPEECH LANGUAGE MODELS (SLMs)

Table 7 shows the list of open-weight and open-source Speech Language Models (SLMs) categorized by the datasets and methods used in training.

Model	Sizes	SFT/Pref. Tuning Langs. [†]	Model Base	SFT	Pref. Tuning
Open-weight SLM					
BAT (Zheng et al., 2024d)	7B	English	Enc-Dec	✓	×
SpeechGPT (Zhang et al., 2023)	13B	English	Dec	✓	✓
Open-source SLM					
Close-weight and Close-source SLM					
Reka (Ormazabal et al., 2024)	7B (Edge), 21B (Flash), UNK (Core)	Multi. (110)	Enc-Dec	✓	✓

Table 7: Pre-trained Speech Language Models. [†]The languages do not include the languages seen by the base model.

Model	Sizes	SFT/Pref. Tuning Langs. [†]	Model Base	SFT	Pref. Tuning
Open-weight VLM					
Falcon 2 VLM	11B [△]	Multi. (11)	Enc-Dec	✓	×
InstructBLIP (Dai et al., 2023b)	7B, 13B (Vicuna)	English	Enc-Dec	✓	×
	3B, 11B (FLAN-T5)	English	Enc-Dec	✓	×
InstructPix2Pix (Brooks et al., 2023)	UNK	English	UNK	✓	×
LLaVA 1.5 (Liu et al., 2024a)	7B, 13B	English	Enc-Dec	✓	×
LLaVA 1.6 (NeXT)	UNK [△]	English	Enc-Dec	✓	×
X-instructblip (Panagopoulou et al., 2023)	7B, 13B	English	Enc-Dec	✓	×
Phi3-Vision (Abdin et al., 2024)	4.2B	English	Enc-Dec	✓	×
Otter (Li et al., 2023)	7B (Dec)	English	Enc-Dec	✓	×
MultiModal-GPT (Gong et al., 2023)	UNK	English	Enc-Dec	✓	×
Stable Diffusion v1.5 (Rombach et al., 2022)	UNK	English	Enc-Dec	✓	×
Video-LLaMA (Zhang et al., 2023)	7B, 13B	English	Dec-Only	✓	×
Open-source VLM					
Close-weight and Close-source SLM					
Reka (Ormazabal et al., 2024)	7B (Edge), 21B (Flash), UNK (Core)	Multi. (110)	Enc-Dec	✓	✓
SORA (Liu et al., 2024)	UNK	UNK	Enc-Dec	✓	✓

Table 8: Pre-trained Vision Language Models. [†]The languages do not include the languages seen by the base model.

3.3.3 VISION LANGUAGE MODELS (VLMs)

Table 8 shows the list of open-weight and open-source Vision Language Models (VLMs) categorized by the datasets and methods used in training.

4. Online Alignment

In this section, we explore into human preference tuning using online methods, where data is continuously sampled. Online preference tuning involves real-time model updates as new data becomes available, enabling the model to dynamically adapt to evolving preferences and new information. This approach allows the alignment process to incorporate new data as it arrives and benefit from online exploration. We discuss the mechanisms of data collection, processing, and real-time model updates, emphasizing the benefits of managing non-stationary environments and enhancing model performance through continuous learning. Various techniques and strategies for implementing especially on-policy tuning are examined to provide a comprehensive understanding of its effective application in human preference tuning. We cover standard RL-based methods (e.g., PPO, which is online and on-policy), online DPO and SFT like algorithms (which can be on-policy or off-policy) and Nash Learning (or self-play) based algorithms.

4.1 Reinforcement Learning Human Feedback (RLHF)

In general, RLHF learns a reward function from human feedback and then optimize that reward function (Christiano et al., 2017). The training for RLHF involves three stages:

- The policy model π_θ interacts with the environment and the parameters of π_θ are updated via RL.
- The pairs of segments are selected from the output produced by the policy model π_θ , and send them to human annotators for comparison.
- The parameters are optimized using reward r to fit the comparisons collected from human.

According to Ziegler, Stiennon, Wu, Brown, Radford, Amodei, Christiano, and Irving (2019), the RLHF pipeline for LMs can be summarized as following:

- **Supervised Fine-Tuning:** A pre-trained LM is instruction-tuned using a dataset consisting of a given instruction prompt, and (typically) a human-written completion. The LM/policy is trained with a cross-entropy loss over the completion only. Often, the SFT model, denoted as `sft` is used to initialize both the reward model and the RLHF policy.
- **Reward Modeling:** RLHF leverages a reward model r_ϕ trained using a dataset of preferences \mathcal{D} . The reward model is trained using the following loss:

$$\text{loss}(r) = \mathbb{E}_{(x, \{y_i\}_i, b) \sim S} \left[\log \frac{e^{r(x, y_b)}}{\sum_i e^{r(x, y_i)}} \right]. \quad (6)$$

or, for pairwise preferences,

$$\mathcal{L}_{\text{RM}}(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_{\text{pref}}} \log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l)). \quad (7)$$

- **Reinforcement Learning:** In this stage, the learned reward model r_{ϕ^*} is used to provide online feedback in the optimization of the policy. In Ziegler et al. (2019), Stiennon et al. (2020), Ouyang et al. (2022), RLHF further maximizes average reward with an extra KL regularization term, i.e.:

$$\mathcal{L}_{\text{RL}}(\phi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(\cdot | x)} [r_{\phi^*}(x, y) - \beta_{\text{reg}} \text{KL}(\pi(\cdot | x) | \pi_{\text{ref}}(\cdot | x))], \quad (8)$$

where $\beta_{\text{reg}} > 0$ is a hyper-parameter controlling the deviation from the reference policy $\pi_{\text{ref}} = \pi^{\text{SFT}}$.

Notably, the reward model is trained by assumption on the Bradley-Terry Model, which leverages pairwise preference datasets — i.e. pairs of preferred and non-preferred responses. There are various methods and variations for training RLHF, primarily categorized into two main approaches: PPO style (with critic) and REINFORCE style (without critic). In the following sections, we will describe these methods in detail.

4.1.1 PROXIMAL POLICY OPTIMIZATION (PPO)

Initially in the original RLHF paper (Ziegler et al., 2019), they use PPO (Schulman et al., 2017) as their optimization strategy. PPO framework is a method for the human preference signals from external reward models with RLHF. The idea is to improve the current state of affairs by introducing an algorithm that attains the data efficiency and reliable performance of TRPO, while using only first-order optimization with a simpler clipped surrogate objective, omitting the expensive second-order optimization presented in TRPO using stochastic gradient ascent. Whereas standard policy gradient methods perform one gradient update per data sample, PPO (Schulman et al., 2017) proposes a novel objective function that enables multiple epochs of minibatch updates. It has some of the benefits of TRPO, but they are much simpler to implement and more efficient. For the optimization, KL-shaped reward (Ahmadian et al., 2024a) is useful as penalty-free optimization of the reward model leads to degradation in the coherence of the model. Optimizing this objective is equivalent to maximizing the following KL-shaped reward in expectation. There are many variants of PPO for RLHF, e.g. P3O (Wu et al., 2023c), or RLHF-V (Yu et al., 2024) for multi-modal models.

Pairwise Proximal Policy Optimization (P3O) P3O (Wu et al., 2023c) is an on-policy RL algorithms that interleaves off-policy updates with on-policy updates. P3O uses the effective sample size between the behavior policy and the target policy to control how far they can be from each other and does not introduce any additional hyper-parameters.

RLHF-V RLHF-V (Yu et al., 2024) enhances MLLM trustworthiness via behavior alignment from fine-grained correctional human feedback. Specifically, RLHF-V collects human preference in the form of segment-level corrections on hallucinations, and performs dense direct preference optimization over the human feedback.

4.1.2 REINFORCE

One notable issue of PPO is its hyperparameter sensitivity and high demand of GPU memory when training the models. Thus many works revisit the REINFORCE (Williams, 1987, 1992) style policy gradient (Sutton et al., 1999) methods to avoid the computation burden raised by training the critic network.

ReMax ReMax (Li et al., 2023c) builds upon the well-known REINFORCE algorithm leveraging three key properties of RLHF: fast simulation, deterministic transitions, and trajectory-level rewards. The name “ReMax” reflects its foundation in REINFORCE and its use of the argmax operator. ReMax modifies the gradient estimation by incorporating a subtractive baseline value as following:

$$\tilde{g}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T [\nabla_{\theta} \log \pi_{\theta}(a_t | x, a_{1:t-1}) \times (r(x^i, a_{1:T}^i) - b_{\theta}(x^i))], \quad (9)$$

where N is the number of prompts, $b_{\theta}(x^i)$ is a baseline value for which ReMax (Li et al., 2023c) chose as

$$b_{\theta}(x^i) = r(x^i, \bar{a}_{1:T}^i), \bar{a}_t^i \in \operatorname{argmax}_{\cdot} \pi_{\theta}(\cdot | x^i, \bar{a}_{1:t-1}^i). \quad (10)$$

This baseline value can be calculated by greedily sampling a response and computing the associated reward value.

REINFORCE Leave One-Out (RLOO) RLOO (Ahmadian et al., 2024a) extends the REINFORCE algorithm by leveraging multiple online samples to achieve unbiased variance reduction. It improves upon REINFORCE in two key ways: (1) The rewards from each sample can serve as a baseline for all other samples, and (2) Policy updates are performed using the average of gradient estimates from each sample, resulting in a variance-reduced multi-sample Monte Carlo (MC) estimate. This is the intuition behind the RLOO estimator, as following:

$$\frac{1}{k} \sum_{i=1}^k \nabla \log \pi(y_{(i)}|x) [R(y_{(i)}, x) - \frac{1}{k-1} \sum_{j \neq i} R(y_{(j)}, x)], \text{ for } y_{(1)}, \dots, y_{(k)} \stackrel{i.i.d.}{\sim} \pi_{\theta}(\cdot|x), \quad (11)$$

where k refers to the number of online samples generated, RLOO_k considers each $y_{(i)}$ individually and uses the remaining $k-1$ samples to create an unbiased estimate of the expected return for the prompt. This approach functions similarly to a parameter-free value function, but it is estimated at each training step.

Group Relative Policy Optimization (GRPO) GRPO proposes in Shao et al. (2024), is also a REINFORCE style algorithm and has been proven to enhance the reasoning capability of LLMs e.g. in training DeepSeek Math and DeepSeek R1 (Guo et al., 2025). For the outcome supervision RL with GRPO, GRPO samples a group of outputs $\{o_1, o_2, \dots, o_G\}$ from the old policy $\pi_{\theta_{old}}$ for each question q and then optimizes the policy model by maximizing the following objective $\mathcal{J}_{GRPO}(\theta) =$:

$$\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[\frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t} | q, o_{i,<t})} \hat{A}_{i,t}, \text{clip} \left(\frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t} | q, o_{i,<t})}, 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{KL}[\pi_{\theta} || \pi_{ref}] \right\},$$

where ε and β are hyper-parameters, and $\hat{A}_{i,t}$ is the advantage calculated based on relative rewards of the outputs inside each group only, which was chosen as:

$$\frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}. \quad (12)$$

Following up work Dr. GRPO (Liu et al., 2025) argue that the $\frac{1}{|o_i|}$ discounting in GRPO and advantage normalization by the variance $\text{std}(r)$ in (12) lead to biased policy gradient estimator and are also empirically unnecessary. Liu et al. (2025) propose the variant named Dr. GRPO by removing these two parts.

4.2 Online Directed Preference Optimization (Online DPO)

4.2.1 ONLINE AI FEEDBACK (OAIF)

OAIF (Guo et al., 2024a) employ a LLM as an annotator during each training iteration. In this process, two responses are sampled from the current model, and the LLM annotator

is prompted to select the preferred response, thereby providing real-time feedback. OAIF aims to gather preferences dynamically for responses generated by the language model being aligned. Given the prohibitive cost of using human feedback, this method leverages an LLM as an online annotator to collect preferences over pairs of responses sampled from the model π_θ during its alignment process. The objective for online DPO yields (please see detailed derivation of DPO in Section 5.1):

$$\mathcal{L}_{\text{OAIF}}(\pi_\theta; \pi_{\text{ref}}) := -\mathbb{E}_{x \sim \mathcal{D}, (y_w, y_l) \sim \pi_{\theta_-}} \left[\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right], \quad (13)$$

in which we note π_{θ_-} to show that preference pairs are generated under π_θ , but we further adopt a stop gradient to prevent it from getting into the loss objective for the gradient computation. The OAIF is illustrated in Algorithm 1 (OAIF algorithm in Guo et al. (2024a)), in which function ℓ can be log-sigmoid (DPO), square (IPO), or ReLU (SLiC) functions.

Algorithm 1 Online AI Feedback (OAIF) for Direct Alignment from Preference (DAP)

- 1: **Input:** Prompt dataset $\mathcal{D}_x = \{x_i\}_{i=1}^N$, an LLM annotator, SFT model π_{θ_0}
 - 2: **for** $t := 0$ to T **do**
 - 3: Sample prompt $x \sim \mathcal{D}_x$
 - 4: Sample response pair $y_1, y_2 \sim \pi_{\theta^t}(\cdot|x)$
 - 5: Use LLM annotator to get preference pair y_w, y_l
 - 6: Update θ^t into θ^{t+1} using $\nabla_{\theta} \ell(x, y_w, y_l, \theta^t)$
 - 7: **end for**
-

4.2.2 ITERATIVE DIRECTED PREFERENCE OPTIMIZATION

Iterative DPO (Xu et al., 2023; Xiong et al., 2024) has been proposed to narrow the gap between the performance offline preference optimization methods like DPO and online methods like RLHF, as RLHF still outperforms offline DPO. Different from DPO that used a fixed offline dataset, iterative DPO proposed to formulate the preference datasets by the generations of the current model and labelers, being either a pretrained reward model or LLM as a judge or the model to be trained itself through specific prompting (Yuan et al., 2024b), thus this pipeline usually appears at the same time with self-rewarding (Yuan et al., 2024b) methods (some paper will even call self-rewarding as iterative DPO methods). For each iteration, if the batch size for preference datasets utilized for policy optimization is only 1, then iterative DPO is essentially the same as online DPO or OAIF, except that the reference policy may be chosen as the last iterated policy instead of always being the SFT policy; otherwise iterative DPO is a hybrid method which combines offline learning in loss function optimization and online sampling in preference data generation. The reference model in the loss objective may differ between different methods, can be fixed SFT model Xiong et al. (2024) or last iterated model (Xu et al., 2023; Yuan et al., 2024b) or some mixtures.

4.2.3 ONLINE PREFERENCE TUNING (OPTUNE)

OPTune (Chen et al., 2024d) is an algorithm for efficient data generation in online RLHF. It improves both generation and training efficiency by selectively regenerating only the lowest-rewarded responses and employing a weighted DPO objective that prioritizes pairs with larger reward gaps. This approach significantly enhances the overall efficiency of the RLHF pipeline, setting the stage for the development of preference-aligned LLMs in a resource-efficient manner. The method enhances both data generation and training efficiency for online preference alignment. To minimize the cost of iterative data regeneration, it employs a straightforward yet effective reward-based prompt selection strategy, updating responses only for prompts with the lowest scores according to the reward model. Additionally, recognizing that converting scalar rewards to binary labels for the online DPO objective results in information loss, the method introduces a weighted DPO loss variant. This variant prioritizes learning from response pairs with larger reward gaps, further boosting online learning efficiency.

4.3 SFT-like

4.3.1 RANK RESPONSES TO ALIGN HUMAN FEEDBACK (RRHF)

RRHF (Yuan et al., 2023) is a method that evaluates sampled responses from various sources using the logarithm of conditional probabilities and aligns these probabilities with human preferences through ranking loss. This approach can utilize responses from multiple origins, including the model’s own outputs, responses from other large language models, and human expert responses, to learn how to rank them effectively. The primary objective is to simplify the complex hyper-parameter tuning and extensive training resources required by PPO. Before training, RRHF samples responses from diverse sources, which can include model-generated responses from the model itself as well as pre-existing human-authored responses of varying quality. During training, RRHF scores these responses based on the log probability provided by the training language model. These scores are then aligned with human preference rankings or labels using ranking loss, ensuring that the model’s outputs are better aligned with human preferences.

4.3.2 REWARD RANKED FINE TUNING (RAFT)

RAFT (Dong et al., 2023) is the combination of ranking samples by rewards and SFT, which iteratively alternates among three steps: 1) The batch is sampled from the generative models; 2) The reward function is used to score the samples and filter them to get a filtered subset of high rewards; and 3) fine-tune the generative models on the filtered subset.

4.3.3 REINFORCED SELF-TRAINING (ReST)

ReST (Gulcehre et al., 2023) is an RLHF algorithm aimed at aligning an LM’s outputs with human preferences. It uses a learned reward function to model human preferences over sequences. In the Markov decision process underlying conditional language modeling, states represent partial sequences, and actions correspond to generated tokens. ReST divides the typical reinforcement learning pipeline into distinct offline stages for dataset growth and policy improvement. Initially, it fine-tunes a model to map input sequences to output

sequences using a dataset of sequence pairs, optimizing with Negative Log-Likelihood (NLL) loss. Then, it creates a new dataset by augmenting the initial training dataset with samples generated by the model. In this phase, conditioning inputs are resampled from the original dataset, similar to self-training, but direct sampling is possible if accessible.

4.3.4 SUPERVISED ITERATIVE LEARNING FROM HUMAN FEEDBACK (SUPERHF)

SuperHF (Mukobi et al., 2023) is an alignment algorithm that enhances data efficiency using a reward model and replaces PPO with a straightforward supervised fine-tuning loss. The core concept involves the language model generating its own training data by sampling a “superbatch” of outputs, filtering these through a reward model, and iteratively fine-tuning on each filtered completion. This method builds upon and unifies previous research by integrating two crucial components: (1) the Kullback-Leibler (KL) divergence penalty and (2) an iterative process of sampling and fine-tuning. Additionally, SuperHF is embedded within a Bayesian inference framework, demonstrating that both RLHF and SuperHF can be understood from a unified theoretical perspective that does not rely on reinforcement learning. This perspective naturally justifies the use of the KL penalty and the iterative approach.

