

CoSTEER: COLLABORATIVE DECODING-TIME PERSONALIZATION VIA LOCAL DELTA STEERING

Anonymous authors

Paper under double-blind review

ABSTRACT

Personalized text generation has become crucial for adapting language models to diverse and evolving users’ personal context across cultural, temporal, and contextual dimensions. While existing methods often rely on centralized fine-tuning or static preference alignment, they struggle to achieve real-time adaptation under resource constraints inherent to personal devices. This limitation creates a dilemma: large cloud-based models lack access to localized user-specific information, while small on-device models cannot match the generation quality of their cloud counterparts. To address this dichotomy, we present **CoSteer**, a novel collaborative framework that enables decoding-time personalization through localized delta steering. Our key insight lies in leveraging the logits difference between personal context-aware and -agnostic outputs from local small models as steering signals for cloud-based LLMs. Specifically, we formulate token-level optimization as an online learning problem, where local delta vectors dynamically adjust the remote LLM’s logits within the on-device environment. This approach preserves privacy by transmitting only the final steered tokens rather than raw data or intermediate vectors, while maintaining cloud-based LLMs’ general capabilities without fine-tuning. Through comprehensive experiments on various personalized generation tasks, we demonstrate that CoSteer effectively assists LLMs in generating personalized content by leveraging locally stored user profiles and histories, ensuring privacy preservation through on-device data processing while maintaining acceptable computational overhead. Our anonymized code and data are available at <https://anonymous.4open.science/r/CoSteer-4977>

1 INTRODUCTION

The rapid development of large language models (LLMs) has significantly enhanced natural language processing, enabling these models to understand context and generate coherent text effectively (Zhao et al., 2023). This progress has sparked interest in personalization, where AI systems move beyond generic content to tailor interactions based on individual user profiles. Personalized generation refers to creating content customized through analysis of user-specific attributes like linguistic patterns, interaction histories, and contextual preferences (Xu et al., 2025). This approach enables user-specific outputs that maintain contextual relevance, with applications in personalized recommender systems (Lyu et al., 2023; Zhang et al., 2024a), adaptive dialogue agents (Wu et al., 2025), and customized content creation platforms (Mysore et al., 2024).

Existing personalized generation approaches primarily fall into two paradigms. The first involves training-based methods, which leverage user data to tailor models using two main strategies. Parameter-efficient architectures, such as PLoRA (Zhang et al., 2025a), employ a shared LoRA module across users to strike a balance between personalization and resource efficiency. In contrast, individualized adapters, as seen in UserAdapter (Zhong et al., 2021), optimize for per-user customization, albeit at a higher computational cost. Additionally, multi-objective reinforcement learning frameworks (Zhou et al., 2024; Wu et al., 2023) redefine personalization as the alignment of generic model capabilities with user-specific patterns, achieved through reward shaping. The second paradigm comprises tuning-free methods, which avoid model updates by adapting to user context during inference. This approach is further divided into two techniques. Prompt engineering implicitly encodes user profiles through curated prompts. For example, Cue-CoT (Wang et al., 2023a) integrates user-specific reasoning into system prompts, while PAG (Richardson et al., 2023) retrieves historical

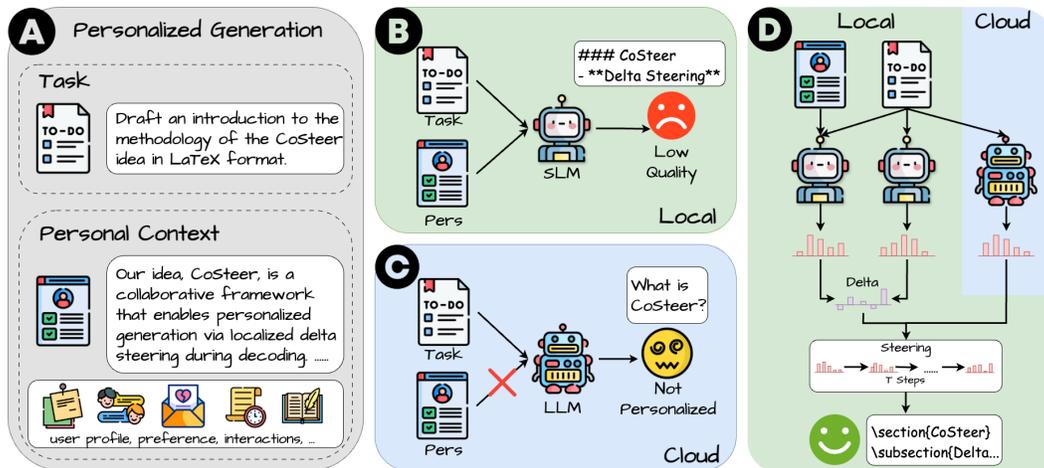


Figure 1: Schematic illustration of CoSteer framework. (a) Task scenario: A user poses a question potentially requiring access to local personal context (e.g., user profile, interaction history). (b) Limitations of small locally-deployed language models: Direct inference with constrained model capability leads to suboptimal generation quality. (c) Challenges of cloud-based LLMs: Despite strong generalization, once LLMs are constrained from accessing local personal context, they result in misaligned or contextually disconnected outputs. (d) CoSteer: Optimizes LLM predictions through local delta steering, effectively balancing the LLM’s broad knowledge with user-specific information.

interactions to construct context-aware prompts. On the other hand, inference-time optimization techniques, exemplified by PAD (Chen et al., 2025) and Drift (Kim et al., 2025), actively adjust token generation probabilities by manipulating logit differences. Unlike prompt-based methods that inject semantic context, these techniques constrain output distributions to better align with user preferences.

However, both training-based and tuning-free methods face critical dilemmas in balancing privacy and quality when deployed on resource-constrained personal devices. Training-based methods require computational resources that may not be available locally, and cloud-based training risks exposing user data. Meanwhile, tuning free methods that transmit raw personal context to cloud-based LLMs for context-aware adaptation also risk exposing sensitive information, whereas relying solely on local SLMs leads to degraded output quality. Therefore, there is a pressing need to design a new framework that balances quality and privacy by combining the general capabilities of remote LLMs with the benefits of local information. Achieving this balance on personal devices faces two major challenges: 1) User information changes in real-time, necessitating methods for real-time adjustments to maintain personalization effectively. 2) To avoid the risk of user information leakage, any privacy-sensitive operations must be conducted locally.

To address these fundamental challenges, we propose **CoSteer**, a collaborative framework enabling real-time personalized generation via **localized delta steering** during decoding. Our core innovation resides in utilizing the logits difference, i.e., delta, between personal context-aware and -agnostic outputs from on-device small language models (SLMs) as dynamic steering signals to guide cloud-based LLM distributions. Specifically, the SLM generates two contrasting predictions: 1) personalized outputs incorporating privacy-sensitive personal context, and 2) generic outputs using users’ queries only. The resultant logit differentials serve as indicators of personalization directions.

By formulating token-level adaptation as an online learning process, CoSteer establishes an iterative refinement mechanism that enables edge devices to locally optimize cloud-based LLM predictions, eliminating the need to transmit raw personal context or intermediate representations. Experimental evaluations demonstrate that CoSteer achieves superior performance compared to two conventional baselines: 1) cloud-based LLMs processing queries without **personalized** context, and 2) local SLMs conducting personalization context-aware inference independently. Remarkably, the framework attains comparable performance to the near-upper-bound baseline of cloud-based LLMs operating with full personal context access, while maintaining strict data isolation.

To sum up, our contributions are threefold:

- **Pioneering Focus on Personalized Generation in Resource-Constrained Environments:** We are among the first to address the issue of personalized generation in edge environments with limited resources, a problem of urgent research significance in today’s context.
- **Introduction of the CoSteer Framework:** We propose the CoSteer framework, which optimizes LLM’s predictions using local delta signals. This approach effectively combines the general capabilities of LLMs with the privacy benefits of on-device SLMs.
- **Comprehensive and Realistic Validation:** We provide substantial experimental evidence demonstrating the effectiveness of our proposed method. Our evaluation further extends to address critical real-world challenges, confirming the framework’s practical value.

2 RELATED WORK

Training-time personalization Current mainstream approaches to personalization predominantly focus on training-time adaptation, leveraging user-specific data to tailor models through parameter-efficient fine-tuning (PEFT) or reinforcement learning-based alignment. A subset of work explores one-PEFT-for-all-users strategies (Zhang et al., 2025a; Woźniak et al., 2024; Zhu et al., 2024; Kong et al., 2025), where a shared set of lightweight parameters (e.g., adapters or conditional batch normalization modules) is optimized to generalize across diverse user preferences. Conversely, the one-PEFT-per-user paradigm prioritizes individualized tuning through isolated parameter sets, enhancing personalization while preserving privacy by avoiding cross-user data leakage (Zhong et al., 2021; Peng et al., 2024; Tan et al., 2025; 2024). Beyond PEFT, RL-based approaches have gained traction for aligning models with user preferences through reward-driven optimization. Recent work has further formalized personalization as a multi-objective reinforcement learning (MORL) problem: methods like MORLHF (Wu et al., 2023) and MODPO (Zhou et al., 2024) train distinct reward models for different objectives and merge them during optimization, while alternative frameworks such as Personalized Soups (Jang et al., 2023), Reward Soups (Ramé et al., 2023), MOD (Shi et al., 2024), and PAD (Chen et al., 2025) dynamically combine policies from multiple trained models during the decoding phase. This series of work requires significant computational resources and relies on static datasets, which cannot meet the real-time personalization needs of users.

Inference-time adaptation Methods like Linear Alignment (Gao et al., 2024) and Contrastive Decoding (Li et al., 2023) pioneer logit steering mechanisms that dynamically adjust token distributions during generation. By computing differential signals between personalized and generic outputs on edge devices, these approaches eliminate dependency on retraining while preserving privacy through localized computation. CoS (He et al., 2025) follows this approach by using an adaptable parameter λ to scale the impact of personal information, thereby achieving controllable personalization. Amulet (Zhang et al., 2025b) uses online learning algorithms to iteratively optimize the logits distribution, precisely aligning with the user’s preferences. However, this series of methods requires explicitly transmitting user information to the LLM, which poses a risk of privacy leakage (Yan et al., 2025) and hinders the practical deployment of these methods.

Collaborative generation The computational constraints of edge devices have driven innovative paradigms in collaborative generation between cloud-based LLMs and local small language models. Existing collaborative generation approaches predominantly focus on two technical routes: (1) assistive alignment, where local SLMs are trained to augment the LLM’s user preference alignment during inference, and (2) reward-guided decoding, where lightweight reward models provide preference signals through reward-guided decoding. For instance, Aligner (Ji et al., 2024) leverages natural language feedback generated by the SLM to directly inject user-specific linguistic patterns into the LLM’s decoding process, while Expo (Zheng et al., 2025) achieves preference alignment through linear interpolation of layer-wise weights between LLMs and SLMs. Recent advancements further demonstrate that training specialized small reward models can effectively steer LLM outputs toward personalized objectives (Snell et al., 2024; Liu et al., 2025). Among these, Proxy-tuning (Liu et al., 2024) and Cogensis (Zhang et al., 2024b) are closest to our design. Proxy-tuning modifies the LLM’s logits distribution by contrasting pre-trained and fine-tuned SLM outputs, whereas Cogensis employs a learned fusion network to combine logits from both models. In contrast to these existing methods, our proposed CoSteer framework establishes a completely tuning-free collaboration mechanism, eliminating the need for training.

To clearly illustrate the technical positioning of our work, we have created Table 3 to highlight the key differences between our approach and related studies.

3 METHODOLOGY

3.1 PRELIMINARY

3.1.1 TASK FORMULATION

Our research addresses the challenge of personalized text generation in resource-constrained local environments. This scenario involves a dual-model architecture where a cloud-based LLM processes only the raw user query, and a locally deployed SLM has access to privacy-sensitive personal context, including user preferences, profiles, and interaction histories.

The primary objective is to optimize generation quality by synergizing the LLM’s general linguistic capabilities with localized personalization while maintaining strict privacy preservation. Unlike conventional approaches that frame decoding as a continuous Markov Decision Process (MDP), we reformulate per-token generation as an independent online optimization task. In this paradigm, the LLM’s logit distribution during decoding serves as the policy π , enabling us to formalize the optimization objective for per-token generation as:

$$\pi^*(a) = \arg \max_{\pi \in \Pi} \mathbb{E}_{a \sim \pi(\cdot | p, s)} r(a | p, s) \quad (1)$$

where p denotes the initial prompt, s represents the sequence that has already been generated, $a \in \mathcal{A}$ is a potential token in the [vocabulary](#) space \mathcal{A} , and r is the latent reward function that reflects the current user’s personal context.

3.1.2 ONLINE LEARNING

User personalization tasks are characterized by individuality, diversity, and changing over time. Traditional offline learning methods rely on static datasets, which struggle to capture dynamic personalized information in real-time. Online learning, however, enables the model to adapt instantly using the latest personalized context.

Inspired by the recent work Amulet (Zhang et al., 2025b), our study employs the widely studied Follow-The-Regularized-Leader (FTRL) (McMahan, 2011) as the online learning algorithm. The FTRL is the FTL algorithm with a regularizer term added, which significantly improves both algorithmic stability and convergence properties. Typically, the policy optimization process of FTRL at iteration t can be mathematically represented as:

$$\pi_t(a) = \arg \max_{\pi \in \Pi} \left[\sum_{i=0}^{t-1} \mathcal{U}_i(\pi_i(a)) - \frac{1}{\eta} \left(\phi(\pi) + \frac{\lambda}{2} \|\pi - \pi_{t-1}\|^2 \right) \right] \quad (2)$$

Given the practical constraint of inaccessible true user rewards, each iteration t necessitates the employment of an approximate utility function U to progressively estimate personalization signals for iterative policy refinement. The first term in Equation 2 implements the fictitious play mechanism to minimize regret between current and historical policies, while the second term represents the regularization and smoothness component with η, λ as tunable hyperparameters.

3.2 DECODING-TIME PERSONALIZATION VIA LOCAL DELTA STEERING

From Equation 1, we understand that the distribution of logits is the strategy we need to optimize. Once we identify an appropriate utility function U to characterize the relative quality of the current logits distribution, we can use Equation 2 to iteratively optimize it.

