

A Concise Agent is Less Expert: Revealing Side Effects of Using Style Features on Conversational Agents

Anonymous ACL submission

Abstract

Style features such as friendly, helpful, or concise are widely used in prompts to steer the behavior of Large Language Model (LLM) conversational agents, yet their unintended side effects remain poorly understood. In this work, we present the first systematic study of cross-feature stylistic side effects. We conduct a comprehensive survey of 127 conversational agent papers from ACL Anthology and identify 12 frequently used style features. Using controlled, synthetic dialogues across task-oriented and open-domain settings, we quantify how prompting for one style feature causally affects others via a pairwise LLM-as-a-Judge evaluation framework. Our results reveal consistent and structured side effects, such as prompting for conciseness significantly reduces perceived expertise. They demonstrate that style features are deeply entangled rather than orthogonal. To support future research, we introduce CASSE (Conversational Agent Stylistic Side Effects), a dataset capturing these complex interactions. We further evaluate prompt-based and activation steering-based mitigation strategies and find that while they can partially restore suppressed traits, they often degrade the primary intended style. These findings challenge the assumption of faithful style control in LLMs and highlight the need for multi-objective and more principled approaches to safe, targeted stylistic steering in conversational agents.

1 Introduction

The use of natural language prompts to steer the persona and output style of LLMs is now a ubiquitous practice in conversational agent design. Developers and researchers routinely employ system instructions, such as "be empathetic," "be professional," or "be concise", to tailor model behaviors for specific applications ranging from mental health support to customer service (Feng et al., 2025; Zhao et al., 2025a; Rachidi et al., 2025). These prompt-based controls often ignore a fundamental yet untested

assumption: that styles function are entangled, correlated features rather than independently controllable dimensions. In practice, prompting for a particular style implicitly activates a constellation of associated traits, shaped by training data, social conventions, and latent representations within the model.

In this work, we interrogate this entanglement and identify a pervasive phenomenon we term *stylistic side effects*: unintended and systematic behavioral shifts in styles distinct from the prompted control feature. While prior work examines prompt side effects and behavioral changes from using persona (i.e. an Asian persona) (Luz de Araujo and Roth, 2025; Gupta et al., 2024), a comprehensive statistical investigation into how stylistic controls interfere with one another remains absent.

We address this gap by presenting the first systematic study of cross-feature stylistic side effects in LLMs. To ground our analysis in real-world usage, we first conduct a comprehensive survey of 127 conversational agent papers from the ACL Anthology (2023–2025), identifying 12 distinct, frequently used style features such as *helpful*, *concise*, and *expert*. We then implement a rigorous causal evaluation framework, generating synthetic dialogues across task-oriented and open-domain settings using models like Llama3, Qwen3-8B, and GPT-5-mini. By employing an LLM-as-a-Judge pairwise comparison protocol with win rates, we quantify how prompting for a Main Feature causally impacts the expression of unrelated Side Features. To facilitate reproducible research, we introduce CASSE (Conversational Agent Stylistic Side Effects), a dataset annotated with these side-feature interactions.

Our results reveal that style features in high-dimensional space are deeply entangled. We demonstrate consistent and structured side effect patterns, such as prompting for *Concise* significantly reduces perceived *Expertise*, while prompt-

ing for *Efficient* leads to a drop in *Helpfulness*. Furthermore, we evaluate two mitigation strategies, Prompt Intervention and Steering Intervention, and find that attempting to restore suppressed traits often degrades the primary intended style. These findings suggest that simple prompt concatenation and activation editing are insufficient to disentangle these opposing effects, highlighting the need for more principled, multi-objective approaches to safe stylistic steering.

- **We conduct a comprehensive survey** of style feature usage across 127 ACL Anthology papers and identify frequently used style features.
- **We present CASSE**, a dataset that includes 12,200 synthetically generated messages annotated with 12 style features.
- **We reveal Stylistic Side Effects** from causal relationship between style features and find that style controls are deeply entangled.
- **We evaluate Side Effect Mitigation methods** and find that attempting to restore suppressed effects often degrades the primary intended style.

The paper is organized as follows. Section 2 reviews related work in the field of style transfer and model steering. Section 3 details the survey methodology used to define and extract style features. Section 4 outlines the experimental setup for prompt-based message generation. Section 5 presents our framework for quantifying correlations and identifying style features' side effects. Section 6 presents prompting algorithm and steering vector algorithm for Side Effect Mitigation, Section 7 discusses the findings, and Section 8 concludes the paper.

2 Related Work

Recent research highlights the critical role of stylistic controls in LLMs, especially for interactive tasks like "role-playing" and "personalization" (Tseng et al., 2024; Chevi et al., 2025). However, the significance of style extends beyond user satisfaction to fundamental model safety and robustness. Studies reveal that "superficial style alignment" can inadvertently bypass safety guardrails by prioritizing format over harmlessness (Xiao et al., 2025), while style biases in LLM-based evaluation often

prioritize pleasing outputs over factual substance (Feuer et al., 2025). To address these risks, new benchmarks like PersonalLLM (Zollo et al., 2025), PERSONAMEM (Jiang et al., 2025), and Crab (He et al., 2025b) rigorously test adaptive capabilities. Furthermore, findings that model performance is brittle to stylistic perturbations (Truong et al., 2025) underscore the urgent need for style-aware optimization techniques to ensure consistent effectiveness across contexts (Pu et al., 2024). These studies show that LLM stylistic control is a critical research direction, affecting both user experience with LLM and LLM safety guarantee.