4.4 Nash Learning

4.4.1 NASH LEARNING FROM HUMAN FEEDBACK (NLHF)

NLHF (Munos et al., 2023) is motivated to address the limitation of reward models (or essentially the Elo ratings) to represent the richness of human preferences as in RLHF. Instead of targeting at maximizing the (regularized) reward, NLHF takes the preference model as the ‘first class citizen’, and pursue ‘a policy that consistently generates responses preferred over those generated by any competing policy’. Thus this policy is the Nash equilibrium of this preference model, the reason the method is named NLHF. Concretely, the (regularized) preference model for two policies π, π' is defined as:

$$\mathcal{P}(\pi > \pi') := \mathbb{E}_{x \sim \rho} \mathbb{E}_{y \sim \pi(\cdot|x), y' \sim \pi'(\cdot|x)} \left[\mathcal{P}(y > y' | x) - \beta_{\text{reg}} \log \frac{\pi(y | x)}{\mu(y | x)} + \beta_{\text{reg}} \log \frac{\pi'(y' | x)}{\mu(y' | x)} \right], \quad (14)$$

and NLHF searches the Nash Equilibrium such that (denote μ as π_{ref} for simplicity here):

$$\pi^* := \arg \max_{\pi} \min_{\pi'} \mathcal{P}(\pi > \pi') - \beta_{\text{reg}} \text{KL}_{\rho}(\pi, \mu) + \beta_{\text{reg}} \text{KL}_{\rho}(\pi', \mu). \quad (15)$$

For optimization, the Nash-MD algorithm proposed in NLHF used a geometric mixture between the current policy π_t and the reference policy μ as the competing policy in the place of π' :

$$\pi_t^{\mu}(y) := \frac{\pi_t(y)^{1-\eta\beta_{\text{reg}}} \mu(y)^{\eta\beta_{\text{reg}}}}{\sum_{y'} \pi_t(y')^{1-\eta\beta_{\text{reg}}} \mu(y')^{\eta\beta_{\text{reg}}}}, \quad (16)$$

where η is a learning rate, and Nash-MD algorithm is a step of mirror descent relative to the regularized policy π_t^{μ} :

$$\pi_{t+1} := \arg \max_{\pi} [\eta \mathcal{P}(\pi > \pi_t^{\mu}) - \text{KL}(\pi, \pi_t^{\mu})], \quad (17)$$

which yields a closed-form solution that:

$$\log \pi_{t+1}(y) = [(1 - \eta\beta_{\text{reg}}) \log \pi_t(y) + \eta\beta_{\text{reg}} \log \mu(y)] + \eta\mathcal{P}(y > \pi_t^\mu) + c, \quad (18)$$

where c is a normalization constant which is independent of y and the algorithm is proved to converge of rate $\frac{1}{T}$ under the tabular setting. For practical concern, when policy is a deep neural network beyond tabular setting, NLHF further proposes Nash-MD-PG motivated by Nash-MD, and the algorithm updates the policy with policy gradient:

$$\nabla_{\theta} \mathcal{P}_{\tau} \left(\pi_{\theta} > \pi'_{\theta_-} \right) = \mathbb{E}_{x \sim \rho, y \sim \pi_{\theta}(\cdot|x), y' \sim \pi'(\cdot|x)} \left[\hat{g}(x, y, y') \right], \quad (19)$$

where π'_{θ_-} denotes a stop-gradient on π'_{θ} with π'_{θ} being a geometric mixture

$$\log \pi'_{\theta}(y | x) := (1 - \lambda) \log (\pi_{\theta}(y | x)) + \lambda \log (\mu(y | x)) + c(x), \quad (20)$$

in which λ is a mixing constant and

$$\hat{g}(x, y, y') := \nabla_{\theta} \log \pi_{\theta}(y | x) \left(\mathcal{P}(y > y' | x) - 1/2 - \beta_{\text{reg}} \text{KL}(\pi_{\theta}(\cdot | x), \mu(\cdot | x)) \right), \quad (21)$$

respectively. NLHF also argues that, Nash equilibrium of the preference model is a solution that better aligns with the diversity of human preferences.

4.4.2 SELF-PLAY PREFERENCE OPTIMIZATION (SPPO)

SPPO (Wu et al., 2024) can be understood as a specific instance of NLHF by taking $\lambda = 0$, i.e., the reference policy is itself. The algorithm can be found in Algorithm 2, given an LLM judge:

Algorithm 2 Self-Play Preference Optimization (SPPO)

- 1: **Input:** base policy π_{θ_0} , preference oracle \mathcal{P} , learning rate η , number of generated samples K
- 2: **for** $t = 0, 1, \dots$ **do**
- 3: Generate synthetic responses by sampling $x \sim \mathcal{D}$ and $y_{1:K} \sim \pi_{\theta_t}(\cdot|x)$
- 4: Annotate the win-rate $\mathcal{P}(y_k \succ y_{k'}|x), \forall k, k' \in [K]$
- 5: Select responses from $y_{1:K}$ to form dataset $D_t = \{(x_i, y_i, \hat{\mathcal{P}}(y_i \succ \pi_{\theta_t}|x_i))\}_{i \in [N]}$
- 6: Optimize $\pi_{\theta_{t+1}}$ according to:

$$\theta_{t+1} \leftarrow \arg \min_{\theta} \mathbb{E}_{(x, y, \hat{\mathcal{P}}(y \succ \pi_{\theta_t}|x)) \sim D_t} \left(\log \left(\frac{\pi_{\theta}(y|x)}{\pi_{\theta_t}(y|x)} \right) - \eta \left(\hat{\mathcal{P}}(y \succ \pi_{\theta_t}|x) - \frac{1}{2} \right)^2 \right).$$

- 7: **end for**
-

4.5 Fine-tuning Diffusion Models

Given the popularity of diffusion based T2I models and its different nature of structural properties to LLM, we have the methods of fine-tuning diffusion models as a separate section of interest. We first briefly review the formulation of text-to-image diffusion generative

models. For a more comprehensive background of diffusion models, we refer the interested readers to existing tutorial and survey papers (Luo, 2022; Cao et al., 2024; Yang et al., 2023; Tang & Zhao, 2024; Chen et al., 2024f; Chan, 2024). DDPM (Sohl-Dickstein et al., 2015; Ho et al., 2020) consider a sequence of positive noise scales $0 < \beta_1, \beta_2, \dots, \beta_N < 1$, and perturb data by gradually adding noise through a stochastic process: for each training data point $x_0 \sim p_{\text{data}}(x)$, a discrete Markov chain $\{x_0, x_1, \dots, x_N\}$ is constructed such that:

$$x_i = \sqrt{1 - \beta_i}x_{i-1} + \sqrt{\beta_i}z_{i-1}, \quad i = 1, \dots, N, \quad (22)$$

where $z_{i-1} \sim \mathcal{N}(0, I)$. For generative modeling, the backward process - a variational Markov chain in the reverse direction - is parameterized with

$$p_\theta(x_{i-1} | x_i) = \mathcal{N}\left(x_{i-1}; \frac{1}{\sqrt{1 - \beta_i}}(x_i + \beta_i s_\theta(i, x_i)), \beta_i I\right), \quad (23)$$

in which $s_\theta(i, x_i)$ is learned by maximizing an evidence lower bound (ELBO). In the context of text-to-image generation, trained s_{θ^*} will also be dependent on an input prompt c for conditional generation, becoming $s_\theta(i, x_i, c)$. For inference process, samples can be generated by starting from pure noise and following the estimated reverse process as:

$$x_{i-1} = \frac{1}{\sqrt{1 - \beta_i}}(x_i + \beta_i s_{\theta^*}(i, x_i, c)) + \sqrt{\beta_i}z_i, \quad i = N, N - 1, \dots, 1. \quad (24)$$

4.5.1 DDPO AND DPOK

We review some key elements in DDPO and DPOK (Black et al., 2024; Fan et al., 2024) to formulate the problem of fine-tuning diffusion models as discrete-time MDPs, and then apply RL algorithms. Note that recent works, Tang (2024), Uehara, Zhao, Biancalani, and Levine (2024a), Uehara, Zhao, Black, Hajiramezanali, Scalia, Diamant, Tseng, Biancalani, and Levine (2024b) extend a continuous-time stochastic control formulation for fine-tuning, which motivates works e.g. Uehara, Zhao, Black, Hajiramezanali, Scalia, Diamant, Tseng, Levine, and Biancalani (2024c), Uehara, Zhao, Hajiramezanali, Scalia, Eraslan, Lal, Levine, and Biancalani (2024d) for more feedback efficient or conservative fine-tuning, Domingo-Enrich, Drozdal, Karrer, and Chen (2024) for a memory-less extension and using adjoint methods for solving the such control problem. Zhao, Chen, Zhang, Yao, and Tang (2025b) propose a continuous-time reinforcement learning framework and a continuous time PPO variant for more robust diffusion models preference learning.

In this paper, however we stick to the discrete time formulation for simplicity and broader community. Consider taking (i, x_i, c) as the state space, and define the action as the next hierarchy x_{i-1} to go to, then Eq. (24) naturally defines a stochastic policy: the stochasticity of the policy comes from $\sqrt{\beta_i}z_i$, thus the policy follows Gaussian with mean determined by $s_{\theta^*}(i, x_i, c)$ with variance β_i :

$$\pi_\theta(x_{i-1} | x_i) \sim \mathcal{N}\left(\frac{1}{\sqrt{1 - \beta_i}}(x_i + \beta_i s_\theta(i, x_i, c)), \beta_i\right), \quad i = N, N - 1, \dots, 1. \quad (25)$$

Given this formulation, Black et al. (2024) directly maximize the expected reward (without regularization) $\mathcal{J}_{\text{DDPO}} = \mathbb{E}_{\theta} [r(x_0, c)]$ by REINFORCE or PPO:

$$\nabla_{\theta} \mathcal{J}_{\text{DDPO}} = \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta} \log p_{\theta}(x_{t-1} | x_t, c) r(x_0, c) \right]. \quad (26)$$

Compare to DDPO, DPOK (Fan et al., 2024) optimize the same regularized reward objective as in Eq. (8):

$$\mathcal{J}_{\text{DPOK}} = \mathbb{E}_{\theta} [r(x_0, c)] - \beta \mathbb{E}_{p(z)} [\text{KL}(p_{\theta}(x_0 | z) \| p_{\text{pre}}(x_0 | z))] \quad (27)$$

They further proposed a clipped gradient algorithm for optimization, motivated by the original PPO clipped objective. In addition, DPOK shows that adding regularization will yield a better generation result compared to the version without regularization.

4.5.2 REWARD FEEDBACK LEARNING (REFL)

ReFL (Xu et al., 2024b) is a supervised fine-tuning method based on its pre-trained reward model ImageReward $r_{\text{IR}}(c, x)$. The objective for ReFL optimization is a linear combination of negative pre-trained loss (for diffusion models) and reward maximization:

$$\mathcal{J}_{\text{ReFL}}(\theta) = \mathcal{J}_{\text{pre}}(\theta) + \lambda \mathbb{E}_{c \sim p_{\mathbf{c}}, x_t \sim p_{\theta}(\cdot | c)} (\phi(r_{\text{IR}}(c, x_t))), \quad (28)$$

in which λ is a scaling constant, ϕ is taken as a ReLU function and $t \in [0, \tilde{T}]$ is a random number for sampling, a technique that Xu et al. (2024b) claim can help stabilize the training instead of always letting t be 0.

4.5.3 DIRECT REWARD FINE-TUNING (DRAFT)

DRaFT (Clark et al., 2023) introduces a straightforward method for fine-tuning diffusion models using differentiable reward functions. The goal is to fine-tune the parameters θ of a pre-trained diffusion model such that images generated by the sampling process maximize a differentiable reward function r :

$$J(\theta) = \mathbb{E}_{\mathbf{c} \sim p_{\mathbf{c}}, \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [r(\text{sample}(\theta, \mathbf{c}, \mathbf{x}_T), \mathbf{c})], \quad (29)$$

where $\text{sample}(\theta, \mathbf{c}, \mathbf{x}_T)$ denotes the sampling process from time $t = T \rightarrow 0$ with context \mathbf{c} . First, DRaFT consider solving Eq. 29 by computing $\nabla_{\theta} r(\text{sample}(\theta, \mathbf{c}, \mathbf{x}_T), \mathbf{c})$ and using gradient ascent. Computing this gradient requires backpropagation through multiple diffusion model calls in the sampling chain, similar to backpropagation through time in a recurrent neural network. To mitigate the memory cost associated with this process, DRaFT employs two strategies: 1) low-rank adaptation (LoRA) (Hu et al., 2021), and 2) gradient checkpointing (Chen et al., 2016).

4.5.4 ALIGNPROP

AlignProp (Prabhudesai et al., 2023) introduces a method that transforms denoising inference within text-to-image diffusion models into a differentiable recurrent policy, effectively

linking conditioning input prompts and sampled noise to generate output images. This approach enables fine-tuning of the denoising model’s weights through end-to-end back-propagation, guided by differentiable reward functions applied to the generated images. The proposed model casts conditional image denoising as a single step MDP with states $\mathcal{S} = \{(x_T, \mathbf{c}), x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\}$, actions are the generated image samples, and the whole DDIM denoising chain corresponds to a differentiable policy that maps states to image samples: $\mathcal{A} = \{x_0 : x_0 \sim \pi_\theta(\cdot | x_T, \mathbf{c}), x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\}$. The reward function is a differentiable function of parameters ϕ that depends only on generated images $R_\phi(x_0), x_0 \in \mathcal{A}$. Given a dataset of prompts input \mathcal{D} , our loss function reads:

$$\mathcal{L}_{\text{align}}(\theta; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{c}^i \in \mathcal{D}} R_\phi(\pi_\theta(x_T, \mathbf{c}^i)). \quad (30)$$

The parameters of the diffusion model using gradient descent on $\mathcal{L}_{\text{align}}$. The policy π is recurrent, and training it is akin to backpropagation through time, a technique commonly used for training recurrent neural networks. The gradient for updating the parameters of the diffusion model with respect to the downstream objective (i.e., the differentiable reward function) is expressed as following:

$$\hat{\nabla}_\theta \mathcal{L}_{\text{align}} = \frac{\partial \mathcal{L}_{\text{align}}}{\partial \theta} + \sum_{t=0}^K \frac{\partial \mathcal{L}_{\text{align}}}{\partial x_t} \cdot \frac{\partial x_t}{\partial \theta}, \quad (31)$$

in which K is uniformly drawn from $[0, T]$ for memory efficiency instead of being T , referred as randomized truncation in Prabhudesai et al. (2023).

4.5.5 PROXIMAL REWARD DIFFERENCE PREDICTION

PRDP (Deng et al., 2024) proposes a loss for matching likelihood difference with reward difference for fine-tuning diffusion models, inspired by DPO. Notice that, (same as derivation in DPO), for any two generations x_0^1 and x_0^2 , the optimal policy (KL-regularized reward) yields:

$$\log \frac{\pi_{\theta^*}(x_0^1 | \mathbf{c})}{\pi_{\text{ref}}(x_0^1 | \mathbf{c})} - \log \frac{\pi_{\theta^*}(x_0^2 | \mathbf{c})}{\pi_{\text{ref}}(x_0^2 | \mathbf{c})} = \frac{r(x_0^1, \mathbf{c}) - r(x_0^2, \mathbf{c})}{\beta_{\text{reg}}} \quad (32)$$

thus PRDP proposes to minimize the MSE error between LHS with θ (replacing θ^*) and RHS. The objective is $\mathcal{L}_{\text{PRDP}}(\pi_\theta; \pi_{\text{ref}}) :=$

$$\mathbb{E}_{\mathbf{c} \sim \mathcal{D}, (x^1, x^2) \sim \pi_\theta(\cdot | \mathbf{c})} \left(\beta_{\text{reg}} \log \frac{\pi_\theta(x_0^1 | x)}{\pi_{\text{ref}}(x_0^1 | x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(x_0^2 | x)}{\pi_{\text{ref}}(x_0^2 | x)} - (r(x_0^1) - r(x_0^2)) \right)^2, \quad (33)$$

Furthermore, they also employ proximal updates (clipping the ratios and optimizing a proximal objective) for stable training of (33), in the same spirit of PPO.

Similar works include Yang, Chen, and Zhou (2024b), which apply the idea of dense reward to DPO-style explicit-reward-free approach on text-to-image diffusion models, so as to suit better to diffusion models’ generation hierarchy.

4.5.6 DIFFUSION LOSS-GUIDED POLICY OPTIMIZATION (DLPO)

DLPO (Chen et al., 2024c) applies online RL to fine-tune TTS diffusion models, where the reward is shaped by the diffusion model’s loss. Incorporating the diffusion model loss into the objective function serves as an additional mechanism to enhance performance and maintain the coherence of the model. The method’s objective is described as following:

$$\mathbb{E}_{c \sim p(c)} \mathbb{E}_{t \sim \mathcal{U}\{1, T\}} \mathbb{E}_{p_\theta(x_{0:T}|c)} [-\alpha r(x_0, c) - \beta \|\tilde{\epsilon}(x_t, t) - \epsilon_\theta(x_t, c, t)\|_2], \quad (34)$$

where α, β are the reward and weights for diffusion model loss, respectively. DLPO uses the following gradient to update the objective:

$$\mathbb{E}_{c \sim p(c)} \mathbb{E}_{t \sim \mathcal{U}\{0, T\}} \mathbb{E}_{p_\theta(x_{1:T}|c)} [- (\alpha r(x_0, c) - \beta \nabla_\theta \|\tilde{\epsilon}(x_t, t) - \epsilon_\theta(x_t, c, t)\|_2) \nabla_\theta \log p_\theta(x_{t-1}|x_t, c)]. \quad (35)$$

The diffusion model objective is incorporated into the reward function as a penalty. This algorithm aligns with the training procedure of TTS diffusion models by integrating the original diffusion model objective $\beta \|\tilde{\epsilon}(x_t, t) - \epsilon_\theta(x_t, c, t)\|_2$ as a penalty in the reward function. This approach effectively prevents model deviation and ensures that the model remains coherent during training.

4.5.7 HUMAN FEEDBACK FOR INSTRUCTIONAL VISUAL EDITING (HIVE)

HIVE (Zhang et al., 2024b) is proposed to improve instruction visual editing models (diffusion models based, e.g., InstructPix2Pix (Brooks et al., 2023)) with human feedback. In instructional supervised training, the stable diffusion model has two conditions $c = [c_I, c_T]$, where c_T is the editing instruction, and c_I is the latent space of the original input image. In the training process, a pre-trained auto-encoder with encoder \mathcal{E} and decoder \mathcal{D} is used to convert between edited image \tilde{x} and its latent representation $z = \mathcal{E}(\tilde{x})$. The diffusion process is composed of an equally weighted sequence of denoising autoencoders $\epsilon_\theta(z_t, t, c)$, $t = 1, \dots, T$, which are trained to predict a denoised variant of their input z_t , which is a noisy version of z . The objective of instructional supervised training is:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(\tilde{x}), c, \epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(z_t, t, c)\|_2^2]. \quad (36)$$

HIVE proposes that optimizing an exponential reward weighted objective for fine-tuning diffusion models:

$$\mathcal{L}_{\text{HIVE}}(\theta) := \mathbb{E}_{\mathcal{E}(\tilde{x}), c, \epsilon \sim \mathcal{N}(0, I), t} [\omega(\tilde{x}, c) \cdot \|\epsilon - \epsilon_\theta(z_t, t, c)\|_2^2], \quad (37)$$

with $\omega(\tilde{x}, c) = \exp(r_\phi(\tilde{x}, c)/\beta)$ being the exponential reward weight for edited image \tilde{x} and condition c , which is motivated by the closed form of optimal solution for RLHF in Eq. (38).