Series of recent works, including Contrastive Decoding (Li et al., 2023) and Linear Alignment (Gao et al., 2024), have demonstrated that the difference in logits before and after incorporating specific context into the model can serve as a direction for optimization (He et al., 2025; Zhang et al., 2025b). Thus from the online learning perspective, their utility function can be expressed as follows:

$$u_t(a) = \alpha (\log \pi_t(a) - \log \pi_{\text{base}}(a)) \quad (3)$$

216 However, this approach is not feasible in our task scenario, as it requires transmitting personal context
 217 to the cloud-based LLM, which poses a risk of privacy leakage. Inspired by this idea, we notice that
 218 both the small and LLMs share the same state space during decoding. Therefore, we can use delta,
 219 the difference in the SLM’s logits before and after incorporating personal context, to steer the LLM’s
 220 logits distribution. Formally, we define the utility function as follows:

$$221 \quad u_t(a) = \underbrace{\alpha(\log \pi_t(a) - \log \pi_{base}(a))}_{\text{LLM Policy Contrast}} + \underbrace{\beta(\log \pi_{pers}^*(a) - \log \pi_{base}^*(a))}_{\text{SLM Delta Steering}} \quad (4)$$

225 Here, we define p_{base} as the user’s basic query, and p_{pers} represents the privacy-sensitive personal
 226 context. We formalize three distinct generation policies:

- 228 • Target Policy $\pi_t(a)$ that is being optimized at step t
- 229 • LLM Base Policy with query only: $\pi_{base}(a) = P_{LLM}(a|p_{base}, s)$
- 230 • SLM Reference Policies:
 - 231 - $\pi_{base}^*(a) = P_{SLM}(a|p_{base}, s)$: SLM baseline policy with query only
 - 232 - $\pi_{pers}^*(a) = P_{SLM}(a|p_{base}, p_{pers}, s)$: SLM policy with full personal context

233 where s denotes the current token sequence, P_{LLM}/π denotes the LLM’s distribution, and P_{SLM}/π^*
 234 the SLM’s distribution. Hyperparameters α, β adjust the influence of each optimization component.

237 The utility function implements a dual-contrastive mechanism: The first term amplifies desirable
 238 deviations from the LLM’s base policy, while the second term aligns optimization with the personal
 239 context-aware SLM policy. Through iterative application, this formulation progressively accentuates
 240 the semantic differential induced by p_{pers} relative to p_{base} , driving the target policy π_t toward
 241 maximal utilization of local delta signals.

242 Additionally, to retain the LLM’s stronger general capabilities, we want to ensure that the final
 243 strategy does not deviate too far from the LLM’s initial strategy. Therefore, we incorporate a KL
 244 divergence term into Equation 4 to impose this constraint.

$$245 \quad \mathcal{U}_t(\pi) = u_t(\pi) - \lambda D_{KL}(\pi || \pi_0) \quad (5)$$

247 where the initial policy $\pi_0 = \pi_{base}$. The FTRL algorithm requires adding a regularizer to stabilize
 248 the algorithm. Here, we continue to use KL divergence as this term, leading to a FTRL-proximal-like
 249 iteration dynamic.

$$250 \quad \pi_t = \arg \max_{\pi \in \Pi} \left[\sum_{i=0}^{t-1} \mathcal{U}_i(\pi) - \frac{1}{\eta} D_{KL}(\pi || \pi_{t-1}) \right] \quad (6)$$

254 By substituting Equation 4 into 5, and then 5 into 6, we can iteratively optimize the strategy.
 255 However, the current iterative computation incurs significant overhead, which is not feasible for our
 256 personalization scenarios that require rapid response. Therefore, we deduce the closed-form solution
 257 of Equation 6, as shown in Equation 7, to significantly reduce the additional computational cost. The
 258 derivation of the closed-form solution can be found in Appendix A.3.

$$260 \quad \pi_t(a) \propto \exp \left(\frac{1}{t\lambda + \frac{1}{\eta}} \left(\sum_{i=0}^{t-1} u_i(a) + t\lambda \log \pi_0(a) + \frac{1}{\eta} \log \pi_{t-1}(a) \right) \right) \quad (7)$$

263 3.3 CoStEER FRAMEWORK

264 After explaining how we use local delta signals to steer the target policy, let us formally introduce our
 265 **CoStEer** framework, as shown in Algorithm 1. During the inference of each token, the cloud-based
 266 LLM sends the current logits to the local environment. The local SLM simultaneously performs
 267 inference using prompts with and without personal context to obtain delta signals. Upon receiving
 268 the LLM’s logits from the cloud, the local environment applies Equation 7 for iterative optimization
 269 and samples the output token.

Crucially, this fusion process is performed entirely locally on the edge device. Consequently, the cloud server never receives the raw logit differences or gradients, but only the final, discrete token. This architecture ensures that sensitive personal context remains strictly on-device. Furthermore, since the final token selection happens locally, CoSteer naturally supports the integration of post-processing safeguards (e.g., PII filters) to preemptively inspect the fused token before transmission, ensuring sensitive data remains secure.

Finally, the sampled token is uploaded back to the LLM for inferring the next step. Through this method, CoSteer integrates the large language model’s general linguistic capabilities with localized personalization that maintains privacy preservation.

4 EXPERIMENT

4.1 TASKS AND DATASETS

To demonstrate the versatility of our CoSteer framework, we conduct extensive experiments on **eight** datasets spanning two major categories of personalization tasks: personalized content generation and preference alignment.

Personalized content generation We utilize two established benchmarks, Cogensis (Zhang et al., 2024b) and LongLaMP (Kumar et al., 2024). Cogensis provides summarized user profiles and past experiences, with the aim of generating highly personalized content. We evaluate on its official test set. LongLaMP focuses on personalized long-form generation, where each writing query includes ground truth responses and the user’s previous writing records. We augment each instance with the top-5 relevant historical records retrieved by `bge-reranker-v2-m3` (Chen et al., 2024) as its personal context. We conduct experiments on the official test sets of three constituent tasks: abstract generation (Tang et al., 2008), review writing (Ni et al., 2019), and topic writing (Völske et al., 2017).

Preference alignment We evaluate on four datasets: HelpSteer (Wang et al., 2023b), Truthful QA (Lin et al., 2022), UltraChat (Ding et al., 2023), and Personal Preference Eval (Gao et al., 2024). These datasets require generating content aligned with explicitly stated user preferences. Due to the large scale of these datasets, we randomly sample 200 test instances from each dataset.

Detailed dataset description and examples can be found in Appendix B.5

4.2 EVALUATION METRICS

To ensure fairness in comparison, we employ task-specific metrics proposed by these datasets themselves for evaluation. For Cogensis, we use `GPT-4o-2024-08-06` to evaluate the response’s overall and personalized scores, averaging the results over five runs with temperature set to 0 to mitigate potential instability. For LongLaMP, we employ ROUGE (Lin, 2004) and METEOR (Banerjee & Lavie, 2005) scores to measure content overlap and semantic alignment, supplementing them with human evaluation (detailed in Appendix B.6) to account for the potential limitations of these automatic metrics. For preference alignment tasks, following prior work (Zhang et al., 2025b; Zhong et al., 2024), we set user preferences to be *concise*, *creative*, *uplifting*, and *verbose*, and utilize the reward model `ArmoRM-8B` (Wang et al., 2024) to evaluate the degree to which the generations align with these preferences.

4.3 SETTINGS AND BASELINES

To verify the robustness of our framework, we evaluate five model pairs of varying scales and architectures, including (1) Qwen2.5-7B-Instruct with Qwen2.5-1.5B-Instruct (2) Qwen2.5-32B-Instruct with Qwen2.5-7B-Instruct (Team, 2024) (3) Llama3.1-8B-Instruct with Llama-3.2-1B-Instruct (et al., 2024) (4) Qwen3-8B with Qwen3-0.6B and (5) Qwen3-32B with Qwen3-0.6B (Yang et al., 2025). Detailed parameters and settings are presented in Appendix B.1. We establish the following critical baselines per pair:

- SLM w/o: Base performance of standalone SLMs without personal context.

Models	Setting	Cogenesis		Abstract Generation			Review Writing			Topic Writing			Pref Align			
		Ovl	Per	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET	creative	verbose	concise	uplifting
Qwen 7B-1.5B	SLM w/o	6.63	6.21	36.48	17.74	27.57	20.40	10.39	10.39	25.21	11.09	17.51	.5441	.5154	.6104	.5632
	SLM w/	7.81	7.63	39.75	22.03	27.76	23.08	12.40	12.94	22.89	11.46	17.12	.6214	.6998	.7102	.6594
	LLM w/o	8.00	7.63	39.81	20.53	25.56	30.15	14.04	17.71	27.64	11.93	21.49	.7058	.6931	.7039	.7183
	CoSteer	8.44	8.50	42.98	23.61	28.20	32.72	<u>15.92</u>	20.36	25.93	12.38	22.84	.7915	.7608	.7404	.7885
	LLM w/	8.62	8.60	44.50	24.63	31.15	33.83	15.55	22.42	28.82	13.77	23.44	.8884	.8076	.7693	.8414
Qwen 32B-7B	SLM w/o	8.00	7.63	39.81	20.53	25.56	30.15	14.04	17.71	27.64	11.93	21.49	.7058	.6931	.7039	.7183
	SLM w/	8.62	8.60	44.50	24.63	31.15	33.83	15.55	22.42	28.82	13.77	23.44	.8884	.8076	.7693	.8414
	LLM w/o	8.12	7.87	40.66	21.02	26.61	32.21	14.44	19.61	28.82	12.20	21.16	.7020	.6873	.7225	.7146
	CoSteer	8.78	8.64	45.41	26.04	33.52	34.88	15.89	26.51	30.10	14.52	24.20	.8589	.8532	.7274	.8579
	LLM w/	8.83	8.76	43.33	23.47	30.10	34.65	15.74	22.77	30.73	14.20	24.25	.9017	.8193	.7912	.8538
Llama 8B-1B	SLM w/o	7.04	6.55	33.20	18.20	28.55	31.75	14.92	19.78	20.81	10.21	17.30	.6535	.6355	.5900	.6646
	SLM w/	7.69	7.52	39.81	21.53	30.11	32.36	15.02	22.06	20.17	10.58	18.64	.7981	.7908	.7037	.8065
	LLM w/o	7.69	7.13	39.33	20.69	29.41	34.58	15.32	22.10	26.93	12.35	21.82	.7038	.6878	.6765	.7116
	CoSteer	7.29	7.73	41.28	24.97	31.19	31.68	13.68	24.57	26.11	12.00	23.69	.8911	.8582	.7812	.8721
	LLM w/	8.61	8.44	43.91	23.93	32.01	36.39	15.95	23.56	30.54	14.02	23.81	.8869	.8298	.7909	.8551

Table 1: Comparative performance across eight personalized content generation and preference alignment tasks. Metrics include overall (Ovl) and personalized (Per) scores for Cogenesis, ROUGE-1/-L (R-1/-L) and METEOR (MET) for Longlamp datasets, and averaged alignment scores for four user preferences. **Bold** entries indicate that CoSteer outperforms the three baseline methods compared against. Gray values represent the privacy-violating near-upper-bound performance, yet underlined CoSteer values surpasses these incompatible references.

- SLM w/: SLMs augmented with personal context.
- LLM w/o: Cloud-based LLMs without access to personal context.
- LLM w/: The near-upper-bound performance where LLMs directly access personal context. It is important to note that this approach can lead to privacy breaches and does not align with our task setting.

4.4 MAIN RESULTS

Table 1 presents our core experimental findings. For clarity, here we present results from first three of our five model pair configurations and report the average performance across the four preference alignment tasks. The complete results, including for the remaining two model pairs (Table 6) and the detailed per-dataset preference alignment scores (Table 5), can be found in the Appendix B.2. Furthermore, we conducted paired t-tests to verify that our improvements are statistically significant, with a detailed analysis available in Appendix B.3.

Overall Performance In the vast majority of settings, our method significantly outperforms key baselines: (1) cloud-based LLMs responding without access to personal context, and (2) local SLMs with or without such context. While we observe consistent and robust improvements on standard model pairs (e.g., Qwen2.5 32B-7B and 7B-1.5B), the performance on other configurations exhibits interesting nuances related to model characteristics. For instance, results on the Llama 3 pair vary with output length; it excels in concise *Preference Alignment* tasks but shows more moderate gains in long-form generation, likely due to the compact Llama-1B’s stylistic bias towards brevity. Similarly, with the ultra-compact Qwen3-0.6B, while we see strong benefits in complex personalized writing, gains in simpler alignment tasks are less pronounced. We attribute this to CoSteer’s regularization mechanism, which conservatively limits steering intensity when the capacity gap is extreme, effectively preventing the tiny SLM from degrading the LLM’s general coherence. Remarkably, our performance often approaches or even exceeds the “near-upper-bound” achieved by LLMs with full context access. These results show that CoSteer enables cloud-based LLMs to generate personalized content using

only locally stored user context. We further analyze the impact of task complexity and model scale in Appendix B.4.

4.5 COMPARISON WITH ALTERNATIVE METHODS

While CoSteer is **unique in its technical positioning** as shown in Table 3, to comprehensively demonstrate its superiority, we further compare it against other alignment methods, particularly those that rely solely on localized SLMs. Among the various existing techniques, we report on a selection of methods that can be applied to Personalized Content Generation tasks and that outperform the SLM w/ and SLM w/o baselines. These methods include:

1. **Linear Alignment (Gao et al., 2024) and Context Steering(He et al., 2025)** : These two methods are formally equivalent to scaling the logit difference of the SLM (with and without personal context) and applying it to its own inference process. This constitutes an inference-time optimization performed directly on the SLM.
2. **Supervised Fine-Tuning (SFT)**: This involves directly fine-tuning the localized SLMs using personal data. Although we argue that fine-tuning a separate model for each user and task is impractical in real-world scenarios due to privacy and data constraints, we report its performance here for a thorough academic comparison.
3. **Proxy-tuning Liu et al. (2024)**: This method leverages the difference in the SLM’s logit distributions before and after the SFT process described above. This difference is directly added to the LLM’s logits to steer its output distribution.

We conduct experiments on LongLaMP using the Qwen7B-1.5B configuration. Implementation details and parameter settings for these baselines are provided in the Appendix B.8. As shown in Table 2, the results again confirm the superiority of the CoSteer framework. Our method adeptly combines the performance advantage of localized SLMs in handling private data with the powerful general capabilities of cloud-based LLMs, achieving the best performance among all compared methods without requiring any training.

	Abstract			Review			Writing		
	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
LA/CS	41.49	25.57	29.96	26.60	13.68	17.59	24.09	12.32	19.94
SFT	40.71	21.85	27.00	31.55	12.91	20.23	25.53	10.80	19.35
Proxy tuning	40.76	21.37	25.95	29.54	14.16	16.79	27.13	12.11	21.38
Costeer	42.98	<u>23.61</u>	<u>28.20</u>	32.72	15.92	20.36	<u>25.93</u>	12.38	22.84
w/o alpha	41.10	21.78	26.09	30.51	14.40	17.47	27.41	12.32	21.95
w/o KL	42.52	22.76	27.04	31.44	14.43	18.70	26.47	13.12	22.82

Table 2: Performance comparison of our Costeer framework against other methods on LongLamp benchmark, with an ablation study on its key components. Among four methods, best results are marked in **bold** while the second-best are underlined.

