

# WARP: ON THE BENEFITS OF WEIGHT AVERAGED REWARDED POLICIES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement learning from human feedback (RLHF) aligns large language models by encouraging their generations to have high rewards, using a reward model trained on human preferences. To prevent forgetting of pre-trained knowledge, RLHF usually incorporates a KL regularization; this forces the policy to remain close to its initialization, though it hinders the reward optimization. To address the trade-off between KL and reward, in this paper we introduce a novel alignment strategy named Weight Averaged Rewarded Policies (*WARP*), merging policies in the weight space at three distinct stages. First, it uses the exponential moving average of the policy as a dynamic anchor in the KL regularization. Second, it applies spherical interpolation to merge independently fine-tuned policies into a new enhanced one. Third, it linearly interpolates between this merged model and the initialization, to recover features from pre-training. This procedure is then applied iteratively, with each iteration’s final model used as an advanced initialization for the next, progressively refining the KL-reward trade-off, achieving superior rewards at fixed KL. Experiments with Gemma policies validate that *WARP* improves their quality and alignment, outperforming open-source models.

## 1 INTRODUCTION

**LLM alignment.** Large language models (LLMs) like Gemini (Gemini Team, 2023) and GPT-4 (OpenAI, 2023), along with their open-weight counterparts (Jiang et al., 2023; Gemma Team et al., 2024), demonstrate remarkable abilities as chatbots, but also for tasks like mathematics and coding (Bubeck et al., 2023). These capabilities largely emerge from pre-training on next-token prediction (Radford et al., 2018; 2019), subsequently refined through supervised fine-tuning (SFT) (Raffel et al., 2020; Wei et al., 2022). As these LLMs become more powerful, aligning them with human values becomes increasingly crucial to ensure safe deployment (Amodei et al., 2016; Hendrycks & Mazeika, 2022). To this end, reinforcement learning from human feedback (RLHF) has become the prominent strategy (Christiano et al., 2017; Ziegler et al., 2019; Stiennon et al., 2020), first learning a reward model (RM) on human preferences, before optimizing the LLM to maximize predicted rewards.

**Challenges in RLHF.** However, RLHF introduces several unresolved challenges (Casper et al., 2023). First, the limited scope of fine-tuning, often restricted to relatively small datasets, can lead to excessive specialization and catastrophic forgetting (French, 1992) of the broad and diverse knowledge acquired during pre-training (Goodfellow et al., 2013; Li & Hoiem, 2017; Kirkpatrick et al., 2017; Kumar et al., 2022). Such *alignment tax* (Ouyang et al., 2022) can degrade the LLM’s reasoning capabilities and performance on NLP benchmarks (Dong et al., 2023a; Lin et al., 2024a). Second, maximizing an imperfect RM presents several issues on its own, as the LLM can learn to exploit loopholes in the RM (Clark & Amodei, 2016; Pan et al., 2022) when it deviates significantly from its initialization (Gao et al., 2023). Such *reward hacking* (Askill et al., 2021; Skalse et al., 2022) can produce outputs that are linguistically flawed (Lewis et al., 2017), excessively verbose (Singhal et al., 2023), or sycophantic (Perez et al., 2022; Sharma et al., 2023), thereby raising misalignment (Taylor et al., 2016; Ngo et al., 2022) and safety (Amodei et al., 2016; Hendrycks & Mazeika, 2022) concerns. Finally, RLHF can reduce the diversity of generations (Kirk et al., 2024), potentially leading to policy collapse (Moalla et al., 2024; Hamilton, 2024). Such *loss of diversity* limits use in creative or exploratory tasks and can result in the LLM systematically refusing to answer. Overall, achieving high rewards based on an imperfect RM on a selected distribution of prompts is insufficient due to potential reward misspecification and distribution shifts upon deployment.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

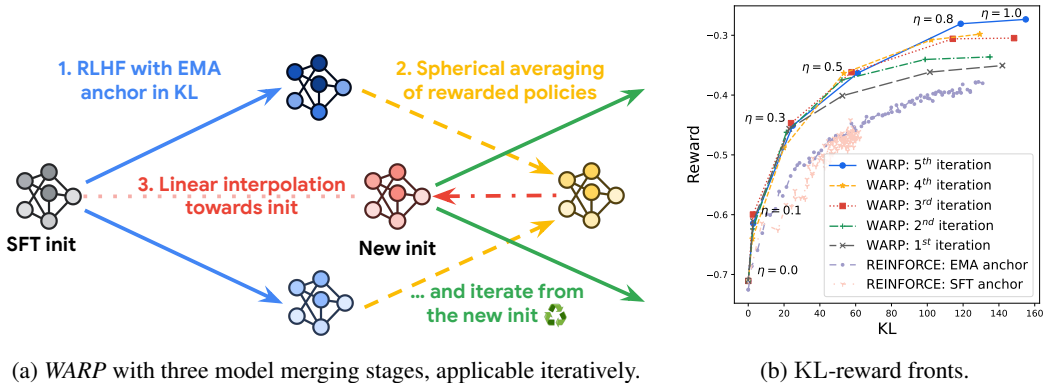


Figure 1: Figure 1(a) illustrates the RLHF alignment process with WARP from a supervised fine-tuned (SFT) LLM. WARP uses model merging by weight averaging at three different stages. First, the exponential moving average (EMA) (Izmailov et al., 2018) of the policy serves as the anchor for KL regularization (Jaques et al., 2017). Second, the independently fine-tuned policies are merged by spherical linear interpolation (SLERP) (Shoemaker, 1985) of their task vectors (Ilharco et al., 2023). Third, we interpolate towards the initialization (LITI) (Wortsman et al., 2022b), revealing a Pareto front of solutions as we slide the interpolating coefficient  $\eta$  from 1 to 0. This results in the “WARP: 1<sup>st</sup> iteration” curve from Figure 1(b) which improves over the REINFORCE (Williams, 1992) fine-tuning trajectories. Critically, iteratively using a point from this front as an advanced initialization for the next episode WARP improves performance. Details in Figure 4(c).

**RL with KL regularization.** To address these issues, previous works constrained the reward optimization by integrating a Kullback-Leibler (KL) regularization (Jaques et al., 2017; Geist et al., 2019), using the SFT initialization as the anchor. As clarified in Section 2, this KL regularization forces the policy to remain close to its initialization (Lazaridou et al., 2020; Lu et al., 2020), mitigating forgetting and reward hacking (Gao et al., 2023). However, employing the SFT model as the anchor may lead to reward underfitting: indeed, there is a fundamental tension between reducing KL and maximizing reward. Thus, different policies should be compared in terms of trade-off between KL-reward as in Figure 1(b), where the  $x$ -axis is the KL and the  $y$ -axis is the reward as estimated by the RM, with the optimal policies located in the top-left of the plot.

**On model merging by weight averaging.** To improve the trade-off between KL and reward during RLHF, we leverage the ability to merge LLMs by weight averaging (WA) (Utans, 1996). WA relies on the linear mode connectivity (Frankle et al., 2020; Neyshabur et al., 2020), an empirical observation revealing linear paths of high performance between models fine-tuned from a shared pre-trained initialization. Model merging was shown to improve robustness under distribution shifts (Izmailov et al., 2018; Wortsman et al., 2022a; Ramé et al., 2022) by promoting generalization and reducing memorization (Ramé et al., 2024), to combine models’ abilities (Ilharco et al., 2023; 2022; Ramé et al., 2023), to reduce forgetting in continual learning (Stojanovski et al., 2022), to enable collaborative (Raffel, 2023) and distributed (Douillard et al., 2023) learning at scale, without computational overheads at inference time. Model merging is increasingly adopted within the open-source community (Goddard et al., 2024; Lambert & Morrison, 2024), leading to state-of-the-art models in specialized domains (Labrak et al., 2024) but also significant advancements on general-purpose benchmarks (Labonne, 2024b;a). In particular, while WA was initially mostly used for discriminative tasks (Wortsman et al., 2022a) such as reward modeling (Ramé et al., 2024), it is now becoming popular for generative tasks (Rofin et al., 2022; Akiba et al., 2024); its use in KL-constrained RLHF has already shown preliminary successes in a few recent works (Ramé et al., 2023; Noukhovitch et al., 2023; Lin et al., 2024a; Liu et al., 2024; Gorbatovski et al., 2024; Munos et al., 2023), further elaborated in Section 5.

**WARP.** In this paper, we propose Weight Averaged Rewarded Policies (WARP), a simple strategy for aligning LLMs, illustrated in Figure 1(a) and detailed in Section 3. WARP is designed to optimize the KL-reward Pareto front of solutions, as demonstrated in Figure 1(b). WARP uses three variants of WA at three different stages of the alignment procedure, for three distinct reasons.

**Stage 1: Exponential Moving Average (EMA).** During RL fine-tuning, instead of regularizing the policy towards the SFT initialization, *WARP* uses the policy’s own exponential moving average (Polyak & Juditsky, 1992) as a dynamic updatable anchor in the KL. This stage enables stable exploration with distillation from an EMA teacher (Tarvainen & Valpola, 2017) and annealed constraint.

**Stage 2: Spherical Linear intERPolation of task vectors (SLERP).** Considering  $M$  policies RL fine-tuned independently with their own *EMA* anchor, we merge them by spherical linear interpolation (Shoemaker, 1985) of their task vectors (Ilharco et al., 2023). This stage creates a merged model with higher reward by combining the strengths of the  $M$  individual policies.

**Stage 3: Linear Interpolation Towards Initialization (LITI).** Considering the merged policy from *SLERP*, *WARP* linearly interpolates towards the initialization, akin to WiSE-FT (Wortsman et al., 2022b). This stage allows to run through an improved Pareto front simply by adjusting the interpolating coefficient  $\eta$  between 1 (high reward but high KL) and 0 (small KL but small reward). Critically, selecting an intermediate value for  $0 < \eta < 1$  offers a balanced model that can serve as a new and improved initialization for subsequent iterations of *WARP*.

**Experiments and discussion.** In Section 4, we validate the efficacy of *WARP* for the fine-tuning of Gemma 7B (Gemma Team et al., 2024). Finally, in Section 6, we discuss the connections between *WARP*, the distributed learning literature (Raffel, 2023; Douillard et al., 2023) and iterated amplification (Christiano et al., 2018), illustrating how *WARP* embodies their principles to enable scaling post-training, for continuous alignment and improvement of LLMs.

## 2 CONTEXT AND NOTATIONS

**RL for LLMs.** We consider a transformer (Vaswani et al., 2017) LLM  $f(\cdot, \theta)$  parameterized by  $\theta$ . Following the foundation model paradigm (Bommasani et al., 2021) and the principles of transfer learning (Oquab et al., 2014), those weights are trained via a three-stage procedure: pre-training through next token prediction, supervised fine-tuning resulting in  $\theta_{\text{sft}}$ , and ultimately, RLHF (Christiano et al., 2017; Ouyang et al., 2022) to optimize a reward  $r$  as determined by a RM trained to reflect human preferences. In this RL stage,  $\theta$  defines a policy  $\pi_\theta(\cdot | \mathbf{x})$  by auto-regressively generating token sequences  $\mathbf{y}$  from the prompt  $\mathbf{x}$ . The primary objective is to find weights maximizing the average reward over a dataset of prompts  $\mathcal{X}$ :  $\arg\max_{\theta} \mathbb{E}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{\mathbf{y} \sim \pi_\theta(\cdot | \mathbf{x})} [r(\mathbf{x}, \mathbf{y})]$ .

**KL vs. reward.** Optimizing solely for  $r$  can (i) forget general abilities from pre-training (French, 1992) as an alignment tax (Ouyang et al., 2022; Lin et al., 2024a), (ii) hack the reward (Askeel et al., 2021; Skalse et al., 2022) leading to potential misalignment, or (iii) reduce the diversity of possible generations (Kirk et al., 2024) (as visible in Appendix F). To mitigate these risks, a KL regularization is usually integrated to balance fidelity to the initialization and high rewards:

$$\arg\max_{\theta} \mathbb{E}_{\mathbf{x} \in \mathcal{X}} \left[ \mathbb{E}_{\mathbf{y} \sim \pi_\theta(\cdot | \mathbf{x})} r(\mathbf{x}, \mathbf{y}) - \beta \text{KL}(\pi_\theta(\cdot | \mathbf{x}) \| \pi_{\theta_{\text{anchor}}}(\cdot | \mathbf{x})) \right], \quad (1)$$

where usually  $\theta_{\text{anchor}} \leftarrow \theta_{\text{sft}}$  and  $\beta$  is an hyperparameter, with high values leading to low KL yet also lower reward. The KL-regularized reward function is then  $r(\mathbf{x}, \mathbf{y}) - \beta \log\left(\frac{\pi_\theta(\mathbf{y} | \mathbf{x})}{\pi_{\theta_{\text{anchor}}}(\mathbf{y} | \mathbf{x})}\right)$ . Our base RL algorithm is a variant of REINFORCE (Williams, 1992). This choice follows recent RLHF works (Roit et al., 2023; Lee et al., 2024a; Ramé et al., 2024) and the findings from Li et al. (2023); Tajwar et al. (2024); Ahmadian et al. (2024) that, in terms of KL-reward trade-off, REINFORCE performs better than the more complex PPO (Schulman et al., 2017) and also better than various offline algorithms such as DPO (Rafailov et al., 2023), IPO (Azar et al., 2023) or RAFT (Dong et al., 2023b). Practitioners then employ early stopping to select an optimal point on the trajectory.

**Weight averaging.** The question of how best to merge models has recently garnered significant attention, driven by the discoveries that deep models can be merged in the weight space (Utans, 1996; Izmailov et al., 2018; Wortsman et al., 2023) instead of in the prediction space, as traditionally done in ensembling (Lakshminarayanan et al., 2017). Specifically, be given two sets of weights  $\theta^1$  and  $\theta^2$ , the different strategies merge them into a new set of weights  $\theta$ , parameterizing the same non-linear

network architecture. For clarity, we collectively refer to them as weight averaging (WA). The most basic one, uniform linear averaging, is also the most common; in this case,  $\theta = \frac{\theta^1 + \theta^2}{2}$ .

### 3 WARP

We introduce a novel alignment framework named Weight Averaged Rewarded Policies (*WARP*), illustrated in Figure 1(a) and described in Algorithm 1 below. *WARP* merges LLMs in the weight space to enhance the KL-reward front of policies. The following Sections 3.1 to 3.3 describe the motivations behind applying three distinct variants of WA at the three different stages of *WARP*. In particular, we summarize the key insights as observations, that will be experimentally validated in Section 4 (and in Appendices C and D), and theoretically motivated in Appendix B when possible. Overall, we observe that *WARP* outperforms other RL alignment strategies, without any memory or inference overhead at test time. However, training *WARP* is costly, requiring multiple RL runs at each iteration: see Section 6 for a detailed discussion on the required compute scaling.