5. Offline Alignment

In this section, we present a detailed explanation for each offline preference tuning method, including SLiC-HF, DPO and its variants. In Table 5, for simplicity, we include representative DPO variants and their final loss objectives. For each DPO variant, we conclude not only the resulting final objective or algorithm, but also both summarize the motivation or the direction the method contributed to for improvement over DPO.

Method	Objective
DPO	$-\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right)$
IPO	$\left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} - \frac{1}{2} \right)^2$
f -DPO	$-\log \sigma \left(\beta_{\text{reg}} f' \left(\frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} \right) - \beta_{\text{reg}} f' \left(\frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right) \right)$
KTO	$-\lambda_w \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - z_{\text{ref}} \right) - \lambda_l \sigma \left(z_{\text{ref}} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right),$ where $z_{\text{ref}} = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\beta_{\text{reg}} \text{KL}(\pi_{\theta}(y x) \pi_{\text{ref}}(y x))]$
ODPO	$-\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} - \Delta_r(x) \right)$
Mallows-DPO	$-\log \sigma \left(\phi(x) \left[\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right] \right)$
R-DPO	$-\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} - (\alpha y_w - \alpha y_l) \right)$
CPO	$-\log p_{\theta}(y_w x) - \log \sigma \left(\beta_{\text{reg}} \log \pi_{\theta}(y_w x) - \beta_{\text{reg}} \log \pi_{\theta}(y_l x) \right)$
ORPO	$-\log p_{\theta}(y_w x) - \lambda \log \sigma \left(\log \frac{p_{\theta}(y_w x)}{1-p_{\theta}(y_w x)} - \log \frac{p_{\theta}(y_l x)}{1-p_{\theta}(y_l x)} \right),$ where $p_{\theta}(y x) = \exp \left(\frac{1}{ y } \log \pi_{\theta}(y x) \right)$
SimPO	$-\log \sigma \left(\frac{\beta_{\text{reg}}}{ y_w } \log \pi_{\theta}(y_w x) - \frac{\beta_{\text{reg}}}{ y_l } \log \pi_{\theta}(y_l x) - \gamma \right)$
RainbowPO	$-\log \sigma \left(\phi(x) \left[\frac{\beta_{\text{reg}}}{ y_w } \log \frac{\pi_{\theta}(y_w x)}{\pi_{\alpha}(y_w x)} - \frac{\beta_{\text{reg}}}{ y_l } \log \frac{\pi_{\theta}(y_l x)}{\pi_{\alpha}(y_l x)} \right] \right)$

Table 9: Various preference optimization DPO objectives. The table is inspired from Meng et al. (2024).

5.1 Offline Directed Preference Optimization (Offline DPO)

One disadvantage of RLHF is that the RL step often requires substantial computational effort (e.g., to carry out the proximal policy optimization). DPO, recently proposed by Rafailov et al. (2024), suggested a possible way to bypass the reward modeling stage and avoid RL, and has attracted great attention. The key idea of DPO is the observation that given a reward function $r(x, y)$, the problem in Eq. (8) has a closed-form solution:

$$\pi_r(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta_{\text{reg}}} r(x, y) \right), \quad (38)$$

where $Z(x) = \sum_y \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta_{\text{reg}}} r(x, y) \right)$ is a normalizing constant. Rearranging the terms, and plug in the ground truth reward r^* with the optimal policy $\pi^* = \pi_{r^*}$ yield:

$$r^*(x, y) = \beta_{\text{reg}} \log \frac{\pi^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta_{\text{reg}} \log Z(x). \quad (39)$$

Through this change of variables, the latent reward $r^*(x, y)$ can be expressed in terms of the optimal policy $\pi^*(y | x)$, the reference policy $\pi_{\text{ref}}(y | x)$ and a constant $Z^*(x)$. Substituting

this r^* expression into Eq. (2) yields:

$$p^*(y_1 \succ y_2 | x) = \sigma \left(\beta_{\text{reg}} \log \frac{\pi^*(y_1 | x)}{\pi_{\text{ref}}(y_1 | x)} - \beta_{\text{reg}} \log \frac{\pi^*(y_2 | x)}{\pi_{\text{ref}}(y_2 | x)} \right), \quad (40)$$

where $Z^*(x)$ cancels out and motivates the DPO objective:

$$\begin{aligned} \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) := \\ - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right], \end{aligned} \quad (41)$$

which is a supervised learning problem, requiring much less computation than the RLHF. To understand the loss objective of DPO, we can further examine its gradient as following:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = \\ - \beta_{\text{reg}} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_\theta(x, y_l) - \hat{r}_\theta(x, y_w))}_{\text{higher weight when estimate is wrong}} \left[\underbrace{\nabla_\theta \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} \right. \right. \\ \left. \left. - \underbrace{\nabla_\theta \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right], \end{aligned} \quad (42)$$

in which

$$\hat{r}_\theta(x, y) = \beta_{\text{reg}} \log \frac{\pi_\theta(y | x)}{\pi_{\text{ref}}(y | x)}, \quad (43)$$

is called the implicit reward model for the policy π_θ in DPO.

5.1.1 IDENTITY PREFERENCE OPTIMIZATION (IPO)

For DPO variants, we first visit IPO, proposed in Azar et al. (2024), motivated to bypass the assumption of Bradley-Terry model in the derivation of DPO (which comes from the reward modeling stage of RLHF). Azar et al. (2024) first propose a generic form of regularized optimization objective as:

$$\max_{\pi} \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi(\cdot | x) \\ y' \sim \mu(\cdot | x)}} [\Psi(p^*(y \succ y' | x))] - \beta_{\text{reg}} D_{\text{KL}}(\pi \| \pi_{\text{ref}}) \quad (44)$$

in which the new introduced function Ψ is non-decreasing. They show that Eq. (44) shares the same optimality as DPO when taking $\Psi(q) = \log(q/(1-q))$ (notably this equivalence still needs the Bradley-Terry model assumption). Furthermore, Azar et al. (2024) show that when $\Psi(x) = x$, i.e., when Ψ is the identity mapping, Eq. (44) is equivalent to:

$$\mathcal{L}_{\text{IPO}}(\pi_\theta; \pi_{\text{ref}}) := \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} - \frac{1}{2} \right)^2, \quad (45)$$

if the offline dataset \mathcal{D} is created by $x \sim \rho$ and $y, y' \sim \mu$. Notice that the derivation of the objective in Eq. (45) does not acquire Bradley-Terry model, thus IPO is *preference model*

free. In Azar et al. (2024), it is also demonstrated through a synthetic bandit experiment that DPO can be prone to overfitting, while IPO could avoid this problem. In addition, also shows that online version of IPO (Calandriello et al., 2024) (see details of online DPO in Section 4.2) is indeed equivalent to Nash-MD proposed in Nash Learning from Human Feedback (Munos et al., 2023).

5.1.2 REJECTION SAMPLING OPTIMIZATION (RSO)

RSO revisits the derivation of DPO and interpret the objective as a maximum likelihood estimator (MLE) of the optimal policy based on Eq. (40) (Liu et al., 2023). However, such a density estimation problem theoretically requires the datasets to be generated from the optimal policy instead of the SFT model in DPO. Thus, RSO algorithm is proposed to generate the datasets from the approximated optimal policy with an aid of a trained reward model r_{ϕ^*} and statistical rejection sampling, see in Algorithm 3. Notice that $\pi_{r_{\phi^*}}$

Algorithm 3 RSO algorithm.

- 1: Start with empty $\mathcal{Y} \leftarrow \{\}$.
 - 2: **while** not enough samples in \mathcal{Y} **do**
 - 3: Generate $y \sim \pi_{\text{sft}}(y | x)$ that is not in \mathcal{Y} .
 - 4: Generate $u \sim U[0, 1]$ and let $M = \min \left\{ m \mid m \geq \frac{\pi_{r_{\phi^*}}(y|x)}{\pi_{\text{sft}}(y|x)} \text{ for all } y \notin \mathcal{Y} \right\}$.
 - 5: **if** $u < \frac{\pi_{r_{\phi^*}}(y|x)}{M\pi_{\text{sft}}(y|x)}$ **then**
 - 6: Accept y and add it to \mathcal{Y} .
 - 7: **else**
 - 8: Reject y .
 - 9: **end if**
 - 10: **end while**
-

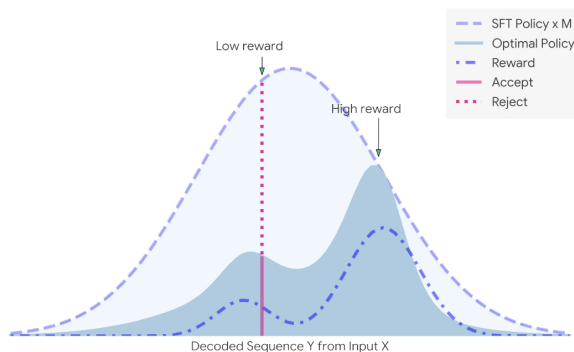


Table 10: RSO illustration in Liu et al. (2023).

is computed by Eq. (38) with the learned reward model r_{ϕ^*} . Liu et al. (2023) show that this *distribution correction* could help improve the performance of DPO by utilizing the resampled preference dataset.

Method	Loss Function	$f(x)$
DPO	log logistic	$-\log \sigma(x)$
IPO	square	$(x - 1)^2$
SLiC-HF	hinge loss	$\max(0, 1 - x)$

Table 11: Unified perspective through loss function in Liu et al. (2023), Tang et al. (2024b).

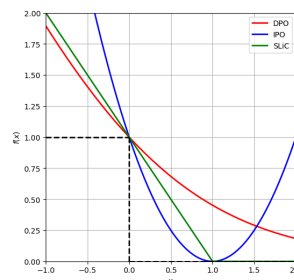


Table 12: Loss function comparison in Tang et al. (2024b).

In addition, RSO also unifies DPO and (normalized) SLiC-HF from the perspective of *loss function*; similar unified perspective also appeared in GPO (Tang et al., 2024b) (see e.g., in Table 1 of it):

$$\mathcal{L}_{\text{GPO}}(\pi_\theta; \pi_{\text{ref}}) := \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[f \left(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right], \quad (46)$$

for any convex function f , like in Table 11. GPO further provides an analysis of this formulation from an policy improvement and policy regularization trade-off. Applying Taylor Expansion of the form above yields:

$$\mathcal{L}_{\text{GPO}}(\pi_\theta; \pi_{\text{ref}}) = f(0) + \beta_{\text{reg}} \underbrace{f'(0)}_{<0} \underbrace{\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\rho_\theta]}_{\text{optimization}} + \frac{1}{2} \beta_{\text{reg}}^2 \underbrace{f''(0)}_{>0} \underbrace{\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\rho_\theta^2]}_{\text{regularization}}, \quad (47)$$

in which $\rho_\theta = \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)}$ denotes ‘implicit reward difference’.

5.1.3 f -DPO

DPO is derived from the RLHF objective which utilized the (reverse) KL divergence to prevent the deviation of new models from old models. f -DPO in Wang, Jiang, Yang, Liu, and Chen (2023) consider extending this statistical distance to general f -divergence. Concretely, for two probability distribution P and Q with probability density function p and q respectively, f -divergence is defined as:

$$D_f(P||Q) = \mathbb{E}_{q(x)} \left[f \left(\frac{p(x)}{q(x)} \right) \right], \quad (48)$$

and reverse KL divergence is a special case when taking $f(x) = x \log(x)$. Wang et al. (2023) first show that through a first order condition of optimality / KKT and the similar change of variable technique in DPO, the RLHF objective with a f -divergence

$$\mathcal{L}_{\text{RLHF}-f}(\phi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(\cdot | x)} [r_{\phi^*}(x, y) - \beta_{\text{reg}} D_f(\pi(\cdot | x) | \pi_{\text{ref}}(\cdot | x))], \quad (49)$$

could yield the f -DPO objective:

$$\begin{aligned} \mathcal{L}_{f\text{-DPO}}(\pi_\theta; \pi_{\text{ref}}) = \\ \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[-\log \sigma \left(\beta_{\text{reg}} f' \left(\frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} \right) - \beta_{\text{reg}} f' \left(\frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right) \right]. \end{aligned} \quad (50)$$

Special cases of Eq. (50) are when taking f divergence as α -divergence and JS-divergence, and Wang et al. (2023) further argue that JS-divergence could possibly yield a better diversity and accuracy tradeoff than reverse KL, through small-scale experiments on e.g., IMDB controllable generation and fine-tuning Pythia 2.8B on Anthropic HH dataset.

5.1.4 KAHNEMAN-TVERSKY OPTIMIZATION (KTO)

KTO (Ethayarajh et al., 2024) is motivated to address the need of pairwise preferences datasets in DPO, which can be scarce and expensive. Instead of maximizing the log-likelihood of preferences in DPO and inspired by Kahneman & Tversky’s prospect theory,

KTO proposes to minimize a human-aware loss function (HALO) that represents the utility of generations and also takes into account the human nature of loss aversion. The resulting KTO objective decouples the pair-preferences into two separate terms that are further linearly combined:

$$\mathcal{L}_{\text{KTO}}(\pi_\theta; \pi_{\text{ref}}) := -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\lambda_w \sigma(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - z_{\text{ref}}) + \lambda_l \sigma(z_{\text{ref}} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}) \right], \quad (51)$$

where $z_{\text{ref}} = \mathbb{E}_{(x, y) \sim \mathcal{D}} [\beta_{\text{reg}} \text{KL}(\pi_\theta(y|x) || \pi_{\text{ref}}(y|x))]$ acts like a subjective value and λ_w, λ_l are additional hyper-parameters to be tuned. If there is only desired/undesired answer, the KTO objective will thus have only one term, which makes it *pairwise preference data free*.

5.1.5 OFFSET DPO (ODPO)

DPO objective cannot reflect the *significance* of the preference pairs i.e., the extent y_w is preferred to y_l , and ODPO (Amini et al., 2024) propose to add a margin to capture this significance; they model this margin, or they call offset Δ_r as a monotonically increasing function $f(\cdot)$ of the difference between the scores associated with the responses:

$$\Delta_r(x, y_w, y_l) = \alpha f(\text{score}(x, y_w) - \text{score}(x, y_l)), \quad (52)$$

where α is a hyper-parameter that controls the extent to which an offset should be enforced. The resulting objective becomes:

$$\mathcal{L}_{\text{ODPO}}(\pi_\theta; \pi_{\text{ref}}) := -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} - \Delta_r(x, y_w, y_l) \right) \right]. \quad (53)$$

5.1.6 MALLOWS-DPO

Mallows-DPO (Chen et al., 2024b) is motivated by DPO’s lack of capability to characterize the diversity of human preferences. Inspired by Mallows Ranking Model (opposed to Bredley-Terry in RLHF and DPO) which has a natural carrier of a dispersion index, Mallows-DPO pays attention to the *dispersion* of the preferences: when human tends to agree about the answer to a certain question, e.g., ‘1 + 1 =?’ , the preference dispersion will be small; however, the dispersion will be large for answer to a general open question. Chen et al. (2024b) propose a contextual scaled objective derived from MLE under Mallows: compared to DPO that puts equal weights on each prompt and preference pairs, the resulting Mallows-DPO adds a contextual scaling factor $\phi(x)$ that represents this dispersion of the preferences of answers to each prompt x :

$$\mathcal{L}_{\text{Mallows-DPO}}(\pi_\theta; \pi_{\text{ref}}) := -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\phi(x) \left[\beta_{\text{reg}} \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta_{\text{reg}} \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right] \right) \right]. \quad (54)$$

To compute this dispersion, Mallows-DPO provided a direct approach by using a normalized predictive entropy of preference pairs $\{y_i^w, y_i^l\}_{i=1,\dots,N}$ with $N = \max(|y^w|, |y^l|)$:

$$\phi(x) = -\log \left(\frac{\frac{1}{2} \sum_{i=1}^{N-1} [H_{\pi_{\text{ref}}}(Y_{i+1} | Y_i = y_i^w) + H_{\pi_{\text{ref}}}(Y_{i+1} | Y_i = y_i^l)]}{\log(n)} \right). \quad (55)$$

To illustrate the effect of this additional term, when dispersion is high: $\phi(x)$ in Eq. (55) will be close to 0, and Mallows-DPO will put less weights on the corresponding preference pairs in the optimization objective to prevent from overfitting; In contrast, when dispersion is low, Mallows-DPO put more weights in the preference optimization objective, for which $\phi(x)$ is large and will lead to stronger effect of alignment.

5.1.7 LR-DPO

LR-DPO (Park et al., 2024), DPO with length regularization is motivated to address the problem of verbosity in the DPO setting. LR-DPO proposed a simple regularization strategy that prevents length exploitation by penalizing the rewards with length of the generation in standard RLHF objective:

$$\mathcal{L}_{\text{LR-RLHF}}(\phi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(\cdot|x)} [r_{\phi^*}(x, y) - \alpha|y| - \beta_{\text{reg}} \text{KL}(\pi(\cdot|x) | \pi_{\text{ref}}(\cdot|x))], \quad (56)$$

in which α is a hyper-parameter that controls the extent of length regularization. Eq. (56) thus similarly yields a supervised learning objective referred as DPO with length regularization:

$$\begin{aligned} & \mathcal{L}_{\text{LR-DPO}}(\pi_{\theta}; \pi_{\text{ref}}) \\ &= - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta_{\text{reg}} \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} - (\alpha|y_w| - \alpha|y_l|) \right). \end{aligned} \quad (57)$$

Park et al. (2024) further show that this can effectively improve model quality by addressing the verbosity issue.