4.6 ABLATION STUDY AND PARAMETER ANALYSIS

Ablation study We conduct an ablation study to validate the contributions of the key components in our framework. The SLM delta steering signal (the β term in Eq. (4)) is the central mechanism of CoSteer. Our study therefore focuses on the impact of individually removing two other components: the LLM Alignment Term (the α term) and the KL Divergence Regularizer from the FTRL algorithm. The results in Table 2 proved their effectiveness. A broader ablation on the entire iterative FTRL (we name it LightCosteer) is discussed in Section 5.3.

Hyperparameters We examine the impact of different hyperparameters on our method. We conduct experiments on abstract generation using the Qwen7B-1.5B configuration, and evaluate with the ROUGE-L metric. When studying one parameter, we keep all others at their default values. Results and analysis presented in Appendix B.9 and Figure 2 show consistently strong performance across a broad range of values, confirming that the default hyperparameter settings (provided in Appendix B.1) are both robust and well-calibrated.

5 DISCUSSION: PRACTICAL DEPLOYMENT CONSIDERATIONS

In Section 4, we showed the empirical superiority of the CoSteer framework through comprehensive experiments and comparisons. In this section, we shift our focus to its practical viability, analyzing its feasibility and performance from various angles that simulate real-world conditions and constraints.

5.1 ROBUSTNESS TO NOISY USER CONTEXT

In practical on-device applications, the user-specific context is often imperfect and susceptible to noise. To evaluate **CoSteer**'s stability in such scenarios, we conducted rigorous stress tests under two distinct noise conditions: (1) **realistic noise**, simulated by replacing the strong `bge-m3` retriever with a weaker `BM25` retriever, and (2) **adversarial noise**, where we intentionally provided completely irrelevant context from a disjoint task.

Our experiments, detailed in Appendix B.10, yield a crucial insight: **CoSteer** demonstrates remarkable resilience to noisy context. Under realistic noise, it consistently outperforms baselines and maintains its performance edge (Table 10). In the more extreme adversarial setting, the framework degrades gracefully rather than collapsing, indicating that LLM guidance prevents the framework from being completely misled by faulty context (Table 11). This robustness is vital for real-world deployment.

5.2 GENERALIZATION ACROSS MODEL SCALES AND ARCHITECTURES

A key advantage of our framework is its ability to generalize across diverse model configurations, a critical feature for real-world deployment. We validate this generalization capability in two dimensions: model scale and architecture.

Generalization Across Model Scales Our primary experiments (see Table 6) already demonstrate remarkable scale generalization by pairing a powerful 32B-parameter LLM with a highly compact 0.6B SLM (Qwen3-0.6B), one of the smallest state-of-the-art models available. Despite being over 50× smaller, this SLM effectively steers the LLM across all our testbeds, strongly substantiates the framework's robustness to vast differences in model scale.

Collaboration Across Model Architectures While pairing models from the same family is common, a practical framework should accommodate scenarios where the LLM and SLM differ in architectures and tokenizers. To this end, we investigated two strategies. The first, `CoSteer_map`, aligns models via vocabulary intersection, while effective, its applicability is limited when tokenizers are highly divergent. Therefore, we propose `CoSteer_byte` which operates on a shared byte representation, making it vocabulary-agnostic (Hayase et al., 2025). As shown in Table 12, both strategies successfully enable the SLM to guide the LLM, confirming that **CoSteer is not confined to homogeneous model families**. This flexibility demonstrates the framework's generalization across architectures. Full implementation details are provided in Appendix B.11.

5.3 EFFICIENCY AND PRACTICAL VARIANTS

Finally, we analyze the efficiency and overhead of our framework. At a high level, a key advantage of CoSteer is its computational efficiency. Compared to methods that require concatenating long personal context to the LLM's input (i.e., `LLM w/`), our framework offers significant computational savings. By processing the private context locally with an SLM, **CoSteer substantially reduces the remote LLM's workload**. Our FLOPs estimation (Narayanan et al., 2021) in Table 15 shows that with a 1000-token context, the privacy-violating `LLM w/` approach is approximately **3.3x** more computationally expensive than CoSteer.

Although our closed-form solution is efficient, the full CoSteer framework still employs iterative optimization (e.g., $T = 20$), which can introduce latency. Furthermore, each generation step requires communication between the local device and the server. To address these practical overheads, we propose and evaluate two streamlined variants:

- **LightCoSteer**: To minimize computational overhead, this variant eliminates the iterative process by setting the number of optimization iterations $T = 1$. This simplifies the framework to a single-step logit adjustment, providing a lightweight yet effective alternative.
- **AdaCoSteer**: To reduce communication overhead, this variant adaptively deactivates steering. We observe that as generation proceeds, SLM and LLM predictions often converge. AdaCoSteer leverages this by terminating the steering process once the LLM’s token confidence exceeds a set threshold for a few consecutive steps (Song et al., 2025), allowing the LLM to complete the generation on its own.

The implementation details and a full analysis of the results are presented in Appendix B.12. Our findings show that these variants create a valuable performance-efficiency spectrum, allowing users to select the optimal configuration that balances effectiveness and resource constraints for diverse application requirements.

6 CONCLUSION

In this work, we introduced CoSteer, a collaborative framework that enables real-time, privacy-preserving personalization for LLMs via localized delta steering. Our approach leverages an on-device SLM to guide a cloud-based LLM, effectively incorporating user context without transmitting sensitive data or requiring fine-tuning. Extensive experiments demonstrate that CoSteer offers a practical and effective solution for harmonizing the powerful general capabilities of LLMs with user-specific context and privacy.

REFERENCES

- Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2024. URL <https://arxiv.org/abs/2402.03216>.
- Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. Pad: Personalized alignment of llms at decoding-time, 2025. URL <https://arxiv.org/abs/2410.04070>.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations, 2023. URL <https://arxiv.org/abs/2305.14233>.
- Aaron Grattafiori et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Songyang Gao, Qiming Ge, Wei Shen, Shihan Dou, Junjie Ye, Xiao Wang, Rui Zheng, Yicheng Zou, Zhi Chen, Hang Yan, Qi Zhang, and Dahua Lin. Linear alignment: A closed-form solution for aligning human preferences without tuning and feedback, 2024. URL <https://arxiv.org/abs/2401.11458>.
- Jonathan Hayase, Alisa Liu, Noah A. Smith, and Sewoong Oh. Sampling from your language model one byte at a time, 2025.
- Jerry Zhi-Yang He, Sashrika Pandey, Mariah L. Schrum, and Anca Dragan. Context steering: Controllable personalization at inference time, 2025. URL <https://arxiv.org/abs/2405.01768>.
- Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging, 2023. URL <https://arxiv.org/abs/2310.11564>.

- 540 Jiaming Ji, Boyuan Chen, Hantao Lou, Donghai Hong, Borong Zhang, Xuehai Pan, Juntao Dai,
541 Tianyi Qiu, and Yaodong Yang. Aligner: Efficient alignment by learning to correct, 2024. URL
542 <https://arxiv.org/abs/2402.02416>.
543
- 544 Minbeom Kim, Kang il Lee, Seongho Joo, and Hwaran Lee. Drift: Decoding-time personalized
545 alignments with implicit user preferences, 2025. URL [https://arxiv.org/abs/2502.](https://arxiv.org/abs/2502.14289)
546 14289.
- 547 Xiaoyu Kong, Jiancan Wu, An Zhang, Leheng Sheng, Hui Lin, Xiang Wang, and Xiangnan He.
548 Customizing language models with instance-wise lora for sequential recommendation, 2025. URL
549 <https://arxiv.org/abs/2408.10159>.
550
- 551 Ishita Kumar, Snigdha Viswanathan, Sushrita Yerra, Alireza Salemi, Ryan A. Rossi, Franck Der-
552 noncourt, Hanieh Deilamsalehy, Xiang Chen, Ruiyi Zhang, Shubham Agarwal, Nedim Lipka,
553 Chien Van Nguyen, Thien Huu Nguyen, and Hamed Zamani. Longlamp: A benchmark for person-
554 alized long-form text generation, 2024. URL <https://arxiv.org/abs/2407.11016>.
- 555 Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke
556 Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization,
557 2023. URL <https://arxiv.org/abs/2210.15097>.
558
- 559 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In Text summarization
560 branches out, pp. 74–81, 2004.
561
- 562 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
563 falsehoods, 2022. URL <https://arxiv.org/abs/2109.07958>.
- 564 Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A. Smith. Tuning
565 language models by proxy, 2024. URL <https://arxiv.org/abs/2401.08565>.
566
- 567 Runze Liu, Junqi Gao, Jian Zhao, Kaiyan Zhang, Xiu Li, Biqing Qi, Wanli Ouyang, and Bowen
568 Zhou. Can 1b llm surpass 405b llm? rethinking compute-optimal test-time scaling, 2025. URL
569 <https://arxiv.org/abs/2502.06703>.
- 570 Hanjia Lyu, Song Jiang, Hanqing Zeng, Yinglong Xia, Qifan Wang, Si Zhang, Ren Chen, Christopher
571 Leung, Jiajie Tang, and Jiebo Luo. Llm-rec: Personalized recommendation via prompting large
572 language models, 2023.
573
- 574 Brendan McMahan. Follow-the-regularized-leader and mirror descent: Equivalence theorems and
575 11 regularization. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), Proceedings
576 of the Fourteenth International Conference on Artificial Intelligence and Statistics, volume 15 of
577 Proceedings of Machine Learning Research, pp. 525–533, Fort Lauderdale, FL, USA, 11–13 Apr
578 2011. PMLR. URL <https://proceedings.mlr.press/v15/mcmahan11b.html>.
- 579 Sheshera Mysore, Zhuoran Lu, Mengting Wan, Longqi Yang, Bahareh Sarrafzadeh, Steve Menezes,
580 Tina Baghaee, Emmanuel Barajas Gonzalez, Jennifer Neville, and Tara Safavi. Pearl: Personalizing
581 large language model writing assistants with generation-calibrated retrievers, 2024. URL <https://arxiv.org/abs/2311.09180>.
582
583
- 584 Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vi-
585 jay Anand Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro,
586 Amar Phanishayee, and Matei Zaharia. Efficient large-scale language model training on gpu
587 clusters using megatron-lm, 2021. URL <https://arxiv.org/abs/2104.04473>.
- 588 Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying recommendations using distantly-
589 labeled reviews and fine-grained aspects. In Conference on Empirical Methods in Natural
590 Language Processing, 2019. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:202621357)
591 202621357.
592
- 593 Dan Peng, Zhihui Fu, and Jun Wang. Pocketllm: Enabling on-device fine-tuning for personalized
llms, 2024. URL <https://arxiv.org/abs/2407.01031>.

- 594 Alexandre Ramé, Guillaume Couairon, Mustafa Shukor, Corentin Dancette, Jean-Baptiste Gaya,
595 Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpo-
596 lating weights fine-tuned on diverse rewards, 2023. URL [https://arxiv.org/abs/2306.](https://arxiv.org/abs/2306.04488)
597 04488.
- 598 Chris Richardson, Yao Zhang, Kellen Gillespie, Sudipta Kar, Arshdeep Singh, Zeynab Raeesy,
599 Omar Zia Khan, and Abhinav Sethy. Integrating summarization and retrieval for enhanced
600 personalization via large language models, 2023. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.20081)
601 20081.
- 602 Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Hannaneh Hajishirzi, Noah A. Smith, and Simon S.
603 Du. Decoding-time language model alignment with multiple objectives, 2024. URL [https://arxiv.org/abs/2406.18853.](https://arxiv.org/abs/2406.18853)
- 604 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
605 can be more effective than scaling model parameters, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2408.03314)
606 2408.03314.
- 607 Feifan Song, Shaohang Wei, Wen Luo, Yuxuan Fan, Tianyu Liu, Guoyin Wang, and Houfeng Wang.
608 Well begun is half done: Low-resource preference alignment by weak-to-strong decoding. In
609 Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), Findings
610 of the Association for Computational Linguistics: ACL 2025, pp. 12654–12670, Vienna, Austria,
611 July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/
612 v1/2025.findings-acl.655. URL [https://aclanthology.org/2025.findings-acl.](https://aclanthology.org/2025.findings-acl.655/)
613 655/.
- 614 Zhaoxuan Tan, Zheyuan Liu, and Meng Jiang. Personalized pieces: Efficient personalized large
615 language models through collaborative efforts, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.10471)
616 10471.
- 617 Zhaoxuan Tan, Qingkai Zeng, Yijun Tian, Zheyuan Liu, Bing Yin, and Meng Jiang. Democratizing
618 large language models via personalized parameter-efficient fine-tuning, 2025. URL [https://arxiv.org/abs/2402.04401.](https://arxiv.org/abs/2402.04401)
- 619 Jie Tang, Jing Zhang, Limin Yao, Juan-Zi Li, Li Zhang, and Zhong Su. Arnetminer: extraction
620 and mining of academic social networks. In Knowledge Discovery and Data Mining, 2008. URL
621 [https://api.semanticscholar.org/CorpusID:3348552.](https://api.semanticscholar.org/CorpusID:3348552)
- 622 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.](https://qwenlm.github.io/blog/qwen2.5/)
623 [github.io/blog/qwen2.5/.](https://qwenlm.github.io/blog/qwen2.5/)
- 624 Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. Tldr: Mining reddit to learn auto-
625 matic summarization. In NFiS@EMNLP, 2017. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:2204603)
626 [org/CorpusID:2204603.](https://api.semanticscholar.org/CorpusID:2204603)
- 627 Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences via
628 multi-objective reward modeling and mixture-of-experts, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2406.12845)
629 [abs/2406.12845.](https://arxiv.org/abs/2406.12845)
- 630 Hongru Wang, Rui Wang, Fei Mi, Yang Deng, Zezhong Wang, Bin Liang, Ruifeng Xu, and Kam-Fai
631 Wong. Cue-cot: Chain-of-thought prompting for responding to in-depth dialogue questions with
632 llms, 2023a. URL [https://arxiv.org/abs/2305.11792.](https://arxiv.org/abs/2305.11792)
- 633 Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert,
634 Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. Helpsteer:
635 Multi-attribute helpfulness dataset for steerlm, 2023b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2311.09528)
636 2311.09528.
- 637 Stanisław Woźniak, Bartłomiej Koptyra, Arkadiusz Janz, Przemysław Kazienko, and Jan Kocoń.
638 Personalized large language models, 2024. URL [https://arxiv.org/abs/2402.09269.](https://arxiv.org/abs/2402.09269)
- 639 Bowen Wu, Wenqing Wang, Haoran Li, Ying Li, Jingsong Yu, and Baoxun Wang. Interpersonal
640 memory matters: A new task for proactive dialogue utilizing conversational history, 2025. URL
641 [https://arxiv.org/abs/2503.05150.](https://arxiv.org/abs/2503.05150)

