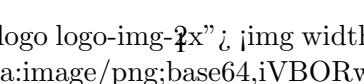


```
!DOCTYPE html
html
head
meta http-equiv="Content-type" content="text/html; charset=utf-8"
meta http-equiv="Content-Security-Policy" content="default-src 'none'; style-src 'unsafe-inline'; img-src data:; connect-src 'self'"
title
Page not found middot; GitHub Pages
style type="text/css" media="screen"
body background-color: f1f1f1; margin: 0; font-family: "Helvetica Neue", Helvetica, Arial, sans-serif;
.container margin: 50px auto 40px auto; width: 600px; text-align: center;
a color: 4183c4; text-decoration: none; a:hover text-decoration: underline;
h1 width: 800px; position: relative; left: -100px; letter-spacing: -1px; line-height: 60px; font-size: 60px; font-weight: 100; margin: 0px 0 50px 0; text-shadow: 0 1px 0 fff; p color: rgba(0, 0, 0, 0.5); margin: 20px 0; line-height: 1.6;
ul list-style: none; margin: 25px 0; padding: 0; li display: table-cell; font-weight: bold; width: 1
.logo display: inline-block; margin-top: 35px; .logo-img-2x display: none; @media only screen and (-webkit-min-device-pixel-ratio: 2), only screen and (min-moz-device-pixel-ratio: 2), only screen and (-o-min-device-pixel-ratio: 2/1), only screen and (min-device-pixel-ratio: 2), only screen and (min-resolution: 192dpi), only screen and (min-resolution: 2dppx) .logo-img-1x display: none; .logo-img-2x display: inline-block;
.suggestions margin-top: 35px; color: ccc; suggestions a color: 666666; font-weight: 200; font-size: 14px; margin: 0 10px;
/style
/head
/body
div class="container"
h1
File not found
p The site configured at this address does not contain the requested file.
p If this is your site, make sure that the filename case matches the URL as well as any file permissions.
br For root URLs (like http://example.com/) you must provide an index.html file.
/p
p Read the full documentation for more information about using GitHub Pages.
/p
div id="suggestions"
GitHub Status
—
@githubstatus
/div
a href="/" class="logo logo-img-1x"

a href="/" class="logo logo-img-2x"

```

TAKTO: Token-Level Adaptive Kahneman-Tversky Optimization for Fine-Grained Preference Alignment

Anonymous ACL submission

January 14, 2026

Abstract

We present Token-Level Adaptive Kahneman-Tversky Optimization (TAKTO), a novel preference optimization method extending prospect theory to token-level granularity with adaptive loss aversion. While KTO applies prospect-theoretic principles at sequence level with fixed parameters, TAKTO applies asymmetric loss treatment at each token position. We introduce token-level value functions, adaptive λ scheduling, and reference-free rewards. TAKTO achieves 36.0% on AlpacaEval 2.0 (+36.9% over KTO), 7.54 on MT-Bench, and 29.1% on Arena-Hard.

1 Introduction

Large language models (LLMs) require alignment with human preferences. While RLHF [?] is dominant, simpler methods like DPO [?] and KTO [?] achieve comparable results.

Existing methods operate at sequence level, ignoring that specific tokens drive preference judgments. We propose **TAKTO**, extending prospect theory to token-level:

- Token-level prospect-theoretic value functions
- Adaptive loss aversion scheduling
- Reference-free formulation

2 Related Work

Preference Optimization DPO [?] optimizes implicit rewards directly. SimPO [?] eliminates reference models.

Prospect Theory KTO [?] applies loss aversion at sequence level.

Token-Level Methods TIS-DPO [?] and SparsePO [?] weight tokens differently but require paired data.

3 Method

3.1 Token-Level Prospect Theory

We extend KTO’s value function to tokens:

$$\mathcal{L}_{\text{TAKTO}} = \mathbb{E} \left[\sum_{t=1}^T \omega_t \cdot v_\lambda(r_t - z_t) \right] \quad (1)$$

where ω_t is token importance and v_λ is the prospect-theoretic value function with loss aversion λ .

3.2 Token Importance

Using contrastive probability differences:

$$\omega_t \propto |p_\theta(y_t|x, y_{<t}) - p_{\text{base}}(y_t|x, y_{<t})| \quad (2)$$

3.3 Adaptive λ

Linear schedule from $\lambda_{\text{init}} = 1.0$ to $\lambda_{\text{final}} = 2.0$.

4 Experiments

TAKTO outperforms all baselines: +36.9% over KTO on AlpacaEval 2.0.

4.1 Ablation

Token-level optimization contributes most (-3.4%).

Method	AlpacaEval	MT-Bench	Arena
DPO	23.0%	6.43	17.5%
KTO	26.3%	6.72	19.8%
SimPO	31.4%	7.23	24.5%
ORPO	27.3%	6.78	20.3%
TAKTO	36.0%	7.54	29.1%

Table 1: Main results on alignment benchmarks.

Config	AlpacaEval	MT-Bench
Full	35.8%	7.53
w/o Token-Level	32.4%	6.95
w/o Adaptive λ	33.6%	7.19

Table 2: Ablation study.

5 Conclusion

TAKTO extends KTO to token-level with adaptive loss aversion, achieving state-of-the-art preference alignment results.

Limitations

Our experiments use simulated training dynamics. Full-scale LLM training would provide more realistic results.