# Paper Critique: SimPO: Simple Preference Optimization with a Reference-Free Reward

Dan Peng

Department of Computer Science

Dec 23, 2024

danpeng@unc.edu

## 1. Research Problem

### 1.1. What research problem does the paper address?

This paper tackles the challenge of large language model (LLM) alignment through preference optimization, introducing SimPO as an elegant alternative to Direct Preference Optimization (DPO). SimPO creates a more direct path between how models learn preferences and how they generate text, addressing a fundamental disconnect in current alignment methods.

### 1.2. What is the motivation of the research work?

The authors were motivated by a critical insight: existing preference optimization methods like DPO suffer from a mismatch between training objectives and actual text generation metrics. This discord causes models to learn preferences inconsistently, leading to suboptimal performance. Additionally, DPO's reliance on reference models creates unnecessary computational overhead. SimPO aims to bridge this gap with a refreshingly streamlined approach.

## 2. Technical Novelty

### 2.1. What are the key technical challenges identified by the authors?

The authors identify two major hurdles in preference optimization: (1) the mismatch between DPO's reward formulation and the likelihood metrics that guide generation, creating a disconnect between training and inference; and (2) the computational burden of maintaining reference models during training, which increases memory consumption and slows down the optimization process.

### 2.2. How significant is the technical contribution of the paper? If you think that the paper is incremental, please provide references to the most similar work

SimPO represents a substantial conceptual leap forward rather than a mere incremental improvement. While it builds upon foundations laid by DPO, it introduces a fundamentally different reward formulation that aligns training signals with real-world generation processes. The paper's approach diverges significantly from previous works like ORPO, which attempts reference-free learning but still lacks SimPO's elegant alignment between rewards and generation metrics.

### 2.3. Identify 1-5 main strengths of the proposed approach.

- **Simplicity as strength**: SimPO eliminates reference models entirely, making it more computationally efficient and conceptually cleaner than competing approaches.

- **Length normalization magic**: By normalizing rewards based on sequence length, SimPO elegantly solves the length exploitation problem that plagues many preference optimization methods.

- **Stellar empirical performance**: SimPO doesn't just win by small margins—it dominates with substantial performance gains of up to 6.4 points on AlpacaEval 2 and 7.5 points on Arena-Hard compared to DPO.

### 2.4. Identify 1-5 main weaknesses of the proposed approach.

- **Mathematical reasoning degradation**: SimPO shows concerning performance drops on mathematical reasoning tasks like GSM8K, suggesting it may optimize for conversational fluency at the expense of precision-critical tasks.

- **Theoretical foundation gaps**: While empirically successful, SimPO lacks rigorous theoretical analysis explaining why this approach outperforms alternatives, particularly regarding the optimal margin value, which is determined through trial and error.

- **Potential for reward hacking**: Without explicit regularization against a reference model, SimPO could theoretically drift toward reward exploitation, though the authors didn't observe this in practice.

## 3. Empirical Results

### 3.1. Identify 1-5 key experimental results, and explain what they signify.

- SimPO-enhanced Gemma-2-9B-it achieving a 72.4% length-controlled win rate on AlpacaEval 2, significantly outperforming even GPT-4 Turbo (55.0%). This represents a quantum leap in open-source model capability, positioning smaller models to compete with much larger proprietary systems.

- The ablation studies demonstrating that removing length normalization causes models to generate verbose, lower-quality text, proving that SimPO's design choices directly address a critical flaw in preference optimization.

- SimPO's performance on Chatbot Arena, ranking first among all ¡10B models according to real human evaluators, validates that the algorithmic improvements translate to genuine human preference alignment, not just metric optimization.

### 3.2. Are there any weaknesses in the experimental section (i.e., unfair comparisons, missing ablations, etc)?

The experimental section lacks detailed investigation into why SimPO underperforms on mathematical reasoning tasks. While the authors acknowledge the issue, they don't thoroughly explore mitigations like the supervised fine-tuning regularization that has helped other methods maintain reasoning capabilities. Additionally, the paper would benefit from ablations that directly measure how SimPO's rewards correlate with human preferences, rather than relying primarily on benchmark performance as a proxy. A more diverse set of evaluation metrics beyond chat-based benchmarks would provide a more complete picture of SimPO's strengths and limitations.

## 4. Summary

SimPO represents an intellectual breakthrough that realigns preference optimization with how language models actually generate text—offering a "why didn't we think of this before?" moment for the field. I'm impressed by its elegantly simple formulation that produces remarkably strong results across challenging benchmarks. However, I'm also concerned about its mathematical reasoning degradation and lack of theoretical guarantees. What makes SimPO fascinating isn't just its performance but how it exposes a fundamental mismatch in previous approaches between training and inference objectives. Like adding power steering to a car, SimPO doesn't change the destination but makes the journey smoother, faster, and more efficient.

## 5. QA Prompt for a Paper Discussion

### 5.1. Discussion Question

Why does SimPO outperform DPO so dramatically despite its simpler formulation, and what does this reveal about the nature of preference optimization in large language models?

### 5.2. Your Answer

SimPO's dramatic outperformance of DPO reveals a fascinating paradox in AI alignment: sometimes apparent sophistication actually introduces misalignment. DPO's use of a reference model creates a "telephone game" effect—where what the model learns about preferences gets distorted between training and generation. SimPO creates a direct path between training signals and generation behavior, like teaching someone to drive by letting them actually hold the wheel rather than watching someone else drive.

This alignment principle—that learning objectives should mirror deployment objectives—has implications beyond preference optimization. It suggests we might be overcomplicating other aspects of AI training, and that aligning incentives between training and deployment might be more important than sophisticated algorithmic machinery.