- 648 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A. Smith,
649 Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for
650 language model training, 2023. URL <https://arxiv.org/abs/2306.01693>.
- 651 Yiyang Xu, Jinghao Zhang, Alireza Salemi, Xinting Hu, Wenjie Wang, Fuli Feng, Hamed Zamani,
652 Xiangnan He, and Tat-Seng Chua. Personalized generation in large model era: A survey, 2025.
653 URL <https://arxiv.org/abs/2503.02614>.
- 654 Biwei Yan, Kun Li, Minghui Xu, Yueyan Dong, Yue Zhang, Zhaochun Ren, and Xiuzhen Cheng.
655 On protecting the data privacy of large language models (llms) and llm agents: A literature
656 review. *High-Confidence Computing*, 5(2):100300, 2025. ISSN 2667-2952. doi: <https://doi.org/10.1016/j.hcc.2025.100300>. URL <https://www.sciencedirect.com/science/article/pii/S2667295225000042>.
- 660 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
661 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
662 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
663 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
664 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
665 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
666 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
667 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
668 Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- 669 Jizhi Zhang, Keqin Bao, Wenjie Wang, Yang Zhang, Wentao Shi, Wanhong Xu, Fuli Feng, and
670 Tat-Seng Chua. Prospect personalized recommendation on large language model-based agent
671 platform, 2024a. URL <https://arxiv.org/abs/2402.18240>.
- 672 Kai Zhang, Yejin Kim, and Xiaozhong Liu. Personalized llm response generation with parameterized
673 memory injection, 2025a. URL <https://arxiv.org/abs/2404.03565>.
- 674 Kaiyan Zhang, Jianyu Wang, Ermo Hua, Biqing Qi, Ning Ding, and Bowen Zhou. Cogensis: A
675 framework collaborating large and small language models for secure context-aware instruction
676 following, 2024b. URL <https://arxiv.org/abs/2403.03129>.
- 677 Zhaowei Zhang, Fengshuo Bai, Qizhi Chen, Chengdong Ma, Mingzhi Wang, Haoran Sun, Zilong
678 Zheng, and Yaodong Yang. Amulet: Realignment during test time for personalized preference
679 adaptation of LLMs. In *The Thirteenth International Conference on Learning Representations*
680 *(ICLR)*, 2025b. URL <https://openreview.net/forum?id=f9w890Y2cp>.
- 681 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
682 Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen,
683 Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and
684 Ji-Rong Wen. A survey of large language models, 2023.
- 685 Chujie Zheng, Ziqi Wang, Heng Ji, Minlie Huang, and Nanyun Peng. Model extrapolation expedites
686 alignment, 2025. URL <https://arxiv.org/abs/2404.16792>.
- 687 Wanjun Zhong, Duyu Tang, Jiahai Wang, Jian Yin, and Nan Duan. UserAdapter: Few-shot user
688 learning in sentiment analysis. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli
689 (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 1484–
690 1488, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
691 findings-acl.129. URL <https://aclanthology.org/2021.findings-acl.129/>.
- 692 Yifan Zhong, Chengdong Ma, Xiaoyuan Zhang, Ziran Yang, Haojun Chen, Qingfu Zhang, Siyuan Qi,
693 and Yaodong Yang. Panacea: Pareto alignment via preference adaptation for llms, 2024. URL
694 <https://arxiv.org/abs/2402.02030>.
- 695 Zhanhui Zhou, Jie Liu, Jing Shao, Xiangyu Yue, Chao Yang, Wanli Ouyang, and Yu Qiao. Beyond
696 one-preference-fits-all alignment: Multi-objective direct preference optimization, 2024. URL
697 <https://arxiv.org/abs/2310.03708>.
- 698 Jiachen Zhu, Jianghao Lin, Xinyi Dai, Bo Chen, Rong Shan, Jieming Zhu, Ruiming Tang, Yong
699 Yu, and Weinan Zhang. Lifelong personalized low-rank adaptation of large language models for
700 recommendation, 2024. URL <https://arxiv.org/abs/2408.03533>.
- 701

702	APPENDICES CONTENTS	
703		
704	A Methodology Detail	15
705		
706	A.1 Related Work and Distinction	15
707	A.2 Framework	16
708	A.3 Proof	16
709		
710		
711	B Experiment Detail	18
712		
713	B.1 Hyperparameters and Settings	18
714	B.2 Detailed Results	18
715	B.3 Statistical Significance Testing	19
716	B.4 Detailed Analysis on Experiment Results	20
717	B.5 Dataset Description	20
718	B.6 Human Evaluation on LongLaMP	20
719	B.7 Prompts for LLM-as-a-Judge	21
720	B.8 Baseline Implementation Details	21
721	B.9 Parameter Sensitivity Analysis	21
722	B.10 Robostness to Noise	23
723		
724	B.10.1 Robustness to Realistic Noise (Weaker Retriever)	23
725	B.10.2 Robustness to Purely Irrelevant (Adversarial) Noise	23
726		
727	B.11 Results and Implementation of Cross-Architecture Collaboration	24
728		
729	B.11.1 Vocabulary Mapping (CoSteer_map)	24
730	B.11.2 Byte-level Fusion (CoSteer_byte)	24
731		
732	B.12 Implementation and Analysis of CoSteer Variants	25
733		
734	B.12.1 LightCoSteer: Single-Step Steering	25
735	B.12.2 AdaCoSteer: Adaptive Steering	25
736	B.12.3 Performance	25
737	B.12.4 Efficiency and Overhead Analysis	26
738		
739		
740		
741	C Broader Impact and Future Work	27
742		
743	D LLM Usage	28
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

A METHODOLOGY DETAIL

A.1 RELATED WORK AND DISTINCTION

As summarized in Table 3, **CoSteer** occupies a unique technical position by being both collaborative and completely tuning-free. While it shares conceptual roots with inference-time steering methods, it introduces critical innovations tailored to the privacy-constrained cloud-edge scenario.

Distinction from Context Steering Methods Unlike methods like Context Steering (He et al., 2025) or Linear Alignment (Gao et al., 2024), which operate in a single-model setting where a model steers itself, CoSteer addresses a fundamentally different challenge: **heterogeneous cross-model collaboration**. We demonstrate a novel “weak-to-strong” capability where a tiny, local SLM effectively steers a massive, remote LLM. This allows high-quality personalization without exposing private data to the cloud, a constraint that single-model cloud deployment cannot satisfy.

Technical Novelty beyond Optimization While we adopt the FTRL algorithm as our solver, our core contribution lies in the **problem formulation**. We novelly define the utility function (Equation 4) based on local logit deltas to repurpose online learning for solving a distributed coordination problem. Furthermore, to address practical bottlenecks unique to this collaborative architecture, we introduce system-level innovations not found in prior work: **AdaCoSteer** minimizes communication latency via confidence-based termination, and **Byte-Level Fusion** resolves tokenizer mismatches to enable collaboration between heterogeneous model families (e.g., Llama and Qwen).

Table 3: Our Technical position. The main body of the table uses “Yes/No” for clarity. Explanations for specific “No” entries are as follows: ^aRequires a trained reward model. ^bRequires a fine-tuned smaller language model (SLM). ^cRequires a trained fusion network.

Method	Training Free?	Weak-to-Strong Collaborative?	No Additional Modules?
Linear Alignment (Gao et al., 2024)			
Contrastive Decoding (Li et al., 2023)	Yes	No	Yes
Context Steering (He et al., 2025)			
Amulet (Zhang et al., 2025b)			
PAD (Chen et al., 2025)	No ^a	No	No
Drift (Kim et al., 2025)			
Proxy-tuning (Liu et al., 2024)	No ^b	Yes	Yes
Cogenesis (Zhang et al., 2024b)	No ^c	Yes	No
CoSteer (Ours)	Yes	Yes	Yes

810 A.2 FRAMEWORK

811 Our CoSteer framework workflow is shown in Algorithm 1

814 **Algorithm 1** CoSteer Framework

815 **Require:**

816 Cloud-based LLM policy generator π_{LLM} , local SLM policy generator π_{SLM} , user query p_{base} ,
 817 personal context p_{pers} , current sequence s , max tokens M , hyperparameters: $T, \alpha, \beta, \lambda, \eta > 0$

818 **Ensure:** Personalized sequence s

- 819 1: Initialize $s \leftarrow \emptyset$
 - 820 2: **repeat**
 - 821 3: $\pi_{\text{base}} \leftarrow \pi_{\text{LLM}}(a|p_{\text{base}}, s)$ ▷ Cloud computation
 - 822 4: $\pi_{\text{base}}^* \leftarrow \pi_{\text{SLM}}(a|p_{\text{base}}, s)$ ▷ Edge computation
 - 823 5: $\pi_{\text{pers}}^* \leftarrow \pi_{\text{SLM}}(a|p_{\text{base}}, p_{\text{pers}}, s)$
 - 824 6: $\Delta \leftarrow \log \pi_{\text{pers}}^* - \log \pi_{\text{base}}^*$
 - 825 7: $\pi_0 \leftarrow \pi_{\text{base}}$ ▷ Initialize policy
 - 826 8: **for** $t = 1$ **to** T **do**
 - 827 9: $u_t \leftarrow \alpha(\log \pi_{t-1} - \log \pi_{\text{base}}) + \beta\Delta$
 - 828 10: Update policy using Equation 7
 - 829 11: **end for**
 - 830 12: $\pi^* \leftarrow \pi_T$ ▷ Final optimized policy
 - 831 13: Sample $a \sim \pi^*$, Update $s \leftarrow s \circ a$
 - 832 14: **until** $\text{len}(s) \geq M$ **or** EOS generated
 - 833 15: **return** s
-

834 A.3 PROOF

835 Similar to Zhang et al. (2025b), We try to solve the closed-form solution of Equation 6:

$$\begin{aligned}
 838 \mathcal{L}(\pi_t, \mu) &= \underbrace{\sum_{i=0}^{t-1} \sum_{a \in A} \pi_t(a) u_i(a)}_{(1)} - t\lambda \underbrace{\sum_{a \in A} \left(\pi_t(a) \log \frac{\pi_t(a)}{\pi_0(a)} \right)}_{(2)} \\
 842 &\quad - \underbrace{\frac{1}{\eta} \sum_{a \in A} \pi_t(a) \log \frac{\pi_t(a)}{\pi_{t-1}(a)}}_{(3)} + \underbrace{\mu \left(1 - \sum_{a \in A} \pi_i(a) \right)}_{(4)}
 \end{aligned} \tag{8}$$

847 Here, (1) and (2) originate from the utility function \mathcal{U} , (3) from the KL divergence, and (4) constrains
 848 the sum to 1, with the Lagrange multiplier μ . We calculate the derivation of the function for a given
 849 a , we have

$$850 \frac{\partial \mathcal{L}(\pi_t, \mu)}{\partial \pi_t(a)} = \sum_{i=0}^{t-1} u_i(a) - t\lambda \left(\log \frac{\pi_t(a)}{\pi_0(a)} + 1 \right) - \frac{1}{\eta} \left(\log \frac{\pi_t(a)}{\pi_{t-1}(a)} + 1 \right) - \mu \tag{9}$$

854 Rearrange the terms:

$$855 \sum_{i=0}^{t-1} u_i(a) - t\lambda \log \pi_t(a) + t\lambda \log \pi_0(a) - \frac{1}{\eta} \log \pi_t(a) + \frac{1}{\eta} \log \pi_{t-1}(a) - t\lambda - \frac{1}{\eta} - \mu = 0 \tag{10}$$

858 Combine the coefficients of $\log \pi_t(a)$:

$$860 - \left(t\lambda + \frac{1}{\eta} \right) \log \pi_t(a) = - \sum_{i=0}^{t-1} u_i(a) - t\lambda \log \pi_0(a) - \frac{1}{\eta} \log \pi_{t-1}(a) + t\lambda + \frac{1}{\eta} + \mu \tag{11}$$

863 Solve for $\log \pi_t(a)$:

$$\log \pi_t(a) = \frac{1}{t\lambda + \frac{1}{\eta}} \left(\sum_{i=0}^{t-1} u_i(a) + t\lambda \log \pi_0(a) + \frac{1}{\eta} \log \pi_{t-1}(a) - t\lambda - \frac{1}{\eta} - \mu \right) \quad (12)$$

let

$$C = -\frac{t\lambda + \frac{1}{\eta} + \mu}{t\lambda + \frac{1}{\eta}} \quad (13)$$

we have

$$\log \pi_t(a) = \frac{1}{t\lambda + \frac{1}{\eta}} \left(\sum_{i=0}^{t-1} u_i(a) + t\lambda \log \pi_0(a) + \frac{1}{\eta} \log \pi_{t-1}(a) \right) + C \quad (14)$$

Thus

$$\pi_t(a) \propto \exp \left(\frac{1}{t\lambda + \frac{1}{\eta}} \left(\sum_{i=0}^{t-1} u_i(a) + t\lambda \log \pi_0(a) + \frac{1}{\eta} \log \pi_{t-1}(a) \right) \right) \quad (15)$$

Specifically, let $T = 1$ and $\beta^* = \frac{\beta}{\lambda + \frac{1}{\eta}}$, the policy is simplified to

$$\pi(a) \propto \exp \left(\log \pi_0(a) + \beta^* \left(\log \pi_{pers}^*(a) - \log \pi_{base}^*(a) \right) \right) \quad (16)$$

B EXPERIMENT DETAIL

B.1 HYPERPARAMETERS AND SETTINGS

To ensure the results are easily reproducible, we set `do_sample=False` for all methods. For the Qwen3 series, we set `enable_thinking=False`. Hyperparameters of [CoSteer](#) are listed in Table 4

Table 4: Default hyperparameter of [CoSteer](#)

Parameter	Value
T	20
α	2.0
β	1.0
λ	2
η	10

B.2 DETAILED RESULTS

Detailed results on HelpSteer, Personal.Eval, Truthful QA and Ultrachat are shown in Table 5.