---

#### Algorithm 1 *WARP* to improve the KL-reward trade-off in alignment

---

**Input:** Weights  $\theta_{\text{sft}}$  pre-trained and supervised fine-tuned  
 Reward model  $r$ , prompt dataset  $\mathcal{X}$ , optimizer Opt  
 $I$  iterations with  $M$  RL runs each for  $T$  training steps  
 $\mu$  EMA update rate,  $\eta$  LITI update rate

- 1: Define  $\theta_{\text{init}} \leftarrow \theta_{\text{sft}}$
- 2: **for** iteration  $i$  from 1 to  $I$  **do**
- 3:   **for** run  $m$  from 1 to  $M$  **do** ▷ Run in parallel
- 4:     Define  $\theta^m, \theta_{\text{ema}}^m \leftarrow \theta_{\text{init}}$
- 5:     **for** step  $t$  from 1 to  $T$  **do**
- 6:       Generate completion  $\mathbf{y} \sim \pi_{\theta^m}(\cdot | \mathbf{x})$  for  $\mathbf{x} \in \mathcal{X}$
- 7:       Compute  $r_\beta(\mathbf{y}) \leftarrow r(\mathbf{x}, \mathbf{y}) - \beta \log \frac{\pi_{\theta^m}(\mathbf{y} | \mathbf{x})}{\pi_{\theta_{\text{ema}}^m}(\mathbf{y} | \mathbf{x})}$  ▷ KL regularized reward
- 8:       Update  $\theta^m \leftarrow \text{Opt}(\theta^m, r_\beta(\mathbf{y}) \nabla_{\theta} [\log \pi_{\theta^m}(\mathbf{y} | \mathbf{x})])$  ▷ Policy gradient
- 9:       Update  $\theta_{\text{ema}}^m \leftarrow (1 - \mu) \cdot \theta_{\text{ema}}^m + \mu \cdot \theta^m$  ▷ Equation (EMA): update anchor
- 10:     **end for**
- 11:   **end for**
- 12:   Define  $\theta_{\text{slerp}}^i \leftarrow \text{slerp}(\theta_{\text{init}}, \{\theta^m\}_{m=1}^M, \lambda = \frac{1}{M})$  ▷ Equation (SLERP): merge  $M$  weights
- 13:   Update  $\theta_{\text{init}} \leftarrow (1 - \eta) \cdot \theta_{\text{init}} + \eta \cdot \theta_{\text{slerp}}^i$  ▷ Equation (LITI): interpolate towards init
- 14: **end for**

**Output:** KL-reward front of weights  $\{(1 - \kappa) \cdot \theta_{\text{sft}} + \kappa \cdot \theta_{\text{slerp}}^I \mid 0 \leq \kappa \leq 1\}$

---

#### 3.1 STAGE 1: EXPONENTIAL MOVING AVERAGE AS A DYNAMIC ANCHOR IN KL REGULARIZATION

**EMA anchor.** RLHF algorithms typically use the SFT initialization as a static anchor (Jaques et al., 2017; Roit et al., 2023) in the KL regularization, but in RL (notably for control tasks) it is common to regularly update the anchor (Schulman et al., 2015; Abdolmaleki et al., 2018). In this spirit, *WARP* uses the policy’s own exponential moving average (EMA) (Polyak & Juditsky, 1992), updated throughout the RL fine-tuning process such as, at each training step with  $\mu = 0.01$ :

$$\theta_{\text{ema}} \leftarrow (1 - \mu) \cdot \theta_{\text{ema}} + \mu \cdot \theta_{\text{policy}}. \quad (\text{EMA})$$

Using  $\theta_{\text{ema}}$  as the anchor  $\theta_{\text{anchor}}$  in Equation (1) provides several benefits, outlined below.

**Observation 1 (EMA).** *Policies trained with an exponential moving average anchor benefit from automatic annealing of the KL regularization and from distillation from a dynamic mean teacher (Tavainen & Valpola, 2017). Empirical evidence in Section 4.1.*

**Benefits from EMA.** Unlike a static SFT anchor, the dynamic nature of an EMA anchor induces a gradual automatic annealing and relaxation of the KL regularization. Specifically, the policy is initially strongly tied to the SFT initialization, and then progressively unleashed, allowing for more aggressive gradient updates later in training, leading to higher rewards. Moreover, by progressively

incorporating knowledge from the training, *EMA* acts as slow weight (Stojanovski et al., 2022; Lee et al., 2024b), and thus performing better than the initialization. But, by also maintaining essential information from the initialization, *EMA* can even perform better than the final policy’s weights; studies (Szegedy et al., 2016; Izmailov et al., 2018; Arpit et al., 2021) (see Morales-Brotons et al. (2024) for a review), and specifically (Kaddour, 2022) within the context of LLMs, indicate that averaging checkpoints over steps improves internal representations and thus predictions. Then, *EMA* guides the policy by KL distillation (Hinton et al., 2015) of high-quality target predictions, akin to a mean teacher (Tarvainen & Valpola, 2017) for self-supervised (Sohn et al., 2020; He et al., 2020; Oquab et al., 2024; Caron et al., 2021; Grill et al., 2020) learning. This also relates to deep RL techniques where *EMA* stabilizes exploration toward a Nash equilibrium (Awheda & Schwartz, 2013; 2016; Gorbatovski et al., 2024; Munos et al., 2023), and approximates mirror descent (Bubeck et al., 2015; Geist et al., 2019; Tomar et al., 2020).

### 3.2 STAGE 2: SPHERICAL LINEAR INTERPOLATION OF INDEPENDENTLY REWARDED POLICIES

**SLERP.** While *EMA* helps for a single RL and a fixed compute budget, it faces limitations due to the similarity of the weights collected along a single fine-tuning (Ramé et al., 2022). In this second stage, we merge  $M$  weights RL fine-tuned independently (each with their own *EMA* anchor). This follows model soups from Wortsman et al. (2022a) and its variants (Ramé et al., 2022; 2023) showing that WA improves generalization, and that task vectors (Ilharco et al., 2023) (the difference between fine-tuned weights and their initialization) can be arithmetically manipulated by linear interpolation (*LERP*) (Utans, 1996). Yet, this time, we use spherical linear interpolation (*SLERP*) (Shoemake, 1985), illustrated in Figure 2 and defined below for  $M = 2$ :

$$\text{slerp}(\theta_{\text{init}}, \theta^1, \theta^2, \lambda) = \theta_{\text{init}} + \frac{\sin[(1 - \lambda)\Omega]}{\sin \Omega} \cdot \delta^1 + \frac{\sin[\lambda\Omega]}{\sin \Omega} \cdot \delta^2, \quad (\text{SLERP})$$

where  $\Omega$  is the angle between the two task vectors  $\delta^1 = \theta^1 - \theta_{\text{init}}$  and  $\delta^2 = \theta^2 - \theta_{\text{init}}$ , and  $\lambda$  the interpolation coefficient. Critically *SLERP* is applied layer by layer, each having a different angle. In Appendix B.3 we clarify how *SLERP* can be used iteratively to merge  $M > 2$  models. To enforce diversity across weights, we simply vary the order in which text prompts  $x$  are given in each run: this was empirically sufficient, though other diversity strategies could help, e.g., varying the hyperparameters or the reward objectives (as explored in Figure 18(c)).

**Benefits from *SLERP* vs. *LERP*.** Merging task vectors, either with *SLERP* or *LERP*, combines their abilities (Ilharco et al., 2023). The difference is that *SLERP* preserves their norms, reaching higher rewards than the base models; this is summarized in Observation 2. In contrast, and as summarized in Observation 3, the more standard *LERP* has less impact on reward, but has the advantage of reducing KL; indeed, as shown in Appendix B, *LERP* tends to pull the merged model towards the initialization, especially as the angle  $\Omega$  between task vectors is near-orthogonal (see Observation 3).

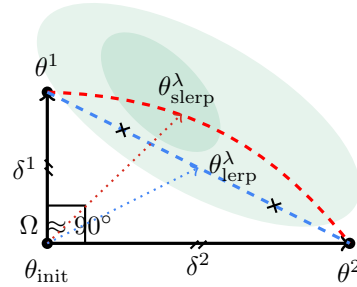


Figure 2: *SLERP* vs. *LERP*.

**Observation 2 (*SLERP*).** Spherical linear interpolation boosts rewards, yet slightly increases KL. Empirical evidence in Section 4.2 and theoretical insights in Lemma 1.

**Observation 3 (*LERP*).** Linear interpolation reduces KL, yet has reduced impact on reward. Empirical evidence in Appendix C.1 and theoretical insights in Lemmas 2 and 3.

**Observation 4 (Task vectors).** Task vectors  $\delta$  are close to orthogonal with  $\Omega \approx 90^\circ$ , while the full weights  $\theta$  are collinear. Empirical evidence in Appendix C.2.

### 3.3 STAGE 3: LINEAR INTERPOLATION TOWARDS INITIALIZATION

**LITI.** In the previous stage, *SLERP* combines multiple policies into one with higher rewards and slightly higher KL. This third stage, inspired by WiSE-FT from Wortsman et al. (2022b), interpolates from the merged model towards the initialization:

$$\theta^\eta \leftarrow (1 - \eta) \cdot \theta_{\text{init}} + \eta \cdot \theta_{\text{slerp}}. \quad (\text{LITI})$$

Adjusting the interpolating coefficient  $\eta \in [0, 1]$  trades off between some newly acquired behaviors leading to high rewards vs. general knowledge from the SFT initialization. Specifically, large values

$\eta \approx 1$  provide high rewards but also high KL, while smaller values  $\eta \approx 0$  lean towards smaller rewards and minimal KL. Fortunately, we observe that the reduction in KL is proportionally greater than the reduction in reward when decreasing  $\eta$ . Then, *LITI* empirically yields Pareto fronts that are noticeably above the “diagonal”, but also above those revealed during the base RL training runs.

**Observation 5** (*LITI*). *Interpolating weights towards the initialization reveals a better Pareto front than the one revealed during RL fine-tuning. Empirical evidence in Figure 1(b) and Section 4.3, and theoretical insights in Lemmas 4 and 5.*

**Benefits from *LITI*.** Previous works tried to understand how weight interpolation can mitigate forgetting while increasing robustness and generalization. Wortsman et al. (2022b) argue that WiSE-FT, akin to *LITI* in supervised learning contexts, recovers generalizable features from pre-training that might be lost during fine-tuning (Kumar et al., 2022), consistently with WA reducing catastrophic forgetting (Stojanovski et al., 2022; Eeckht et al., 2022) in continual learning. Then in the context of RL, Lin et al. (2024a) argue that *LITI* increases feature diversity, efficiently balancing between generality and task specificity. Finally, Jang et al. (2024) argues that the geometric projection of the ideal weights is located between the merged model and the initialization.

### 3.4 ITERATIVE WARP

**Iterative training.** The model merging strategies previously described not only establish an improved Pareto front of solutions, but also set the stage for iterative improvements. Indeed, if the computational budget is sufficient, we can apply those three stages iteratively, using  $\theta^\eta$  from previous Pareto front (usually with  $\eta = 0.3$ , choice ablated in Appendix D.3) as the initialization  $\theta_{\text{init}}$  for the next iteration, following the model recycling (Don-Yehiya et al., 2023; Ramé et al., 2023) strategies. Then, the entire training procedure is made of multiple iterations, each consisting of those three stages, where the final weight from a given iteration serves as an improved initialization for the next one.

**Observation 6** (Iterative *WARP*). *The application of WARP iteratively progressively refines the Pareto front. Empirical evidence in Sections 4.4 and 4.5.*

## 4 EXPERIMENTS: ON THE BENEFITS OF WARP

**Setup.** We consider the Gemma 7B (Gemma Team et al., 2024) LLM, which we seek to fine-tune with RLHF into a better conversational agent. We use REINFORCE (Williams, 1992) policy gradient to optimize the KL-regularized reward. The dataset  $\mathcal{X}$  contains conversation prompts. We generate on-policy samples with temperature 0.9, batch size of 128, Adam (Kingma & Ba, 2015) optimizer with learning rate  $10^{-6}$  and warmup of 100 steps. *SLERP* is applied independently to the 28 layers. Except when stated otherwise, we train for  $T = 9k$  steps, with KL strength  $\beta = 0.1$ , *EMA* update rate  $\mu = 0.01$ , merging  $M = 2$  policies uniformly  $\lambda = 0.5$ , and *LITI* update rate  $\eta = 0.3$ ; we analyze those values in Appendix D. We rely on a high capacity reward model, whose architecture is an order of magnitude larger than our policy.

**Summary.** In our experiments, we analyze the KL to the SFT policy (reflecting the forgetting of pre-trained knowledge) and the reward (evaluating alignment to the RM). In Section 4.1, we first show the benefits of using an *EMA* anchor; then in Section 4.2, we show that merging policies trained independently helps. Moreover, in Section 4.3, we show that *LITI* improves the KL-reward Pareto front; critically, repeating those three *WARP* stages can iteratively improve performances in Section 4.4. A limitation is that our RM accurately approximates true human preferences only in low KL region (Gao et al., 2023). Therefore, we finally report other metrics in Section 4.5, specifically comparing against open-source baselines such as Mixtral (Jiang et al., 2024), and reporting performance on standard benchmarks such as MMLU (Hendrycks et al., 2020).

