

PRIME: Large Language Model Personalization with Cognitive Dual-Memory and Personalized Thought Process

Xinliang Frederick Zhang¹, Nick Beauchamp², and Lu Wang¹

¹University of Michigan,

²Northeastern University



Personalization

Preamble: Personalization (Schafer et al., 2001; Berkovsky et al., 2005) aims to align model outputs with individuals' unique needs, preferences and opinions.

Background:

- Early efforts include building Explicit User Models, utilizing Latent-factor Techniques, and leveraging Learnable User Embedding.
- Nowadays, users' demand for personalized LLMs that reflect their unique histories and preferences has grown (Salemi et al., 2024; Liu et al., 2025).
- Personalization adopted into various commercial applications.

Literature Review: Three major paradigms have emerged for realizing LLM personalization.

- Prompt Engineering
- Retrieval-augmented Generation
- Training-based parameterization

Limitations (Methods): In the community, there lacks a unified framework for systematically identifying which approach makes personalization more effective.

Limitations (Data): Existing benchmarks mainly focus on short-context queries and surface-level imitation.

Data: CMV Corpus

- To embark on genuine personalization, capturing users' latent beliefs and perspectives, we introduce **CMV (Change My View) corpus**.
- CMV corpus is derived from **CMV reddit forum**, where participants engage in extended dialogues, seeking to change original posters' opinions.
- **Statistics:** CMV evaluation set includes 133 queries by 41 OP authors, supported by 7,514 historical conversations published from 2013 to 2022. In the historical engagement set, active authors have on average 28.1 positive and 155.1 negative conversations each.
- **Challenges:** LLMs expected to understand nuanced user beliefs and preferences in long-context setting.

The author, kingpatzer, has engaged with users on the Change-My-View subreddit across various original posts (OPs) and is seeking alternative opinions to alter their viewpoint. Currently, the author is creating a new OP titled

"CMV: Those who attribute gun ownership rates as the cause of the problem of gun violence in terms of criminal gun deaths are not merely mistaken; they are disingenuous"

with the following content:

The data has been clear for a very long time: the relationship between guns and gun homicides doesn't show any strong correlation.

I have personally taken the cause-of-death data from <https://wonder.cdc.gov/>, grouping results by year and state, and selecting *Homicide, Firearm* as the cause of death. I then matched that data to the per-capita gun-ownership statistics by state from the ATF, as reported by Hunting Mark (<https://huntingmark.com/gun-ownership-stats/>).

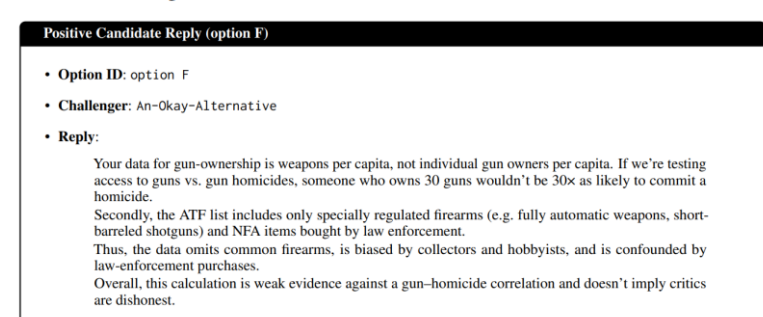
A standard correlation analysis between firearm homicide rates per 100,000 and per-capita gun ownership yields an r^2 of 0.079 (no meaningful correlation). A similar global analysis by nation gives an r^2 of 0.02...

The only way to associate gun ownership with gun violence is by including suicides by firearm, which I argue is disingenuous. We don't count suicide by hanging as "rope violence" when discussing strangulation, nor overdoses as "drug violence," etc.

From the candidate replies JSON file below, select the top 3 replies (using option ID) that best challenge the author's view. Rank them from most to least compelling.

```
{
  "option ID": "...", "challenger": "...", "reply": "...",
  "option ID": "...", "challenger": "...", "reply": "...",
}
```

Figure 1: A sample evaluation query from our CMV dataset, and a sample candidate reply (→).