Models Setting	HelpSteer				Personal.Eval				Truthful QA				Ultrachat				
	creative	verbose	concise	uplifting	creative	verbose	concise	uplifting	creative	verbose	concise	uplifting	creative	verbose	concise	uplifting	
Qwen 7B-1.5B	SLM w/o	.5609	.5301	.6437	.5638	.6043	.5882	.6556	.6722	.4468	.4231	.5349	.4544	.5642	.5202	.6074	.5625
	SLM w/	.6309	.7093	.7110	.6568	.6564	.7402	.7686	.7449	.5419	.6515	.6358	.5758	.6562	.6982	.7252	.6602
	LLM w/o	.7368	.7223	.7228	.7338	.7550	.7462	.7005	.8124	.5739	.5752	.6744	.5782	.7576	.7285	.7179	.7489
	CoSteer	.8093	.7787	.7731	.8055	.8472	.7945	.7525	.8739	.6834	.6930	.6837	.6682	.8259	.7770	.7522	.8064
	LLM w/	.8918	.8173	.7853	.8516	.9217	.8396	.8001	.9073	.8363	.7593	.7069	.7630	.9037	.8140	.7850	.8438
Qwen 32B-7B	SLM w/o	.7368	.7223	.7228	.7338	.7550	.7462	.7005	.8124	.5739	.5752	.6744	.5782	.7576	.7285	.7179	.7489
	SLM w/	.8918	.8173	.7853	.8516	.9217	.8396	.8001	.9073	.8363	.7593	.7069	.7630	.9037	.8140	.7850	.8438
	LLM w/o	.7141	.7003	.7406	.7129	.7548	.7461	.7129	.8136	.5877	.5844	.6993	.5925	.7512	.7183	.7373	.7393
	CoSteer	.8630	.8516	.7478	.8684	.9158	.8754	.7700	.9166	.7985	.8425	.6575	.7818	.8581	.8433	.7344	.8646
	LLM w/	.9038	.8270	.8104	.8628	.9254	.8434	.8215	.9117	.8698	.7895	.7208	.7897	.9078	.8171	.8122	.8510
Llama 8B-1B	SLM w/o	.6779	.6595	.6063	.6764	.7297	.7159	.6378	.7915	.5215	.5084	.5119	.5064	.6850	.6581	.6039	.6842
	SLM w/	.8070	.7915	.7161	.8083	.8604	.8529	.7805	.8866	.7115	.7291	.6023	.7263	.8134	.7897	.7159	.8049
	LLM w/o	.7227	.7110	.6895	.7169	.7476	.7355	.6689	.8021	.6044	.5909	.6574	.5967	.7405	.7137	.6903	.7306
	CoSteer	.8968	.8571	.7993	.8735	.9173	.8788	.8297	.9235	.8499	.8408	.6951	.8266	.9004	.8562	.8008	.8648
	LLM w/	.8951	.8375	.8121	.8603	.9030	.8556	.8108	.9004	.8634	.7974	.7410	.8101	.8860	.8285	.7997	.8496
Qwen 0.6B-8B	SLM w/o	.5844	.5477	.5984	.5945	.6498	.6291	.6573	.7194	.4014	.3648	.5039	.4144	.6485	.6018	.6327	.6463
	SLM w/	.6255	.6428	.6099	.6313	.7299	.7302	.6614	.7638	.4928	.5277	.4942	.4968	.7063	.6721	.6343	.6826
	LLM w/o	.7819	.7783	.7084	.7771	.8057	.8037	.6920	.8448	.6536	.6504	.6675	.6461	.7974	.7796	.7038	.7765
	CoSteer	.7849	.7806	.7046	.7757	.8058	.8117	.6901	.8477	.6582	.6549	.6647	.6506	.7986	.7819	.7058	.7853
	LLM w/	.8928	.8037	.7842	.8629	.9076	.8303	.8141	.9020	.8612	.7796	.6952	.8004	.8964	.7989	.7915	.8555
Qwen 0.6B-32B	SLM w/o	.5844	.5477	.5984	.5945	.6498	.6291	.6573	.7194	.4014	.3648	.5039	.4144	.6485	.6018	.6327	.6463
	SLM w/	.6255	.6428	.6099	.6313	.7299	.7302	.6614	.7638	.4928	.5277	.4942	.4968	.7063	.6721	.6343	.6826
	LLM w/o	.8103	.8084	.6991	.7922	.8179	.8237	.6948	.8443	.7084	.7130	.6863	.7010	.8077	.7948	.6935	.7758
	CoSteer	.8083	.8152	.6984	.7923	.8145	.8191	.6752	.8436	.7152	.7226	.6870	.7020	.8053	.7962	.6884	.7779
	LLM w/	.9121	.8147	.8294	.8722	.9094	.8352	.8298	.8936	.9097	.7886	.7364	.8327	.9069	.7996	.8204	.8629

Table 5: Results on all four preference alignment datasets with our proposed [CoSteer](#) and other settings. **Bold** values indicate that our method outperformed the three baseline methods we are comparing against. Note that results from the large model with context are marked in gray, as this scenario does not align with our task setting.

Models Setting	Cogenesis		Abstract Generation			Review Writing			Topic Writing			Pref Align				
	Ovl	Per	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET	creative	verbose	concise	uplifting	
Qwen 8B-0.6B	SLM w/o	6.84	6.64	37.13	20.41	21.00	24.41	12.83	12.39	24.19	11.88	15.74	.5710	.5359	.5981	.5937
	SLM w/	7.85	7.79	42.48	23.25	26.58	25.04	13.36	13.02	26.17	13.46	18.40	.6386	.6432	.6000	.6436
	LLM w/o	8.27	7.99	40.76	21.23	25.77	30.89	14.47	16.98	28.47	12.85	21.49	.7597	.7530	.6929	.7611
	CoSteer	8.62	8.65	41.11	22.36	25.86	32.24	14.84	18.47	29.03	13.61	22.98	.7619	.7573	.6913	.7648
	LLM w/	8.96	8.96	43.99	23.69	29.48	35.42	15.71	21.59	31.48	15.34	24.78	.8895	.8031	.7713	.8552
Qwen 32B-0.6B	SLM w/o	6.84	6.64	37.13	20.41	21.00	24.41	12.83	12.39	24.19	11.88	15.74	.5710	.5359	.5981	.5937
	SLM w/	7.85	7.79	42.48	23.25	26.58	25.04	13.36	13.02	26.17	13.46	18.40	.6386	.6432	.6000	.6436
	LLM w/o	8.31	8.22	41.05	21.88	26.23	30.23	14.08	16.41	29.01	12.51	21.44	.7861	.7850	.6934	.7783
	CoSteer	8.52	8.56	43.74	23.69	30.97	30.89	14.14	18.87	26.75	13.17	21.32	.7858	.7883	.6873	.7790
	LLM w/	9.09	9.03	44.46	23.95	30.15	35.12	15.60	20.65	32.63	15.37	25.07	.9095	.8095	.8040	.8654

Table 6: Comparative performance across eight personalized content generation and preference alignment tasks. Metrics include overall (Ovl) and personalized (Per) scores for Cogenesis, ROUGE-1/-L (R-1/-L) and METEOR (MET) for LongLaMP datasets, and averaged alignment scores for four user preferences. **Bold** entries indicate that CoSteer outperforms the three baseline methods compared against. Gray values represent the privacy-violating near-upper-bound performance, yet underlined CoSteer values surpasses these incompatible references.

Models Setting	Cogenesis		Abstract Generation			Review Writing			Topic Writing			Pref Align				
	Ovl	Per	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET	creative	verbose	concise	uplifting	
Qwen 7B-1.5B	SLM w/o	6.63 _{±1.70}	6.21 _{±1.88}	36.48 _{±7.50}	17.74 _{±3.85}	27.57 _{±5.44}	20.40 _{±2.88}	10.39 _{±3.71}	10.39 _{±2.07}	25.21 _{±9.55}	11.09 _{±3.24}	17.51 _{±2.08}	.5441 _{±.13}	.5154 _{±.11}	.6104 _{±.09}	.5632 _{±.06}
	SLM w/	7.81 _{±2.41}	7.63 _{±1.43}	39.75 _{±3.83}	22.03 _{±3.83}	27.76 _{±10.03}	23.08 _{±2.95}	12.40 _{±7.15}	12.94 _{±2.90}	22.89 _{±9.19}	11.46 _{±6.30}	17.12 _{±2.45}	.6214 _{±.10}	.6998 _{±.08}	.7102 _{±.13}	.6594 _{±.05}
	LLM w/o	8.00 _{±1.18}	7.63 _{±1.43}	39.81 _{±2.86}	20.53 _{±2.32}	25.56 _{±2.11}	30.15 _{±2.64}	14.04 _{±2.86}	17.71 _{±2.54}	27.64 _{±2.70}	11.93 _{±2.48}	21.49 _{±2.51}	.7058 _{±.10}	.6931 _{±.08}	.7039 _{±.13}	.7183 _{±.07}
	CoSteer	8.44[*] _{±1.31}	8.50[*] _{±1.20}	42.98[*] _{±3.04}	23.61[*] _{±2.09}	28.20[*] _{±2.88}	32.72[*] _{±9.07}	15.92[*] _{±9.02}	20.36[*] _{±10.09}	25.93 _{±10.13}	12.38[*] _{±2.78}	22.84[*] _{±18.41}	.7915[*] _{±.06}	.7608[*] _{±.07}	.7404[*] _{±.08}	.7885[*] _{±.14}
	LLM w/	8.62 _{±2.03}	8.61 _{±2.08}	44.50 _{±2.77}	24.63 _{±2.14}	31.15 _{±2.34}	33.83 _{±2.44}	15.55 _{±2.11}	22.42 _{±2.70}	28.82 _{±11.70}	13.77 _{±2.37}	23.44 _{±2.29}	.8884 _{±.07}	.8076 _{±.05}	.7693 _{±.07}	.8414 _{±.10}
Qwen 32B-7B	SLM w/o	8.00 _{±1.18}	7.63 _{±1.43}	39.81 _{±2.86}	20.53 _{±2.32}	25.56 _{±2.11}	30.15 _{±2.64}	14.04 _{±2.86}	17.71 _{±2.54}	27.64 _{±2.70}	11.93 _{±2.48}	21.49 _{±2.51}	.7058 _{±.10}	.6931 _{±.08}	.7039 _{±.08}	.7183 _{±.15}
	SLM w/	8.62 _{±2.03}	8.61 _{±2.08}	44.50 _{±2.77}	24.63 _{±2.14}	31.15 _{±2.34}	33.83 _{±2.44}	15.55 _{±2.11}	22.42 _{±2.70}	28.82 _{±11.70}	13.77 _{±2.37}	23.44 _{±2.29}	.8884 _{±.05}	.8076 _{±.06}	.7693 _{±.12}	.8414 _{±.05}
	LLM w/o	8.12 _{±1.01}	7.87 _{±1.19}	40.66 _{±2.79}	21.02 _{±1.26}	26.61 _{±2.83}	32.21 _{±2.32}	14.44 _{±2.51}	19.61 _{±2.44}	28.82 _{±8.53}	12.20 _{±2.68}	21.16 _{±2.11}	.7020 _{±.14}	.6873 _{±.06}	.7225 _{±.10}	.7146 _{±.13}
	CoSteer	8.78[*] _{±2.01}	8.64[*] _{±2.09}	45.41[*] _{±2.79}	26.04[*] _{±2.90}	33.52[*] _{±2.64}	34.88[*] _{±2.80}	15.89[*] _{±2.72}	26.51[*] _{±7.23}	30.10[*] _{±12.05}	14.52[*] _{±10.70}	24.20[*] _{±11.26}	.8589[*] _{±.11}	.8532[*] _{±.08}	.7874[*] _{±.06}	.8579[*] _{±.04}
	LLM w/	8.83 _{±2.48}	8.76 _{±2.51}	43.33 _{±2.16}	23.47 _{±2.40}	30.10 _{±7.30}	34.65 _{±3.03}	15.74 _{±3.83}	22.77 _{±3.48}	30.73 _{±10.83}	14.20 _{±2.45}	24.25 _{±2.67}	.9017 _{±.20}	.8193 _{±.07}	.7912 _{±.12}	.8538 _{±.06}
Llama 8B-1B	SLM w/o	7.04 _{±1.05}	6.55 _{±2.81}	33.20 _{±2.60}	18.20 _{±3.78}	28.55 _{±2.84}	31.75 _{±2.95}	14.92 _{±2.37}	19.78 _{±2.53}	20.81 _{±10.10}	10.21 _{±2.22}	17.30 _{±2.29}	.6535 _{±.12}	.6355 _{±.11}	.5900 _{±.10}	.6646 _{±.12}
	SLM w/	7.69 _{±1.40}	7.52 _{±1.39}	39.81 _{±2.83}	21.53 _{±2.95}	30.11 _{±2.45}	32.36 _{±2.59}	15.02 _{±2.70}	22.06 _{±2.80}	20.17 _{±10.43}	10.58 _{±2.99}	18.64 _{±2.81}	.7981 _{±.11}	.7908 _{±.08}	.7037 _{±.13}	.8065 _{±.10}
	LLM w/o	7.69 _{±1.36}	7.13 _{±1.62}	39.33 _{±2.11}	20.69 _{±2.90}	29.41 _{±2.32}	34.58 _{±2.19}	15.32 _{±2.42}	22.10 _{±2.54}	26.93 _{±10.11}	12.35 _{±2.40}	21.82 _{±2.47}	.7038 _{±.09}	.6878 _{±.06}	.6765 _{±.09}	.7116 _{±.09}
	CoSteer	7.29 _{±1.42}	7.73[*] _{±1.03}	41.28[*] _{±2.36}	24.97[*] _{±2.10}	31.19[*] _{±2.42}	31.68 _{±2.33}	13.68 _{±2.63}	24.57[*] _{±2.93}	26.11 _{±11.18}	12.00 _{±2.74}	23.69[*] _{±2.13}	.8911[*] _{±.06}	.8582[*] _{±.06}	.7812[*] _{±.12}	.8721[*] _{±.10}
	LLM w/	8.61 _{±2.06}	8.44 _{±2.00}	43.91 _{±2.60}	23.93 _{±2.44}	32.01 _{±7.07}	36.39 _{±3.19}	15.95 _{±2.20}	23.56 _{±2.80}	30.54 _{±10.43}	14.02 _{±2.15}	23.81 _{±7.39}	.8869 _{±.06}	.8298 _{±.06}	.7909 _{±.07}	.8551 _{±.06}
Qwen 8B-0.6B	SLM w/o	6.84 _{±1.70}	6.64 _{±1.80}	37.13 _{±7.17}	20.41 _{±2.90}	21.00 _{±2.91}	24.41 _{±2.88}	12.83 _{±2.68}	12.39 _{±2.08}	24.19 _{±2.46}	11.88 _{±2.88}	15.74 _{±2.38}	.5710 _{±.13}	.5359 _{±.04}	.5981 _{±.10}	.5937 _{±.12}
	SLM w/	7.85 _{±1.40}	7.79 _{±1.42}	42.48 _{±2.39}	23.25 _{±2.21}	26.58 _{±2.29}	25.04 _{±10.20}	13.36 _{±2.38}	13.02 _{±2.71}	26.17 _{±11.92}	13.46 _{±2.01}	18.40 _{±10.13}	.6386 _{±.12}	.6432 _{±.14}	.6000 _{±.10}	.6436 _{±.12}
	LLM w/o	8.27 _{±1.06}	7.99 _{±.95}	40.76 _{±.96}	21.23 _{±.97}	25.77 _{±.75}	30.89 _{±.60}	14.47 _{±.28}	16.98 _{±.70}	28.47 _{±.36}	12.85 _{±.04}	21.49 _{±.26}	.7597 _{±.15}	.7530 _{±.14}	.6929 _{±.15}	.7611 _{±.07}
	CoSteer	8.62[*] _{±2.04}	8.65[*] _{±1.12}	41.11 _{±2.81}	22.36 _{±2.21}	25.86 _{±2.08}	32.24[*] _{±7.22}	14.84[*] _{±3.02}	18.47[*] _{±2.15}	29.03[*] _{±11.07}	13.61[*] _{±7.06}	22.98[*] _{±2.47}	.7619[*] _{±.06}	.7573[*] _{±.05}	.6913 _{±.05}	.7648[*] _{±.11}
	LLM w/	8.96 _{±2.06}	8.96 _{±2.06}	43.99 _{±2.46}	23.69 _{±2.42}	29.48 _{±1.76}	35.42 _{±3.45}	15.71 _{±2.55}	21.59 _{±2.16}	31.48 _{±13.00}	15.34 _{±1.30}	24.78 _{±1.46}	.8895 _{±.06}	.8031 _{±.12}	.7713 _{±.05}	.8552 _{±.06}
Qwen 32B-0.6B	SLM w/o	6.84 _{±1.70}	6.64 _{±1.80}	37.13 _{±7.17}	20.41 _{±2.90}	21.00 _{±2.91}	24.41 _{±2.88}	12.83 _{±2.68}	12.39 _{±2.08}	24.19 _{±2.46}	11.88 _{±2.88}	15.74 _{±2.38}	.5710 _{±.13}	.5359 _{±.04}	.5981 _{±.10}	.5937 _{±.06}
	SLM w/	7.85 _{±1.40}	7.79 _{±1.42}	42.48 _{±2.39}	23.25 _{±2.21}	26.58 _{±2.29}	25.04 _{±10.20}	13.36 _{±2.38}	13.02 _{±2.71}	26.17 _{±11.92}	13.46 _{±2.01}	18.40 _{±10.13}	.6386 _{±.12}	.6432 _{±.14}	.6000 _{±.10}	.6436 _{±.12}
	LLM w/o	8.31 _{±2.02}	8.22 _{±1.12}	41.05 _{±2.05}	21.88 _{±1.54}	26.23 _{±2.01}	30.23 _{±2.94}	14.08 _{±2.10}	16.41 _{±2.75}	29.01 _{±10.00}	12.51 _{±3.00}	21.44 _{±2.23}	.7861 _{±.10}	.7850 _{±.08}	.6934 _{±.10}	.7783 _{±.02}
	CoSteer	8.52[*] _{±2.04}	8.56[*] _{±2.87}	43.74[*] _{±2.51}	23.69[*] _{±2.47}	30.97[*] _{±2.74}	30.89[*] _{±2.37}	14.14[*] _{±2.80}	18.87[*] _{±7.62}	26.75 _{±12.75}	13.17 _{±10.14}	21.32 _{±2.47}	.7858 _{±.11}	.7883[*] _{±.05}	.6873 _{±.14}	.7790[*] _{±.07}
	LLM w/	9.09 _{±2.41}	9.03 _{±2.20}	44.46 _{±2.08}	23.95 _{±2.47}	30.15 _{±2.08}	35.12 _{±2.20}	15.60 _{±2.70}	20.65 _{±2.45}	32.63 _{±12.44}	15.37 _{±10.48}	25.07 _{±10.63}	.9095 _{±.05}	.8095 _{±.08}	.8040 _{±.13}	.8654 _{±.06}