### 4.1 STAGE 1: EXPONENTIAL MOVING AVERAGE AS A DYNAMIC ANCHOR IN KL REGULARIZATION

In Figures 3(a) and 3(b), we compare the training trajectories of different REINFORCE variants, where the changes lie in the choice of the anchor in the KL regularization and of the hyperparameter  $\beta$  controlling its strength. Results are computed every 100 training steps. In our proposed version, the anchor is the *EMA* of the trained policy with  $\beta = 0.1$  and an *EMA* update rate  $\mu = 0.1$  (other

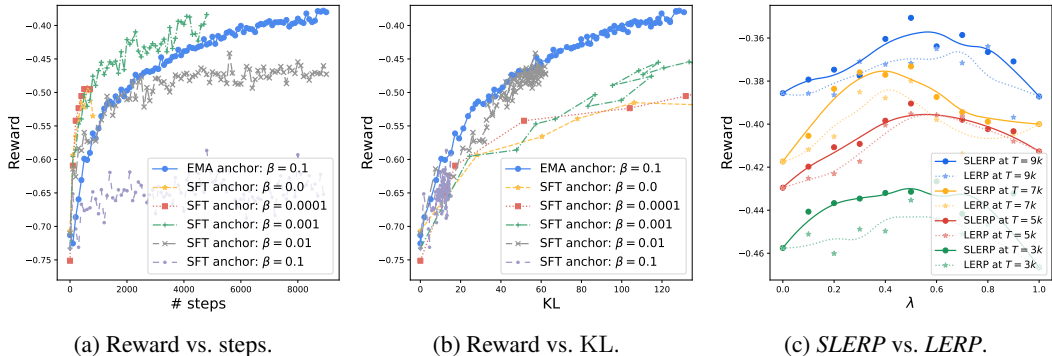


Figure 3: **EMA and SLERP experiments.** We first compare RL runs with different anchors and strengths  $\beta$  in the KL regularization. We show their results along training in Figure 3(a), and their KL-reward Pareto fronts in Figure 3(b). We perform evaluation every 100 steps, and train them for  $T = 9k$  steps, though we stopped the trainings if they ever reach a KL of 200 (e.g., after  $T = 1k$  training steps when  $\beta = 0.0$ ). Figure 3(c) plots the reward obtained when merging two policies (trained independently after  $T$  RL steps with their own *EMA* anchor) with interpolating coefficient  $\lambda$ ; highest rewards are with *SLERP* for  $\lambda = 0.5$  and  $T = 9k$  steps.

values are ablated in Figure 15). As the Pareto front for our strategy is above and to the left in Figure 3(b), this confirms the superiority of using such an adaptive anchor. The baseline variants all use the SFT as the anchor, with different values of  $\beta$ . The lack of regularization ( $\beta = 0.0$ ) leads to very fast optimization of the reward in Figure 3(a), but largely through hacking, as visible by the KL exploding in just a few training steps in Figure 3(b). In contrast, higher values such as  $\beta = 0.1$  fail to optimize the reward as regularization is too strong, causing a quick reward saturation around  $-0.62$  in Figure 3(a). Higher values such as  $\beta = 0.01$  can match our *EMA* anchor in low KL regime, but saturates around a reward of  $-0.46$ . In contrast, as argued in Observation 1, the dynamic *EMA* anchor progressively moves away from the SFT initialization, causing implicit annealing of the regularization. In conclusion, relaxing the anchor with *EMA* updates improves the trade-off between KL and reward, at any given KL level, for a fixed compute budget. We refer the interested reader to additional experiments in Figure 14 from Appendix D.2 where we compare the trained policies with their online *EMA* version.

#### 4.2 STAGE 2: SPHERICAL LINEAR INTERPOLATION OF INDEPENDENTLY REWARDED POLICIES

In Figure 3(c), we plot  $\lambda \rightarrow r(\text{slerp}(\theta_{\text{init}}, \theta^1, \theta^2, \lambda))$  showing reward convexity when interpolating policies via *SLERP*, validating Observation 2. This mirrors the linear mode connectivity (Frankle et al., 2020) property across weights fine-tuned from a shared initialization, i.e., the fact that interpolated weights perform better than the initial models (recovered for  $\lambda = 0$  or  $\lambda = 1$ ). Moreover, *SLERP* consistently obtains higher rewards than *LERP*; yet, this is at slightly higher KL, as further detailed in Appendices B and C.1, where we analyze respectively their theoretical and empirical differences.

#### 4.3 STAGE 3: LINEAR INTERPOLATION TOWARDS INITIALIZATION

In Figure 4(a), we merge policies trained for  $T$  steps, and then apply the *LITI* procedure. Critically, sliding the interpolating coefficient  $\eta \in \{0, 0.1, 0.3, 0.5, 0.8, 1.0\}$  reveals various Pareto fronts, consistently above the training trajectories obtained during the two independent RL fine-tunings. Interestingly, longer fine-tunings improve performances, at high KL, but also at lower KL, simply by using a smaller  $\eta$  afterwards. Then in Figure 4(b), we report the Pareto fronts when merging up to  $M = 5$  weights. We note that all Pareto fronts revealed when applying *LITI* are consistently above the ones from RL fine-tunings, validating Observation 5. More precisely, best results are achieved by merging an higher number of policies  $M$ , suggesting a promising scaling direction.

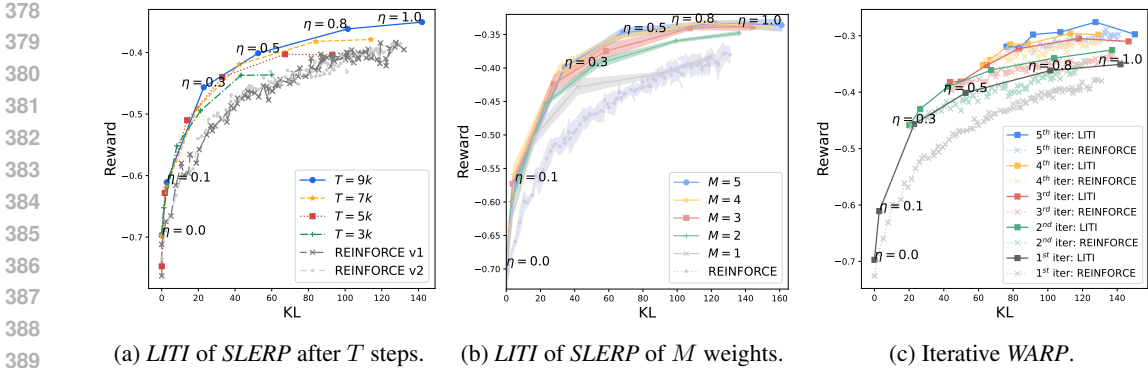


Figure 4: **LITI and iterative experiments.** Figure 4(a) considers the LITI of the SLERP of  $M = 2$  policies after  $T$  steps with  $\lambda = 0.5$ , interpolating towards their SFT init as we slide  $\eta$ , revealing Pareto fronts above the  $M = 2$  REINFORCE training trajectories. Then Figure 4(b) plots the LITI of the SLERP of  $M$  weights with  $\lambda = \frac{1}{M}$  after  $T = 9k$  steps: light-colored areas show standard deviations across 5 experiments. The iterative WARP procedure is illustrated in Figure 4(c); we fine-tune  $M = 2$  policies with their own EMA as the anchor, merge them with SLERP, interpolate towards their init with LITI, and iteratively leverage the weights obtained with  $\eta = 0.3$  as the new initialization for the next iteration.

#### 4.4 ITERATIVE WARP

In Figure 4(c), we apply the iterative procedure described in Section 3.4. At each of the  $I = 5$  iterations we train  $M = 2$  policies for  $T$  steps, with  $T = 9k$  for the first iteration, and  $T = 7k$  for iterations 2 and 3, and then  $T = 5k$  for computational reasons. The LITI curves interpolate towards their own initialization (while Figure 1(b) interpolated towards the SFT initialization, see Appendix D.4 for a comparison). We systematically observe that LITI curves are above the RL training trajectories used to obtain the inits. Results get better at every iteration, validating Observation 6, although with reduced returns after a few iterations.

#### 4.5 COMPARISONS AND BENCHMARKS

**Side by side comparisons.** To conclude our experiments, we compare our trained policies against Mistral (Jiang et al., 2023) and Mixtral (Jiang et al., 2024) LLMs. Each policy generates a candidate answer on an held-out collection of prompts. Then we compute side by side preference rates (Zheng et al., 2023) with “much better”, “better” and “slightly better” receiving scores of  $\pm 1.5$ ,  $\pm 1$ , and  $\pm 0.5$  respectively (and ties receiving a score of 0). The reported score is then the average over all prompts. A positive score represents better policies. The results in Table 1 validate the efficiency of WARP, as our policies are preferred over Mistral variants, and also outperform the two open-sourced Gemma 7B sharing the same pre-training. However, we note that results stagnate after the 3<sup>rd</sup> iteration.

Table 1: Side by side comparisons.

Methods	Mistral 7B v1	Mistral 7B v2	Mixtral 8x7B
Gemma 7B 1.0	0.24	-0.01	-0.08
Gemma 7B 1.1	0.37	0.16	0.08
REINFORCE EMA anchor	0.37	0.16	0.07
WARP: 1 <sup>st</sup> iter	0.42	0.23	0.13
WARP: 2 <sup>nd</sup> iter	<b>0.45</b>	0.25	0.16
<b>WARP: 3<sup>rd</sup> iter</b>	<b>0.45</b>	<b>0.26</b>	<b>0.18</b>
WARP: 4 <sup>th</sup> iter	<b>0.45</b>	0.25	0.16
WARP: 5 <sup>th</sup> iter	<b>0.45</b>	0.24	0.17



Table 2: **Benchmark results.**

Methods	MBPP	MMLU	GSM8K	MATH	HumanEval	BBH
Gemma 7B 1.1	39.0	56.4	55.6	25.6	46.9	53.1
<b>WARP</b>	<b>45.4</b>	<b>57.6</b>	<b>66.8</b>	<b>31.0</b>	<b>50.0</b>	<b>58.8</b>

**Benchmarks.** Table 2 compares *WARP* (3<sup>rd</sup> iter) and Gemma 7B 1.1 (Gemma Team et al., 2024) on popular benchmarks in the zero-shot setup: MBPP (Austin et al., 2021) and HumanEval (Chen et al., 2021) benchmarking coding capabilities, MMLU (Hendrycks et al., 2020) assessing STEM knowledge, the GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) benchmarks targeting reasoning abilities, and the Big Bench Hard (BBH) (Suzgun et al., 2022) benchmark evaluating general capabilities through questions that were deemed difficult for frontier LLMs. *WARP* has particularly strong results on mathematics benchmarks, suggesting higher analytical capabilities. Finally, we refer the reader to Appendix G, where we show how the performances in terms of reward-KL are transposed in terms of SxS and accuracy on real-world benchmarks.

## 5 RELATED WORK

**How to merge models.** The most common model merging strategy is *LERP*, initially used to average checkpoints collected along a single run, uniformly (Szegedy et al., 2016; Izmailov et al., 2018) or with an exponential moving average (*EMA*) (Polyak & Juditsky, 1992). Following the linear model connectivity (Frankle et al., 2020) observation, the model soups variants (Wortsman et al., 2022a; Ilharco et al., 2023; Ramé et al., 2023) linearly interpolate from different fine-tunings; this relies on the shared pre-training, limiting divergence (Neysshabur et al., 2020) such as models remain in constrained weight regions (Gueta et al., 2023), which also suggests that pre-training mitigates the need to explicitly enforce trust regions in gradient updates (Schulman et al., 2015; 2017). Subsequent works such as TIES merging (Yadav et al., 2023) and DARE (Yu et al., 2023) reduce interferences in multi-task setups with sparse task vectors (Ilharco et al., 2023). In contrast, we use *SLERP*, introduced in Shoemake (1985), increasingly popular in the open-source community (Goddard et al., 2024) but relatively underexplored in the academic literature, with limited studies such as Kim et al. (2024). Some tried to align weights trained from scratch (Entezari et al., 2022; Ainsworth et al., 2022) or with different architectures (Wan et al., 2024); yet, the methods are complex, less robust, and usually require additional training.

**Benefits of model merging.** WA boosts generalization by reducing variance (Wortsman et al., 2022a; Ramé et al., 2022), decreasing memorization (Lin et al., 2024b; Zaman et al., 2023; Ramé et al., 2024) and flattening the loss landscape (Cha et al., 2021). Additionally, merging weights combines their strengths (Ilharco et al., 2023), which helps in multi-task setups (Ilharco et al., 2022; Ramé et al., 2023), to tackle forgetting (Stojanovski et al., 2022; Eecka et al., 2022; Alexandrov et al., 2024) or to provide better initializations (Don-Yehiya et al., 2023), as explored in Jain et al. (2023); Jang et al. (2024); Huang et al. (2024) for iterative procedures in classification tasks. In particular, we considered using the geometric insights from Eq. 2 in Jang et al. (2024); yet, as our task vectors are nearly orthogonal  $\Omega \approx 90^\circ$  (see Appendix C.2), using the update rule  $\eta \rightarrow \frac{2 \cos \Omega}{1 + \cos \Omega}$  failed. WA is now also used in RL setups (Nikishin et al., 2018; Gaya et al., 2022; Lawson & Qureshi, 2023); for example, *WARM* (Ramé et al., 2024) merges reward models to boost their efficiency, robustness and reliability. Actually, *WARP* is conceived as a response to *WARM*, demonstrating that model merging can tackle two key RLHF challenges; policy learning in *WARP* and reward design in *WARM*. The most similar works are the following, which also explore how WA can improve policy learning. Noukhovitch et al. (2023) propose an iterative approach with the *EMA* as a new initialization for subsequent iterations. Gorbatovski et al. (2024) and Munos et al. (2023) use *EMA* as the reference, but only for direct preference optimization. Ramé et al. (2023) employ *LERP* to improve alignment in multi-objective RLHF when dealing with different objectives; similarly, Xiao et al. (2023) target multi-task setups with *LERP*. Finally, Lin et al. (2024a) and Fu et al. (2024) use model merging to reduce the alignment tax, although without incorporating *EMA* during training, without merging multiple rewarded policies and not iteratively. Critically, none of these works focus on KL as a measure of forgetting, use *EMA* as the anchor in KL, apply *SLERP* or use *LITI* as the initialization

486 for subsequent RL iterations. In contrast, *WARP* integrates all those elements, collectively leading to  
487 an LLM outperforming Mixtral (Jiang et al., 2024).  
488

## 489 6 DISCUSSION 490

491 **Distributed learning for parallelization and open-source.** *WARP* addresses a crucial challenge:  
492 aligning LLMs with human values and societal norms, while preserving the capabilities that emerged  
493 from pre-training. To this end, we leverage a (perhaps surprising) ability: policies trained in parallel  
494 can combine their strengths within a single policy by weight averaging. Then, the distributed nature  
495 of *WARP* makes it flexible and scalable, as it is easily parallelizable by enabling intermittent weight  
496 sharing across workers. Actually, iterative *WARP* shares similarities with DiLoCo (Douillard et al.,  
497 2023): by analogy, the first stage performs inner optimization on multiple workers independently;  
498 the second stage merges gradients from different workers; the third stage performs SGD outer  
499 optimization with a learning rate equal to  $\eta$ . More generally, *WARP* could facilitate open-source  
500 (Goddard et al., 2024) collaborative training of policies (Raffel, 2023), optimizing resource and  
501 supporting privacy in federated learning (McMahan et al., 2017) scenarios; collaborators could  
502 train and share their LLMs, while keeping their data and RMs private. In particular, we show in  
503 Appendix E that *WARP* can handle diverse objectives.

504 **Iterated amplification.** *WARP* improves LLM alignment by leveraging the principles of iterated  
505 amplification (Christiano et al., 2018) and progressive collaboration of multiple agents. By analogy,  
506 model merging via WA acts as an effective alternative to debate (Irving et al., 2018), with agents  
507 communicating within the weight space instead of the token space, ensuring that only essential  
508 information is retained (Ramé et al., 2024). Then, *WARP* refines the training signal by combining  
509 insights and exploration from diverse models, iteratively achieving higher rewards through self-  
510 distillation (Tarvainen & Valpola, 2017), surpassing the capabilities of any single agent. If this is the  
511 way forward, then an iterative safety assessment would be required to detect and mitigate potential  
512 risks early, ensuring that the development remains aligned with safety standards.