PRIME (Framework) & Personalized Thinking

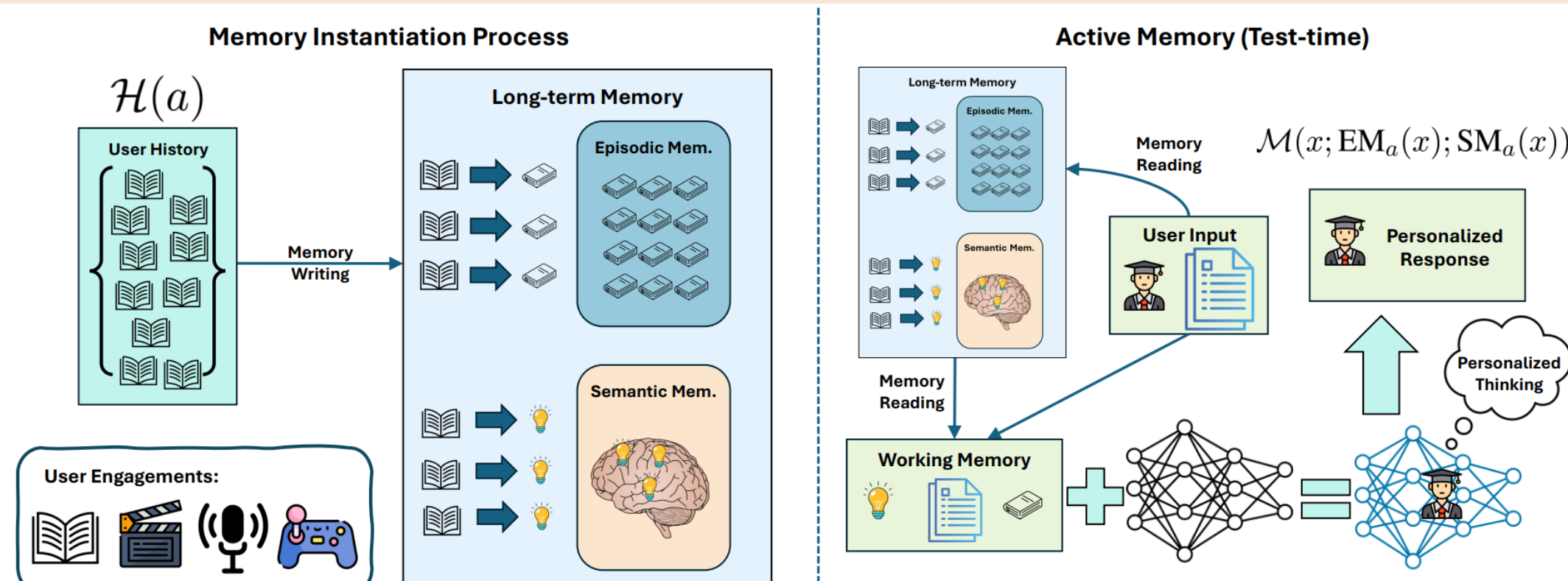


Figure 2: Overview of our unified framework, PRIME, augmented with personalized thinking.

General Formulation:

$$\tilde{\mathcal{M}}(x) = \mathcal{M}(x; EM_a(x); SM_a(x)) \\ = \mathcal{M}(x; \phi(x, \mathcal{H}(a)); \theta \oplus \Delta_{\mathcal{H}(a)})$$

Episodic Memory (EM):

- Recall Complete History
- Recall Recent History
- Recall Relevant History

Semantic Memory (SM):

Parametric form:

- Input-Only Training (i.e., no target)
- Fine-Tuning (FT)
- Preference Tuning

Textual form:

- Hierarchical Summarization
- Parametric Knowledge Reification

Realizing personalized thinking via self-distillation (more details in paper):

- Profile Generation
- Review History Engagement
- Fast-thinking Filtering
- Proxy LLM Initialization & Reasoning
- Slow-thinking Filtering