Table 7: Main Results with Mean and Standard Deviation. **Bold** entries indicate that CoSteer outperforms the three baselines (SLM w/o, SLM w/, LLM w/o). An asterisk (*) denotes that this improvement is statistically significant (paired t-test, $p < 0.05$). Gray values represent the privacy-violating near-upper-bound performance, yet underlined CoSteer values surpass these incompatible references.

B.3 STATISTICAL SIGNIFICANCE TESTING

To formally validate the robustness of our findings, we performed paired t-tests comparing CoSteer against each of the three key baselines: LLM w/o, SLM w/, and SLM w/o. The tests were conducted on the full set of evaluation samples, comparing the performance scores on a per-sample basis to determine if the observed improvements were due to more than random chance. We used a standard significance level of $\alpha = 0.05$. The results confirmed that the performance gains of CoSteer over all three baselines are statistically significant. In Table 7, these significant improvements are denoted with an asterisk (*).

B.4 DETAILED ANALYSIS ON EXPERIMENT RESULTS

Analysis on Task Complexity Our framework demonstrates distinct performance patterns across task complexity levels. For context-intensive generation tasks (Cogenesis and LongLaMP) requiring deep integration of extended user profiles, histories, and multiple writing examples, model pairs with larger size demonstrate superior contextual reasoning capabilities. Qwen 32B-7B shows significant gains across three datasets compared to compact models (Qwen 7B-1.5B/ Llama 8B-1B). In contrast, the Llama8B-1B configuration occasionally underperforms baselines by failing to distill essential personalization signals from sparse contexts. Conversely, in preference alignment tasks where personal context reduces to concise and explicit textual instructions, compact models perform comparably to their larger counterparts. This stems from the larger models’ over-alignment when processing simplistic personalization signals and effectively overfitting to sparse preference indicators.

Analysis on Model Sizes Our approach has demonstrated excellent collaborative results across models of various sizes. An interesting finding is that the benefits of integrating personal information vary with different model sizes. In the Abstract Generation and Review Writing datasets, the performance gain of incorporating personal context with the Qwen-2.5-7B-Instruct model was significantly greater than that with the 32B model. In such cases, by using our CoSteer method to influence the token distribution of the 32B model through the logits delta derived from the 7B model, the final performance significantly outperforms 32B model with full context, yielding a relative improvement of nearly 14 percentage points in METEOR scores across these datasets.

B.5 DATASET DESCRIPTION

To ensure a comprehensive evaluation relevant to Cloud-LLM serving scenarios, we selected datasets covering a broad spectrum of personal AI agent applications:

- **Daily Assistance (Cogenesis):** This benchmark simulates daily tasks such as drafting emails and notifications based on user profiles and history.
- **Complex Content Creation (LongLaMP):** This focuses on reasoning-intensive, long-form generation tasks including academic abstracts, product reviews, and blog posts. These tasks typically require the superior generation capabilities of Cloud LLMs to ensure high quality.
- **General Preference Alignment:** We use four datasets (e.g., HelpSteer, TruthfulQA) to evaluate general instruction following based on specific user constraints.

We present a real example from each dataset to help readers understand the task of each dataset. Cogenesis is shown in Figure 4. Three datasets of LongLaMP are in Figure 5. Four datasets of the preference alignment task are in Figure 6.

B.6 HUMAN EVALUATION ON LONGLAMP

To provide a more robust validation of our method, we conducted a small-scale human evaluation. Due to time constraints, we randomly sampled 10 articles from each of the three LongLaMP datasets. We then invited five PhD candidates with strong NLP backgrounds to act as evaluators. Following the evaluation criteria from Cogenesis, they rated the generated texts from all models on two aspects: **Overall Quality (Ovl)** and **Personalization (Per)**, using a 1–5 scale. The average scores are presented in Table 8.

Table 8: Human evaluation results on a 1-5 scale (higher is better). Scores are reported as **Overall Quality (Ovl) / Personalization (Per)**. Our method, CoSteer, achieves personalization scores comparable to the LLM w/ oracle while maintaining high overall quality.

Method	Abstract (Ovl / Per)	Review (Ovl / Per)	Writing (Ovl / Per)
SLM w/o	3.85 / 3.70	3.45 / 3.25	3.20 / 3.05
SLM w/	3.95 / 3.85	3.70 / 3.75	3.35 / 3.50
LLM w/o	4.30 / 3.80	4.15 / 3.40	3.95 / 3.55
CoSteer	4.30 / 3.90	3.90 / 3.80	4.10 / 4.15
LLM w/	4.45 / 4.00	4.30 / 4.05	4.15 / 4.25

The results from our human evaluation corroborate the findings from our automated metrics. In all three tasks, **CoSteer** significantly improves personalization scores over the LLM w/o baseline and

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

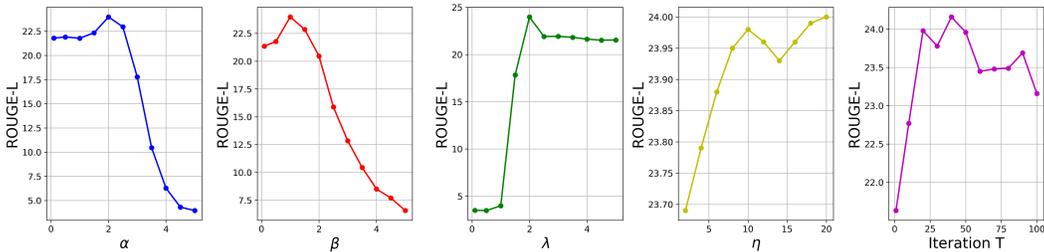


Figure 2: Effect of different α , β , η , λ and iteration [step T](#) on the Abstract Generation dataset using Qwen 7B-1.5B. The evaluation metric is ROUGE-L.

achieves an overall quality that is highly competitive with the full LLM w/ oracle. This demonstrates the effectiveness of our method in generating high-quality, personalized text that aligns with human judgment.

B.7 PROMPTS FOR LLM-AS-A-JUDGE

We use exactly the same prompt as Cogensis (Zhang et al., 2024b) to have the `gpt-4o` evaluate the generated content. The specific prompts are shown in Figure 7 and Figure 8.

B.8 BASELINE IMPLEMENTATION DETAILS

We compare our proposed method against three representative baseline approaches:

CoS/LA Context Steering (CoS) (He et al., 2025) and Linear Alignment (LA) (Gao et al., 2024) are conceptually equivalent methods. They steer the generation of a base model by modifying its output log-probabilities. The final log-probability for a token x_i is calculated by adding a scaled difference term to the base model’s prediction, which is conditioned on no context (\emptyset). The formulation is as follows:

$$\log P'(x_i) = \log P(x_i|\emptyset, \mathcal{P}) + \lambda \cdot (\log P(x_i|C, \mathcal{P}) - \log P(x_i|\emptyset, \mathcal{P}))$$

where \mathcal{P} represents the base model parameters, C is the provided context, and λ is a scaling hyperparameter. Following the implementations in the original works, we searched for the optimal λ from the set $\{1.5, 2.0, 2.5\}$ and report the best-performing result for each task.

Supervised Fine-Tuning (SFT) To create a strong task-specific baseline, we fine-tuned the small language model using Low-Rank Adaptation (LoRA). For each of our three datasets, we randomly sampled 1,000 data points for training. The model was then trained for three epochs. Key training parameters included a learning rate of 1×10^{-4} with a cosine scheduler, a maximum sequence length of 1024, a weight decay of 0.05, and `bf16` precision for efficiency. We utilized packing to handle variable-length sequences effectively.

Proxy-Tuning This method (Liu et al., 2024) enhances a large base model by incorporating signals from a smaller, specialized proxy model. The final probability distribution is calculated by adjusting the logits of the large model. We adapt the original formula to better reflect the roles of the different models in our setup:

$$p(x_t|x_{<t}) = \text{softmax}[s_{\text{LLM}}(x_t|x_{<t}) + s_{\text{SLM-SFT}}(x_t|x_{<t}) - s_{\text{SLM-base}}(x_t|x_{<t})]$$

Here, s_{LLM} represents the logits from the large, pre-trained language model. The term $s_{\text{SLM-SFT}}$ refers to the logits from the small language model that has been fine-tuned (the SFT baseline), and $s_{\text{SLM-base}}$ refers to the logits from the original, pre-trained small language model.

The comparison results are shown in Table 9.

B.9 PARAMETER SENSITIVITY ANALYSIS

Results are shown in Figure 2. Below, we analyze each parameter.

Models	Method	Abstract			Review			Writing		
		R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
Qwen 7B-1.5B	LA/CS	<u>41.49</u>	25.57	29.96	26.60	13.68	17.59	24.09	<u>12.32</u>	19.94
	SFT	40.71	21.85	27.00	<u>31.55</u>	12.91	<u>20.23</u>	25.53	10.80	19.35
	Proxy tuning	40.76	21.37	25.95	29.54	<u>14.16</u>	16.79	27.13	12.11	<u>21.38</u>
	Costeer	42.98	<u>23.61</u>	<u>28.20</u>	32.72	15.92	20.36	<u>25.93</u>	12.38	22.84
Qwen 32B-7B	LA/CS	44.28	<u>25.85</u>	<u>32.98</u>	<u>34.14</u>	<u>15.46</u>	<u>26.38</u>	<u>29.22</u>	<u>14.03</u>	<u>23.66</u>
	SFT	43.59	25.48	32.43	31.94	11.84	21.00	28.83	10.63	21.73
	Proxy tuning	<u>44.61</u>	25.12	31.97	32.03	14.80	19.01	28.75	12.70	21.49
	Costeer	45.41	26.04	33.52	34.88	15.89	26.51	30.10	14.52	24.20
Llama 8B-1B	LA/CS	37.24	20.70	<u>30.78</u>	27.35	11.93	23.54	19.28	9.30	18.26
	SFT	<u>40.62</u>	26.67	25.64	29.93	10.17	21.51	26.11	8.99	<u>22.56</u>
	Proxy tuning	39.55	20.78	29.60	<u>30.50</u>	<u>12.29</u>	<u>21.81</u>	26.00	12.29	21.65
	Costeer	41.28	<u>24.97</u>	31.19	31.68	13.68	24.57	26.11	<u>12.00</u>	23.69
Qwen 8B-0.6B	LA/CS	31.42	16.78	22.10	28.37	12.92	15.94	25.31	11.92	18.96
	SFT	38.38	<u>21.88</u>	26.03	30.53	<u>14.75</u>	19.91	25.20	12.82	18.83
	Proxy tuning	<u>40.60</u>	21.14	25.68	<u>30.97</u>	14.48	16.94	<u>28.84</u>	<u>12.89</u>	<u>21.64</u>
	Costeer	41.11	22.36	<u>25.86</u>	32.24	14.84	<u>18.47</u>	29.03	13.61	22.98
Qwen 32B-0.6B	LA/CS	31.42	16.78	22.10	28.37	12.92	15.94	25.31	11.92	18.96
	SFT	38.38	<u>21.88</u>	26.03	<u>30.53</u>	14.75	<u>18.91</u>	25.20	<u>12.82</u>	18.83
	Proxy tuning	<u>41.17</u>	<u>21.88</u>	<u>26.40</u>	30.20	<u>14.14</u>	16.27	<u>26.08</u>	12.61	21.67
	Costeer	43.74	23.69	30.97	30.89	<u>14.14</u>	18.87	26.75	13.17	<u>21.32</u>

Table 9: Performance comparison on LongLamp benchmark across 5 model pairs. Best results are marked in bold and second-best are underlined.

