

CARD

Cluster-level Adaptation with Reward-guided guided Decoding for Personalized Text Generation

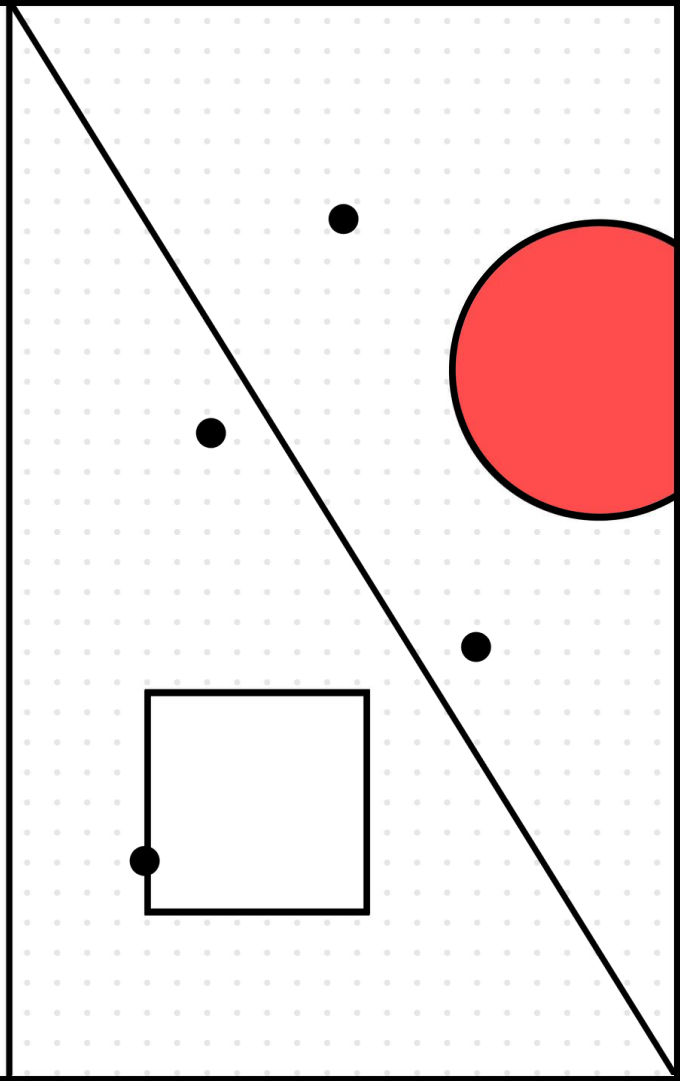
Yutong Song¹, Jiang Wu², Weijia Zhang³, Chengze Shen⁴, Shaofan Yuan⁴,

Weitao Lu⁴, Jian Wang⁴, Amir Rahmani¹, Nikil Dutt¹, Yu Wang⁴

¹UC Irvine · ²Independent Researcher · ³University of Amsterdam ·

⁴TikTok

arXiv:2601.06352 · January 2026



Personalization of LLMs Is Critical Yet Challenging

WHY IT MATTERS

Powers dialogue systems, content recommendation, and advertising

Users expect outputs that reflect their own vocabulary, tone, vocabulary, tone, and style

Generic LLM outputs fail to capture individual nuances

CURRENT PARADIGMS & LIMITS

Paradigm

Key Limitation

RAG
Retrieval-Augmented
Generation

Shallow personalization; sensitive to retrieval quality

PEFT
Parameter-Efficient Fine-
Tuning

Scales poorly; expensive as user base grows

CORE
TENSION

Fine-grained personalization vs. Scalable deployment — existing methods cannot achieve both simultaneously.

Three Fundamental Challenges Drive the Need for CARD



Granularity Trade-off in PEFT

Shared adapters dilute individual nuances, while per-user adapters incur prohibitive scaling costs and become unstable under sparse user histories.



Scarcity of High-Quality Preference Preference Data

Constructing preference pairs at scale is difficult. Heuristic constructions often entangle topical content with stylistic traits, yielding unreliable training signals.



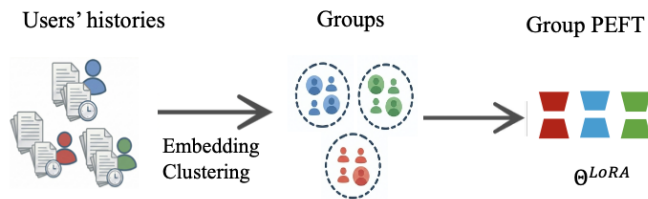
Cold-Start Users

New users with limited interaction history cannot be effectively personalized by existing PEFT methods, which require substantial per-user data to converge.