513 **Scaling alignment.** The *WARP* procedure increases the compute training cost by performing multiple  
514 fine-tunings at each iteration. Yet, this should be viewed as “a feature rather than a bug”. Specifically,  
515 by preventing memorization and forgetting, we see *WARP* as a fine-tuning method that can transform  
516 additional compute allocated to alignment into enhanced capabilities and safety. This would allow  
517 scaling (the traditionally cheap) post-training alignment, in the same way pre-training has been  
518 scaled (Hoffmann et al., 2022). Indeed, historically, pre-training efforts have benefited much more  
519 from compute scaling, fine-tuning efforts remaining significantly cheaper. Critically for large-scale  
520 deployment, the acquired knowledge is within a single (merged) model, thus without inference  
521 or memory overhead, in contrast to “more agents” approaches (Li et al., 2024; Wang et al., 2024).  
522 Finally, although *WARP* improves policy optimization, it is important to recognize that *WARP* does not  
523 address other critical challenges in RLHF (Casper et al., 2023): to mitigate the safety risks (Amodei  
524 et al., 2016; Hendrycks & Mazeika, 2022; Hendrycks, 2023) from misalignment (Taylor et al., 2016;  
525 Ngo et al., 2022), *WARP* should be part of a broader responsible AI framework.

## 526 7 CONCLUSION 527

528 We introduce Weight Averaged Rewarded Policies (*WARP*), a novel RLHF strategy to align LLMs  
529 with three distinct stages of model merging: exponential moving average as a dynamic anchor during  
530 RL, spherical interpolation to combine multiple policies rewarded independently, and interpolation  
531 towards the shared initialization. This iterative application of *WARP* improves the KL-reward Pareto  
532 front, aligning the LLMs while protecting the knowledge from pre-training, and compares favorably  
533 against state-of-the-art baselines. We hope *WARP* could contribute to safe and powerful AI systems  
534 by scaling alignment, and spur further exploration of the magic behind model merging.  
535  
536  
537  
538  
539

## REFERENCES

- 540  
541  
542 Abbas Abdomaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin  
543 Riedmiller. Maximum a posteriori policy optimisation. *ICLR*, 2018. (p. 4)
- 544  
545 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Ahmet Üstün, and  
546 Sara Hooker. Back to basics: Revisiting REINFORCE style optimization for learning from human  
547 feedback in LLMs. *arXiv preprint*, 2024. (p. 3)
- 548  
549 Samuel K. Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models  
550 modulo permutation symmetries. In *ICLR*, 2022. (p. 9)
- 551  
552 Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, and David Ha. Evolutionary optimization of  
553 model merging recipes. *arXiv preprint*, 2024. (p. 2)
- 554  
555 Anton Alexandrov, Veselin Raychev, Mark Niklas Müller, Ce Zhang, Martin Vechev, and Kristina  
556 Toutanova. Mitigating catastrophic forgetting in language transfer via model merging. *arXiv  
557 preprint*, 2024. (p. 9)
- 558  
559 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.  
560 Concrete problems in AI safety. *arXiv preprint*, 2016. (pp. 1 and 10)
- 561  
562 Devansh Arpit, Huan Wang, Yingbo Zhou, and Caiming Xiong. Ensemble of averages: Improving  
563 model selection and boosting performance in domain generalization. In *NeurIPS*, 2021. (p. 5)
- 564  
565 Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones,  
566 Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez,  
567 Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark,  
568 Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for  
569 alignment. *arXiv preprint*, 2021. (pp. 1 and 3)
- 570  
571 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
572 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
573 models. *arXiv preprint*, 2021. (p. 9)
- 574  
575 Mostafa D Awheda and Howard M Schwartz. Exponential moving average Q-learning algorithm. In  
576 *ADPRL*, 2013. (p. 5)
- 577  
578 Mostafa D Awheda and Howard M Schwartz. Exponential moving average based multiagent rein-  
579 forcement learning algorithms. *Artificial Intelligence Review*, 2016. (p. 5)
- 580  
581 Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal  
582 Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human  
583 preferences. *arXiv preprint*, 2023. (p. 3)
- 584  
585 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,  
586 Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportuni-  
587 ties and risks of foundation models. *arXiv preprint*, 2021. (p. 3)
- 588  
589 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,  
590 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:  
591 Early experiments with gpt-4. *arXiv preprint*, 2023. (p. 1)
- 592  
593 Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends  
594 in ML*, 2015. (p. 5)
- 595  
596 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
597 Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. (p. 5)
- 598  
599 Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier  
600 Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems  
601 and fundamental limitations of reinforcement learning from human feedback. *TMLR*, 2023. (pp. 1  
602 and 10)

- 594 Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and  
595 Sungrae Park. SWAD: Domain generalization by seeking flat minima. In *NeurIPS*, 2021. (p. 9)  
596
- 597 Mark Chen, Jerry Tworek, Heewoo Jun, et al. Evaluating large language models trained on code.  
598 *arXiv preprint*, 2021. (p. 9)
- 599 Paul Christiano, Buck Shlegeris, and Dario Amodei. Supervising strong learners by amplifying weak  
600 experts. *arXiv preprint*, 2018. (pp. 3 and 10)  
601
- 602 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
603 reinforcement learning from human preferences. In *NeurIPS*, 2017. (pp. 1 and 3)  
604
- 605 Jack Clark and Dario Amodei. Faulty Reward Functions in the Wild. [https://openai.com/r](https://openai.com/research/faulty-reward-functions)  
606 [esearch/faulty-reward-functions](https://openai.com/research/faulty-reward-functions), 2016. (p. 1)
- 607 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
608 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
609 math word problems. *arXiv preprint arXiv:2110.14168*, 2021. (p. 9)  
610
- 611 Imre Csiszár. I-divergence geometry of probability distributions and minimization problems. *The*  
612 *annals of probability*, pp. 146–158, 1975. (p. 22)
- 613 Shachar Don-Yehiya, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen.  
614 CoLD fusion: Collaborative descent for distributed multitask finetuning. In *ACL*, 2023. (pp. 6 and 9)  
615
- 616 Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang,  
617 Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected  
618 by supervised fine-tuning data composition. *arXiv preprint*, 2023a. (p. 1)
- 619 Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao,  
620 Jipeng Zhang, KaShun SHUM, and Tong Zhang. RAFT: Reward ranked finetuning for generative  
621 foundation model alignment. *TMLR*, 2023b. (p. 3)
- 622
- 623 Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna  
624 Kuncoro, Marc’Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-  
625 communication training of language models. *arXiv preprint*, 2023. (pp. 2, 3, and 10)  
626
- 627 Steven Vander Eeck et al. Weight averaging: A simple yet effective method to overcome catastrophic  
628 forgetting in automatic speech recognition. *arXiv preprint*, 2022. (pp. 6 and 9)
- 629 Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation  
630 invariance in linear mode connectivity of neural networks. In *ICLR*, 2022. (p. 9)  
631
- 632 Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear mode  
633 connectivity and the lottery ticket hypothesis. In *ICML*, 2020. (pp. 2, 7, 9, and 22)  
634
- 635 Robert M French. Semi-distributed representations and catastrophic forgetting in connectionist  
636 networks. *Connection Science*, 1992. (pp. 1 and 3)
- 637 Tingchen Fu, Deng Cai, Lemaoy Liu, Shuming Shi, and Rui Yan. Disperse-then-merge: Pushing the  
638 limits of instruction tuning via alignment tax reduction. *arXiv preprint*, 2024. (p. 9)  
639
- 640 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In  
641 *ICML*, 2023. (pp. 1, 2, and 6)
- 642 Jean-Baptiste Gaya, Laure Soulier, and Ludovic Denoyer. Learning a subspace of policies for online  
643 adaptation in reinforcement learning. In *ICLR*, 2022. (p. 9)  
644
- 645 Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A theory of regularized markov decision  
646 processes. In *ICML*, 2019. (pp. 2 and 5)  
647
- Google Gemini Team. Gemini: A family of highly capable multimodal models. 2023. (p. 1)

- 648 Google Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
649 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open  
650 models based on gemini research and technology. *arXiv preprint*, 2024. (pp. 1, 3, 6, and 9)  
651
- 652 Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian  
653 Benedict, Mark McQuade, and Jacob Solawetz. Arcee’s mergekit: A toolkit for merging large  
654 language models. *arXiv preprint*, 2024. (pp. 2, 9, and 10)
- 655 Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical  
656 investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint*, 2013.  
657 (p. 1)  
658
- 659 Alexey Gorbatoevski, Boris Shaposhnikov, Alexey Malakhov, Nikita Surnachev, Yaroslav Aksenov,  
660 Ian Maksimov, Nikita Balagansky, and Daniil Gavrilov. Learn your reference model for real good  
661 alignment. *arXiv preprint*, 2024. (pp. 2, 5, and 9)
- 662 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
663 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
664 et al. Bootstrap your own latent-a new approach to self-supervised learning. *NeurIPS*, 2020. (p. 5)  
665
- 666 Almog Gueta, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen. Knowl-  
667 edge is a region in weight space for fine-tuned language models. In *EMNLP*, 2023. (p. 9)
- 668 Sil Hamilton. Detecting mode collapse in language models via narration. *arXiv preprint*, 2024. (pp. 1  
669 and 30)  
670
- 671 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
672 unsupervised visual representation learning. In *CVPR*, 2020. (p. 5)
- 673 Dan Hendrycks. Natural selection favors AIs over humans. *arXiv preprint*, 2023. (p. 10)  
674
- 675 Dan Hendrycks and Mantas Mazeika. X-risk analysis for AI research. *arXiv preprint*, 2022. (pp. 1  
676 and 10)  
677
- 678 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
679 Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*, 2020. (pp. 6  
680 and 9)
- 681 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
682 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In  
683 *NeurIPS*, 2021. (p. 9)
- 684 Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In  
685 *NeurIPS*, 2015. (p. 5)  
686
- 687 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
688 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
689 Training compute-optimal large language models. In *NeurIPS*, 2022. (p. 10)
- 690 Zitong Huang, Ze Chen, Bowen Dong, Chaoqi Liang, Erjin Zhou, and Wangmeng Zuo. Imwa:  
691 Iterative model weight averaging benefits class-imbalanced learning tasks. *arXiv preprint*, 2024.  
692 (p. 9)  
693
- 694 Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon  
695 Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by interpolating  
696 weights. In *NeurIPS*, 2022. (pp. 2 and 9)
- 697 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt,  
698 Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *ICLR*, 2023. (pp. 2,  
699 3, 5, 9, 19, and 20)  
700
- 701 Geoffrey Irving, Paul Christiano, and Dario Amodei. Ai safety via debate. *arXiv preprint*, 2018.  
(p. 10)

- 702 Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson.  
703 Averaging weights leads to wider optima and better generalization. In *UAI*, 2018. (pp. 2, 3, 5, 9, 19,  
704 and 27)
- 705 Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Convergence and  
706 generalization in neural networks. In *NeurIPS*, 2018. (p. 21)
- 708 Samyak Jain, Sravanti Addepalli, Pawan Kumar Sahu, Priyam Dey, and R Venkatesh Babu. Dart:  
709 Diversify-aggregate-repeat training improves generalization of neural networks. In *CVPR*, 2023.  
710 (p. 9)
- 711 Dong-Hwan Jang, Sangdoon Yun, and Dongyoon Han. Model stock: All we need is just a few  
712 fine-tuned models. *arXiv preprint*, 2024. (pp. 6, 9, 20, 21, and 25)
- 714 Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E Turner,  
715 and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with  
716 kl-control. In *ICML*, 2017. (pp. 2, 4, and 19)
- 717 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
718 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
719 Mistral 7b. *arXiv preprint*, 2023. (pp. 1 and 8)
- 720 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
721 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
722 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. (pp. 6, 8, and 10)
- 723 Jean Kaddour. Stop wasting my time! saving days of ImageNet and BERT training with latest weight  
724 averaging. In *NeurIPS Workshop*, 2022. (p. 5)
- 726 Minchul Kim, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Token fusion: Bridging  
727 the gap between token pruning and token merging. In *WACV*, 2024. (p. 9)
- 728 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.  
729 (p. 6)
- 731 Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward  
732 Grefenstette, and Roberta Raileanu. Understanding the effects of RLHF on LLM generalisation  
733 and diversity. In *ICLR*, 2024. (pp. 1, 3, and 30)
- 734 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A  
735 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming  
736 catastrophic forgetting in neural networks. In *PNAS*, 2017. (p. 1)
- 737 Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-  
738 tuning can distort pretrained features and underperform out-of-distribution. In *ICLR*, 2022. (pp. 1  
739 and 6)
- 740 Maxime Labonne. Merge Large Language Models with mergekit, 2024a. URL <https://huggingface.co/blog/mlabonne/merge-models>. (p. 2)
- 741 Maxime Labonne. NeuralBeagle14-7B. <https://huggingface.co/mlabonne/NeuralBeagle14-7B-GGUF>, 2024b. (p. 2)
- 742 Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and  
743 Richard Dufour. Biomistral: A collection of open-source pretrained large language models for  
744 medical domains. *arXiv preprint*, 2024. (p. 2)
- 745 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
746 uncertainty estimation using deep ensembles. In *NeurIPS*, 2017. (p. 3)
- 747 Nathan Lambert and Jacob Morrison. Model merging lessons in The Waifu Research Department,  
748 2024. URL <https://www.interconnects.ai/p/model-merging>. (p. 2)
- 749 Daniel Lawson and Ahmed H Qureshi. Merging decision transformers: Weight averaging for forming  
750 multi-task policies. In *ICLR RRL Workshop*, 2023. (p. 9)