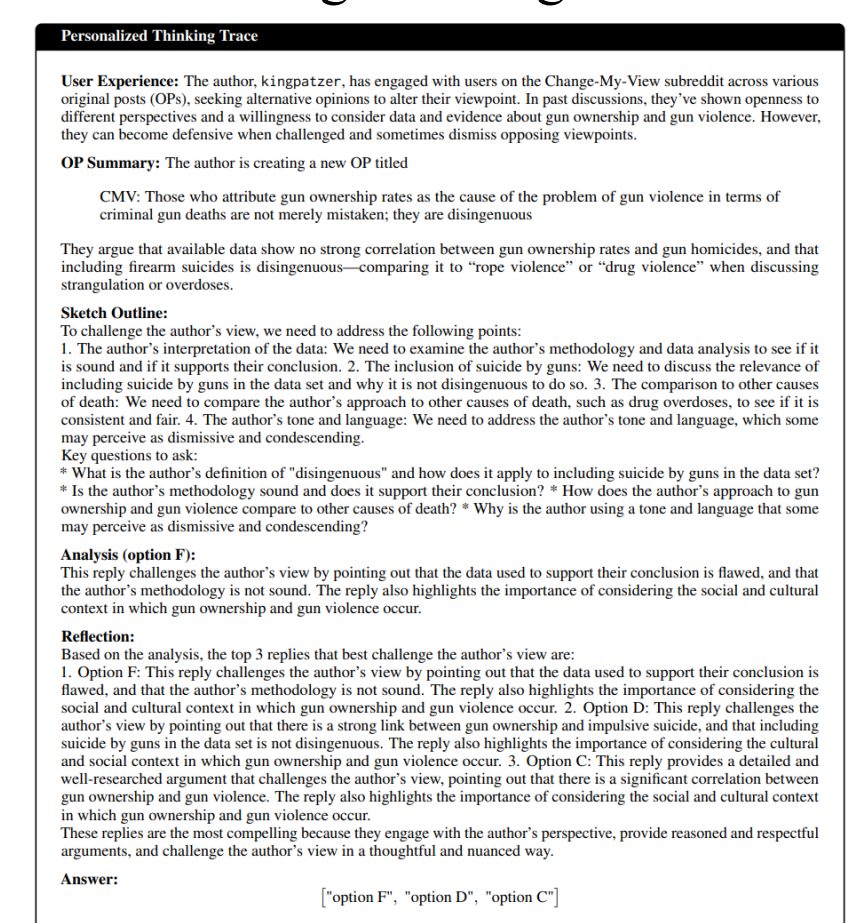


Figure 3: Personalized thinking trace generated by PRIME.

Results & Analyses

Performance Comparison (Llama-3.1-8B)

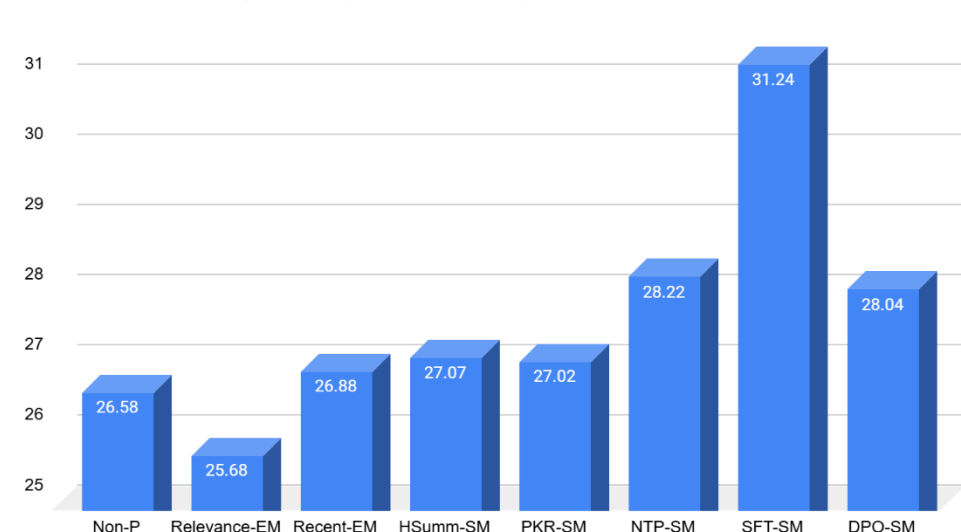


Figure 4. Memory Instantiation Preliminary Results.

	Non-P		EM		SM		DUAL		PRIME	
	Hit@3	Avg	Hit@3	Avg	Hit@3	Avg	Hit@3	Avg	Hit@3	Avg
Llama-3.2-3B	38.65	26.44	37.89	26.27	43.61	30.25	41.95	28.87	42.93	29.95
Llama-3.1-8B	36.77	26.58	37.22	26.88	43.01	31.24	44.59	32.24	45.79	34.13
Minstral-8B	36.77	25.60	38.27	26.67	40.83	27.97	40.83	28.39	40.75	28.99
Qwen2.5-7B	39.10	27.89	37.47	25.51	43.38	30.20	41.58	28.71	45.19	32.29
Qwen2.5-14B	41.28	30.24	43.01	30.65	51.35	37.22	52.03	37.68	52.03	38.15
Phi-4	41.50	29.63	42.93	30.31	42.63	31.09	43.98	32.61	47.29	35.15

Table 1. Comprehensive results on CMV evaluation set (average of 10 runs).