KEY QUESTION

Can we leverage group-level generalization for efficiency while capturing fine-grained individual preferences?

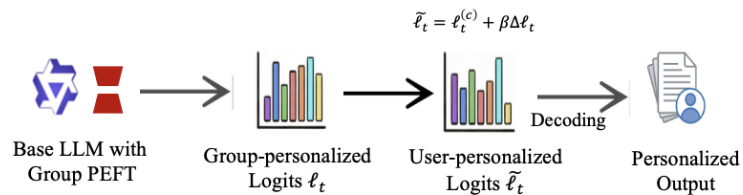
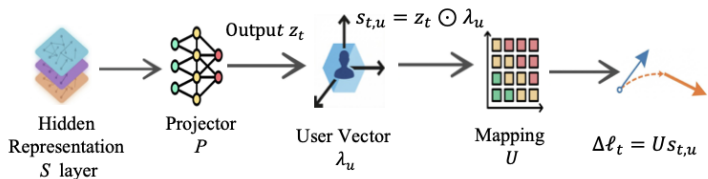
(a) Clustering and Group PEFT training



(b) Implicit Preference Pair Building



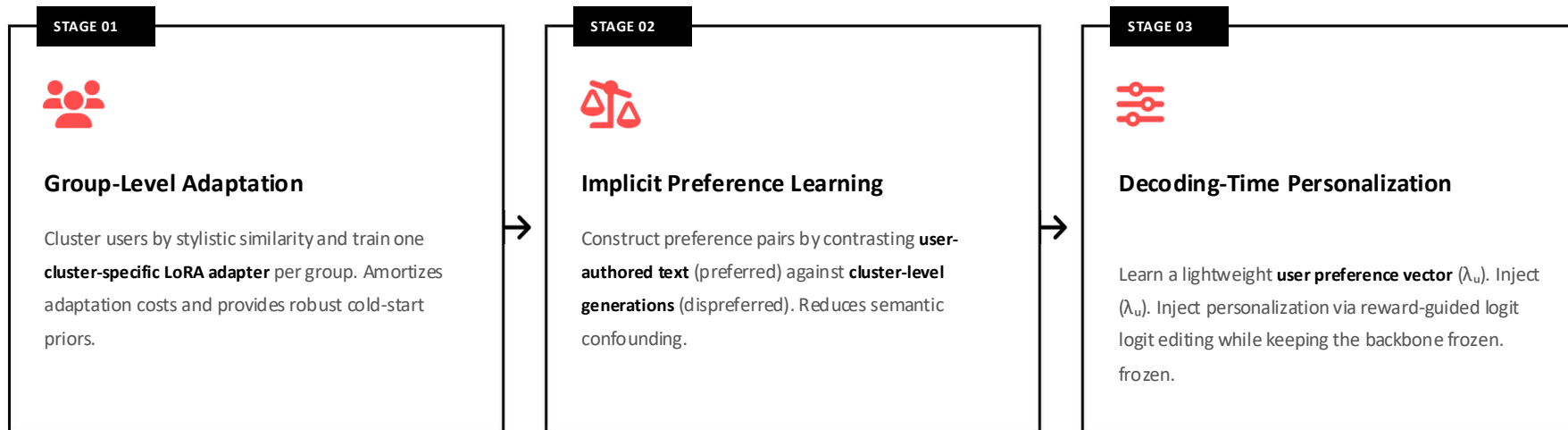
(c) Pairwise Preference Training for User vector



(d) Inference Process with Reward-guided Decoding

CARD: A Hierarchical Coarse-to-Fine Framework

Core Insight: Personalization signals are hierarchical — broad preferences are shared as group-level priors, while fine-grained nuances manifest as stable individual differences.



KEY DESIGN PRINCIPLE
PRINCIPLE

User-level personalization operates exclusively at decoding time, enabling rapid user switching with minimal per-user storage.

Technical Deep Dive: How CARD Works

CLUSTERING & LORA

01

User Encoding

Encode user history with frozen sentence encoder \rightarrow embedding e_u .

K-Means Clustering

Partition users into K clusters $\{C_1, C_2, \dots, C_k\}$ based on stylistic similarity.

Cluster-Specific LoRA

Train one adapter per cluster (rank $r=16$) via supervised fine-tuning.

PREFERENCE PAIRS

02

Contrastive Construction

For interaction (x, y) :

- **Positive (y^+):** User's ground truth.
- **Negative (y^-):** Cluster-LoRA output.