- 756 Angeliki Lazaridou, Anna Potapenko, and Olivier Tieleman. Multi-agent communication meets  
757 natural language: Synergies between functional and structural language learning. In *ACL*, 2020.  
758 (p. 2)
- 759 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton  
760 Bishop, Victor Carbune, and Abhinav Rastogi. RLAIIF: Scaling reinforcement learning from  
761 human feedback with AI feedback. In *ICML*, 2024a. (p. 3)
- 762 Jaerin Lee, Bong Gyun Kang, Kihoon Kim, and Kyoung Mu Lee. Grokfast: Accelerated grokking by  
763 amplifying slow gradients. *arXiv preprint*, 2024b. (p. 5)
- 764 Mike Lewis, Denis Yarats, Yann N Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal?  
765 end-to-end learning for negotiation dialogues. *arXiv preprint*, 2017. (p. 1)
- 766 Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. More agents is all you need. *arXiv*  
767 *preprint*, 2024. (p. 10)
- 770 Zhizhong Li and Derek Hoiem. Learning without forgetting. *TPAMI*, 2017. (p. 1)
- 771 Ziniu Li, Tian Xu, Yushun Zhang, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple,  
772 effective, and efficient reinforcement learning method for aligning large language models. *arXiv*  
773 *preprint*, 2023. (p. 3)
- 774 Yong Lin, Hangyu Lin, Wei Xiong, Shizhe Diao, Jianmeng Liu, Jipeng Zhang, Rui Pan, Haoxiang  
775 Wang, Wenbin Hu, Hanning Zhang, Hanze Dong, Renjie Pi, Han Zhao, Nan Jiang, Heng Ji, Yuan  
776 Yao, and Tong Zhang. Mitigating the alignment tax of rlhf. *arXiv preprint*, 2024a. (pp. 1, 2, 3, 6,  
777 and 9)
- 778 Yong Lin, Lu Tan, Yifan Hao, Honam Wong, Hanze Dong, Weizhong Zhang, Yujiu Yang, and Tong  
779 Zhang. Spurious feature diversification improves out-of-distribution generalization. In *ICLR*,  
780 2024b. (p. 9)
- 781 Tianlin Liu, Shangmin Guo, Leonardo Bianco, Daniele Calandriello, Quentin Berthet, Felipe Llinares,  
782 Jessica Hoffmann, Lucas Dixon, Michal Valko, and Mathieu Blondel. Decoding-time realignment  
783 of language models. In *ICML*, 2024. (p. 2)
- 784 Yuchen Lu, Soumye Singhal, Florian Strub, Aaron Courville, and Olivier Pietquin. Countering  
785 language drift with seeded iterated learning. In *ICML*, 2020. (p. 2)
- 786 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
787 Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.  
788 (p. 10)
- 789 Skander Moalla, Andrea Miele, Razvan Pascanu, and Caglar Gulcehre. No representation, no trust:  
790 Connecting representation, collapse, and trust issues in ppo. *arXiv preprint*, 2024. (pp. 1 and 30)
- 791 Daniel Morales-Brotons, Thijs Vogels, and Hadrien Hendrikx. Exponential moving average of  
792 weights in deep learning: Dynamics and benefits. *TMLR*, 2024. (p. 5)
- 793 Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland,  
794 Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. Nash  
795 learning from human feedback. *arXiv preprint*, 2023. (pp. 2, 5, and 9)
- 796 Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning?  
797 In *NeurIPS*, 2020. (pp. 2, 9, 21, and 25)
- 798 Richard Ngo, Lawrence Chan, and Soren Mindermann. The alignment problem from a deep learning  
799 perspective. *arXiv preprint*, 2022. (pp. 1 and 10)
- 800 Evgenii Nikishin, Pavel Izmailov, Ben Athiwaratkun, Dmitrii Podoprikin, Timur Garipov, Pavel  
801 Shvechikov, Dmitry Vetrov, and Andrew Gordon Wilson. Improving stability in deep reinforcement  
802 learning with weight averaging. In *UDL*, 2018. (p. 9)
- 803 Michael Noukhovitch, Samuel Lavoie, Florian Strub, and Aaron Courville. Language model align-  
804 ment with elastic reset. In *NeurIPS*, 2023. (pp. 2 and 9)

- 810 OpenAI. Gpt-4 technical report. 2023. (p. 1)  
811
- 812 Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level  
813 image representations using convolutional neural networks. In *CVPR*, 2014. (p. 3)
- 814 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
815 Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas  
816 Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael  
817 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Ar-  
818 mand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision.  
819 *TMLR*, 2024. (p. 5)
- 820 Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent  
821 space: Improved editing of pre-trained models. In *NeurIPS*, 2023. (p. 21)  
822
- 823 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
824 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
825 instructions with human feedback. *NeurIPS*, 2022. (pp. 1 and 3)
- 826 Alexander Pan, Kush Bhatia, and Jacob Steinhardt. The effects of reward misspecification: Mapping  
827 and mitigating misaligned models. In *ICLR*, 2022. (p. 1)  
828
- 829 Ethan Perez, Sam Ringer, Kamilé Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit,  
830 Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors  
831 with model-written evaluations. *arXiv preprint*, 2022. (p. 1)
- 832 Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM*,  
833 1992. (pp. 3, 4, and 9)  
834
- 835 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under-  
836 standing by generative pre-training. 2018. (p. 1)
- 837 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
838 models are unsupervised multitask learners. 2019. (p. 1)  
839
- 840 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea  
841 Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv*  
842 *preprint*, 2023. (p. 3)
- 843 Colin Raffel. Building Machine Learning Models Like Open Source Software. *ACM*, 2023. (pp. 2, 3,  
844 and 10)
- 845 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
846 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text  
847 transformer. *JMLR*, 2020. (p. 1)  
848
- 849 Alexandre Ramé, Matthieu Kirchmeyer, Thibaud Rahier, Alain Rakotomamonjy, Patrick Gallinari,  
850 and Matthieu Cord. Diverse weight averaging for out-of-distribution generalization. In *NeurIPS*,  
851 2022. (pp. 2, 5, 9, 21, 27, and 29)
- 852 Alexandre Ramé, Kartik Ahuja, Jianyu Zhang, Matthieu Cord, Léon Bottou, and David Lopez-Paz.  
853 Model ratatouille: Recycling diverse models for out-of-distribution generalization. In *ICML*, 2023.  
854 (pp. 5, 6, and 9)  
855
- 856 Alexandre Ramé, Nino Vieillard, Léonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier  
857 Bachem, and Johan Ferret. WARM: On the benefits of weight averaged reward models. In *ICML*,  
858 2024. (pp. 2, 3, 9, and 10)
- 859 Alexandre Ramé, Guillaume Couairon, Mustafa Shukor, Corentin Dancette, Jean-Baptiste Gaya,  
860 Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpo-  
861 lating weights fine-tuned on diverse rewards. In *NeurIPS*, 2023. (pp. 2 and 9)  
862
- 863 Mark Rofin, Nikita Balagansky, and Daniil Gavrilov. Linear interpolation in parameter space is good  
enough for fine-tuned language models. *arXiv preprint*, 2022. (p. 2)



- 864 Paul Roit, Johan Ferret, Lior Shani, Roei Aharoni, Geoffrey Cideron, Robert Dadashi, Matthieu  
865 Geist, Sertan Girgin, Léonard Hussenot, Orgad Keller, et al. Factually consistent summarization  
866 via reinforcement learning with textual entailment feedback. In *ACL*, 2023. (pp. 3 and 4)  
867
- 868 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region  
869 policy optimization. In *ICML*, 2015. (pp. 4 and 9)
- 870 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
871 optimization algorithms. *arXiv preprint*, 2017. (pp. 3 and 9)  
872
- 873 Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text  
874 generation. In *ACL*, 2020. (p. 30)
- 875 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman,  
876 Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. Towards understanding  
877 sycophancy in language models. *arXiv preprint*, 2023. (p. 1)  
878
- 879 Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang.  
880 Loose lips sink ships: Mitigating length bias in reinforcement learning from human feedback. In  
881 *ACL*, 2023. (p. 29)
- 882 Ken Shoemake. Animating rotation with quaternion curves. In *SIGGRAPH*, 1985. (pp. 2, 3, 5, 9,  
883 and 19)  
884
- 885 Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating  
886 length correlations in rlhf. *arXiv preprint*, 2023. (pp. 1 and 29)  
887
- 888 Joar Max Viktor Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining  
889 and characterizing reward gaming. In *NeurIPS*, 2022. (pp. 1 and 3)
- 890 Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk,  
891 Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with  
892 consistency and confidence. In *NeurIPS*, 2020. (p. 5)  
893
- 894 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
895 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *NeurIPS*,  
896 2020. (p. 1)
- 897 Zafir Stojanovski, Karsten Roth, and Zeynep Akata. Momentum-based weight interpolation of strong  
898 zero-shot models for continual learning. In *NeurIPS Workshop*, 2022. (pp. 2, 5, 6, and 9)  
899
- 900 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
901 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. Challenging  
902 big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint*, 2022. (p. 9)
- 903 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking  
904 the inception architecture for computer vision. In *CVPR*, 2016. (pp. 5 and 9)  
905
- 906 Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano  
907 Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal,  
908 on-policy data. *arXiv preprint*, 2024. (p. 3)
- 909 Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency  
910 targets improve semi-supervised deep learning results. *NeurIPS*, 2017. (pp. 3, 4, 5, 10, and 27)  
911
- 912 Jessica Taylor, Eliezer Yudkowsky, Patrick LaVictoire, and Andrew Critch. Alignment for advanced  
913 machine learning systems. *Ethics of AI*, 2016. (pp. 1 and 10)
- 914 Manan Tomar, Lior Shani, Yonathan Efroni, and Mohammad Ghavamzadeh. Mirror descent policy  
915 optimization. *arXiv preprint*, 2020. (p. 5)  
916
- 917 Joachim Utans. Weight averaging for neural networks and local resampling schemes. In *AAAI*, 1996.  
(pp. 2, 3, 5, and 20)

- 918 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz  
919 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. (p. 3)  
920
- 921 Fanqi Wan, Xinting Huang, Deng Cai, Xiaojun Quan, Wei Bi, and Shuming Shi. Knowledge fusion  
922 of large language models. *arXiv preprint*, 2024. (p. 9)
- 923 Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances  
924 large language model capabilities. *arXiv preprint*, 2024. (p. 10)  
925
- 926 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,  
927 Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *ICLR*,  
928 2022. (p. 1)
- 929 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
930 learning. *Reinforcement learning*, 1992. (pp. 2, 3, 6, and 19)  
931
- 932 Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes,  
933 Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig  
934 Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy  
935 without increasing inference time. In *ICML*, 2022a. (pp. 2, 5, 9, 20, and 21)
- 936 Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Hanna Hajishirzi, Ali Farhadi,  
937 Hongseok Namkoong, and Ludwig Schmidt. Robust fine-tuning of zero-shot models. In *CVPR*,  
938 2022b. (pp. 2, 3, 5, 6, 19, and 21)
- 939 Mitchell Wortsman, Suchin Gururangan, Shen Li, Ali Farhadi, Ludwig Schmidt, Michael Rabbat,  
940 and Ari S. Morcos. lo-fi: distributed fine-tuning without communication. *TMLR*, 2023. (p. 3)  
941
- 942 Shitao Xiao, Zheng Liu, Peitian Zhang, and Xingrun Xing. LM-cocktail: Resilient tuning of language  
943 models via model merging. *arXiv preprint*, 2023. (p. 9)
- 944 Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging:  
945 Resolving interference when merging models. In *NeurIPS*, 2023. (p. 9)  
946
- 947 Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are super mario:  
948 Absorbing abilities from homologous models as a free lunch. *arXiv preprint*, 2023. (p. 9)
- 949 Kerem Zaman, Leshem Choshen, and Shashank Srivastava. Fuse to forget: Bias reduction and  
950 selective memorization through model fusion. *arXiv preprint*, 2023. (p. 9)  
951
- 952 Chujie Zheng, Ziqi Wang, Heng Ji, Minlie Huang, and Nanyun Peng. Weak-to-strong extrapolation  
953 expedites alignment. *arXiv preprint*, 2024. (pp. 24 and 25)
- 954 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
955 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
956 Judging LLM-as-a-judge with MT-bench and chatbot arena. In *NeurIPS*, 2023. (p. 8)  
957
- 958 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
959 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*  
960 *preprint*, 2019. (p. 1)  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

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

## WARP: On the Benefits of Weight Averaged Rewarded Policies

### Supplementary material

This supplementary material is organized as follows:

- Appendix A provides additional illustration of the *WARP* procedure.
- Appendix B details theoretical insights on task vectors, *SLERP*, *LERP* and *LITI*.
- Appendix C details empirical insights on task vectors, *SLERP*, *LERP* and *LITI*.
- Appendix D shows the impact of different design choices in *WARP*.
- Appendix E investigates a potential length bias in *WARP*, and how to fix it.
- Appendix F explores the relationship between KL and diversity in generations.
- Appendix G provides additional SxS and benchmark results.

### A STRATEGY ILLUSTRATION

In Figure 5, we propose an alternative illustration of *WARP*, where the different stages are more detailed than in Figure 1(a). Then in Figure 6, we also refine our illustration showcasing the similarity and difference between *SLERP* and *LERP*.

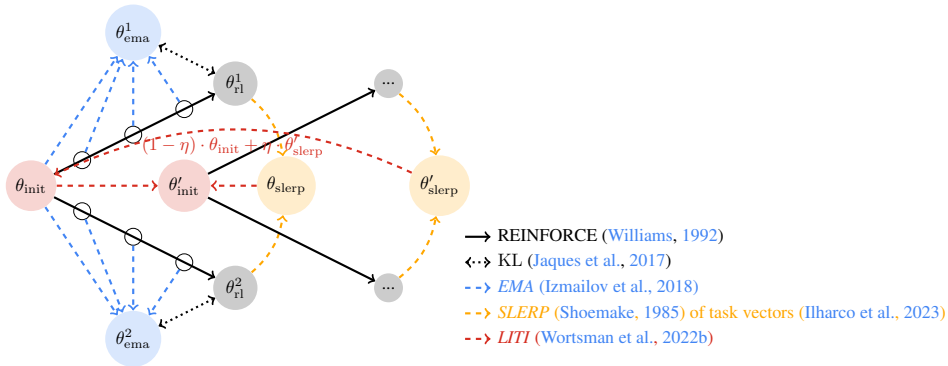


Figure 5: **Detailed illustration of the *WARP* strategy.** From a (pre-trained and supervised fine-tuned) LLM  $\theta_{\text{init}}$ , we launch  $M = 2$  fine-tunings (black arrows  $\longrightarrow$ ). The innovation of *WARP* lies in the use of model merging by weight averaging at three different stages. First, the exponential moving averages (*EMA*, blue dashed arrows  $\dashrightarrow$ ) of the policy (collected at different training steps) serves as the anchor for the KL regularization (black double-headed dotted arrows  $\langle \cdot \cdot \rangle$ ). The fine-tuned networks are weight averaged using spherical linear interpolation of task vectors (*SLERP*, yellow dashed arrows  $\dashrightarrow$ ). Third, we interpolate towards the initialization (*LITI*, red dashed arrows  $\dashrightarrow$ ). This obtained model  $\theta'_{\text{init}}$  serves as an updated initialization for the next iteration, progressively refining the model’s capabilities and alignment. Overall, the final model  $\theta'_{\text{slerp}}$  has high reward but also high KL. Then, by interpolation towards the SFT init, we reveal a KL-reward Pareto front of solutions:  $\{(1 - \eta) \cdot \theta_{\text{sft}} + \eta \cdot \theta_{\text{slerp}}^I \mid 0 \leq \eta \leq 1\}$ .