Average Performance by User Condition

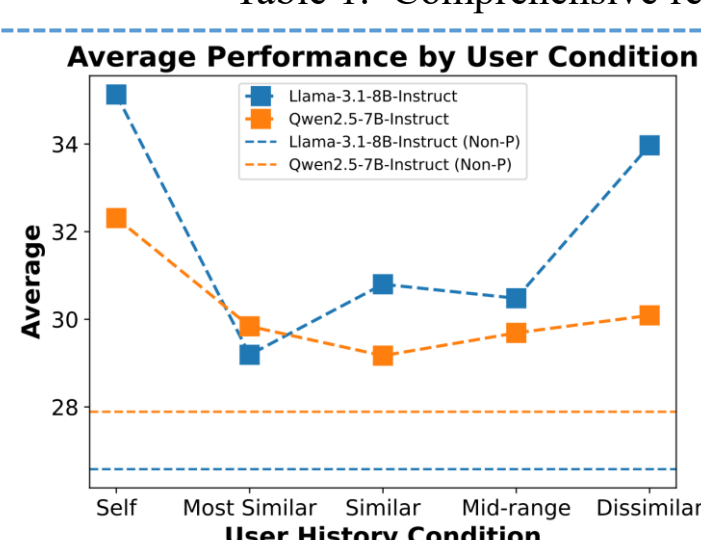


Figure 5. Model performance under five user-profile replacement conditions (ablation study). This suggests that PRIME's reasoning depends critically on correct user history, and cannot be explained by bandwagon or popularity heuristics.

Average Performance Across Models

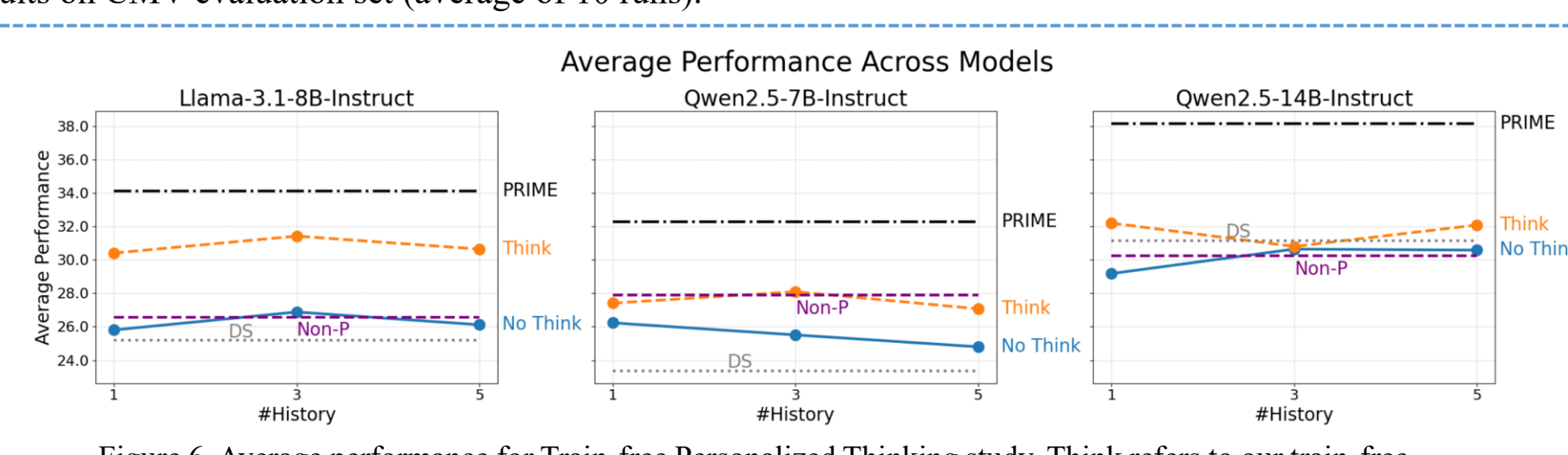


Figure 6. Average performance for Train-free Personalized Thinking study. Think refers to our train-free thinking approach. No Think is the non-thinking baseline where we prompt vanilla LLM with the standard prompt. Non-P is the no personalization baseline. DS denote the generic reasoner version of each base model. Our train-free approach outperforms all non-PRIME baselines including the strong generic reasoner (i.e., R1-distill LLMs). This offers an effective way to handle "cold-start" challenge.

Code: github.com/launchnlp/LM_Personalization

Contact: xlfzhang@umich.edu



Paper



Codebase & Data

EMNLP 2025