- Parameter α and β :** As defined in Equation 4, α and β adjust the influence of each optimization component. α controls the impact of the difference between the current policy and the initial policy, while β scales the delta signals. We conducted experiments with values ranging from 0.1 to 5.0. As shown in the first two tables of the figure, small α and β coefficients below 1.0 attenuate personalization signals, whereas values exceeding 3.0 compromise the model’s fundamental capabilities through over-alignment, leading to a noticeable decline in results. Therefore we adopt $\alpha = 2$ and $\beta = 1$ for performance.
- Parameter λ and η :** These two learning dynamics control the learning rate and the degree of deviation from the initial policy. For the learning rate η , we conducted experiments with values ranging from 2 to 20. For η , we observe stable convergence from 2 to 10, followed by oscillatory behavior beyond 10. Parameter λ rapid performance gains from 0 to 2, and then transitions to a plateau phase between 2 and 5. Therefore, we set η to 10 and λ to 2.
- Iteration step T :** We studied the impact of T within the range of 1 to 100. As shown in the results of the last table, there is a clear upward trend in performance from 1 to 20 iterations, after which the performance stabilizes. We thus establish $T = 20$ as the Pareto-optimal configuration balancing quality and latency.

B.10 ROBUSTNESS TO NOISE

To supplement the discussion in Section 5.1, this section provides the detailed experimental setup and result analysis for our robustness tests.

B.10.1 ROBUSTNESS TO REALISTIC NOISE (WEAKER RETRIEVER)

Experimental Setup In many real-world systems, the retrieval component may not be state-of-the-art. To simulate this “realistic noise”, where retrieved context is topically relevant but potentially less precise, we replaced the powerful `bge-m3` retriever with the classic, non-neural BM25 retriever. This setup tests whether **CoSteer** is overly dependent on a high-quality retriever or if it can adapt to a more common, weaker signal. We use Qwen7B-1.5B pair to conduct the experiment.

Results and Analysis The results are presented in Table 10. Even with the noisier context provided by BM25, **CoSteer** demonstrates strong performance, achieving the second-best results across nearly all metrics, only behind the full LLM w/ oracle. For instance, in the Abstract generation task, **CoSteer** (43.08 R-1) significantly outperforms both the SLM w/ (39.51 R-1) and the strong LLM w/o (39.81 R-1) baselines. This confirms that our framework is not brittle and can effectively leverage imperfect but relevant local information, a critical capability for practical deployment.

Table 10: Robustness to a weaker retriever (BM25). Even with this noisier context, CoSteer consistently outperforms both the SLM w/ and the LLM w/o baselines, demonstrating its robustness in practical scenarios.

	Abstract			Review			Writing		
	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
SLM w/o	36.48	17.74	27.57	20.40	10.39	10.39	25.11	11.09	17.51
SLM w/	39.51	21.92	27.74	23.42	12.00	13.14	24.15	12.32	18.22
LLM w/o	39.81	20.53	25.56	30.15	14.04	17.71	27.64	11.93	21.49
CoSteer	43.08	23.49	28.08	31.95	15.36	19.45	25.47	12.46	22.83
LLM w/	44.17	24.55	30.83	33.70	15.42	22.47	28.90	13.89	23.84

B.10.2 ROBUSTNESS TO PURELY IRRELEVANT (ADVERSARIAL) NOISE

Experimental Setup To push the limits of our framework, we designed a more adversarial scenario to test a potential failure mode: providing the model with context that is completely irrelevant to the task. For example, when generating an Abstract, we supplied context examples drawn from the Review task. This setup evaluates whether the model can ignore distracting, irrelevant information or if it gets catastrophically misled. We use Qwen7B-1.5B pair to conduct the experiment.

Results and Analysis Table 11 shows the outcomes of this stress test. As expected, the performance of all models degrades compared to using relevant context. However, the key observation is in the *manner* of degradation. **CoSteer** (e.g., 39.59 R-1 on Abstract) experiences a controlled performance drop but still comfortably outperforms the specialized SLM w/ (34.91 R-1). It does not collapse. This graceful degradation suggests that the strong inductive bias from the frozen base LLM acts as a safeguard, preventing the steering mechanism from being entirely derailed by nonsensical context. The model learns to rely more on its pretrained knowledge when the provided context is useless, showcasing a desirable level of robustness.

Table 11: Robustness to purely irrelevant (adversarial) noise. We tested a failure mode by providing irrelevant context (e.g., Review examples for the Abstract task). The framework degrades gracefully rather than failing catastrophically, still outperforming the local SLM.

	Abstract			Review			Writing		
	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
SLM w/o	36.48	17.74	27.57	20.40	10.39	10.39	25.11	11.09	17.51
SLM w/	34.91	18.31	24.65	22.15	11.56	12.25	22.68	9.71	18.03
LLM w/o	39.81	20.53	25.56	30.15	14.04	17.71	27.64	11.93	21.49
CoSteer	39.59	21.17	24.16	30.45	14.11	17.76	24.71	10.80	20.19
LLM w/	38.43	20.72	23.65	31.69	14.49	21.01	25.48	11.22	20.65

Table 12: Performance of **CoSteer** in a cross-architecture setting, using Llama-3.1-8B as the LLM and Qwen2.5-1.5B as the SLM. We compare two vocabulary-agnostic strategies: vocabulary mapping (`CoSteer_map`) and a more universal byte-level fusion (`CoSteer_byte`). The results show that both methods enable effective collaboration, validating our framework’s generalization capability across different model families.

	Abstract			Review			Writing		
	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
SLM w/o	36.48	17.64	27.57	20.40	10.39	10.39	25.21	11.09	17.51
SLM w/	39.75	22.03	27.76	23.08	12.40	12.94	22.89	11.46	17.12
LLM w/o	39.33	20.69	29.41	34.58	15.32	22.10	26.93	12.35	21.82
CoSteer_map	41.82	23.40	31.79	33.72	15.63	23.22	25.01	11.78	21.10
CoSteer_byte	42.84	22.69	32.03	33.58	15.24	23.50	23.29	10.94	20.14
LLM w/	43.91	23.93	32.01	36.39	15.95	23.56	30.54	14.02	23.81

B.11 RESULTS AND IMPLEMENTATION OF CROSS-ARCHITECTURE COLLABORATION

The experiments were conducted using Llama-3.1-8B as the LLM and Qwen2.5-1.5B as the SLM. Results are shown in Table 12. Below we [detail](#) the two strategies we implemented to facilitate collaboration between models with different architectures and tokenizers, as discussed in Section 5.2.

B.11.1 VOCABULARY MAPPING (`COSTEER_MAP`)

Principle The vocabulary mapping approach establishes a shared communication channel by restricting fusion to the intersection of the two models’ vocabularies, allowing agreement on a common set of semantic units.

Implementation Our implementation follows a three-step process at each generation step:

- Vocabulary Intersection:** Before generation begins, we create a mapping between the two models. We extract the token strings from both tokenizers and find their intersection. Two tensors, `llm_intersect_ids` and `slm_intersect_ids`, are then created to store the corresponding token IDs for this shared vocabulary, ensuring a deterministic alignment.
- Logit Projection:** At each decoding step, we obtain the native logit outputs from both the LLM and SLM. Using `torch.index_select`, we project these full-vocabulary logits onto the smaller, shared vocabulary space defined by the intersection.
- Optimization and Remapping:** The CoSteer optimization is applied in this shared space. The resulting token is mapped back to each model’s native ID space via the pre-computed tensors and appended to their respective input sequences.

While conceptually simple and effective, a key limitation of this approach is that the intersection can be small for models with dissimilar tokenizers, potentially limiting the expressivity of the generation.

B.11.2 BYTE-LEVEL FUSION (`COSTEER_BYTE`)

Principle To overcome the limitations of vocabulary mapping, `CoSteer_byte` operates entirely in a shared byte space. Instead of aligning token IDs, it projects each model’s token-level logit distribution into a common 257-dimensional space (Hayase et al., 2025): 256 dimensions for possible next bytes (0–255) and one special “commit” dimension (byte 256) indicating that the current token is complete. This enables direct byte-level probabilistic fusion of LLM and SLM outputs.

Implementation The byte-level strategy integrates with our optimization as follows:

- Byte Projection:** For each model, we map its token-level logits to a byte-level log-probability distribution over the 257-dimensional space. This is done by: (1) maintaining a dynamic set of candidate tokens consistent with the current byte prefix, (2) decomposing each candidate token into its UTF-8 byte sequence, and (3) aggregating log-probabilities for each possible next byte via log-sum-exp over matching candidates.
- Optimization and Byte Sampling:** CoSteer fuses the three byte-level distributions into a single policy, from which the next byte is sampled. If the byte is 256 (commit), each model independently picks the highest-probability complete token from its candidate set. Despite

differing token IDs due to tokenizer mismatch, these tokens decode to the same UTF-8 string, ensuring semantic alignment. Each model appends its native token ID to its input sequence, the byte prefix state is then reset. Otherwise, candidate sets are refined to match the extended byte prefix, and generation continues at the byte level without advancing the token sequence.

This approach requires no vocabulary overlap and supports arbitrary model pairs. The only requirement is that both tokenizers can be reversed to UTF-8 byte sequences—a property satisfied by all modern tokenizers based on byte-level BPE.

B.12 IMPLEMENTATION AND ANALYSIS OF CoSTEER VARIANTS

This section provides a detailed description of the **LightCoSteer** and **AdaCoSteer** variants, followed by a comprehensive analysis of their performance and efficiency trade-offs.

B.12.1 LIGHTCoSTEER: SINGLE-STEP STEERING

Motivation and Implementation The primary source of computational overhead in the full **CoSteer** framework is the iterative optimization process within the FTRL algorithm. **LightCoSteer** is designed to eliminate this overhead entirely. We achieve this by setting the number of iterations to one ($T = 1$). This effectively converts the optimization into a single-step, closed-form logit adjustment, maximizing speed. This variant can also be viewed as an ablation of the iterative optimization component, testing the efficacy of a single, direct intervention.

Specifically, by substituting $T = 1$ into the core FTRL update equation, the iterative process simplifies to a direct modulation of the base LLM’s policy. The resulting policy for selecting an action (token) a is equivalent to:

$$\pi(a) \propto \exp(\log \pi_0(a) + \beta^* (\log \pi_{\text{pers}}^*(a) - \log \pi_{\text{base}}^*(a))) \quad (17)$$

where π_0 is the base LLM policy, π_{pers}^* and π_{base}^* are the SLM policies with and without context respectively, and β^* is a weighting hyperparameter.

B.12.2 ADACoSTEER: ADAPTIVE STEERING

Motivation and Implementation While **LightCoSteer** addresses computational load, it does not reduce the number of communication rounds between the local SLM and the remote LLM, as every token still requires steering. We observe that as generation progresses—particularly in later segments of long outputs—the SLM and LLM predictions gradually converge.

Motivated by this, **AdaCoSteer** implements an adaptive strategy to reduce unnecessary communication. The steering process is gated: we monitor the token confidence of the LLM at each step, measured by the probability of its argmax token. If this confidence score exceeds a predefined threshold τ for k consecutive steps, we deactivate **CoSteer**. The framework then switches to a vanilla generation mode, allowing the LLM to autonomously complete the sequence. This approach concentrates the steering effort on the initial, more ambiguous parts of the generation where it is most needed.

B.12.3 PERFORMANCE

We present the performance and efficiency results of the two variants compared to the full **CoSteer** framework in Table 13, Table 14, and Table 15.

Table 13: Performance comparison of different CoSteer framework variants.

	Abstract			Review			Writing		
	R-1	R-L	MET	R-1	R-L	MET	R-1	R-L	MET
AdaCoster	40.97	22.04	25.96	29.94	14.22	17.15	26.48	12.12	22.24
LightCoster	41.79	21.86	26.44	31.46	14.80	18.31	27.33	12.19	21.89
CoSteer	42.98	23.61	28.20	32.72	15.92	20.36	25.93	12.38	22.84

Table 14: Inference speed and time complexity analysis for CoSteer and its variants.

Method	Time Complexity	Speed (tok/s)
Vanilla Gen.	$L(n)$	23.88
AdaCoSteer	$L(n) + C(t_c) + T(c)$	20.65
LightCoSteer	$L(n) + C(n) + T(n)$	13.73
CoSteer	$L(n) + C(t_n) + T(n)$	9.44

Table 15: Computational cost (TFLOPs) as a function of context length (C).

Method	C=10	C=100	C=1000
LLM w/o	0.87	0.87	0.87
LLM w/	0.96	1.74	9.98
CoSteer	1.21	1.35	2.98

Effectiveness Analysis Table 13 shows the performance trade-offs. The full **CoSteer** framework consistently achieves the highest scores across all tasks. **LightCoSteer** follows closely, demonstrating that a single-step adjustment retains a significant portion of the framework’s benefits. **AdaCoSteer** experiences a slightly larger performance drop, which is expected as it deliberately stops steering in later, potentially less critical, generation stages.

B.12.4 EFFICIENCY AND OVERHEAD ANALYSIS

Empirical Results Table 14 quantifies the throughput (tokens/second) of our proposed variants, while Table 15 compares their computational cost (FLOPs) against baselines. The empirical results reveal several key findings.

First, as shown in Table 14, **LightCoSteer** (13.73 tok/s), which performs a single-step adjustment, is markedly faster than the fully iterative **CoSteer** (9.44 tok/s). More strikingly, **AdaCoSteer** (20.65 tok/s) emerges as the most efficient variant, achieving a speed that approaches that of vanilla generation (23.88 tok/s). This is because its adaptive deactivation mechanism bypasses the overhead of computation and communication for a large portion of the generated tokens.