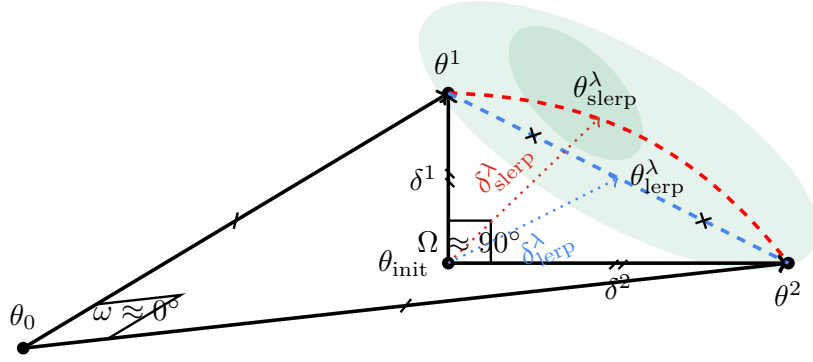


Figure 6: Illustration of the difference between the full weights  $\theta^m$  and their task vectors  $\delta^m = \theta^m - \theta_{\text{init}}$ , where darker areas are of better performance. We found in Appendix C.2 that  $\Omega \approx 90^\circ$  where  $\Omega$  is the angle between task vectors such as  $\cos \Omega = \frac{\delta^1 \cdot \delta^2}{\|\delta^1\| \|\delta^2\|}$ , while  $\omega$  the angle between the full weights such as  $\cos \omega = \frac{\theta^1 \cdot \theta^2}{\|\theta^1\| \|\theta^2\|}$  satisfies  $\omega \approx 0^\circ$ .

## B THEORETICAL INSIGHTS ON TASK VECTORS, *SLERP*, *LERP* AND *LITI*

Based on the insights from Ilharco et al. (2023) that task vectors (the differences between a fine-tuned model and its initialization) are semantically manipulable and interpretable units in the weight space, we compare *SLERP* and *LERP* merging operations by analyzing their task vectors.

**Background.** Linear interpolation (*LERP*) (Utans, 1996) is the simplest merging strategy, notably used in the model soups variants (Wortsman et al., 2022a), and defined as:

$$\text{lerp}(\theta^1, \theta^2, \lambda) = (1 - \lambda) \cdot \theta^1 + \lambda \cdot \theta^2. \quad (\text{LERP})$$

Then, as illustrated in Figure 6, the task vector for *LERP* with interpolating coefficient  $\lambda$  is given by:  $\delta_{\text{lerp}}^\lambda = \text{lerp}(\theta^1, \theta^2, \lambda) - \theta_{\text{init}} = (1 - \lambda) \cdot \delta^1 + \lambda \cdot \delta^2$ . Similarly, we define  $\delta_{\text{slerp}}^\lambda = \text{slerp}(\theta_{\text{init}}, \theta^1, \theta^2, \lambda) - \theta_{\text{init}}$  where *slerp* is defined in Equation (*SLERP*).

### B.1 THEORETICAL INSIGHTS ON THE *SLERP* AND *LERP* TASK VECTORS

We denote  $\Omega$  the angle between the the task vectors  $\delta^1$  and  $\delta^2$ :

$$\cos \Omega = \frac{\delta^1 \cdot \delta^2}{\|\delta^1\| \|\delta^2\|}. \quad (2)$$

Based on the empirical observations from Jang et al. (2024), confirmed in our Figure 11(c), we introduce the following Assumption 1 for simplicity.

**Assumption 1** (Task vectors of equal norm). *Independently fine-tuned task vectors have a same norm  $l$ :*

$$\|\delta^1\| = \|\delta^2\| = l. \quad (3)$$

**Lemma 1** (*SLERP* task vector). *Under Assumption 1, SLERP preserves the norm of the task vector:*

$$\|\delta_{\text{slerp}}^\lambda\| = l. \quad (4)$$

*Proof.* By definition,

$$\delta_{\text{slerp}}^\lambda = \frac{\sin[(1 - \lambda)\Omega]}{\sin \Omega} \cdot \delta^1 + \frac{\sin[\lambda\Omega]}{\sin \Omega} \cdot \delta^2 \quad (5)$$

Then, as  $\delta^1 \cdot \delta^2 = l^2 \cos \Omega$ ,

$$\frac{\|\delta_{\text{slerp}}^\lambda\|^2}{l^2} = \left( \frac{\sin[(1-\lambda)\Omega]}{\sin \Omega} \right)^2 + 2 \frac{\sin[(1-\lambda)\Omega]}{\sin \Omega} \frac{\sin[\lambda\Omega]}{\sin \Omega} \cos(\Omega) + \left( \frac{\sin[\lambda\Omega]}{\sin \Omega} \right)^2 \quad (6)$$

$$= \frac{\sin^2[(1-\lambda)\Omega] + 2 \sin[(1-\lambda)\Omega] \sin[\lambda\Omega] \cos(\Omega) + \sin^2[\lambda\Omega]}{\sin^2 \Omega} \quad (7)$$

$$= \frac{\sin^2 \Omega}{\sin^2 \Omega} \quad (8)$$

$$= 1 \quad (9)$$

using trigonometric identities, proving Lemma 1.  $\square$

**Lemma 2** (LERP task vector). *Under Assumption 1, LERP reduces the norm of the task vector:*

$$\|\delta_{\text{lerp}}^\lambda\| = l \sqrt{1 - 2(1 - \cos \Omega)(\lambda - \lambda^2)}. \quad (10)$$

We recover that averaging weights with  $\lambda = 0.5$  tends to reduce the norm of the task vectors, as previously highlighted in Jang et al. (2024).

*Proof.* By definition:

$$\delta_{\text{lerp}}^\lambda = (1 - \lambda) \cdot \delta^1 + \lambda \cdot \delta^2. \quad (11)$$

Then, as  $\delta^1 \cdot \delta^2 = l^2 \cos \Omega$ ,

$$\frac{\|\delta_{\text{slerp}}^\lambda\|^2}{l^2} = (1 - \lambda)^2 + 2\lambda(1 - \lambda) \cos \Omega + \lambda^2 \quad (12)$$

$$= 1 - 2\lambda(1 - \cos \Omega) + 2\lambda^2(1 - \cos \Omega) \quad (13)$$

$$= 1 - 2(1 - \cos \Omega)(\lambda - \lambda^2), \quad (14)$$

proving Lemma 2 when  $0 < \lambda < 1$ .  $\square$

## B.2 THEORETICAL INSIGHTS ON THE KL

### B.2.1 LINEAR REGIME

**Assumption 2** (Linear regime (Wortsman et al., 2022b)). *We assume that the predictions of a model  $f$ , with weights initialized from  $\theta_0$  and fine-tuned into  $\theta$ , can be approximated by first-order Taylor expansion:  $\forall \mathbf{x}$ ,*

$$f(\mathbf{x}, \theta) \approx f(\mathbf{x}, \theta_0) + (\theta - \theta_0) \cdot \nabla_{\theta} f(\mathbf{x}, \theta_0). \quad (15)$$

Assumption 2 defines a neural tangent (Jacot et al., 2018) space in which the relationship between weights and functions is linear. As previously argued in Wortsman et al. (2022a); Ramé et al. (2022), this Taylor expansion is reasonable partly because weights remain close during fine-tunings (Neyshabur et al., 2020), as confirmed in Figure 11 where they have equal norms and a cosine of one. Yet, please note that Ortiz-Jimenez et al. (2023) highlighted some limitations.

### B.2.2 KL VARIATIONS FOR LERP

We consider  $\theta^1$  and  $\theta^2$  weights fine-tuned from a shared SFT initialization  $\theta_{\text{sft}}$ . Then in the linear regime from Assumption 2, weight and prediction ensembling behaves similarly:

$$f(\mathbf{x}, (1 - \lambda) \cdot \theta^1 + \lambda \cdot \theta^2) \approx (1 - \lambda) \cdot f(\mathbf{x}, \theta^1) + \lambda \cdot f(\mathbf{x}, \theta^2). \quad (16)$$

This similarity enables to prove the following Lemma 3.

**Lemma 3** (LERP reduces KL). *For an interpolating coefficient  $0 \leq \lambda \leq 1$ , denoting  $\pi_\lambda$  the LERP policy from weight interpolation  $(1 - \lambda) \cdot \theta^1 + \lambda \cdot \theta^2$ , and  $\hat{\pi}_\lambda$  the ensembling policy from prediction interpolation  $(1 - \lambda) \cdot \pi_{\theta^1} + \lambda \cdot \pi_{\theta^2}$ , then under Assumption 2,*

$$\text{KL}(\pi_\lambda \| \pi_{\theta_{\text{sft}}}) \approx \text{KL}(\hat{\pi}_\lambda \| \pi_{\theta_{\text{sft}}}) \leq (1 - \lambda) \text{KL}(\pi_{\theta^1} \| \pi_{\theta_{\text{sft}}}) + \lambda \text{KL}(\pi_{\theta^2} \| \pi_{\theta_{\text{sft}}}), \quad (17)$$

*i.e., the KL for LERP is lower than the interpolated KL.*

1134 *Proof.* The following proof applies the linear assumption and properties of the KL divergence.

1135  
1136 **Approximation of KL.** The first approximate equality is a direct application of Assumption 2 to  $\pi_\lambda$ .  
1137 Precisely, applying Equation (16) to the definition of  $\pi_\lambda = \pi_{(1-\lambda)\theta_1 + \lambda\theta_2}$  yields that  $\pi_\lambda \approx \hat{\pi}_\lambda$ .

1138 **Upper bound of the KL.** The KL divergence is convex in both its arguments (Csiszár, 1975), thus  
1139 we directly have that

$$1140 \text{KL}((1-\lambda) \cdot \pi_{\theta^1} + \lambda \cdot \pi_{\theta^2} || \pi_{\theta_{\text{sft}}}) \leq (1-\lambda)\text{KL}(\pi_{\theta^1} || \pi_{\theta_{\text{sft}}}) + \lambda\text{KL}(\pi_{\theta^2} || \pi_{\theta_{\text{sft}}}), \quad (18)$$

1141 which completes the proof.

1142

1143

1144  $\square$

1145 **Remark 1.** Lemma 3 shows that the LERP  $\pi_\lambda$  is closer in KL to the original SFT initialization. This  
1146 relates to Lemma 2, where we show that the linear interpolation reduces the norm to the initialization.  
1147 As the interpolation brings the weights of the models closer, it is natural that it would also bring the  
1148 resulting policies closer.

1149

### 1150 B.2.3 KL AND REWARD VARIATION FOR LITI

1151 We now consider a given weight  $\theta$  (in practice either obtained from LERP or SLERP of multiple  
1152 fine-tuned weights) and its associated task vector  $\delta = \theta - \theta_{\text{sft}}$ . In the linear regime from Assumption 2,  
1153 for each  $\eta \in [0, 1]$ , we have the following:

$$1154 f(\mathbf{x}, \theta_{\text{sft}} + \eta \cdot \delta) - f(\mathbf{x}, \theta_{\text{sft}}) \approx \eta \cdot (f(\mathbf{x}, \theta_{\text{sft}} + \delta) - f(\mathbf{x}, \theta_{\text{sft}})). \quad (19)$$

1155 We try to show that:

$$1156 \text{KL}(\pi_{\theta_{\text{sft}} + \eta \cdot \delta} || \pi_{\theta_{\text{sft}}}) \leq \eta \cdot \text{KL}(\pi_{\theta_{\text{sft}} + \delta} || \pi_{\theta_{\text{sft}}}). \quad (20)$$

1157 **Lemma 4** (KL upper bound for interpolated distributions). For an interpolating coefficient  $0 \leq \eta \leq 1$ ,  
1158 denoting  $\pi_\eta$  the LITI policy from weight interpolation  $\theta_{\text{sft}} + \eta \cdot \delta$ , and  $\hat{\pi}_\eta$  the ensembling policy from  
1159 prediction interpolation  $(1-\eta) \cdot \pi_{\theta_{\text{sft}}} + \eta \cdot \pi_{\theta_{\text{sft}} + \delta}$ , then under Assumption 2,

$$1160 \text{KL}(\pi_\eta || \pi_{\theta_{\text{sft}}}) \approx \text{KL}(\hat{\pi}_\eta || \pi_{\theta_{\text{sft}}}) \leq \eta \text{KL}(\pi_{\theta_{\text{sft}} + \delta} || \pi_{\theta_{\text{sft}}}). \quad (21)$$

1161

1162

1163 *Proof.* The following proof uses the same method as the one of Lemma 3. We use Assumption 2 to  
1164 link the policy with the interpolation of policies, and the inequality is a result of the KL convexity.

1165 **Approximation of KL.** The first approximate equality is a direct application of Assumption 2 to  $\pi_\eta$ .  
1166 Precisely, applying Equation (19) to the definition of  $\pi_\eta = \pi_{\theta_{\text{sft}} + \eta \cdot \delta}$  yields that  $\pi_\eta \approx \hat{\pi}_\eta$ .

1167 **Upper bound of the KL.** Using the fact that the KL is convex, we have

$$1171 \text{KL}(\eta \cdot \pi_{\theta_{\text{sft}} + \delta} + (1-\eta) \cdot \pi_{\theta_{\text{sft}}} || \pi_{\theta_{\text{sft}}}) \leq \eta \text{KL}(\pi_{\theta_{\text{sft}} + \delta} || \pi_{\theta_{\text{sft}}}). \quad (22)$$

1172

1173  $\square$

1174 **Assumption 3** (LITI reward is above the expected reward). The rewards for the LITI interpolated  
1175 weights are above the interpolated rewards:

$$1176 r(\pi_0 + \eta \cdot (\pi - \pi_{\theta_{\text{sft}}})) \geq \eta r(\pi) + (1-\eta)r(\pi_{\theta_{\text{sft}}}), \quad (23)$$

1177

1178 This Assumption 3 is based on observations from Figure 9(b), and extends to a reward maximization  
1179 setup the notion of linear mode connectivity (Frankle et al., 2020), usually defined w.r.t. the accuracy  
1180 in supervised learning.

1181 **Lemma 5** (LITI for KL-reward trade-off). Be given the convexity of the KL from Lemma 4 and  
1182 the concavity of the reward  $r$  in Assumption 3, then the reward vs. KL front of LITI is above the  
1183 diagonal. Illustration in Figure 7.