5.1.8 CONTRASTIVE PREFERENCE OPTIMIZATION (CPO)

CPO (Xu et al., 2024a) is motivated to improve the memory and speed efficiency of DPO by neglecting the reference policy, further accompanied by a SFT loss term:

$$\mathcal{L}_{\text{CPO}}(\pi_{\theta}) := -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log p_{\theta}(y_w|x) + \log \sigma \left(\beta_{\text{reg}} \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\theta}(y_l|x)} \right) \right]. \quad (58)$$

5.1.9 ODDS RATIO PREFERENCE OPTIMIZATION (ORPO)

Opposed to maximizing the likelihood ratios of winning and losing answers in the preference pair in DPO, ORPO (Hong et al., 2024) propose that odds ratio can be a more sensible choice.

$$\begin{aligned} & \mathcal{L}_{\text{ORPO}}(\pi_{\theta}) := \\ & - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log p_{\theta}(y_w|x) + \lambda \log \sigma \left(\log \frac{p_{\theta}(y_w|x)}{1 - p_{\theta}(y_w|x)} - \log \frac{p_{\theta}(y_l|x)}{1 - p_{\theta}(y_l|x)} \right) \right]. \end{aligned} \quad (59)$$

where $p_\theta(y|x) = \exp\left(\frac{1}{|y|} \log \pi_\theta(y|x)\right)$. ORPO is similar to CPO in the sense that it is also reference model free and combined with a SFT loss; in addition, notably that ORPO also adopts a form of length regularization by normalizing the likelihoods with respect to the length, as in the definition of $p_\theta(y|x)$; finally, they compute odds ratio instead of the original likelihood ratio.

5.1.10 SIMPO

SimPO (Meng et al., 2024) proposes a simple yet effective objective that is claimed to match or even outperform the performance of DPO:

$$\mathcal{L}_{\text{SimPO}}(\pi_\theta) := -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\frac{\beta_{\text{reg}}}{|y_w|} \log \pi_\theta(y_w|x) - \frac{\beta_{\text{reg}}}{|y_l|} \log \pi_\theta(y_l|x) - \gamma \right) \right], \quad (60)$$

where γ is introduced as a target reward margin to help separating the winning and losing responses. SimPO is similar to CPO in the sense of being reference model free; it also adopted the length normalization for the likelihoods as in ORPO; finally, it additionally introduced a constant margin to be tuned that could help to further improve the performance by encouraging a larger difference between the normalized likelihoods.

5.1.11 RAINBOWPO

Inspired by the paper Rainbow on improving DQN for better performance, RainbowPO (Zhao et al., 2024a) demystifies the effectiveness of existing DPO variants by categorizing their key components into several broad directions, and integrate the identified effective components into a single cohesive objective:

$$\mathcal{L}_{\text{RainbowPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} f \left[\phi(x) \left(\frac{\beta}{|y_w|^\eta} \log \frac{\pi_\theta(y_w | x)}{\pi_\alpha(y_w | x)} - \frac{\beta}{|y_l|^\eta} \log \frac{\pi_\theta(y_l | x)}{\pi_\alpha(y_l | x)} \right) \right], \quad (61)$$

in which $\eta \in \{0, 1\}$, and π_α is referred to a mixing policy mechanism they propose for formulating a better reference policy by mixing policy π_{ref} and π_γ , defined as:

$$\pi_\alpha(y | x) \propto \pi_{\text{ref}}^\alpha(y | x) \cdot \pi_\gamma^{1-\alpha}(y | x), \quad (62)$$

and π_γ is a policy which assumes to exist (which can be understood as the reference policy taken by SimPO (Meng et al., 2024)), such that the model is perfect at distinguishing the preference pairs in the dataset:

$$\pi_\gamma(y_w | x)^{1/|y_w|} / \pi_\gamma(y_l | x)^{1/|y_l|} = \exp(\gamma), \quad (63)$$

for any prompt x . Zhao et al. (2024a) show that optimizing such generic objective can yield the best performance on downstream task of tuning Llama3-8B-Instruct for instruction-following capabilities, benefiting from composition of effective elements.

5.2 Diffusion Models or Multi-Modal Models

5.2.1 DIFFUSION-DPO

Diffusion-DPO (Wallace et al., 2024) is adapting DPO to diffusion models. It uses a fixed dataset and each example contains a prompt and a pairs of images generated from a reference model with human label. Similar to RL for diffusion, the goal is still to align the base diffusion models to human preferences. The derivation is similar to RL framework for diffusion in DDPO and DPOK, and also DPO for Language Models:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x_0^w, x_0^l) \sim \mathcal{D}} \log \sigma(\beta_{\text{reg}} \mathbb{E}_{\substack{x_{1:T}^w \sim p_\theta(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim p_\theta(x_{1:T}^l | x_0^l)}} \left[\log \frac{p_\theta(x_{0:T}^w)}{p_{\text{ref}}(x_{0:T}^w)} - \log \frac{p_\theta(x_{0:T}^l)}{p_{\text{ref}}(x_{0:T}^l)} \right]). \quad (64)$$

However, the main concern left is that the likelihood term of the generations $p_\theta(x_{0:T})$ is not tractable if only given generation x_0 . Wallace et al. (2024) further propose to use the forward process $q(x_{1:T} | x_0)$ of diffusion to match the distribution of backward process $p_\theta(x_{1:T} | x_0)$, and yield the final DPO-Diffusion objective:

$$\begin{aligned} L_{\text{DPO-diffusion}}(\theta) = & -\mathbb{E}_{(x_0^w, x_0^l) \sim \mathcal{D}, t \sim \mathcal{U}[0, T], x_t^w \sim q(x_t^w | x_0^w), x_t^l \sim q(x_t^l | x_0^l)} \log \sigma(-\beta_{\text{reg}} T \\ & [\text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| p_\theta(x_{t-1}^w | x_t^w)) - \text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| p_{\text{ref}}(x_{t-1}^w | x_t^w)) \\ & - \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| p_\theta(x_{t-1}^l | x_t^l)) + \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| p_{\text{ref}}(x_{t-1}^l | x_t^l))]), \quad (65) \end{aligned}$$

with each term can be readily computed. Other variants of Diffusion-DPO includes Diffusion-KTO (Li et al., 2024a) and Diffusion-NPO (Wang et al., 2025), which are mainly motivated by related preference learning algorithm in LLMs and adopt these successful methods for diffusion models. Motivated by Rich Feedback (Liang et al., 2024) paper, Rich Preference Optimization (RPO, Zhao, Chen, Guo, Winata, Ou, Huang, Yao, and Tang (2025a)) explores the promise of curating preference pairs from rich feedback. They utilize VLMs for generating concrete textual feedback signals on misalignment parts of prompt and diffusion models generated images, which further yield concrete and concise editing instructions. They utilize these instructions to edit the original image and get better images, which reveals the necessity of informative preference pairs.

5.2.2 POVID

POVID (Zhou et al., 2024b) proposes a method for performing preference optimization in visual language models (VLM) with synthetically generated preferences. This is mainly aimed at attenuating the hallucination problems in VLMs that arises due to lack of alignment between the language and vision modalities. Specifically, the authors use the ground-truth instructions as the preferred response and employ a two-stage approach to generate dis-preferred responses: first, use GPT-4V to inject hallucinatory texts into the preferred responses, and second, add diffusion noise to the image to trigger the inherent hallucination behavior of the VLM by making the image difficult for the VLM to understand. Both the

strategies are merged together in an reformulation of the DPO loss as:

$$\begin{aligned} \mathcal{L}_{\text{POVID}}(\pi_{\theta}; \pi_{ref}) = & -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\alpha \log \frac{\pi_{\theta}(y_w | x)}{\pi_{ref}(y_w | x)} \right. \right. \\ & \left. \left. - \left(\beta_{reg_1} \log \frac{\pi_{\theta}(y_l^t | x)}{\pi_{ref}(y_l^t | x)} + \beta_{reg_2} \log \frac{\pi_{\theta}(y_l^n | x^n)}{\pi_{ref}(y_l^n | x^n)} \right) \right) \right], \end{aligned} \quad (66)$$

where $\alpha, \beta_{reg_1}, \beta_{reg_2}$ are coefficients for balancing preferred responses (y_w) and dispreferred responses (y_l^t, y_l^n). y_l^t indicates the dispreferred response generated using GPT-4V, and y_l^n denotes the dispreferred response generated using the noisy image x^n .

5.3 Sequence Likelihood Calibration (SLiC-HF)

SLiC-HF (Zhao et al., 2023) uses a sequence level contrastive learning training method to align the model’s sequence likelihood over the decoded sequences by measuring their similarity with given reference sequences. The main reason to use a contrastive objective is to put more loss on negative sequence compared to positive sequences such that model puts more probability mass on generating positive sequences. Further, this specific formulation allows the use of human preference for ranking directly by using offline policy preference data \mathcal{D} or by training a separate predictive ranking model on offline data. SLiC-HF obtains a supervised fine-tuned model $\pi_{\theta_{ref}}(y | x)$, which we denote as the reference model for consistency with RLHF pipelines on a reference dataset $(x, y_{target}) \sim \mathcal{D}$. The preference datasets $\{y_w, y_l\}_m$ is formulated by uniformly drawing answer pairs from $\pi_{\theta_{ref}}(\cdot | x)$ and ranking them by their similarity (from a score computed by a pre-trained model denoted as $s(y, y_{ref}; x)$) to the target answer y_{ref} . The step after is to align the SFT model’s sequence likelihood using the SLiC loss (Zhao et al., 2022):

$$\mathcal{L}_{\text{SLiC}}(\pi_{\theta}; \pi_{ref}) = \sum L^{\text{cal}}(\theta, x, y_{target}, \{y_w, y_l\}_m) + \lambda L^{\text{reg}}(\theta, \theta_{ref}; x, y_{target}), \quad (67)$$

in which L^{cal} is the calibration loss from SLiC and L^{reg} is the regularization loss to prevent the aligned model stray away from the SFT model. Taking a special case of L^{cal} and L^{reg} to be a rank calibration loss and cross entropy loss respectively, Eq. (67) becomes:

$$\mathcal{L}_{\text{SLiC}}(\pi_{\theta}; \pi_{ref}) = \underbrace{\max(0, \delta - \log \pi_{\theta}(y_w | x) + \log \pi_{\theta}(y_l | x))}_{\text{rank calibration loss}} - \underbrace{\lambda \log \pi_{\theta}(y_{ref} | x)}_{\text{regularization}}, \quad (68)$$

where, in the first term of calibration loss, we are maximizing the likelihood corresponding to the positive sequence y_w and minimizing negative sequence y_l and the margin δ is a hyper-parameter represents which can be a constant or prompt dependent score/rank difference; the second term is just standard SFT loss. As a remark, one can use a secondary reward model, opposed to the similarity function in SLiC, trained on human preference data to classify positive or negative pairs (y_w, y_l) .

6. Combined Policies and Sampling-Agnostic Alignment

In this section, we explore some other directions proposed in literature for improving the effectiveness of human preference tuning. We discuss ExPO (Zheng et al., 2024a), which

proposed that combining two aligned models by extrapolating from their weights could enhance the alignment quality of the model; we discuss P3O (Fakoor et al., 2020), which utilized both on-policy and off-policy sampling; we also introduce sampling-agnostic alignment methods that can be applied to both off-policy and on-policy approaches.

6.1 ExPO

ExPO (Zheng et al., 2024a) provides a simple and training-free method for enhancing the alignment of large language models (LLMs) with human preferences. The core insight behind ExPO is that a model trained with DPO/RLHF can be viewed as an interpolation between two models with differing strengths. By leveraging this concept, one can potentially extrapolate a stronger model if the other two models are available. Specifically, if we denote the model π_{ExPO} as the interpolation of two other models, π_a and π_b , which may be trained using different alignment methods and datasets, ExPO assumes that combining these models will yield improved alignment. The stronger, better-aligned model π_{ExPO} is then obtained by extrapolating from the weights of these two relatively weaker models (which is reminiscent to Model Soups (Wortsman et al., 2022)), as formulated below:

$$\pi_{\text{ExPO}} = (1 + \alpha)\pi_a - \alpha\pi_b = \pi_a + \alpha(\pi_a - \pi_b) = \pi_a + \alpha\Delta\pi. \quad (69)$$

This method is shown to work when π_a and π_b are a stronger model from a combination of SFT model and a model further preference trained on top of it respectively. However in naive cases of choosing arbitrary π_a and π_b , it has shown to cause model collapse or degradation. Nevertheless broader applicability of this approach requires further research.

6.2 Policy-on Policy-off Policy Optimization (P3O)

P3O (Fakoor et al., 2020) is a simple and effective algorithm that uses the effective sample size to automatically manage the combination of on-policy and off-policy optimization. It performs gradient ascent using the gradient. Fakoor et al. (2020) describe how P3O integrates the following on-policy update with the off-policy update:

$$\nabla_{\theta}^{\text{on}} J(\pi_{\theta}) = \mathbb{E}_{s \sim d^{\pi_{\theta}}, a \sim \pi_{\theta}} [g(\pi_{\theta})], \quad (70)$$

$$\nabla_{\theta}^{\text{off}} J(\pi_{\theta}) = \mathbb{E}_{s \sim d_{\text{reg}}^{\beta}, a \sim \beta_{\text{reg}}} [\bar{\rho}_c g(\pi_{\theta})], \quad (71)$$

where π_{θ} denotes a policy that is parameterized by parameters $\theta \in \mathbb{R}^n$, and $q^{\pi_{\theta}}$ and $v^{\pi_{\theta}}$ denote a parameterization of the state-action and state-only value functions, respectively. It is also denoted the baselined policy gradient integrand in short by following:

$$g(\pi_{\theta}) = \hat{A}^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s), \quad (72)$$

$$\hat{A}^{\pi_{\theta}}(s, a) = \hat{q}^{\pi_{\theta}}(s, a) - \hat{v}^{\pi_{\theta}}(s). \quad (73)$$

It forms a unified policy optimization as following:

$$\mathbb{E}_{s \sim d^{\pi_{\theta}}, a \sim \pi_{\theta}} [g(\pi_{\theta})] + \mathbb{E}_{s \sim d_{\text{reg}}^{\beta}, a \sim \beta_{\text{reg}}} [\bar{\rho}_c g(\pi_{\theta})] - \lambda \nabla_{\theta} \mathbb{E}_{s \sim d_{\text{reg}}^{\beta}, a \sim \beta_{\text{reg}}} \text{KL}(\beta_{\text{reg}}(\cdot | s) \| \pi_{\theta}(\cdot | s)). \quad (74)$$

The first term above is the standard on-policy gradient. The second term is the off-policy policy gradient with truncation of the IS ratio using a constant c while the third term allows explicit control of the deviation of the target policy π_θ from β_{reg} . Further, the KL-divergence term can be rewritten as $\mathbb{E}_{s \sim d^{\pi_\theta}, a \sim \pi_\theta} [\log \rho]$ and therefore minimizes the importance ratio ρ over the entire replay buffer β_{reg} . There are two hyper-parameters in the P3O gradient: the IS ratio threshold c and the KL regularization co-efficient λ .

6.3 Reinforced Token Optimization (RTO)

Standard RLHF and DPO’s reward models are all based on the whole generation, thus the whole pipeline is in some sense closer to bandit instead of classical MDP based RL. Inspired by that nature of auto-regressive models is next token prediction, RTO (Zhong et al., 2024) derives a *token-wise* reward function from preference data and conducts policy optimization using this learned reward signal. Broadly, RTO formulates the optimization problem as an MDP and involves two primary steps: (1) learning a token-wise reward from preference data, and (2) optimizing this reward through RL training methods like PPO.

Theoretical Version. Consider the offline setting by assuming that we have an offline dataset $\mathcal{D} = \{(\tau^w, \tau^l)\}$ that contains several trajectory pairs, where $\tau^w = \{(s_h^w, a_h^w)\}_{h=1}^H$ is preferred over $\tau^l = \{(s_h^l, a_h^l)\}_{h=1}^H$. Each pair of trajectories shares the same initial state (i.e., $s_1^w = s_1^l$), but differs in the subsequent tokens. RTO computes the maximum likelihood estimator θ_{mle} based on \mathcal{D} by maximizing the log likelihood and calculates the pessimistic reward \hat{r} via token-wise reward learning. The output of the algorithm is policy $\hat{\pi}$.