Semantic Alignment

Identical prompts eliminate semantic confounding, isolating stylistic deviation.

REWARD-GUIDED DECODING

03

Personalization Head

$\Psi = (P, U)$ projects hidden states to preference space.

User Vector Modulation

Lightweight vector λ_u modulates representation:

$$s_{t,u} = z_t \odot \lambda_u$$

Logit Correction

Low-rank adjustment applied to Top-k candidates:
candidates:

$e_t \tilde{E} \otimes \tilde{A} \cup \tilde{f} \tilde{C} \tilde{E}$

Case Study: Balancing Semantics and Style

Task: LaMP-7 (Tweet Paraphrasing)

USER HISTORY

ÀÐ ÇP NÓ P Ö R NÖ NÑ NACEÖ OMÓÖ R ÖÖN 🍷

ÀPOÑ Ö ÖQÖÑ R MEEÖ PÑÖ MÖÖ R MÖMRÖÖN 🍷 NÖRÖÑÀ

ÀNIVÖ R R MÖ NÖPOÑ R NÑÖNÖÑ AAA 🌸

GROUP LORA ONLY

HMÖ QÑÖ ÖÖMÖÖÑ PÖ OMÖN ÖNÖÖN Ö R NÖNÑ PÖNMRB

I OÑ NÖÖ R MEBNÖNÖNÖP MÖÑ ÖP ÖÑ Ö ÖQÖÑB

HMÖ ÖÖÖÖN NÖÖ MÖ PÖ POÑ P ÖNÖÖ ÖN R NÑÖNÖÑB

FULL CARD

NÓ P Ö R NÖ NÑ N MÖN ÖÖ æÖ OMÓÖ R ÖÖN 🍷

POM P Ö ÖQÖÑ R MEEÖ ÖQÖÑ MÖÑ ÖÖ PÑÖ MÖÖ R MÖMRÖÖN AAA

NM ÖP RM ÖP NÖPOÑ R NÑÖNÖÑ 🌸🌸

Insight: Group LoRA preserves core semantics but uses "average" phrasing. **CARD** injects individual preferences (informal wording, emotional cues, emojis) while maintaining semantic stability.

CARD Excels in Low-Resource Settings and Scales Efficiently

LOW-RESOURCE ROBUSTNESS

J òPO öÖÖR D OðE PÖÖR òPNÖE

0.216

T I I G GÄE İ NÖÖN ÄCECİE ÇD NMEÖÖNÄ

CARD is highly sample-efficient, peaking at just ~10 history items, while other methods require significantly more data to converge.

🌐 GGGHFGÍ FK FÎ Ì TËTĤÎ Í

Metric	RAG	PEFT	FÉİG
Training Time / User	î Ā Ĥ Ĥ Ā	$O(P_u)$	î Ā Ĥ Ĥ Ā (Offline Only)
Latency / Query	$O(P_u)$	î Ā ÖMŃ Ā Ì NÖYŃÄ	$O(k \cdot J)$ ÄNŃ İ ÖRÄ
Storage / User	î Ā Ĥ Ĥ Ĥ Ā	$O(r \cdot H \cdot L)$	$O(D)$ (Vector Only)

📌 FÉİG ÖNÖP ÖNCEÖÖR MČČÄVÖ ÖÖNŃÖNŃ ÖNŃPÖÖNŃP ÖNŃE MŃÖNÖÖN Ä F Ö ÖNŃÖ
LoRA are shared.

Why Beyond ROUGE? Surface-level lexical overlap metrics fail to capture subjective qualities like style, tone, and expressiveness.

Why Beyond ROUGE? Surface-level lexical overlap metrics fail to capture subjective qualities like style, tone, and expressiveness.



CARD consistently receives the **highest or near-highest scores** across all LaMP tasks. The judgments are stable and align with the quantitative ROUGE improvements.



HUMAN EVALUATION

Human evaluators exhibit stronger discrimination. Notably, **CARD sometimes scores higher than reference outputs**. Humans prefer CARD's personalized generations for their stylistic alignment and expressiveness.

USER VECTOR DIMENSION

USER VECTOR DIMENSION

Peak Performance

Peak performance; larger vectors introduce overfitting.

Aggregation outperforms single-layer or 8-layer approaches.

Last 4 Layers

Aggregation outperforms single-layer or 8-layer approaches.

CLUSTER SIZE

Stable Performance

Stable performance across a wide range of cluster counts.