1184

1185 *Proof.* We obtain a policy  $\pi_\theta$  fine-tuned from  $\pi_{\theta_{\text{sft}}}$ . The LITI policy for  $\theta_\eta = (1-\eta) \cdot \theta_{\text{sft}} + \eta \cdot \theta$  is  
1186 noted  $\pi_\eta$ . Combining the approximation from Lemma 4 and Assumption 3, we have that

$$1187 r(\pi_\eta) \geq (1-\eta)r(\pi_{\theta_{\text{sft}}}) + \eta r(\pi_\theta). \quad (24)$$

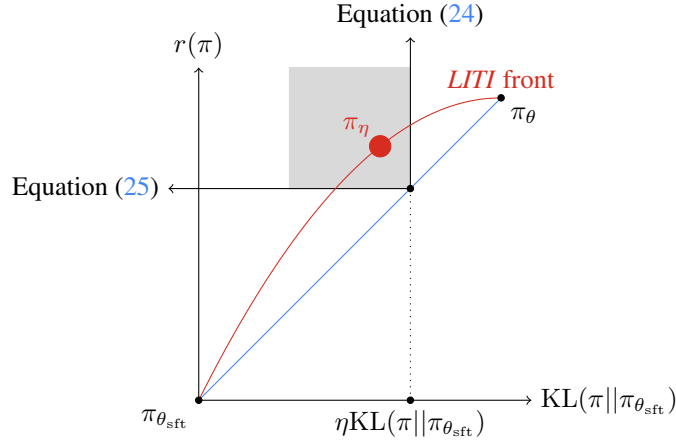


Figure 7: Illustration of Lemma 5. Based on experimental observation and theoretical insights, we see that the **Pareto front of the LITI policy** is better than the identity. It highlights how Equations (24) and (25) place LITI policies on the KL-reward plane.

And, from Lemma 4, we also have that

$$\text{KL}(\pi_\eta || \pi_{\theta_{\text{sft}}}) \leq \eta \text{KL}(\pi_\theta || \pi_{\theta_{\text{sft}}}). \quad (25)$$

This means that for every LITI coefficient  $\eta$ , the LITI policy has a higher reward than the interpolated reward at a lower KL. Geometrically, this means that each point on the Reward-KL front from LITI is on the top left quadrant of the plane according to the corresponding point on the diagonal.  $\square$

### B.3 UNIFORMLY AVERAGING $M > 2$ WEIGHTS WITH SLERP

The SLERP merging formula from Equation (SLERP) is only defined for  $M = 2$  weights. We trivially (and certainly suboptimally) generalize this to  $M > 2$  weights in the uniform averaging setup, thus giving an equal coefficient to each of them, i.e.,  $\lambda = \frac{1}{M}$ . In that setup, removing the dependency of  $\theta_{\text{init}}$  that is assumed shared, we generalize SLERP to merge  $M$  weights uniformly through the iterative procedure defined below:

$$\text{slerpm}(\{\theta^m\}_{m=1}^M) = \text{slerp}\left(\text{slerpm}(\{\theta^m\}_{m=1}^{M-1}), \theta^M, \lambda = \frac{1}{M}\right). \quad (26)$$

Though these operations are not associative, the standard deviations in performances are small, as indicated by the shaded areas in Figure 4(b).

## C EMPIRICAL INSIGHTS ON TASK VECTORS, SLERP, LERP AND LITI

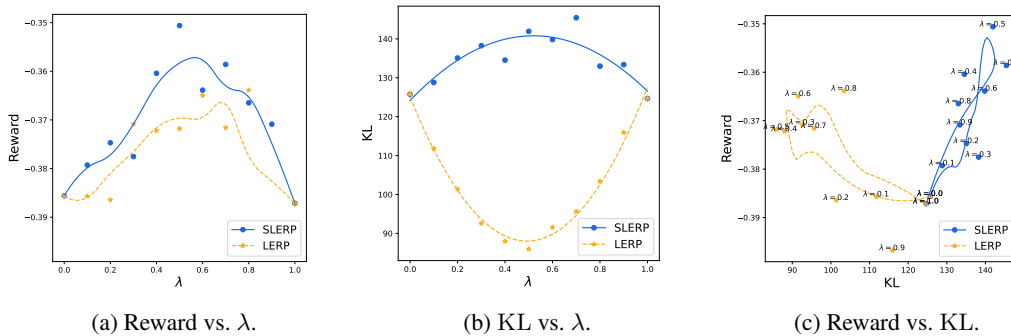
### C.1 EMPIRICAL INSIGHTS ON THE DIFFERENCE BETWEEN SLERP AND LERP

We now empirically investigate how those theoretical differences between SLERP and LERP affect the performance of the merged policies.

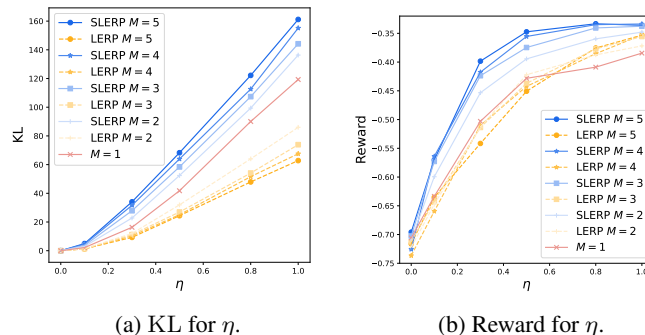
**SLERP vs. LERP.** In Figure 8 we adjust the interpolating coefficient  $\lambda$ , highlighting distinct behaviors for SLERP and LERP. SLERP consistently enhances rewards more than LERP, as depicted in Figures 3(c) and 8(a). However, a comprehensive evaluation must consider both KL and reward. As shown in Figure 8(b), LERP consistently reduces KL, corroborating with Lemma 2 that LERP reduces the norm of updates (while SLERP preserves it). When plotting these metrics together in Figure 8(c), we observe that SLERP and LERP target different regions on the Pareto front: SLERP achieves higher rewards at the expense of increased KL, while the main impact of LERP is to lower KL. This is consistent with Lemmas 2 and 3, be given the orthogonal angles between task vectors  $\Omega \approx 90^\circ$  (as shown in Figure 11(a)).

1242 **Combining *SLERP* and *LERP* with *LITI*.** We also compare the behaviours of *SLERP* and *LERP*  
 1243 when we apply *LITI*, as we adjust the interpolating coefficient  $\eta$ . Figure 9(a) and Figure 9(b) validate  
 1244 that KL is convex with regard to  $\eta$  while the reward is concave with regard to  $\eta$ , for different values  
 1245 of  $M$ . This is also highlighted in Figure 10(a), which reproduces the results from Figure 4(b) (and  
 1246 maintaining the same axis limits), replacing *SLERP* by *LERP*: this leads to critical changes in the  
 1247 Pareto fronts. Indeed, increasing  $M$  now tends to decrease KL for *LERP*, while it used to increase  
 1248 reward with *SLERP*. In Figure 10(b), we explore the extrapolation strategies from (Zheng et al.,  
 1249 2024), using  $0 \leq \eta \leq 2$  to compare the full extrapolated fronts from *LERP* and *SLERP*. While both  
 1250 perform similarly on low KL, our results suggest that *SLERP* perform better in high KL regions.

1251 **Conclusion.** *SLERP* demonstrates some key advantages. In particular, it reveals the full Pareto front  
 1252 of solutions, while *LERP* only exposes a portion; extrapolation Figure 10(b) with  $\eta > 1$  can partially  
 1253 mitigate this but as our experiments suggest, *LERP* curves consistently lag behind *SLERP* curves in  
 1254 high-reward regions. Moreover, from a practical perspective, *SLERP* scales the choice of  $\eta$  effectively,  
 1255 where 1 represents full updates and a fixed value of 0.3 always corresponds to the same operational  
 1256 region, optimizing for high reward and KL.

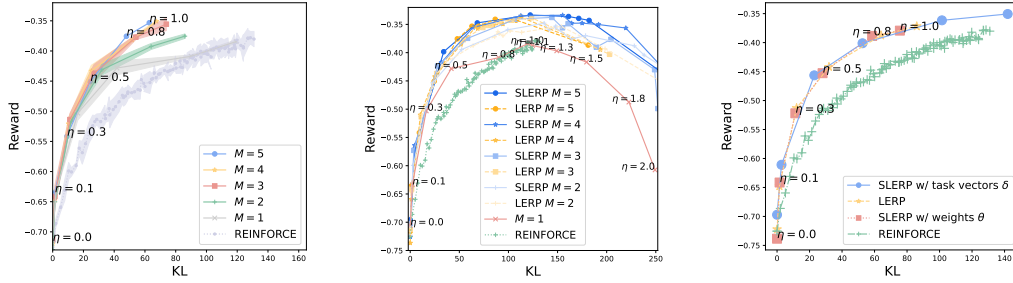


1269 Figure 8: ***SLERP* vs. *LERP* when sliding the interpolating coefficient  $\lambda$ .** Considering  $M = 2$   
 1270 weights after  $T = 9k$  RL steps, we merge them using either *SLERP* or *LERP*, while sliding the  
 1271 interpolating coefficient  $\lambda$  between 0 and 1. We then evaluate the merged checkpoints. Figure 8(a)  
 1272 shows that *SLERP* leads to higher reward than *LERP*, as previously in Figure 3(c). Figure 8(b) shows  
 1273 that *LERP* significantly reduces the KL (consistently with Lemma 3) while *SLERP* slightly increases  
 1274 it. Figure 8(c) shows how this impact the KL-reward Pareto front, where larger markers/darker colors  
 1275 indicate higher values of  $\lambda$ ; while *SLERP* covers high KL-high reward regions, *LERP* tends to cover  
 1276 regions of lower KL and thus also lower rewards.



1278 Figure 9: ***SLERP* vs. *LERP* when sliding the interpolating coefficient  $\eta$  of *LITI*.** In Figure 9(a)  
 1279 we show that the KL is convex (and almost linear) with regard to  $\eta$ , consistently with Lemma 4. In  
 1280 contrast, Figure 9(b) shows that the reward is concave, validating Assumption 3.





(a) LITI of LERP of  $M$  weights. (b) Extrapolation with  $0 \leq \eta \leq 2$ . (c) SLERP w/o task vectors.

Figure 10: **SLERP vs. LERP when sliding the interpolating coefficient  $\eta$  of LITI.** Figure 10(a) merges  $M$  policies with LERP and  $\lambda = \frac{1}{M}$  (the endpoints on the top right of the solid lines), and then interpolates towards their SFT init, where light-colored areas show standard deviations across 5 experiments, and with  $0 \leq \eta \leq 1$ . In contrast, in Figure 10(b) we investigate extrapolation (Zheng et al., 2024), using  $0 \leq \eta \leq 2$  enabling to compare the full fronts of solutions with LERP and SLERP. Finally, Figure 10(c) confirms that applying SLERP on the full weights  $\theta$  rather than on the task vectors  $\delta$  perform very similarly to LERP.

## C.2 EMPIRICAL INSIGHTS ON THE ROLE OF TASK VECTORS

We now explore the effectiveness of applying SLERP on task vectors  $\delta$  vs. full weights  $\theta$ , as illustrated in Figure 6. To this end, in Figure 11 we draw inspiration from Jang et al. (2024) and plot the angles  $\Omega$  and  $\omega$  and norms of  $\delta$  and  $\theta$ .

**Angles of task vectors  $\Omega \approx 90^\circ$ .** Figure 11(a) shows that the task vectors are typically orthogonal ( $\Omega \approx 90^\circ$ ), highlighting the diverse trajectories of the different RL fine-tunings. This is in contrast with (Jang et al., 2024) for supervised fine-tunings, where  $\Omega$  typically range between  $40^\circ$  and  $80^\circ$ . We suspect that this is related to the underlying differences between reinforcement and supervised learning; in RL the policies are trained on their own generations, creating more orthogonal task vectors, whereas in supervised learning the LLM try to imitate the groundtruth labels, leading to more similar task vectors. The orthogonality of our task vectors prevents the use of the update rule  $\eta \rightarrow \frac{2 \cos \Omega}{1 + \cos \Omega}$  suggested from Eq. 2 in Jang et al. (2024), as it would lead to  $\eta \approx 0$ , deleting any potential update.

**Angles of full weights  $\omega \approx 0^\circ$ .** In contrast, Figure 11(b) show that full weights remain collinear ( $\omega \approx 0^\circ$ ). This explains the empirical results from Figure 10(c), where applying SLERP directly to full weights results in behaviors similar to LERP. Indeed, as the angles  $\omega \approx 0^\circ$ , spherical interpolation effect is minimal because  $\sin(x) \approx x + \mathcal{O}(x^3)$ , and thus  $\frac{\sin[\lambda\omega]}{\sin \omega} \approx \frac{\lambda\omega}{\omega} \approx \lambda$ .

**Norms consistency.** Figure 11(c) confirms the consistency in the norms of different task vectors, supporting our Assumption 1. This uniformity is aligned with previous research (Jang et al., 2024). As a side note, this consistency extends to full weights  $\theta$ , confirming that fine-tuning typically results in minimal changes to the overall weight (Neysshabur et al., 2020).

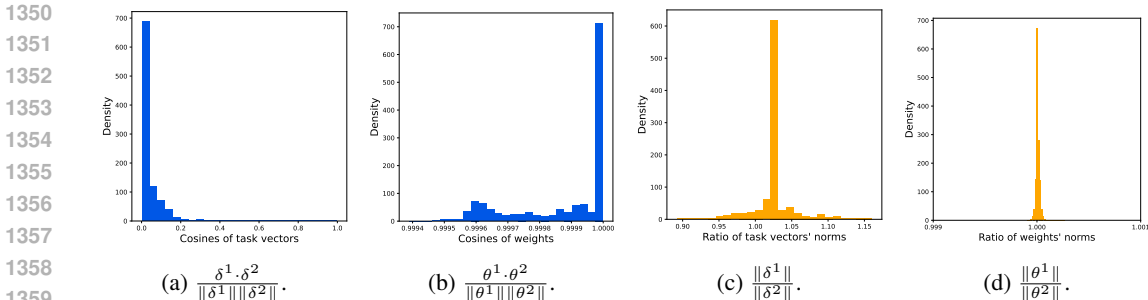


Figure 11: **Angles and norms of (full) weights  $\theta^m$  and their task vectors  $\delta^m = \theta^m - \theta_{\text{init}}$ .** The histograms are across the 28 layers of the Gemma 7B architecture. Figure 11(a) plots the histograms of task vector cosines. Figure 11(b) plots the histograms of weights cosines. Figure 11(c) plots the histograms of task vector norms ratio. Figure 11(d) plots the histograms of weights norms ratio.

### D EMPIRICAL INVESTIGATION OF SEVERAL DESIGN CHOICES

We include several experiments showcasing the robustness of *WARP* to different design choices, while further demonstrating its superiority in terms of KL-reward trade-off. Specifically, Appendix D.1 analyzes the performances along training at different steps  $T$ ; Appendix D.2 provides results with different values for the hyperparameters  $\mu$  and  $\beta$ ; Appendix D.3 shows the impact of the update rate  $\eta$  to provide an improved initialization for the 2<sup>nd</sup> iteration of *WARP*; finally, Appendix D.4 shows that in iterative *WARP*, interpolating towards the episode initialization or the SFT initialization both perform similarly.