Practical Version Similar to learning the reward model in RLHF, the key challenge left is to learn the token-wise reward from the offline data. For sentence level reward, popular frameworks outlined in InstructGPT (Ouyang et al., 2022), Claude (Bai et al., 2022a), and LLaMA2 (Touvron et al., 2023a) replace the last layer of the LLM with a linear layer for a scalar output and maximize the log-likelihood, which thus cannot be naively used for token-level reward. RTO observes that, given a trajectory $\tau = \{(s_h, a_h)\}_{h=1}^H$, denoting $\pi_{\beta_{\text{reg}}}^*(a|s) = \exp\{(Q_{\beta_{\text{reg}}}^*(s, a) - V_{\beta_{\text{reg}}}^*(s))/\beta_{\text{reg}}\}$ as the optimal policy, the KL regularization can be rewritten as:

$$\sum_{h=1}^H \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(a_h|s_h)}{\pi_{\text{ref}}(a_h|s_h)} = \sum_{h=1}^H (Q_{\beta_{\text{reg}}}^*(s_h, a_h) - V_{\beta_{\text{reg}}}^*(s_h) - \log \pi_{\text{ref}}(a_h|s_h))$$

$$= \sum_{h=1}^H r(s_h, a_h) - V_{\beta_{\text{reg}}}^*(s_1) \tag{75}$$

$$+ \underbrace{\sum_{h=1}^{H-1} (\mathbb{E}_{s' \sim \mathcal{P}(\cdot|s_h, a_h)} [V_{\beta_{\text{reg}}}^*(s')] - V_{\beta_{\text{reg}}}^*(s_{h+1}))}_{(\star)}, \tag{76}$$

in which the second equality follows from the fact that:

$$Q_{\beta_{\text{reg}}}^\pi(s, a) = r_{\beta_{\text{reg}}}(s, a) + \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} [V_{\beta_{\text{reg}}^\pi}(s')], \tag{77}$$

with $r_{\beta_{\text{reg}}}(s, a) = r(s, a) + \beta_{\text{reg}} \log \pi_{\text{ref}}(a|s)$. RTO focuses on the typical LLM generation scenario where the transition kernel is deterministic. Then, $(\star) = 0$, yielding that

$$\sum_{h=1}^H r(s_h, a_h) = \sum_{h=1}^H \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(a_h|s_h)}{\pi_{\text{ref}}(a_h|s_h)} + V_{\beta_{\text{reg}}}^*(s_1). \quad (78)$$

Building upon this result and combining it with the definition of the BT model, for any trajectory pair $\{\tau^j = \{(s_h^j, a_h^j)\}_{h=1}^H\}_{j=1}^2$ satisfying $s_1^1 = s_1^2$, we have:

$$\begin{aligned} \mathbb{P}(\tau^1 \succ \tau^2) &= \sigma \left(\sum_{h=1}^H r(s_h^1, a_h^1) - \sum_{h=1}^H r(s_h^2, a_h^2) \right) \\ &= \sigma \left(\sum_{h=1}^H \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(a_h^1|s_h^1)}{\pi_{\text{ref}}(a_h^1|s_h^1)} - \sum_{h=1}^H \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(a_h^2|s_h^2)}{\pi_{\text{ref}}(a_h^2|s_h^2)} \right). \end{aligned} \quad (79)$$

Similar to the bandit setting where the learning objective is equivalent to a BT model with sentence-wise reward $r^*(x, y) = \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(y|x)}{\pi_{\text{ref}}(y|x)}$ (Rafailov et al., 2024), it shows that the learning objective in token-wise MDP equivalents to a BT model with a token-wise reward function

$$r^*(s_h = (x, y_{1:h-1}), a_h = y_h) = \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(a_h|s_h)}{\pi_{\text{ref}}(a_h|s_h)} = \beta_{\text{reg}} \log \frac{\pi_{\beta_{\text{reg}}}^*(y_h|x, y_{1:h-1})}{\pi_{\text{ref}}(y_h|x, y_{1:h-1})}, \quad (80)$$

where x is the prompt, $y_{1:h-1}$ is the tokens generated so far, and y_h is the token chosen at the current step. RTO assigns the defined token-wise reward function to each step. Formally, for any h , it is defined as following:

$$\begin{aligned} &\beta_{\text{reg}}^1 \log \frac{\pi_{\beta_{\text{reg}}}^*(y_h|x, y_{1:h-1})}{\pi_{\text{ref}}(y_h|x, y_{1:h-1})} - \beta_{\text{reg}}^2 \log \frac{\pi(y_h|x, y_{1:h-1})}{\pi_{\text{ref}}(y_h|x, y_{1:h-1})} \\ &\approx \beta_{\text{reg}}^1 \log \frac{\pi_{\text{dpo}}(y_h|x, y_{1:h-1})}{\pi_{\text{ref}}(y_h|x, y_{1:h-1})} - \beta_{\text{reg}}^2 \log \frac{\pi(y_h|x, y_{1:h-1})}{\pi_{\text{ref}}(y_h|x, y_{1:h-1})} := r_{\text{rto}}((x, y_{1:h-1}), y_h), \end{aligned} \quad (81)$$

as the token-wise reward used by RTO, where β_{reg}^1 and β_{reg}^2 are tuning parameters, and π is the current policy to be updated. In the last step, RTO uses π_{dpo} , the policy learned by DPO, as a proxy for the unknown optimal $\pi_{\beta_{\text{reg}}}^*$. Then RTO employs PPO to optimize the model with respect to the token-wise reward r_{rto} . The idea of transformation from sequence level preferences to token level guidance also appeared in an earlier work by Yang, Zhang, Xia, Feng, Xiong, and Zhou (2024c).

7. Training-Free Optimization

Another area currently being explored in the literature is training-free preference optimization or inference-time alignment. The main goal of this approach is to use inference-time techniques to cater to user-specific or task-specific preferences at runtime, instead of relying on explicit training, while preserving the original model's integrity. Inference-time techniques can include a broad range of methods, such as, in-context learning, model merging,

chat vectors, decode time intervention, and the use of external models during decoding. Existing work in these categories implicitly depends on test-time scaling to achieve performance. Since there is significant overlap in this area, this section clusters the broad types of techniques and details exemplar works for each cluster.

7.1 In-Context Alignment

This group of work builds on vanilla in-context learning to achieve alignment without changing model weights. They aim to steer the model’s behavior towards a preferred outcome, either through changes in the system prompt or in-context examples. Lin et al. (2024) utilize a specially crafted system prompt and as few as three stylistic examples, enabling the model to adapt effectively to response styles. Lee, Park, Kim, and Seo (2024) generalize the specialized system prompts, allowing users to specify behavioral preferences in system prompts, thereby influencing model behavior. Although alignment is training-free, the authors did need to train models to interpret a vast array of system prompts. Song, Fan, Zhang, Wang, and Wang (2024a) use a contrastive scoring mechanism to assess the difference in log probabilities between the model’s samples for each response before and after in-context learning, and then selects the response that aligns with a human preference score. Cheng, Liu, Zheng, Ke, Wang, Dong, Tang, and Huang (2024) optimize (iteratively updates) user prompts to better match the LLM’s input processing without changing LLM weights, effectively training a small surrogate ‘reward’ model from pairwise human preferences. This is similar in essence to RLHF, but the changes are applied to the prompts themselves, viewed as a form of iterative prompt tuning.

7.2 Reward Model-Based Alignment

Several methods use feedback from a separate reward model or LLM to iteratively improve the model’s responses. These techniques generally rely on selecting the best response by using a reward model or LLM-As-A-Judge, which acts as a proxy for human preference, and refining the output through iterative feedback. Li et al. (2025) introduce Test-Time Preference Optimization (TPO), which utilizes a reference model to sample multiple responses to a prompt. A reward model then scores these responses and selects a winner and a rejected response. A critique LLM is then used to provide feedback on why the winner outperforms the rejected response. Another call to a critique LLM uses these feedback to define how to improve current response, and then based on the instructions generate sample responses for the next round. This refinement continues until the final response is obtained, demonstrating improvements over DPO-baselined models.

Other works focus on Best-of- N sampling-based methods, where N candidate responses are generated for a prompt and evaluated using a reward model to select the highest-scoring response. The granularity of rewards can vary between using response level rewards, segment level rewards and token level rewards during decoding. Heuristics for selecting the best response include Best-of N (Lightman et al., 2024; Huang et al., 2025), Monte Carlo tree search (Zhang et al., 2024; Zhao et al., 2024b), reward model based exploitation and exploration (Hung et al., 2024) or rejection sampling during decoding the response (Khanov et al., 2024; Li et al., 2024).

7.3 Vector-Based Alignment

This body of work suggests that encoded representations of preference dimensions can be learned and extracted, similar to task arithmetic (Ilharco et al., 2023). Multiple such representations, corresponding to different domain-specific alignments, can be extracted through linear operations between the base and aligned models, and then applied to the base model in various ways to achieve different alignment goals. Huang et al. (2024) lead the way of using ‘alignment vectors’ obtained from the difference between base models and some trained version of the same model (e.g., a continued pre-trained version), which regains alignment characteristics. Shahriar, Qi, Pappas, Doss, Sunkara, Halder, Mager, and Benajiba (2024) extract these ‘alignment vectors’ and achieves multi-domain alignment on the fly by performing weighted mixing, then applying them to the base model. Shirafuji, Takenaka, and Taguchi (2025) focus on debiasing alignment, specifically training a biased model and subtracting the base LLM to obtain the ‘alignment vector’ that exemplifies the direction of bias. They then use linear operations with the ‘alignment vector’ to de-bias other pre-trained language models.

7.4 Decoding Time Alignment

Methods in this section influence decoding at the token level directly. For example, Khanov et al. (2024), Chen, Zhang, Luo, Chai, and Liu (2025), Kuang, Sun, McFaddin, Ma, and Ettl (2024) modify the conditional probability of token generation by using a general-purpose or personalized reward model on the vocabulary during decoding, thus generating tokens that align with preferences. Similarly, Liu et al. (2024) adjust the conditional probability of token generation by interpolating between an aligned model and a reference model during decoding, allowing dynamic control over alignment without retraining. While these methods are training-free from a preference optimization perspective, the reward model itself still needs training unless an off-the-shelf model can be repurposed. Zhu, Liu, Zhang, Guo, and Mao (2025) take a slightly different approach by dynamically adapting the model’s behavior at the token level using surrogate reward signals instead of actual reward models during decoding, based on preference principles such as ethical guidelines or specific stylistic preferences. Zhang, Bai, Chen, Ma, Wang, Sun, Zheng, and Yang (2025) similarly treat each token generation as an online learning problem, adjusting the model’s outputs to align with user-specific preferences provided through user or system prompts. Li, Wei, Zhao, Zhang, and Zhang (2023b), on the other hand, uses self-reflection on the reference model itself alongside preference principles, allowing the model to adjust its token-level decisions, i.e., select token or rewind generation during the decoding phase. Gao, Ge, Shen, Dou, Ye, Wang, Zheng, Zou, Chen, Yan, et al. (2024) also work via influencing the decoding by directly estimating aligned responses through a method called self-contrastive decoding. Finally, Li, Patel, Viégas, Pfister, and Wattenberg (2023b) adopt a different strategy by identifying specific attention heads within the model strongly correlated with truthful responses. During inference, the activations of these attention heads are adjusted in directions that promote truthfulness, achieving alignment.

7.5 Model Merging Based Alignment

Recent work in model merging has shown that performing some type of weighted linear interpolation of model weights has compositional effects (Wortsman et al., 2022). This has been demonstrated in other domains such as federated learning (Kairouz et al., 2021; Das & Brunschwiler, 2019) in IID or non-IID domains. Zheng et al. (2024a) perform linear interpolation between a base SFT-ed model and an aligned model to produce a new model with improved alignment with human preferences. Jang, Kim, Lin, Wang, Hessel, Zettlemoyer, Hajishirzi, Choi, and Ammanabrolu (2024) involve independently training policy models on distinct preference dimensions and subsequently merging their parameters post-training to optimize different directions. This approach reduces computational complexity from exponential to linear concerning the number of preferences and allows for efficient integration of new preferences without retraining the entire model.

7.6 Comparative Analysis: Inference-Time vs. Training-Based Alignment

There are several key differences between inference-time and training-based alignment approaches:

1. **Alignment Speed:** Inference-time methods are proposed as fast alternatives to their training-based counterparts, such as RLHF or DPO, which can take days. However, inference costs cannot be discounted entirely, as multiple LLM calls are often necessary for even a single response in inference-time alignment.
2. **Complexity of alignment:** Inference based optimizations are frequently somewhat limited to the performance of the context and the models attention to the context, and decoding strategies and reward models instead of explicit training performance.
3. **On the Fly Domain Transfer:** training-free approaches can adapt to changing preferences or domains on the fly, whereas training-based approaches typically require expensive retraining.
4. **Data Efficiency:** Training or retraining requires extensive collection of preference data, which can be costly. Inference-time algorithms often perform well with vastly smaller amounts of in-context data.
5. **Latency Impact:** Training-based approaches have fixed inference latency after training, whereas the latency of inference-time alignment approaches is typically higher due to test-time scaling.

It is important to note that there is no conclusive evidence to suggest that either training-free or training-based alignment techniques consistently outperform the other. The choice between the two should depend on resource constraints and the specific requirements of the use case.

8. Evaluation

In the context of aligning models with human preferences it is crucial to evaluate not just the models' core abilities but also how aligned they are with human preferences. In Section 4 -

Section 7 we have explored different types of preference optimization techniques, but without standardized and comprehensive evaluation metrics, it is difficult to assess the practical performance of these aligned models or compare different optimization approaches. In this section, we survey key evaluation techniques across speech, language, vision, and reward models, enabling practitioners to make informed decisions about model development and deployment. Each of the evaluation methods presented in this section, from AlpacaEval for language models to SpeechLMScore for speech models, can be applied across all optimization techniques within their respective modalities.

8.1 LLM As A Judge for LLM Evaluation

Human evaluation is both costly and time-consuming. Developing an automatic evaluation method that closely aligns with human assessments can significantly reduce evaluation time and accelerate research progress. In this context, we outline the benchmarks employed for automatic evaluation using LLMs.

8.1.1 ALPACAEVAL

AlpacaEval (Dubois et al., 2024) win rate (against GPT4) is an LLM-based automatic evaluation that has high-level agreement to human. To further improve the fairness of the evaluation and address the verbosity issue of GPT4 as a judge, Dubois et al. (2024) introduce a length-controlled version of AlpacaEval that aims to conduct measurement with outputs with similar lengths. The metric is used in AlpacaEval calculates win-rates for models across a variety of NLP tasks to measure of model capabilities compared to a baseline by using an LLM judge. **AlpacaEval 2.0:** The judge uses GPT4-Turbo to replace GPT-3 based model “text-davinci-003” in the 1.0 version, which makes it more challenging and have a metric that better reflects the current SOTA model.

8.1.2 CHATBOTARENA

ChatbotArena (Chiang et al., 2024) is a benchmarking platform for Large Language Models (LLMs) that conducts anonymous, randomized ‘battles’ in a crowdsourced environment. On this platform, users can pose questions and receive responses from two anonymous LLMs. After reviewing the answers, users vote for the response they prefer, with the identities of the models revealed only after voting. This crowdsourcing approach effectively gathers a diverse array of user prompts, accurately reflecting real-world LLM applications. Utilizing this data, they apply a range of advanced statistical techniques, from the Bradley-Terry model (Bradley & Terry, 1952) to the E-values framework (Vovk & Wang, 2021), to estimate model rankings as reliably and efficiently as possible.

8.1.3 MT-BENCH

MT-bench (Zheng et al., 2024b) is a series of open-ended questions designed to evaluate a chatbot’s multi-turn conversational and instruction-following abilities. It is used in the platform that assesses these capabilities in a crowdsourced battle format. This platform is particularly useful for evaluating the quality of LLM-generated responses, utilizing judges like GPT-4. Consequently, employing LLM as a judge provides a scalable and explainable method to approximate human preferences, which would otherwise be very costly to obtain.

8.1.4 HELM

HELM (Liang et al., 2022) is a large-scale reproducible and transparent framework for evaluating LLM models to enhance the transparency of language models. The framework has seven metrics, such as accuracy, calibration, robustness, fairness, bias, toxicity, and efficiency.

8.2 Vision Language Model Evaluation

8.2.1 VHELM

VHELM¹ is an extension of the HELM framework (Liang et al., 2022) with the adaptation methods to assess the performance of VLMs by scoring the winning rates against the GPT-4V model.

8.2.2 MMSTAR

MMStar (Chen et al., 2024e) is a multi-modal benchmark consisting of 1,500 samples meticulously curated by human experts. It evaluates six core capabilities and 18 specific criteria to assess the multi-modal capacities of LVLMS. The samples are selected from existing benchmarks using an automated process, followed by human review to ensure each sample demonstrates visual dependency, minimal data leakage, and requires advanced multi-modal skills.

8.3 Speech Language Model Evaluation

8.3.1 SPEECHLMScore

SpeechLMScore (Maiti et al., 2023) calculates the average log-probability of a speech signal by converting it into discrete tokens and assessing the average probability of generating the token sequence. Formally, $\text{SpeechLMScore}(\mathbf{x}|\theta)$ is defined as:

$$\text{SpeechLMScore}(\mathbf{d}|\theta) = \frac{1}{T} \sum_{i=1}^T \log p(d_i|d_{<i}, \theta), \quad (82)$$

where θ is an LM used to generate the score. Specifically, to compute SpeechLMScore, the process involves: i) encoding the speech into discrete tokens $\mathbf{d} = d_1 \cdots d_T$, and ii) iteratively calculating the log probability of each token d_i given all preceding tokens $d_1 \cdots d_{i-1}$ using θ , i.e., $\log p(d_i|d_{<i}, \theta)$. SpeechLMScore thus measures the average log-probability of a sequence of speech tokens. This metric is closely related to the perplexity of a speech sample, essentially indicating how perplexed a speech language model is when presented with a set of discrete tokens from speech \mathbf{x} .

8.3.2 SPEECHBERTScore

SpeechBERTScore (Saeki et al., 2024) evaluates the BERTScore for self-supervised dense speech features derived from both generated and reference speech, even when these sequences differ in length. This method utilizes BERTScore as a metric to assess the quality

1. <https://crfm.stanford.edu/helm/vhelm/latest/>.

of speech generation. By computing the BERTScore for SSL feature sequences from both the generated and reference speech, SpeechBERTScore effectively captures their semantic alignment.

8.4 Reward Model Evaluation

To assess the quality of reward models, it is crucial to evaluate their performance using appropriate benchmarks. Zhu, Frick, Wu, Zhu, Ganesan, Chiang, Zhang, and Jiao (2024), Jiang, Ren, and Lin (2023) propose using validation sets from previous RLHF training processes, such as Anthropic’s Helpful and Harmless data (Bai et al., 2022a) or OpenAI’s Learning to Summarize (Stiennon et al., 2020). Additionally, newly released preference data, aimed at expanding the diversity of preference training datasets, such as UltraFeedback (Cui et al., 2023), UltraInteract (Yuan et al., 2024a), and Nectar (Zhu et al., 2024), lack test sets, necessitating a new style of evaluation for reward models. RewardBench is a benchmark dataset and codebase designed for this purpose (Lambert et al., 2024). The dataset comprises a collection of prompt-chosen-rejected triplets that span various domains, including chat, reasoning, and safety. This allows for a comprehensive evaluation of how reward models perform on challenging, structured, and out-of-distribution queries. Winata et al. (2024) propose METAMETRICS, a new method to construct a meta-metric that is aligned with human preferences by calibrating multiple metrics by using Bayesian optimization and boosting methods, which has been further applied to machine translation (Anugraha et al., 2024).

9. Discussion and Research Directions

In this section, we describe topics related to human preferences that are either underexplored or still in their early stages. We also discuss potential future research areas that could be highly beneficial for advancing the field.

9.1 Discussion

9.1.1 EFFECTIVENESS OF OPTIMIZATION COMPONENTS

In the literature on preference tuning, the comparative performance of different methods remains unclear, particularly when comparisons are not conducted under fair conditions. This is largely because RL is highly sensitive to changes in hyper-parameters, and running multiple hyper-parameter configurations is very costly. For instance, when a new method is proposed, the baseline may not be fully optimized, resulting in weaker baselines. Another issue in automatic evaluation using LLMs as judges is the bias introduced by the pre-training data. A model might prefer predictions generated by a similar type of model. For example, a GPT-4 model may favor outputs from its own model family over those from other models, such as Llama. Additionally, judge models may have a preference for longer sequences or text in certain positions (Zheng et al., 2024b). Therefore, finding a less biased model is crucial during evaluation. Consequently, the effectiveness of each method, along with their optimized components and the models used in automatic evaluation, needs further investigation and careful consideration.