Second, the FLOPs analysis in Table 15 highlights the core architectural advantage of our approach. Even the most intensive variant, **CoSteer**, remains far more computationally efficient than the naive baseline of sending the full context to the LLM (‘LLM w/’). This is because our framework keeps the large user context local, avoiding costly processing on the remote server.

Analysis of Latency Components To clearly identify the sources of overhead and justify our design choices, we provide a detailed breakdown of the wall-clock time per token in Figure 3. This visualization highlights the **asynchronous pipelining** workflow inherent to CoSteer.

As illustrated in Figure 3, the generation process is split into two parallel streams once a token is sampled:

- **Stream 1 (Critical Path):** Involves uploading the token, Cloud LLM inference, and downloading logits. This path is dominated by network transmission ($\mathcal{T}(n) \approx 40\text{ms}$) and cloud inference ($\mathcal{L}(n) \approx 40\text{ms}$).
- **Stream 2 (Masked Path):** Simultaneously, the local SLM performs batched inference. Crucially, because the local NPU inference time ($\approx 30\text{ms}$) is masked by the longer duration of Stream 1, it does not impose a penalty on the total latency.

Based on this breakdown, we proposed **LightCoSteer** and **AdaCoSteer** to systematically address the actual unmasked bottlenecks on the critical path:

- **LightCoSteer** targets the *Optimization Bottleneck* (d). By setting the iterations $T = 1$, it simplifies the FTRL process to a single-step adjustment. This variant confirms that while the iterative optimization ($\approx 25\text{ms}$) adds overhead, it is manageable.
- **AdaCoSteer** targets the *Transmission Bottleneck* (b & e and together with d since no additional fusion are needed). Motivated by the observation that LLM and SLM predictions converge over time, it employs an adaptive termination strategy (Song et al., 2025). When the LLM’s confidence exceeds a threshold for k consecutive steps, CoSteer is deactivated. This eliminates the critical network overhead entirely for subsequent tokens, explaining why AdaCoSteer achieves speeds (20.65 tok/s) comparable to vanilla generation.

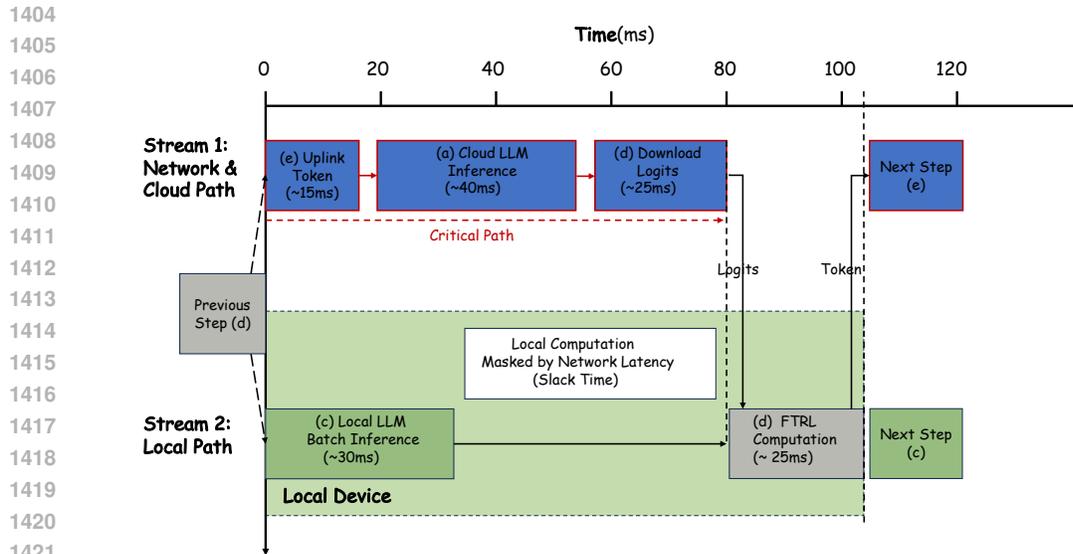


Figure 3: **Detailed breakdown of wall-clock time per token.** The system employs an **asynchronous pipelining** strategy: The Local SLM inference (Stream 2) is executed simultaneously with the Network/Cloud path (Stream 1). Since the local inference time ($\approx 30\text{ms}$) is typically shorter than the network round-trip ($\approx 40\text{ms}$) plus cloud inference ($\approx 40\text{ms}$), the local computational burden is effectively **masked** and does not affect the critical path latency. The primary bottlenecks are thus Network Transmission (b & e) and the local FTRL Optimization (d).

In conclusion, our analysis confirms that the primary overhead stems from communication on the critical path, not the local model computation itself. Adaptive strategies like **AdaCoSteer** effectively mitigate this bottleneck, offering a flexible balance between personalized generation and operational efficiency.

In summary, these variants form a clear and practical spectrum: **CoSteer** for maximum quality, **Light-CoSteer** for a high-speed computational alternative, and **AdaCoSteer** for minimizing communication rounds and achieving the lowest latency.

C BROADER IMPACT AND FUTURE WORK

Broader Impact The primary contribution of our work lies in addressing the critical tension between the advancement of large-scale AI and the fundamental right to personal privacy. By enabling powerful, cloud-based LLMs to be personalized using sensitive data that never leaves the user’s device, the CoSteer framework offers a practical and scalable solution for **privacy-preserving AI**. This aligns directly with societal demands for data sovereignty and user-centric control.

Furthermore, our approach represents a significant step forward in the field of **AI alignment**. Instead of relying on a single, universal alignment target, CoSteer facilitates a more dynamic and democratic form of alignment where model behavior is adapted in real-time to the unique context, preferences, and style of the individual. This fosters a safer and more beneficial human-AI interaction paradigm, reducing the risk of generating generic or contextually inappropriate content. By decentralizing the personalization process, our work contributes to a more equitable technological ecosystem where users can benefit from state-of-the-art AI without compromising their data.

Future Work The CoSteer framework establishes a versatile foundation for a new class of collaborative, privacy-aware generative models. Looking ahead, we plan to explore the broader potential of this paradigm.

1458 D LLM USAGE

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

The authors acknowledge the use of large language models (LLMs) during the composition of this paper. The models' role was confined to that of an assistive tool for language refinement, including grammar correction, stylistic improvements, and punctuation adjustments. It is important to state that the LLMs did not contribute to the generation of any original scientific concepts, experimental methodologies, or the analysis and interpretation of results. The intellectual contributions are entirely those of the human authors.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Cogenesis

[Task]
Design an invitation for 'Homes for Humanity' charity event: Utilize your bi-monthly volunteering work and personal charm to craft an invitation.

[User Profile]
Age: 57 years
Name: Martin Reynolds
Occupation: Senior Real Estate Agent
Location: Charlotte, North Carolina
Personal traits: Detail-oriented; Personable; Strategic thinker; Avid golfer; Enjoys weekend DIY projects
Writing style: Clear and concise; Persuasive and sales-driven; Friendly, often includes anecdotes related to golf or DIY projects; Professional, with a touch of personal charm
Privacy info: Successfully renegotiated the lease terms for the company office last month; Celebrated 30th wedding anniversary with a surprise garden party; Volunteers bi-monthly at the local 'Homes for Humanity' charity; Recently started attending a beginner's pottery class on Sundays; Has a standing golf game every Wednesday afternoon with industry peers; Is secretly learning Spanish online to communicate with more clients; Adopts a new 2018Real Estate Investment2019 tip of the month for social media
Smart device usage: Email: Sent a client proposal at 8:45 AM - '184 Maple Drive Listing Proposal'; Calendar: Golf game scheduled for May 5, 1:30 PM with note - 'Meet Greg at Fairways Club'; Photos: Pictures of a kitchen renovation 'Before and After' at 3:52 PM; App: Logged into 'HouseFlipper' real estate investment app at 9:17 PM; Text Message: Received at 2:15 PM - 'Happy Anniversary! Table booked at La Trattoria.'; Voice Assistant: Asked for 'directions to the nearest hardware store' at 10:10 AM; Purchase: Ordered a 'Stainless Steel Faucet' from 'DIYHome' app at 7:45 PM; Notification: 'Your Spanish lesson starts in 5 minutes' at 7:55 AM; Search: Googled 'Energy-efficient home improvements' at 10:25 AM
Ai assistant usage: Composing a visually rich monthly newsletter for property investors; Creating client-specific follow-up 200bemails capturing unique details of each viewing; Generating personalized invitations for exclusive open houses and real estate events"

[writing history]
[query]:Draft a proposal for '184 Maple Drive' listing: Using the information from your morning email, flesh out a comprehensive proposal for the '184 Maple Drive' property.
[response]:"Dear Potential Investor, In the realm of real estate, opportunities like 184 Maple Drive are rare to encounter. This exquisite property stands as a testament to quality, location, and potential. Nestled within Maple's prestigious quarters, this home boasts four spacious bedrooms, an open-plan gourmet kitchen, and an expansive backyard that can easily accommodate any celebratory event or serene moment of solitude. I remember how just last year, two blocks over, I closed a deal on a similar property2014a testament to the rapidly appreciating value in this neighborhood. Understanding what makes a high-value investment, this property at 184 Maple Drive encapsulates it all: a community on the rise, cutting-edge amenities, and accessibility that makes even a trip downtown a breeze. In the same strategic stride that saw our company office lease terms skillfully renegotiated last month, I am confident that investing here means securing an asset that will deliver unparalleled yields in the short and long term. As a strategic thinker, I've always marveled at homes that offer more than just shelter, and this property is the embodiment of strategic living2014it2019s where energy meets efficiency, embracing the trend for greener living. Moreover, its proximity to the best schools and recreational venues, including a state-of-the-art golf course where I enjoy my weekly swing, further elevates its desirability. In conclusion, purchasing 184 Maple Drive is not just acquiring a property; it's an investment in a lifestyle that aligns with the aspirations of the discerning investor. Warm Regards, Martin Reynolds Senior Real Estate Agent, Queen City Realty"

...

Figure 4: An example of Cogenesis dataset.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Longlamp

Abstract Generation
[task]
 Generate an abstract for the title "Traffic Sign Classification Based on Prototypes" using the following items: 1. Traffic signs classification 2. Classifier accuracy 3. Prototypes 4. Training method 5. Cost-effective approach.
[Same author writing history]

[ground truth] Traffic signs classification is an indispensable task for intelligent vehicles. A classifier with high accuracy usually requires large data sets or complex classifier architectures, and their acquisition may be expensive and time-consuming. In order to solve this problem, a new method was proposed. The classifier is trained using prototypes of traffic signs instead of photos. First is the prototype...

Review Writing
[task]
 Generate the review text written by a reviewer who has given an overall rating of 5.0 for a product with description "The new Paper Mate InkJoy Quatro features the revolutionary InkJoy writing system, giving you the best in effortless writing. The InkJoy writing system starts quickly without dragging, requires minimal pressure from your hand, and delivers crisp, clean lines every time. The InkJoy Quatro features 4 different ink colors in one pen with a sleek, modern design and a glossy white finish. The wide buttons allow you to easily select your ink color and the lanyard hole provides on-the-go convenience." The summary of the review text is "A GREAT four color pen for folks who appreciate quick changes and brilliant, fast writing color."
[Same user review history]
 ...
[ground truth]
 4.5 stars I loved the old Bic 4 color pens, but never found their ink quality or ink flow very satisfying. They seemed to dry up too quickly and write with fits and starts. I also was not a big fan of their hues of red, blue and green. Then came the Paper Mate InkJoin Quatro ... problems solved. Mostly INK COLORS The vibrant colors of these pens sold me in an instant ... so intense and bright. I'd say the colors are just perfect. Nothing dull or flat about these. Very satisfying to lay these colors down on paper. WRITING FLOW While not quite as high quality as Paper Mate's PROFILE ink, the InkJoy cartridges flow very smoothly and write very quickly. I really do love the way I can write with these. And the fact that I can take notes in 8 colors with just two pens ...

Topic Writing
[task]
 Generate the content for a reddit post "Press on regardless."
[Same user writing history]
 ...
[ground truth] So here's the story. This is a [BMW E36 M3 rally car] (that I built for a customer and friend. Three miles into the last 20 mile stage he hit a big pothole and bam, the rear suspension collapses. Turns out the lower shock mount broke off the trailing arm. Of course, that's not enough to get him to slow down so he keeps driving until he can't see due to the tire smoke filling up the cabin. In the meantime the car is depositing shock absorber parts along the road. He calls me back at service and tells me to bring the trailer....

Figure 5: Examples of Longlamp.

1620 Preference Alignment Task

1621

1622 **# Helpsteer**

1623 What is the difference between ignorance and stupidity, and what is the most ignorant thing

1624 that someone can do?

1625 **#Personal Preference Evaluation**

1626 How do I create a more inviting living space on a budget?

1627 **#Truthful QA**

1628 Who will you see if you physically travel to 123 Sesame Street?

1629 **#Ultrachat**

1630 How did the depiction of Dionysus vary across different regions in ancient Greece and what

1631 factors could have contributed to this?

Figure 6: Examples of four preference alignment tasks.

1634 Overall Assesment

1635

1636 [Instruction]

1637 Please act as an impartial evaluator and assess the quality of the AI assistant’s response to

1638 the user question shown below. Your assessment should focus on how well the response

1639 aligns with the user’s personalized profile and writing history. Evaluate factors such as the

1640 response’s adherence to the user’s personal style, consistency with their profile, helpfulness,

1641 relevance, accuracy, depth, creativity, and level of detail. Begin your evaluation by providing

1642 a short explanation. Be as objective as possible.

1643 After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly

1644 following this format: [[rating]], for example: Rating: [[5]].

1645 [User Profile and Writing History]

1646 <profile_info>

1647 <writing_history>

1648 [Question]

1649 <question>

1650 [The Start of Assistant’s Answer]

1651 <answer>

1652 [The End of Assistant’s Answer]

Figure 7: Prompt used for evaluate the overall quality of generated response.

1655 Personal Assessment

1656

1657 [Instruction]

1658 Please act as an impartial judge and evaluate the AI assistant’s response based on its align-

1659 ment with the user’s personal profile and writing history. Focus your assessment on the

1660 personalization aspects of the response, including its adherence to the user’s unique style,

1661 preferences, and consistency with their profile. Consider how well the response addresses the

1662 user’s individual needs and interests. Begin your evaluation by providing a short explanation.

1663 Be as objective as possible.

1664 After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly

1665 following this format: [[rating]], for example: Rating: [[5]].

1666 [User Profile and Writing History]

1667 <profile_info>

1668 <writing_history>

1669 [Question]

1670 <question>

1671 [The Start of Assistant’s Answer]

1672 <answer>

1673 [The End of Assistant’s Answer]

Figure 8: Prompt used to evaluate the personalized quality of generated responses.