#### D.1 ANALYZING THE NUMBER OF TRAINING STEPS

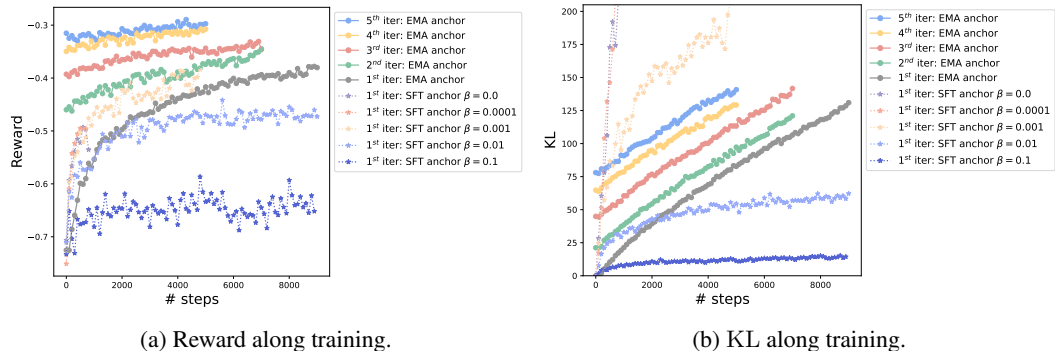


Figure 12: **Rewards and KL at different number of training steps  $T$ .** Figures 12(a) and 12(b) complement Figure 3(b) and Figure 4(c), this time plotting rewards and KL separately as a function of the number of training steps  $T$ . Regarding iterative *WARP*, we observe that each iteration has higher rewards but also higher KL (by starting at training step 0 from a new initialization). Regarding the baseline (REINFORCE with SFT anchor), we observe that low values of  $\beta$  lead to very fast hacking of the reward, as visible by the KL exploding, while high values of  $\beta$  slow down the training.

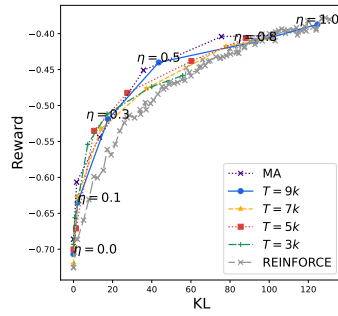


Figure 13: *LITI* with  $M = 1$  at different number of training steps  $T$ . The reward gain is significantly reduced compared to Figure 4(a) where we first merged  $M = 2$  policies before applying *LITI*. We also try to perform moving average (MA) (Izmailov et al., 2018) before applying *LITI*, averaging checkpoints collected along a single RL fine-tuning at steps  $\{6k, 7k, 8k, 9k\}$ ; this does not improve performances, suggesting the need to merge weights from independent fine-tunings to have enough diversity (Ramé et al., 2022).

## D.2 ANALYZING THE VALUES OF $\mu$ AND $\beta$

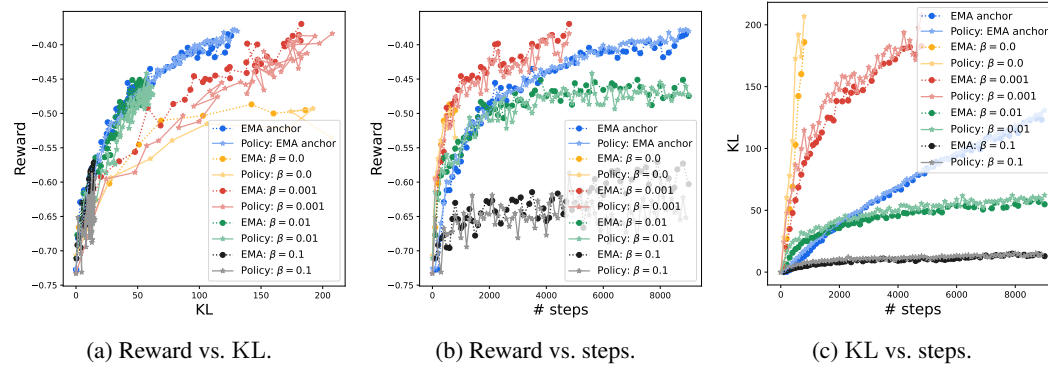


Figure 14: *EMA vs. their base policies*, extending Figures 3(a) and 3(b). Figure 14(a) shows that the *EMA* of all variants (with SFT anchor) perform similarly or better than their base policies in KL-reward. As a reminder, we perform evaluation every 100 steps, and train them for  $T = 9k$  steps, though we stopped the trainings if the base policy ever reaches a KL of 200. This confirms Observation 1; the benefits of our variant with *EMA* anchor is partly explained by distillation from an improved mean teacher (Tarvainen & Valpola, 2017).

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

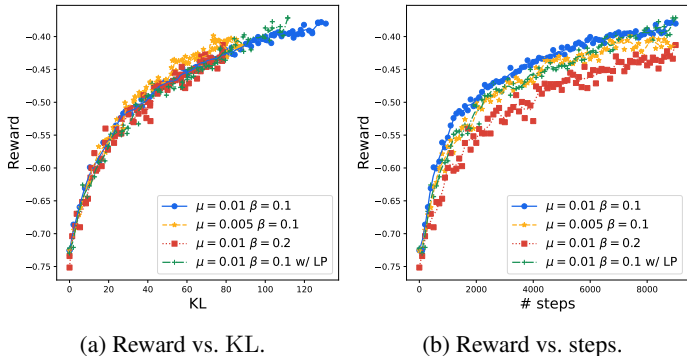


Figure 15: **Experiments ablating the values for the EMA update rate  $\mu$  and the KL regularization strength  $\beta$ .** So far we have systematically used  $\mu = 0.01$  and  $\beta = 0.1$  for all EMA-based runs, including in the iterative WARP. These hyperparameters were chosen at the project’s onset and have remained unchanged. In Figures 15(a) and 15(b) we increase regularization with  $\mu = 0.005$  and  $\beta = 0.2$ . Our results indicate that reducing  $\mu$  or increasing  $\beta$  behaves similarly, marginally improving the KL-reward Pareto front but slowing down training. Additionally, we include the training trajectory when using a length penalty (LP), as detailed in Appendix E.

D.3 ANALYZING THE VALUES OF  $\eta$

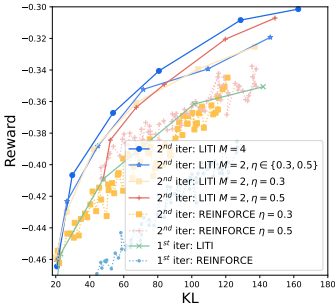


Figure 16: **Experiments ablating the LITI update rate  $\eta$ .** As we initiate the 2<sup>nd</sup> iteration of WARP, selecting an appropriate value for  $\eta$  is key, as it determines the starting point  $\theta^\eta$  and functions similarly to an outer learning rate (see Section 6). We usually set  $\eta = 0.3$ , but now provide results with an increased  $\eta = 0.5$ , starting the 2<sup>nd</sup> iteration from a more “advanced” position on the previous Pareto front. We run and average  $M = 2$  fine-tunings from each of those two initializations for  $T = 7k$  steps, before applying LITI. Our results indicate that a higher  $\eta$  (0.5) performs better in regions of high KL, whereas a lower  $\eta$  (0.3) helps in regions with KL below 65. This suggests that the optimal choice for  $\eta$  is compute-dependent; a lower  $\eta$  is appropriate if further iterations can explore high KL regions, whereas a limited compute budget might benefit from a higher  $\eta$ . This resembles the learning rate trade-off in optimization, where lower rates improve results but require more training steps. As a final note, we can also use different  $\eta$  for the different fine-tunings; notably, we observe that merging all those  $M = 4$  RLs perform better (though it doubles the compute).

D.4 INTERPOLATE TOWARDS THE INITIALIZATION? OR TOWARDS THE SFT?

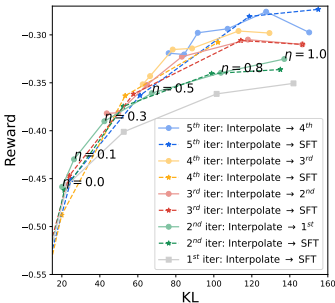
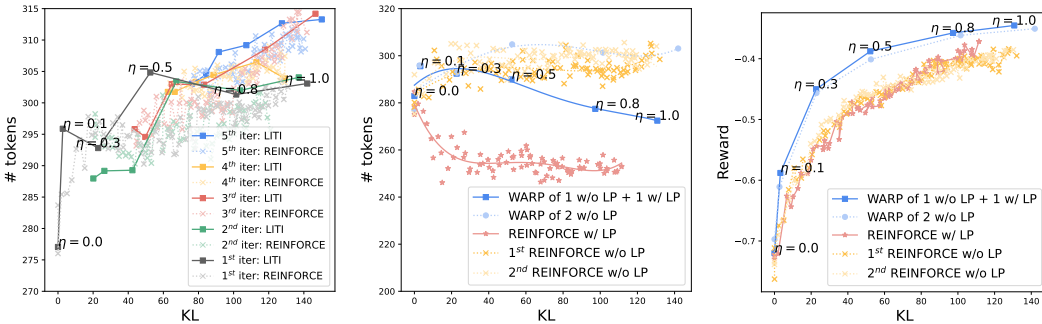


Figure 17: **Experiments ablating the initialization in LITI.** We compare LITI either towards the episode-specific initialization (the  $\theta^\eta$  selected from previous iteration) or towards the SFT (the initialization of the 1<sup>st</sup> episode). The two resulting fronts are similar. However, in our iterative experiments we interpolate towards the episode-specific initialization as it allows maintaining a constant  $\eta$  at each WARP iteration, enabling a smooth progression towards the high KL regions.

E ADDRESSING LENGTH BIAS IN WARP



(a) Length bias in iterative WARP. (b) Adding length penalty (LP). (c) Benefits of diversity.

Figure 18: **Addressing length bias in WARP.** Figure 18(a) explores how length and KL change in successive WARP iterations. Figure 18(b) demonstrates the effectiveness of length penalty (LP) in reducing output length, and how such policies can merge with others trained without LP. Finally, Figure 18(c) shows that merging policies trained with different objectives further improves the KL-reward trade-off.

**Problem: length bias.** We investigate a potential length bias in WARP. LLMs after RLHF tend to be unnecessarily verbose (Shen et al., 2023) because RMs often prefer longer generations to shorter ones, leading to this form of reward hacking. We confirm such a phenomenon in Figure 18(a), where the length of the generation increases with higher KL values. This trend is even more pronounced in iterative WARP, where the 3<sup>rd</sup> iteration generates longer sentences than the 1<sup>st</sup> iteration at same KL.

**Mitigation strategy: length penalty.** To mitigate this length bias, we integrate a length penalty (LP) into the reward:  $-0.0005 \times \text{len}(y)$ , following Singhal et al. (2023). From SFT, we launch one RL fine-tuning run with LP, highlighted with red stars in Figure 18(b). This LP leads to shorter outputs as KL increases along training, in contrast to policies trained without LP.

**SLERP with different configurations.** Figure 18(b) displays the generation lengths from a SLERP of two policies, one trained with the LP and the other without. Critically, merging policies from diverse training configurations not only mitigates the length bias but also improves the Pareto front, as illustrated in Figure 18(c). This improvement is likely due to the increased diversity across policies, which appears beneficial for generalization, as shown in supervised learning (Ramé et al., 2022).

**Conclusion.** Those experiments highlight the possibility to fix the length bias, and also the benefits of merging policies trained with diverse rewards.

## F DIVERSITY IN PREDICTIONS

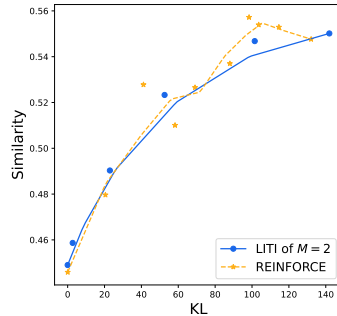
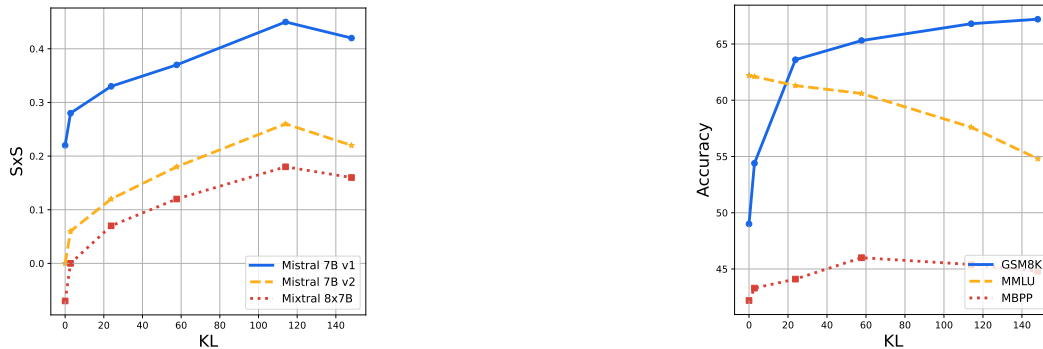


Figure 19: **Confirming diversity loss in RLHF.** The  $x$ -axis is the KL compared to the SFT initialization; the  $y$ -axis is the similarity across two generations from a given policy when decoding with temperature 0.9.

Finally, we investigate the loss in diversity across generations when aligning LLMs, as reported in Kirk et al. (2024). This could have negative consequences for creative or exploratory tasks, or even lead to policy collapse (Moalla et al., 2024; Hamilton, 2024). In Figure 19 we plot the BLEURT similarity (Sellam et al., 2020) across generations, during REINFORCE, or in LITI (as we interpolate back towards the SFT initialization). We observe that KL is strongly positively correlated with similarity across generations, confirming that RLHF induces a loss of diversity across generations. This experiment confirms that protecting the KL enables to trade-off between alignment and other benefits from pre-training, such as diversity in generations.

## G SxS AND BENCHMARK SCORES AT DIFFERENT KL



(a) SxS.

(b) Benchmark.

Figure 20: **SxS and benchmark scores at different KL.** The different checkpoints were obtained by LITI between the SFT and the SLERP at the end of the 3<sup>rd</sup> iter, with coefficients  $\eta \in \{0, 0.1, 0.3, 0.5, 0.8, 1.0\}$ . In particular, the one with  $\eta = 0.8$  was highlighted in Tables 1 and 2. In terms of SxS, hacking appears around around 110 of KL. In terms of accuracies, the alignment tax is benchmark dependent; while GSM8K seems to benefit from RLHF, scores on MMLU significantly reduces while they are stable on MBPP.