9.1.2 OFFLINE VS. ONLINE ALGORITHMS

Through theoretical and experimental analysis, Xu, Fu, Gao, Ye, Liu, Mei, Wang, Yu, and Wu (2024a) explore the limitations of DPO and find that DPO is sensitive to distribution shifts between the base model outputs and preference data. They suggest that iterative DPO, which involves continuous updating, is more effective than training on static data. However, they also find that DPO fails to improve performance on challenging tasks such as code generation. From a different perspective, Tang, Guo, Zheng, Calandriello, Cao, Tarassov, Munos, Pires, Valko, and Cheng (2024a) clarify the confusion surrounding the limitations of offline algorithms’ performance, often attributed to the bounded performance of offline algorithms. The paper discusses that the dichotomy between online and offline algorithms is frequently inaccurate in practice. An offline algorithm that continuously updates its data stream effectively functions as an online algorithm. Consequently, the shortcomings identified in offline learning can be mitigated by adopting a more careful approach to the data generation process.

9.1.3 LLM AS A JUDGE CAN BE UNRELIABLE

It has widely been aware that LLM-as-a-judge can be susceptible to length bias: Wang, Ivison, Dasigi, Hessel, Khot, Chandu, Wadden, MacMillan, Smith, Beltagy, et al. (2023a) have noticed that when evaluating 13B parameter models in head-to-head comparisons with the Davinci-003 model, win rates have a strong correlation (0.96) with the average number of unique tokens in the model’s response. Zheng, Pang, Du, Liu, Jiang, and Lin (2024c) find that popular benchmarks of LLM-as-a-judge can be cheated by even a null model to achieve impossible high winning rate. These possible bias of LLM-as-a-judge calls for further research on how to improve the existing usage or more reliable evaluation methods while still being efficient.

9.2 Research Directions

Here, we explore potential research directions that offer significant opportunities for further investigation and development. These avenues hold promise for both academic researchers and industry practitioners, providing ground for innovative studies and practical applications. We summarize key ideas and topics that could drive future advancements in the field, highlighting areas where there is ample room for exploration and growth.

9.2.1 MULTILINGUAL, MULTICULTURAL, AND PLURALISTIC PREFERENCE TUNING

While significant resources have been allocated to enhance the safety of LLMs for deployment, the safety of multilingual LLMs remains underexplored. Ahmadian, Ermis, Goldfarb-Tarrant, Kreutzer, Fadaee, and Hooker (2024b) is one of the pioneering works pushing the boundaries of aligning language models by optimizing for both general and safety performance simultaneously in a multilingual setting using Distributional DPO. Similarly, Li et al. (2024b) propose exploring DPO training to reduce toxicity in multilingual open-ended generations. Another line of research focuses on using multilingual alignment based on human preferences to improve reasoning abilities, aiming to align reasoning processes in other languages with those in the dominant language (She et al., 2024). There is still ample

room for exploration in the multilingual space, particularly in examining the cultural aspects of multilingualism (Adilazuarda et al., 2024; AlKhamissi et al., 2024) and improving the alignment of LLM for generation (Winata et al., 2021b). It is crucial to cover more diverse languages, including regional languages, different dialects (Aji et al., 2022), and code-switching (Winata et al., 2021a), which are common phenomena in bilingual and multilingual communities (Winata et al., 2024). Additionally, the exploration of multilingual topics in vision-language and speech tasks remains open for further investigation.

9.2.2 MULTI-MODALITY

While alignment in LLMs has been extensively studied, alignment for multi-modal models has not yet been investigated to the same extent. Sun et al. (2023) and Zhou et al. (2024b) align LLaVA (Liu et al., 2024a) using PPO and DPO, respectively. Similarly, Li, Xie, Li, Chen, Wang, Chen, Yang, Wang, and Kong (2023c) and Yu, Hu, Yao, Zhang, Zhao, Wang, Wang, Pan, Xue, and Li (2023) employ DPO and its variations to align the Qwen-VL (Bai et al., 2023) and Muffin (Yu et al., 2023) models. Notably, in addition to different alignment strategies and base models, all these works introduce novel preference datasets for alignment, varying in size, collection methods, and generation schemes. Consequently, while each of these studies offers valuable insights into alignment for multi-modal LLMs, it can sometimes be challenging to attribute reported improvements to specific proposed choices. Furthermore, Amirloo, Fauconnier, Roesmann, Kerl, Boney, Qian, Wang, Dehghan, Yang, and Gan (2024) examine each component of multi-modal alignment independently, involving sampling from the model during policy optimization.

9.2.3 SPEECH APPLICATIONS

The application of preference tuning in speech technology is in its early stages, offering many opportunities for future exploration. As research advances, preference tuning is expected to enhance various speech-related technologies, including TTS and speech recognition systems, by incorporating human preferences to improve performance and user satisfaction. In TTS, it can help select the most natural and pleasing synthetic voices (Zhang et al., 2024), while in speech recognition, it can ensure more accurate and contextually appropriate transcriptions. Additionally, preference tuning can benefit voice assistants, automated customer service systems, and language learning tools by creating more intuitive and effective interfaces. Ongoing research and experimentation will be essential to fully realize the potential of preference tuning in speech technology, aiming to develop systems that are both technically proficient and closely aligned with human communication and preferences.

9.2.4 UNLEARNING

Yao, Xu, and Liu (2023b), Zhang, Lin, Bai, and Mei (2024a) propose an alignment technique for unlearning by utilizing negative examples, which are easier and cheaper to collect than the positive examples needed for preference tuning. This method is considered computationally efficient, with costs comparable to light supervised finetuning. They demonstrate that unlearning is particularly appealing when resources are limited and the priority is to stop generating undesirable outputs. Despite using only negative samples, unlearning can achieve better alignment performance than RLHF. The unlearning method can be very

useful in removing harmful responses, erasing copyright-protected content, and reducing hallucinations. This approach is promising and has potential for further exploration in future work.

9.2.5 BENCHMARKING PREFERENCE TUNING METHODS

Developing a comprehensive benchmark for various preference tuning methods is essential for gaining a clearer understanding of their individual effectiveness. Currently, the effectiveness of each method is somewhat unclear, making it difficult to fully appreciate their value. By creating a benchmark, we can systematically assess and compare these methods, thereby clarifying their strengths and weaknesses. This effort to elucidate the usefulness of each approach is vital for advancing our knowledge and improving the application of preference tuning techniques. Such a benchmark would not only enable more informed decisions when selecting the most suitable method for specific tasks but also stimulate innovation by identifying areas that require further refinement and development. Ultimately, this initiative aims to enhance the overall effectiveness and reliability of preference tuning methods across various applications.

9.2.6 MECHANISTIC UNDERSTANDING OF PREFERENCE TUNING METHODS

Despite the popularity of preference tuning methods for LLM alignment, explanations for their underlying mechanisms in terms of how models become “aligned” still lack, thus making it difficult to explain phenomena like jailbreaks (Chao et al., 2023). Taking toxicity reduction as the task and applying DPO on GPT2-medium, Lee, Bai, Pres, Wattenberg, Kummerfeld, and Mihalcea (2024) suggest that capabilities may be rather bypassed instead of removed. Thus, current preference tuning methods may still be vulnerable to reverse-engineering and the models tuned are easy to be unaligned again. More interpretation of preference tuning methods could possibly address these concerns by ensuring that models not only meet alignment objectives more reliably but also provide clearer insights into how these objectives are achieved; it could also possibly help lead to better preference tuning methods that can mitigate issues such as jailbreaking and other forms of misalignment, where models exhibit undesirable behaviors despite appearing aligned during training.

9.2.7 PRIVACY AND SAFETY

There has also been concern raised about the private and safety aspects about existing alignment techniques. Wu, Inan, Backurs, Chandrasekaran, Kulkarni, and Sim (2023a) propose to combine DP-SGD with RLHF using PPO for a privacy preserving alignment. For safety, Dai, Pan, Sun, Ji, Xu, Liu, Wang, and Yang (2023a) propose safe RLHF by introducing a safety constraint to ensure the aligned model to be both helpful and less harmless. Recent study in Qi, Panda, Lyu, Ma, Roy, Beirami, Mittal, and Henderson (2024) introduce the concept of “shallow safety alignment” to uncover a fundamental vulnerability in current safety alignment approaches and proposes “deep safety alignment” as a promising defense. These works pin down several future research directions in the privacy and safety aspects of RLHF and general alignment algorithms.

Acknowledgments

Genta Indra Winata, Hanyang Zhao, and Anirban Das contributed equally on this paper. Wenpin Tang and Hanyang Zhao are supported by NSF grant DMS-2206038, a start-up grant at Columbia University, and the Columbia Innovation Hub grant. The works of Hanyang Zhao and David D. Yao are part of a Columbia-CityU/HK collaborative project that is supported by InnotHK Initiative, The Government of the HKSAR and the AIFT Lab.

References

- Abdin, M., Jacobs, S. A., Awan, A. A., Aneja, J., Awadallah, A., Awadalla, H., Bach, N., Bahree, A., Bakhtiari, A., & Behl, H. (2024). Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., & Anadkat, S. (2023). Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Adilazuarda, M. F., Mukherjee, S., Lavania, P., Singh, S., Dwivedi, A., Aji, A. F., O’Neill, J., Modi, A., & Choudhury, M. (2024). Towards measuring and modeling” culture” in llms: A survey. *arXiv preprint arXiv:2403.15412*.
- Adler, B., Agarwal, N., Aithal, A., Anh, D. H., Bhattacharya, P., Brundyn, A., Casper, J., Catanzaro, B., Clay, S., Cohen, J., et al. (2024). Nemotron-4 340b technical report. *arXiv preprint arXiv:2406.11704*.
- Ahmadian, A., Cremer, C., Gallé, M., Fadaee, M., Kreutzer, J., Üstün, A., & Hooker, S. (2024a). Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*.
- Ahmadian, A., Ermis, B., Goldfarb-Tarrant, S., Kreutzer, J., Fadaee, M., & Hooker, S. (2024b). The multilingual alignment prism: Aligning global and local preferences to reduce harm. *arXiv preprint arXiv:2406.18682*.
- Aji, A. F., Winata, G. I., Koto, F., Cahyawijaya, S., Romadhony, A., Mahendra, R., Kurniawan, K., Moeljadi, D., Prasojo, R. E., & Baldwin, T. (2022). One country, 700+ languages: Nlp challenges for underrepresented languages and dialects in indonesia. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7226–7249.
- AlKhamissi, B., ElNokrashy, M., AlKhamissi, M., & Diab, M. (2024). Investigating cultural alignment of large language models. *arXiv preprint arXiv:2402.13231*.
- Allal, L. B., Li, R., Kocetkov, D., Mou, C., Akiki, C., Ferrandis, C. M., Muennighoff, N., Mishra, M., Gu, A., & Dey, M. (2023). Santacoder: don’t reach for the stars!. *arXiv preprint arXiv:2301.03988*.
- Almazrouei, E., Alobeidli, H., Alshamsi, A., Cappelli, A., Cojocaru, R., Debbah, M., Goffinet, E., Hesslow, D., Launay, J., & Malartic, Q. (2023). The falcon series of open language models. *arXiv preprint arXiv:2311.16867*.

- Amini, A., Vieira, T., & Cotterell, R. (2024). Direct preference optimization with an offset. *arXiv preprint arXiv:2402.10571*.
- Amirloo, E., Fauconnier, J.-P., Roesmann, C., Kerl, C., Boney, R., Qian, Y., Wang, Z., Dehghan, A., Yang, Y., & Gan, Z. (2024). Understanding alignment in multimodal llms: A comprehensive study. *arXiv preprint arXiv:2407.02477*.
- Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., & Chen, Z. (2023). Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Anthropic, A. (2024). The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card, 1*.
- Anugraha, D., Kuwanto, G., Susanto, L., Wijaya, D. T., & Winata, G. I. (2024). Metametrics-mt: Tuning meta-metrics for machine translation via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation, USA. Association for Computational Linguistics*.
- Aryabumi, V., Dang, J., Talupuru, D., Dash, S., Cairuz, D., Lin, H., Venkitesh, B., Smith, M., Marchisio, K., & Ruder, S. (2024). Aya 23: Open weight releases to further multilingual progress. *arXiv preprint arXiv:2405.15032*.
- Azar, M. G., Guo, Z. D., Piot, B., Munos, R., Rowland, M., Valko, M., & Calandriello, D. (2024). A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR.
- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., & Huang, F. (2023). Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., & Henighan, T. (2022a). Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., & McKinnon, C. (2022b). Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Bellagente, M., Tow, J., Mahan, D., Phung, D., Zhuravinskyi, M., Adithyan, R., Baicoianu, J., Brooks, B., Cooper, N., & Datta, A. (2024). Stable lm 2 1.6 b technical report. *arXiv preprint arXiv:2402.17834*.
- Biderman, S., Schoelkopf, H., Anthony, Q. G., Bradley, H., O’Brien, K., Hallahan, E., Khan, M. A., Purohit, S., Prashanth, U. S., & Raff, E. (2023). Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR.
- Black, K., Janner, M., Du, Y., Kostrikov, I., & Levine, S. (2024). Training diffusion models with reinforcement learning. In *The Twelfth International Conference on Learning Representations*.
- Bradley, R. A., & Terry, M. E. (1952). Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4), 324–345.

- Brooks, T., Holynski, A., & Efros, A. A. (2023). Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18392–18402.
- Brown, H., Lee, K., Mireshghallah, F., Shokri, R., & Tramèr, F. (2022). What does it mean for a language model to preserve privacy?. In *Proceedings of the 2022 ACM conference on fairness, accountability, and transparency*, pp. 2280–2292.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., & Askell, A. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877–1901.
- Cahyawijaya, S., Lovenia, H., Aji, A. F., Winata, G., Wilie, B., Koto, F., Mahendra, R., Wibisono, C., Romadhony, A., & Vincentio, K. (2023). Nusacrowd: Open source initiative for indonesian nlp resources. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13745–13818.
- Cahyawijaya, S., Lovenia, H., Koto, F., Putri, R. A., Dave, E., Lee, J., Shadieq, N., Cenggoro, W., Akbar, S. M., & Mahendra, M. I. (2024). Cendol: Open instruction-tuned generative large language models for indonesian languages. *arXiv preprint arXiv:2404.06138*.
- Cai, Z., Cao, M., Chen, H., Chen, K., Chen, K., Chen, X., Chen, X., Chen, Z., Chen, Z., & Chu, P. (2024). Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.
- Calandriello, D., Guo, D., Munos, R., Rowland, M., Tang, Y., Pires, B. A., Richemond, P. H., Lan, C. L., Valko, M., & Liu, T. (2024). Human alignment of large language models through online preference optimisation. *arXiv preprint arXiv:2403.08635*.
- Cao, H., Tan, C., Gao, Z., Xu, Y., Chen, G., Heng, P.-A., & Li, S. Z. (2024). A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering*.
- Chakraborty, S., Qiu, J., Yuan, H., Koppel, A., Huang, F., Manocha, D., Bedi, A. S., & Wang, M. (2024). Maxmin-rlhf: Towards equitable alignment of large language models with diverse human preferences. *arXiv preprint arXiv:2402.08925*.
- Chan, S. H. (2024). Tutorial on diffusion models for imaging and vision. *arXiv preprint arXiv:2403.18103*.
- Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., & Wong, E. (2023). Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*.
- Chen, C., & Shu, K. (2023). Can llm-generated misinformation be detected?. *arXiv preprint arXiv:2309.13788*.
- Chen, C., Hu, Y., Wu, W., Wang, H., Chng, E. S., & Zhang, C. (2024a). Enhancing zero-shot text-to-speech synthesis with human feedback. *arXiv preprint arXiv:2406.00654*.
- Chen, H., Zhao, H., Lam, H., Yao, D., & Tang, W. (2024b). Mallows-dpo: Fine-tune your llm with preference dispersions. *arXiv preprint arXiv:2405.14953*.
- Chen, J., Byun, J.-S., Elsner, M., & Perrault, A. (2024c). Reinforcement learning for fine-tuning text-to-speech diffusion models. *arXiv preprint arXiv:2405.14632*.

- Chen, L., Chen, J., Liu, C., Kirchenbauer, J., Soselia, D., Zhu, C., Goldstein, T., Zhou, T., & Huang, H. (2024d). Optune: Efficient online preference tuning. *arXiv preprint arXiv:2406.07657*.
- Chen, L., Li, J., Dong, X., Zhang, P., Zang, Y., Chen, Z., Duan, H., Wang, J., Qiao, Y., & Lin, D. (2024e). Are we on the right way for evaluating large vision-language models?. *arXiv preprint arXiv:2403.20330*.
- Chen, M., Mei, S., Fan, J., & Wang, M. (2024f). An overview of diffusion models: Applications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*.
- Chen, R., Zhang, X., Luo, M., Chai, W., & Liu, Z. (2025). PAD: Personalized alignment at decoding-time. In *The Thirteenth International Conference on Learning Representations*.
- Chen, T., Xu, B., Zhang, C., & Guestrin, C. (2016). Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*.
- Chen, Z., Zhou, K., Zhao, W. X., Wang, J., & Wen, J.-R. (2024). Low-redundant optimization for large language model alignment. *arXiv preprint arXiv:2406.12606*.
- Cheng, J., Liu, X., Zheng, K., Ke, P., Wang, H., Dong, Y., Tang, J., & Huang, M. (2024). Black-box prompt optimization: Aligning large language models without model training. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3201–3219.
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., & Gonzalez, J. E. (2023). Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3), 6.
- Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhang, H., Zhu, B., Jordan, M., & Gonzalez, J. E. (2024). Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., & Gehrmann, S. (2023). Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240), 1–113.
- Chowdhury, S. R., Kini, A., & Natarajan, N. (2024). Provably robust dpo: Aligning language models with noisy feedback. *arXiv preprint arXiv:2403.00409*.
- Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Clark, K., Vicol, P., Swersky, K., & Fleet, D. J. (2023). Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*.
- Conover, M., Hayes, M., Mathur, A., Xie, J., Wan, J., Shah, S., Ghodsi, A., Wendell, P., Zaharia, M., & Xin, R. (2023). Free dolly: Introducing the world’s first truly open instruction-tuned llm. *Company Blog of Databricks*.

- Costa-jussà, M. R., Cross, J., Çelebi, O., Elbayad, M., Heafield, K., Heffernan, K., Kalbassi, E., Lam, J., Licht, D., & Maillard, J. (2022). No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., Xie, G., Liu, Z., & Sun, M. (2023). Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*.
- Dai, J., Pan, X., Sun, R., Ji, J., Xu, X., Liu, M., Wang, Y., & Yang, Y. (2023a). Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.
- Dai, W., Li, J., Li, D., Tiong, A., Zhao, J., Wang, W., Li, B., Fung, P., & Hoi, S. (2023b). Instructblip: Towards general-purpose vision-language models with instruction tuning. arxiv 2023. *arXiv preprint arXiv:2305.06500*, 2.
- Dang, J., Ahmadian, A., Marchisio, K., Kreutzer, J., Üstün, A., & Hooker, S. (2024). Rlhf can speak many languages: Unlocking multilingual preference optimization for llms. *arXiv preprint arXiv:2407.02552*.
- Das, A., & Brunschwiler, T. (2019). Privacy is what we care about: Experimental investigation of federated learning on edge devices. In *International Workshop on Challenges in Artificial Intelligence and Machine Learning for Internet of Things held in conjunction with ACM SenSys*. Association for Computing Machinery, Inc.
- Deng, F., Wang, Q., Wei, W., Hou, T., & Grundmann, M. (2024). Prdp: Proximal reward difference prediction for large-scale reward finetuning of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7423–7433.
- Ding, N., Chen, Y., Xu, B., Qin, Y., Hu, S., Liu, Z., Sun, M., & Zhou, B. (2023). Enhancing chat language models by scaling high-quality instructional conversations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3029–3051.
- Domingo-Enrich, C., Drozdal, M., Karrer, B., & Chen, R. T. (2024). Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control. *arXiv preprint arXiv:2409.08861*.
- Dong, H., Xiong, W., Goyal, D., Zhang, Y., Chow, W., Pan, R., Diao, S., Zhang, J., Shum, K., & Zhang, T. (2023). Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*.
- Dong, H., Xiong, W., Pang, B., Wang, H., Zhao, H., Zhou, Y., Jiang, N., Sahoo, D., Xiong, C., & Zhang, T. (2024). Rlhf workflow: From reward modeling to online rlhf. *arXiv preprint arXiv:2405.07863*.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., & Fan, A. (2024). The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Dubois, Y., Galambosi, B., Liang, P., & Hashimoto, T. B. (2024). Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.

- Ethayarajh, K., Xu, W., Muennighoff, N., Jurafsky, D., & Kiela, D. (2024). Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*.
- Fakoor, R., Chaudhari, P., & Smola, A. J. (2020). P3o: Policy-on policy-off policy optimization. In *Uncertainty in artificial intelligence*, pp. 1017–1027. PMLR.
- Fan, A., Bhosale, S., Schwenk, H., Ma, Z., El-Kishky, A., Goyal, S., Baines, M., Celebi, O., Wenzek, G., & Chaudhary, V. (2021). Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*, 22(107), 1–48.
- Fan, Y., Watkins, O., Du, Y., Liu, H., Ryu, M., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Lee, K., & Lee, K. (2024). Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36.
- Feng, W., He, X., Fu, T.-J., Jampani, V., Akula, A., Narayana, P., Basu, S., Wang, X. E., & Wang, W. Y. (2022). Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*.
- Gao, S., Ge, Q., Shen, W., Dou, S., Ye, J., Wang, X., Zheng, R., Zou, Y., Chen, Z., Yan, H., et al. (2024). Linear alignment: A closed-form solution for aligning human preferences without tuning and feedback. *arXiv preprint arXiv:2401.11458*.
- Gong, T., Lyu, C., Zhang, S., Wang, Y., Zheng, M., Zhao, Q., Liu, K., Zhang, W., Luo, P., & Chen, K. (2023). Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*.
- Groeneveld, D., Beltagy, I., Walsh, P., Bhagia, A., Kinney, R., Tafjord, O., Jha, A. H., Ivison, H., Magnusson, I., & Wang, Y. (2024). Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*.
- Gulcehre, C., Paine, T. L., Srinivasan, S., Konyushkova, K., Weerts, L., Sharma, A., Sidhant, A., Ahern, A., Wang, M., & Gu, C. (2023). Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*.
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Del Giorno, A., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., & Saarikivi, O. (2023). Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.
- Guo, B., Zhang, X., Wang, Z., Jiang, M., Nie, J., Ding, Y., Yue, J., & Wu, Y. (2023). How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv:2301.07597*.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. (2025). Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Guo, S., Zhang, B., Liu, T., Liu, T., Khalman, M., Llinares, F., Rame, A., Mesnard, T., Zhao, Y., & Piot, B. (2024a). Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*.
- Guo, Y., Cui, G., Yuan, L., Ding, N., Wang, J., Chen, H., Sun, B., Xie, R., Zhou, J., & Lin, Y. (2024b). Controllable preference optimization: Toward controllable multi-objective alignment. *arXiv preprint arXiv:2402.19085*.

- Hartvigsen, T., Gabriel, S., Palangi, H., Sap, M., Ray, D., & Kamar, E. (2022). Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3309–3326.
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. In *Neurips*, Vol. 33, pp. 6840–6851.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., de Las Casas, D., Hendricks, L. A., Welbl, J., & Clark, A. (2022). Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, pp. 30016–30030.
- Hong, J., Lee, N., & Thorne, J. (2024). Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*, 2(4), 5.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Huang, A., Block, A., Liu, Q., Jiang, N., Foster, D. J., & Krishnamurthy, A. (2025). Is best-of-n the best of them? coverage, scaling, and optimality in inference-time alignment. *arXiv preprint arXiv:2503.21878*.
- Huang, S.-C., Li, P.-Z., Hsu, Y.-c., Chen, K.-M., Lin, Y. T., Hsiao, S.-K., Tsai, R., & Lee, H.-Y. (2024). Chat vector: A simple approach to equip llms with instruction following and model alignment in new languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 10943–10959.
- Hung, C.-Y., Majumder, N., Mehrish, A., & Poria, S. (2024). Inference time alignment with reward-guided tree search. *arXiv preprint arXiv:2406.15193*.
- Ilharco, G., Ribeiro, M. T., Wortsman, M., Schmidt, L., Hajishirzi, H., & Farhadi, A. (2023). Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Jang, J., Kim, S., Lin, B. Y., Wang, Y., Hessel, J., Zettlemoyer, L., Hajishirzi, H., Choi, Y., & Ammanabrolu, P. (2024). Personalized soups: Personalized large language model alignment via post-hoc parameter merging. In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*.
- Ji, J., Liu, M., Dai, J., Pan, X., Zhang, C., Bian, C., Chen, B., Sun, R., Wang, Y., & Yang, Y. (2024). Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36.
- Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y. J., Madotto, A., & Fung, P. (2023a). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1–38.
- Ji, Z., Yu, T., Xu, Y., Lee, N., Ishii, E., & Fung, P. (2023b). Towards mitigating llm hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 1827–1843.

- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., & Saulnier, L. (2023). Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., Casas, D. d. l., Hanna, E. B., & Bressand, F. (2024). Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Jiang, D., Ren, X., & Lin, B. Y. (2023). Llm-blender: Ensembling large language models with pairwise comparison and generative fusion. In *Proceedings of the 61th Annual Meeting of the Association for Computational Linguistics (ACL 2023)*.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. (2021). Advances and open problems in federated learning. *Foundations and trends® in machine learning*, 14(1–2), 1–210.
- Khanov, M., Burapachee, J., & Li, Y. (2024). ARGS: Alignment as reward-guided search. In *The Twelfth International Conference on Learning Representations*.
- Kim, D., Kim, Y., Song, W., Kim, H., Kim, Y., Kim, S., & Park, C. (2024). sdpo: Don't use your data all at once. *arXiv preprint arXiv:2403.19270*.
- Kim, H., Yu, Y., Jiang, L., Lu, X., Khashabi, D., Kim, G., Choi, Y., & Sap, M. (2022). Prosocialdialog: A prosocial backbone for conversational agents. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 4005–4029.
- Kim, S., Suk, J., Longpre, S., Lin, B. Y., Shin, J., Welleck, S., Neubig, G., Lee, M., Lee, K., & Seo, M. (2024). Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*.
- Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., & Levy, O. (2023). Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 36652–36663.
- Kool, W., van Hoof, H., & Welling, M. (2019). Buy 4 reinforce samples, get a baseline for free!. *Deep RL Meets Structured Prediction*.
- Köpf, A., Kilcher, Y., von Rütte, D., Anagnostidis, S., Tam, Z. R., Stevens, K., Barhoum, A., Nguyen, D., Stanley, O., & Nagyfi, R. (2024). Openassistant conversations-democratizing large language model alignment. *Advances in Neural Information Processing Systems*, 36.
- Kuang, N. L., Sun, W., McFaddin, S., Ma, Y., & Ettl, M. (2024). Towards personalized language models via inference-time human preference optimization. In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*.
- Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Devlin, J., & Lee, K. (2019). Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7, 453–466.

- Lai, V., Nguyen, C., Ngo, N., Nguyen, T., Deroncourt, F., Rossi, R., & Nguyen, T. (2023). Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 318–327.
- Lambert, N., Pyatkin, V., Morrison, J., Miranda, L., Lin, B. Y., Chandu, K., Dziri, N., Kumar, S., Zick, T., & Choi, Y. (2024). Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*.
- Laurencon, H., Saulnier, L., Wang, T., Akiki, C., Villanova del Moral, A., Le Scao, T., Von Werra, L., Mou, C., González Ponferrada, E., Nguyen, H., et al. (2022). The bigscience roots corpus: A 1.6 tb composite multilingual dataset. *Advances in Neural Information Processing Systems*, 35, 31809–31826.
- Le Scao, T., Fan, A., Akiki, C., Pavlick, E., Ilić, S., Hesslow, D., Castagné, R., Luccioni, A. S., Yvon, F., & Gallé, M. (2023). Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Lee, A., Bai, X., Pres, I., Wattenberg, M., Kummerfeld, J. K., & Mihalcea, R. (2024). A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. *arXiv preprint arXiv:2401.01967*.
- Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., & Gu, S. S. (2023). Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*.
- Lee, S., Park, S. H., Kim, S., & Seo, M. (2024). Aligning to thousands of preferences via system message generalization. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., & Zettlemoyer, L. (2020). Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880.
- Li, B., Zhang, Y., Chen, L., Wang, J., Pu, F., Yang, J., Li, C., & Liu, Z. (2023). Mimic-it: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*.
- Li, B., Wang, Y., Grama, A., & Zhang, R. (2024). Cascade reward sampling for efficient decoding-time alignment. *arXiv preprint arXiv:2406.16306*.
- Li, H., Koto, F., Wu, M., Aji, A. F., & Baldwin, T. (2023a). Bactrian-x: Multilingual replicable instruction-following models with low-rank adaptation. *arXiv preprint arXiv:2305.15011*.
- Li, K., Patel, O., Viégas, F., Pfister, H., & Wattenberg, M. (2023b). Inference-time intervention: Eliciting truthful answers from a language model. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Li, L., Xie, Z., Li, M., Chen, S., Wang, P., Chen, L., Yang, Y., Wang, B., & Kong, L. (2023c). Silk: Preference distillation for large visual language models. *arXiv preprint arXiv:2312.10665*.

- Li, R., Allal, L. B., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., Li, J., & Chim, J. (2023d). Starcoder: may the source be with you!. *arXiv preprint arXiv:2305.06161*.
- Li, S., Kallidromitis, K., Gokul, A., Kato, Y., & Kozuka, K. (2024a). Aligning diffusion models by optimizing human utility. *arXiv preprint arXiv:2404.04465*.
- Li, X., Yong, Z.-X., & Bach, S. H. (2024b). Preference tuning for toxicity mitigation generalizes across languages. *arXiv preprint arXiv:2406.16235*.
- Li, Y., Hu, X., Qu, X., Li, L., & Cheng, Y. (2025). Test-time preference optimization: On-the-fly alignment via iterative textual feedback. *arXiv preprint arXiv:2501.12895*.
- Li, Y., Bubeck, S., Eldan, R., Del Giorno, A., Gunasekar, S., & Lee, Y. T. (2023a). Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*.
- Li, Y., Wei, F., Zhao, J., Zhang, C., & Zhang, H. (2023b). Rain: Your language models can align themselves without finetuning. *arXiv preprint arXiv:2309.07124*.
- Li, Z., Xu, T., Zhang, Y., Lin, Z., Yu, Y., Sun, R., & Luo, Z.-Q. (2023c). Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. In *Forty-first International Conference on Machine Learning*.
- Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., & Kumar, A. (2022). Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.
- Liang, Y., He, J., Li, G., Li, P., Klimovskiy, A., Carolan, N., Sun, J., Pont-Tuset, J., Young, S., & Yang, F. (2024). Rich human feedback for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19401–19411.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., & Cobbe, K. (2024). Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Lin, B. Y., Ravichander, A., Lu, X., Dziri, N., Sclar, M., Chandu, K., Bhagavatula, C., & Choi, Y. (2024). The unlocking spell on base LLMs: Rethinking alignment via in-context learning. In *The Twelfth International Conference on Learning Representations*.
- Liu, H., Li, C., Li, Y., & Lee, Y. J. (2024a). Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26296–26306.
- Liu, H., Li, C., Wu, Q., & Lee, Y. J. (2024b). Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Liu, N., Li, S., Du, Y., Torralba, A., & Tenenbaum, J. B. (2022). Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pp. 423–439. Springer.
- Liu, R., Garrette, D., Saharia, C., Chan, W., Roberts, A., Narang, S., Blok, I., Mical, R., Norouzi, M., & Constant, N. (2023). Character-aware models improve visual

- text rendering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16270–16297.
- Liu, T., Guo, S., Bianco, L., Calandriello, D., Berthet, Q., Llinares-López, F., Hoffmann, J., Dixon, L., Valko, M., & Blondel, M. (2024). Decoding-time realignment of language models. In *Forty-first International Conference on Machine Learning*.
- Liu, T., Zhao, Y., Joshi, R., Khalman, M., Saleh, M., Liu, P. J., & Liu, J. (2023). Statistical rejection sampling improves preference optimization. *arXiv preprint arXiv:2309.06657*.
- Liu, Y. (2020). Multilingual denoising pre-training for neural machine translation. *arXiv preprint arXiv:2001.08210*.
- Liu, Y., Zhang, K., Li, Y., Yan, Z., Gao, C., Chen, R., Yuan, Z., Huang, Y., Sun, H., & Gao, J. (2024). Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., & Lin, M. (2025). Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.
- Longpre, S., Hou, L., Vu, T., Webson, A., Chung, H. W., Tay, Y., Zhou, D., Le, Q. V., Zoph, B., & Wei, J. (2023). The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pp. 22631–22648. PMLR.
- Lovenia, H., Mahendra, R., Akbar, S. M., Miranda, L. J. V., Santoso, J., Aco, E., Fadhilah, A., Mansurov, J., Imperial, J. M., & Kampman, O. P. (2024). Seacrowd: A multilingual multimodal data hub and benchmark suite for southeast asian languages. *arXiv preprint arXiv:2406.10118*.
- Lu, Y., Zhu, W., Li, L., Qiao, Y., & Yuan, F. (2024). Llamax: Scaling linguistic horizons of llm by enhancing translation capabilities beyond 100 languages. *arXiv preprint arXiv:2407.05975*.
- Luo, C. (2022). Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*.
- Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., & Jiang, D. (2023). Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*.
- Maheshwary, R., Yadav, V., Nguyen, H., Mahajan, K., & Madhusudhan, S. T. (2024). M2lingual: Enhancing multilingual, multi-turn instruction alignment in large language models. *arXiv preprint arXiv:2406.16783*.
- Maiti, S., Peng, Y., Saeki, T., & Watanabe, S. (2023). Speechlmscore: Evaluating speech generation using speech language model. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE.
- Manyika, J., & Hsiao, S. (2023). An overview of bard: an early experiment with generative ai. *AI. Google Static Documents*, 2.
- Meng, Y., Xia, M., & Chen, D. (2024). Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.

- Muennighoff, N., Wang, T., Sutawika, L., Roberts, A., Biderman, S., Le Scao, T., Bari, M. S., Shen, S., Yong, Z. X., & Schoelkopf, H. (2023). Crosslingual generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15991–16111.
- Mukobi, G., Chatain, P., Fong, S., Windesheim, R., Kutyniok, G., Bhatia, K., & Alberti, S. (2023). Superhf: Supervised iterative learning from human feedback. *arXiv preprint arXiv:2310.16763*.
- Munos, R., Valko, M., Calandriello, D., Azar, M. G., Rowland, M., Guo, Z. D., Tang, Y., Geist, M., Mesnard, T., & Michi, A. (2023). Nash learning from human feedback. *arXiv preprint arXiv:2312.00886*.
- Nguyen, X.-P., Zhang, W., Li, X., Aljunied, M., Tan, Q., Cheng, L., Chen, G., Deng, Y., Yang, S., & Liu, C. (2023). Seallms—large language models for southeast asia. *arXiv preprint arXiv:2312.00738*.
- Ni, J., Xue, F., Deng, Y., Phang, J., Jain, K., Shah, M. H., Zheng, Z., & You, Y. (2023). Instruction in the wild: A user-based instruction dataset. *GitHub repository*.
- Ormazabal, A., Zheng, C., d’Autume, C. d. M., Yogatama, D., Fu, D., Ong, D., Chen, E., Lamprecht, E., Pham, H., & Ong, I. (2024). Reka core, flash, and edge: A series of powerful multimodal language models. *arXiv preprint arXiv:2404.12387*.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., & Ray, A. (2022). Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35, 27730–27744.
- Panagopoulou, A., Xue, L., Yu, N., Li, J., Li, D., Joty, S., Xu, R., Savarese, S., Xiong, C., & Niebles, J. C. (2023). X-instructblip: A framework for aligning x-modal instruction-aware representations to llms and emergent cross-modal reasoning. *arXiv preprint arXiv:2311.18799*.
- Park, R., Rafailov, R., Ermon, S., & Finn, C. (2024). Disentangling length from quality in direct preference optimization. *arXiv preprint arXiv:2403.19159*.
- Parmar, J., Prabhumoye, S., Jennings, J., Patwary, M., Subramanian, S., Su, D., Zhu, C., Narayanan, D., Jhunjhunwala, A., & Dattagupta, A. (2024). Nemotron-4 15b technical report. *arXiv preprint arXiv:2402.16819*.
- Prabhudesai, M., Goyal, A., Pathak, D., & Fragkiadaki, K. (2023). Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*.
- Qi, X., Panda, A., Lyu, K., Ma, X., Roy, S., Beirami, A., Mittal, P., & Henderson, P. (2024). Safety alignment should be made more than just a few tokens deep. *arXiv preprint arXiv:2406.05946*.
- Qiao, Y., Li, K., Lin, J., Wei, R., Jiang, C., Luo, Y., & Yang, H. (2024). Robust domain generalization for multi-modal object recognition. *arXiv preprint arXiv:2408.05831*.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., & Clark, J. (2021). Learning transferable visual models from natural

- language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR.
- Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., & Young, S. (2021). Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., & Finn, C. (2024). Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140), 1–67.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2), 3.
- Rando, J., & Tramèr, F. (2023). Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*.
- Reid, M., Savinov, N., Teplyashin, D., Lepikhin, D., Lillicrap, T., Alayrac, J.-b., Soricut, R., Lazaridou, A., Firat, O., & Schrittwieser, J. (2024). Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695.
- Saeki, T., Maiti, S., Takamichi, S., Watanabe, S., & Saruwatari, H. (2024). Speechbertscore: Reference-aware automatic evaluation of speech generation leveraging nlp evaluation metrics. *arXiv preprint arXiv:2401.16812*.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., & Salimans, T. (2022). Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35, 36479–36494.
- Sanh, V., Webson, A., Raffel, C., Bach, S. H., Sutawika, L., Alyafeai, Z., Chaffin, A., Stiegler, A., Scio, T. L., & Raja, A. (2021). Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Shahriar, S., Qi, Z., Pappas, N., Doss, S., Sunkara, M., Halder, K., Mager, M., & Benajiba, Y. (2024). Inference time llm alignment in single and multidomain preference spectrum. *arXiv preprint arXiv:2410.19206*.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al. (2024). Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Shazeer, N. (2020). Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*.

- She, S., Huang, S., Zou, W., Zhu, W., Liu, X., Geng, X., & Chen, J. (2024). Mapo: Advancing multilingual reasoning through multilingual alignment-as-preference optimization. *arXiv preprint arXiv:2401.06838*.
- Shirafuji, D., Takenaka, M., & Taguchi, S. (2025). Bias vector: Mitigating biases in language models with task arithmetic approach. In *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 2799–2813.
- Shoeybi, M., Patwary, M., Puri, R., LeGresley, P., Casper, J., & Catanzaro, B. (2019). Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*.
- Singh, S., Vargus, F., Dsouza, D., Karlsson, B. F., Mahendiran, A., Ko, W.-Y., Shandilya, H., Patel, J., Mataciunas, D., & OMahony, L. (2024). Aya dataset: An open-access collection for multilingual instruction tuning. *arXiv preprint arXiv:2402.06619*.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. PMLR.
- Soltan, S., Ananthakrishnan, S., FitzGerald, J., Gupta, R., Hamza, W., Khan, H., Peris, C., Rawls, S., Rosenbaum, A., & Rumshisky, A. (2022). Alexatm 20b: Few-shot learning using a large-scale multilingual seq2seq model. *arXiv preprint arXiv:2208.01448*.
- Song, F., Fan, Y., Zhang, X., Wang, P., & Wang, H. (2024a). Icdpo: Effectively borrowing alignment capability of others via in-context direct preference optimization. *arXiv preprint arXiv:2402.09320*.
- Song, F., Yu, B., Li, M., Yu, H., Huang, F., Li, Y., & Wang, H. (2024b). Preference ranking optimization for human alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38, pp. 18990–18998.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., & Christiano, P. F. (2020). Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33, 3008–3021.
- Sun, Y., He, J., Cui, L., Lei, S., & Lu, C.-T. (2024). Exploring the deceptive power of llm-generated fake news: A study of real-world detection challenges. *arXiv preprint arXiv:2403.18249*.
- Sun, Z., Shen, S., Cao, S., Liu, H., Li, C., Shen, Y., Gan, C., Gui, L.-Y., Wang, Y.-X., & Yang, Y. (2023). Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*.
- Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12.
- Tang, W. (2024). Fine-tuning of diffusion models via stochastic control: entropy regularization and beyond. *arXiv preprint arXiv:2403.06279*.
- Tang, W., & Zhao, H. (2024). Score-based diffusion models via stochastic differential equations—a technical tutorial. *arXiv preprint arXiv:2402.07487*.

- Tang, Y., Guo, D. Z., Zheng, Z., Calandriello, D., Cao, Y., Tarassov, E., Munos, R., Pires, B. Á., Valko, M., & Cheng, Y. (2024a). Understanding the performance gap between online and offline alignment algorithms. *arXiv preprint arXiv:2405.08448*.
- Tang, Y., Guo, Z. D., Zheng, Z., Calandriello, D., Munos, R., Rowland, M., Richemond, P. H., Valko, M., Pires, B. Á., & Piot, B. (2024b). Generalized preference optimization: A unified approach to offline alignment. *arXiv preprint arXiv:2402.05749*.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., & Hashimoto, T. B. (2023). Stanford alpaca: an instruction-following llama model (2023). *URL https://github.com/tatsu-lab/stanford_alpaca, 1(9)*.
- Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., & Hauth, A. (2023). Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Team, G., Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M. S., & Love, J. (2024). Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Tonmoy, S., Zaman, S., Jain, V., Rani, A., Rawte, V., Chadha, A., & Das, A. (2024). A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313*.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., & Azhar, F. (2023a). Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., & Bhosale, S. (2023b). Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Uehara, M., Zhao, Y., Biancalani, T., & Levine, S. (2024a). Understanding reinforcement learning-based fine-tuning of diffusion models: A tutorial and review. *arXiv preprint arXiv:2407.13734*.
- Uehara, M., Zhao, Y., Black, K., Hajiramezanali, E., Scalia, G., Diamant, N. L., Tseng, A. M., Biancalani, T., & Levine, S. (2024b). Fine-tuning of continuous-time diffusion models as entropy-regularized control. *arXiv preprint arXiv:2402.15194*.
- Uehara, M., Zhao, Y., Black, K., Hajiramezanali, E., Scalia, G., Diamant, N. L., Tseng, A. M., Levine, S., & Biancalani, T. (2024c). Feedback efficient online fine-tuning of diffusion models. *arXiv preprint arXiv:2402.16359*.
- Uehara, M., Zhao, Y., Hajiramezanali, E., Scalia, G., Eraslan, G., Lal, A., Levine, S., & Biancalani, T. (2024d). Bridging model-based optimization and generative modeling via conservative fine-tuning of diffusion models. *Advances in Neural Information Processing Systems*, 37, 127511–127535.
- Üstün, A., Aryabumi, V., Yong, Z.-X., Ko, W.-Y., D’souza, D., Onilude, G., Bhandari, N., Singh, S., Ooi, H.-L., & Kayid, A. (2024). Aya model: An instruction finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*.

- Vovk, V., & Wang, R. (2021). E-values: Calibration, combination and applications. *The Annals of Statistics*, 49(3), 1736–1754.
- Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A., Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., & Naik, N. (2024). Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8228–8238.
- Wang, C., Jiang, Y., Yang, C., Liu, H., & Chen, Y. (2023). Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints. *arXiv preprint arXiv:2309.16240*.
- Wang, F.-Y., Shui, Y., Piao, J., Sun, K., & Li, H. (2025). Diffusion-npo: Negative preference optimization for better preference aligned generation of diffusion models. In *The Thirteenth International Conference on Learning Representations*.
- Wang, H., Xiong, W., Xie, T., Zhao, H., & Zhang, T. (2024). Interpretable preferences via multi-objective reward modeling and mixture-of-experts. *arXiv preprint arXiv:2406.12845*.
- Wang, Y., Ivison, H., Dasigi, P., Hessel, J., Khot, T., Chandu, K., Wadden, D., MacMillan, K., Smith, N. A., Beltagy, I., et al. (2023a). How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36, 74764–74786.
- Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A., Khashabi, D., & Hajishirzi, H. (2023b). Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–13508.
- Wang, Y., Mishra, S., Alipoormolabashi, P., Kordi, Y., Mirzaei, A., Naik, A., Ashok, A., Dhanasekaran, A. S., Arunkumar, A., & Stap, D. (2022). Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5085–5109.
- Wang, Z., Dong, Y., Delalleau, O., Zeng, J., Shen, G., Egert, D., Zhang, J. J., Sreedhar, M. N., & Kuchaiev, O. (2024). Helpsteer2: Open-source dataset for training top-performing reward models. *arXiv preprint arXiv:2406.08673*.
- Wei, A., Haghtalab, N., & Steinhardt, J. (2024). Jailbroken: How does LLM safety training fail?. *Advances in Neural Information Processing Systems*, 36.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35, 24824–24837.
- Weyssow, M., Kamanda, A., & Sahraoui, H. (2024). Codeultrafeedback: An llm-as-a-judge dataset for aligning large language models to coding preferences. *arXiv preprint arXiv:2403.09032*.
- Williams, R. J. (1987). *Reinforcement-learning connectionist systems*. College of Computer Science, Northeastern University.

- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8, 229–256.
- Winata, G. I., Anugraha, D., Susanto, L., Kuwanto, G., & Wijaya, D. T. (2024). Metametrics: Calibrating metrics for generation tasks using human preferences. *arXiv preprint arXiv:2410.02381*.
- Winata, G. I., Cahyawijaya, S., Liu, Z., Lin, Z., Madotto, A., & Fung, P. (2021a). Are multilingual models effective in code-switching?. In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pp. 142–153.
- Winata, G. I., Madotto, A., Lin, Z., Liu, R., Yosinski, J., & Fung, P. (2021b). Language models are few-shot multilingual learners. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pp. 1–15.
- Winata, G. I., Zhang, R., & Adelani, D. I. (2024). Miners: Multilingual language models as semantic retrievers. *arXiv preprint arXiv:2406.07424*.
- Wortsman, M., Ilharco, G., Gadre, S. Y., Roelofs, R., Gontijo-Lopes, R., Morcos, A. S., Namkoong, H., Farhadi, A., Carmon, Y., Kornblith, S., et al. (2022). Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pp. 23965–23998. PMLR.
- Wu, F., Inan, H. A., Backurs, A., Chandrasekaran, V., Kulkarni, J., & Sim, R. (2023a). Privately aligning language models with reinforcement learning. *arXiv preprint arXiv:2310.16960*.
- Wu, S., Irsoy, O., Lu, S., Dabrovolski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D., & Mann, G. (2023b). Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*.
- Wu, T., Zhu, B., Zhang, R., Wen, Z., Ramchandran, K., & Jiao, J. (2023c). Pairwise proximal policy optimization: Harnessing relative feedback for llm alignment. *arXiv preprint arXiv:2310.00212*.
- Wu, X., Sun, K., Zhu, F., Zhao, R., & Li, H. (2023d). Human preference score: Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2096–2105.
- Wu, Y., Sun, Z., Yuan, H., Ji, K., Yang, Y., & Gu, Q. (2024). Self-play preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*.
- Xiong, W., Dong, H., Ye, C., Wang, Z., Zhong, H., Ji, H., Jiang, N., & Zhang, T. (2024). Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*.
- Xu, C., Sun, Q., Zheng, K., Geng, X., Zhao, P., Feng, J., Tao, C., & Jiang, D. (2023). Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.
- Xu, H., Sharaf, A., Chen, Y., Tan, W., Shen, L., Durme, B. V., Murray, K., & Kim, Y. J. (2024a). Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In *Forty-first International Conference on Machine Learning*.

- Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., & Dong, Y. (2024b). Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36.
- Xu, J., Lee, A., Sukhbaatar, S., & Weston, J. (2023). Some things are more cringe than others: Preference optimization with the pairwise cringe loss. *arXiv preprint arXiv:2312.16682*.
- Xu, S., Fu, W., Gao, J., Ye, W., Liu, W., Mei, Z., Wang, G., Yu, C., & Wu, Y. (2024a). Is dpo superior to ppo for llm alignment? a comprehensive study. *arXiv preprint arXiv:2404.10719*.
- Xu, Z., Jiang, F., Niu, L., Deng, Y., Poovendran, R., Choi, Y., & Lin, B. Y. (2024b). Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *arXiv preprint arXiv:2406.08464*.
- Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Zhang, W., Cui, B., & Yang, M.-H. (2023). Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4), 1–39.
- Yang, R., Ding, R., Lin, Y., Zhang, H., & Zhang, T. (2024a). Regularizing hidden states enables learning generalizable reward model for llms. *arXiv preprint arXiv:2406.10216*.
- Yang, S., Chen, T., & Zhou, M. (2024b). A dense reward view on aligning text-to-image diffusion with preference. *arXiv preprint arXiv:2402.08265*.
- Yang, S., Zhang, S., Xia, C., Feng, Y., Xiong, C., & Zhou, M. (2024c). Preference-grounded token-level guidance for language model fine-tuning. *Advances in Neural Information Processing Systems*, 36.
- Yao, J.-Y., Ning, K.-P., Liu, Z.-H., Ning, M.-N., & Yuan, L. (2023a). Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*.
- Yao, Y., Xu, X., & Liu, Y. (2023b). Large language model unlearning. *arXiv preprint arXiv:2310.10683*.
- Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., & Wu, Y. (2022a). Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*.
- Yu, J., Xu, Y., Koh, J. Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., & Ayan, B. K. (2022b). Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3), 5.
- Yu, T., Hu, J., Yao, Y., Zhang, H., Zhao, Y., Wang, C., Wang, S., Pan, Y., Xue, J., & Li, D. (2023). Reformulating vision-language foundation models and datasets towards universal multimodal assistants. *arXiv preprint arXiv:2310.00653*.
- Yu, T., Yao, Y., Zhang, H., He, T., Han, Y., Cui, G., Hu, J., Liu, Z., Zheng, H.-T., & Sun, M. (2024). Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13807–13816.

- Yuan, L., Cui, G., Wang, H., Ding, N., Wang, X., Deng, J., Shan, B., Chen, H., Xie, R., & Lin, Y. (2024a). Advancing llm reasoning generalists with preference trees. *arXiv preprint arXiv:2404.02078*.
- Yuan, W., Pang, R. Y., Cho, K., Sukhbaatar, S., Xu, J., & Weston, J. (2024b). Self-rewarding language models. *arXiv preprint arXiv:2401.10020*.
- Yuan, Z., Yuan, H., Tan, C., Wang, W., Huang, S., & Huang, F. (2023). Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*.
- Zhang, D., Zhoubian, S., Hu, Z., Yue, Y., Dong, Y., & Tang, J. (2024). ReST-MCTS*: LLM self-training via process reward guided tree search. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Zhang, D., Li, S., Zhang, X., Zhan, J., Wang, P., Zhou, Y., & Qiu, X. (2023). Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv preprint arXiv:2305.11000*.
- Zhang, D., Li, Z., Li, S., Zhang, X., Wang, P., Zhou, Y., & Qiu, X. (2024). Speechalign: Aligning speech generation to human preferences. *arXiv preprint arXiv:2404.05600*.
- Zhang, H., Li, X., & Bing, L. (2023). Video-llama: An instruction-tuned audio-visual language model for video understanding. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 543–553.
- Zhang, R., Lin, L., Bai, Y., & Mei, S. (2024a). Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*.
- Zhang, S., Yang, X., Feng, Y., Qin, C., Chen, C.-C., Yu, N., Chen, Z., Wang, H., Savarese, S., & Ermon, S. (2024b). Hive: Harnessing human feedback for instructional visual editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9026–9036.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., & Lin, X. V. (2022). Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Zhang, Z., Bai, F., Chen, Q., Ma, C., Wang, M., Sun, H., Zheng, Z., & Yang, Y. (2025). Amulet: Realignment during test time for personalized preference adaptation of LLMs. In *The Thirteenth International Conference on Learning Representations*.
- Zhao, H., Chen, H., Guo, Y., Winata, G. I., Ou, T., Huang, Z., Yao, D. D., & Tang, W. (2025a). Fine-tuning diffusion generative models via rich preference optimization. *arXiv preprint arXiv:2503.11720*.
- Zhao, H., Chen, H., Zhang, J., Yao, D. D., & Tang, W. (2025b). Score as action: Fine-tuning diffusion generative models by continuous-time reinforcement learning. *arXiv preprint arXiv:2502.01819*.
- Zhao, H., Winata, G. I., Das, A., Zhang, S.-X., Yao, D. D., Tang, W., & Sahu, S. (2024a). Rainbowpo: A unified framework for combining improvements in preference optimization. *arXiv preprint arXiv:2410.04203*.

- Zhao, S., Brekelmans, R., Makhzani, A., & Grosse, R. (2024b). Probabilistic inference in language models via twisted sequential monte carlo. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Zhao, W., Ren, X., Hessel, J., Cardie, C., Choi, Y., & Deng, Y. (2024c). Wildchat: 1m chatgpt interaction logs in the wild. *arXiv preprint arXiv:2405.01470*.
- Zhao, Y., Joshi, R., Liu, T., Khalman, M., Saleh, M., & Liu, P. J. (2023). Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.
- Zhao, Y., Khalman, M., Joshi, R., Narayan, S., Saleh, M., & Liu, P. J. (2022). Calibrating sequence likelihood improves conditional language generation. In *The eleventh international conference on learning representations*.
- Zheng, C., Wang, Z., Ji, H., Huang, M., & Peng, N. (2024a). Weak-to-strong extrapolation expedites alignment. *arXiv preprint arXiv:2404.16792*.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., & Xing, E. (2024b). Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Zheng, X., Pang, T., Du, C., Liu, Q., Jiang, J., & Lin, M. (2024c). Cheating automatic llm benchmarks: Null models achieve high win rates. *arXiv preprint arXiv:2410.07137*.
- Zheng, Z., Peng, P., Ma, Z., Chen, X., Choi, E., & Harwath, D. (2024d). Bat: Learning to reason about spatial sounds with large language models. *arXiv preprint arXiv:2402.01591*.
- Zhong, H., Feng, G., Xiong, W., Zhao, L., He, D., Bian, J., & Wang, L. (2024). Dpo meets ppo: Reinforced token optimization for rlhf. *arXiv preprint arXiv:2404.18922*.
- Zhou, C., Liu, P., Xu, P., Iyer, S., Sun, J., Mao, Y., Ma, X., Efrat, A., Yu, P., & Yu, L. (2024a). Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.
- Zhou, Y., Cui, C., Rafailov, R., Finn, C., & Yao, H. (2024b). Aligning modalities in vision large language models via preference fine-tuning. *arXiv preprint arXiv:2402.11411*.
- Zhu, B., Frick, E., Wu, T., Zhu, H., Ganesan, K., Chiang, W.-L., Zhang, J., & Jiao, J. (2024). Starling-7b: Improving helpfulness and harmlessness with rlaf. In *First Conference on Language Modeling*.
- Zhu, M., Liu, Y., Zhang, L., Guo, J., & Mao, Z. (2025). On-the-fly preference alignment via principle-guided decoding. In *The Thirteenth International Conference on Learning Representations*.
- Zhuang, Z., Lei, K., Liu, J., Wang, D., & Guo, Y. (2023). Behavior proximal policy optimization. *arXiv preprint arXiv:2302.11312*.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., & Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.