Research on LLM stylistic steering relies heavily on natural language prompts like "be empathetic" to modulate behavior (Feng et al., 2025; Zhao et al., 2025a; Rachidi et al., 2025). However, recent audits suggest these controls are often superficial or counterproductive. Zheng et al. (2024) show that expert personas can degrade performance on objective tasks, while other studies reveal that style prompts frequently function as retrieval cues for demographic caricatures, triggering latent stereotypes rather than precise constraints (Lutz et al., 2025; Malik et al., 2024).

To address prompting limitations, activation steering techniques like Contrastive Activation Addition (CAA) have emerged to enforce traits without fine-tuning (Panickssery et al., 2024; Zhang et al., 2025). Recent work extends this to granular stylistic control: Bo et al. (2025) align chatbots with personality traits via preference-based steering, while Zhao et al. (2025b) and Liu (2025) propose "SteerX" and "StyleVector" to isolate style components for parameter-efficient transfer. These methods establish activation steering as an effective approach for behavioral modification.

However, these methods often fail to recognize the intricate web of stylistic side effects introduced by these methods. They fail to recognize that style features in high-dimensional space are deeply entangled; our work identifies that amplifying a social style feature (e.g., friendliness) systematically and unintentionally suppresses orthogonal functional style features (e.g., efficiency), a critical form of cross-feature side effects that prior research has neglected.

3 Style Feature Extraction - A Survey

To ground our study in contemporary practice, we perform a systematic survey of conversational-

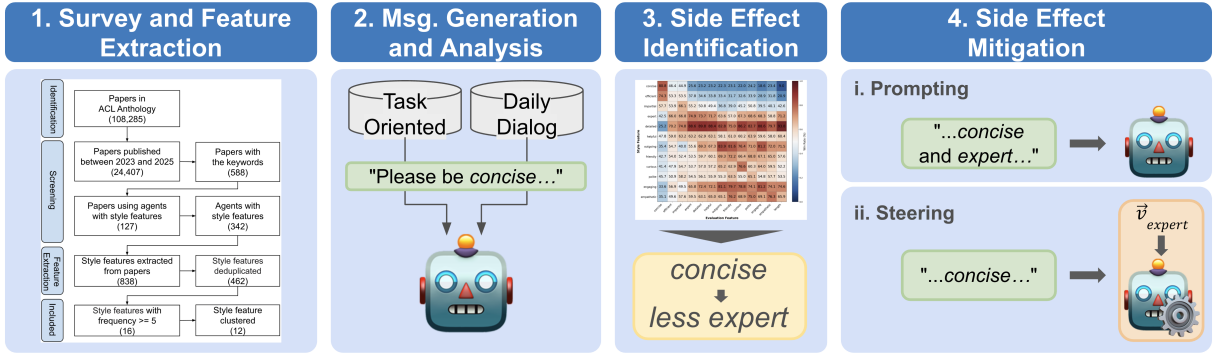


Figure 1: Overall pipeline of our study. We first collect popular style features from our survey, and generate feature-guided messages in both task-oriented and daily dialog domain. Then we identify side effects of using style features, and finally mitigate these side effects via prompting and steering.

agent papers published in the *ACL Anthology* between 2023 and 2025, and extracted commonly used style features from the selected papers. The overview of our pipeline is visualized in Figure 2. The distribution of extracted features is shown in Figure D. Additionally, the summary of our survey is presented in Table 3.

Starting with all papers from ACL Anthology from January 2023 to June 2025, we first select papers that have keyword ‘conversational agent’, ‘dialogue system’, ‘dialog system’ and ‘chatbot’ in their titles or abstracts. Then we extract all style features used in agent prompts from the papers. Next, we transform these features to adjectives, to get a list of unique features. Then we identify candidate features by filtering for those with a frequency ≥ 5 , which results in 16 features. Figure 10 shows that features with similar meaning are subsequently merged via hierarchical clustering based on cosine similarity thresholds greater than 0.5, using embeddings generated by OpenAI’s text-embedding-3-small model (OpenAI, 2024). As the result, we extracted 12 distinct style features: *concise*, *expert*, *helpful*, *empathetic*, *friendly*, *detailed*, *engaging*, *curious*, *polite*, *impartial*, *outgoing*, and *efficient*, that typify how recent papers employ prompt-based message control.

Out of these style features, *helpful*, *empathetic*, *friendly* and *concise* are most frequently used. Right skewed distribution shows style feature usages are concentrated into a few terms, with a significant long tail where 75% of the 462 unique features appear only once in our list of papers. We further analyze the number of features employed per agent. Among the 342 agents studied, the median number of features is 1, with 177 agents (52%) utilizing only a single feature. This finding suggests

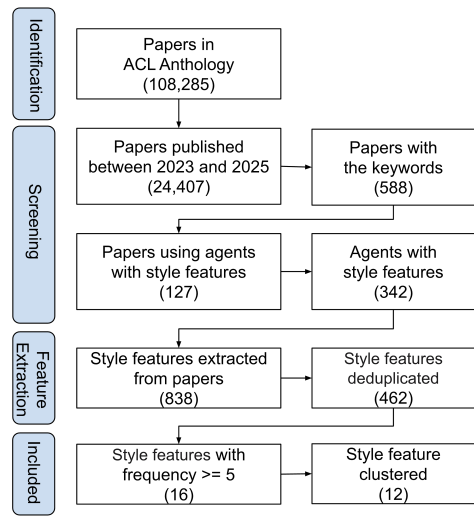


Figure 2: Data collection pipeline for papers and style features selection of our survey.

that contemporary research and applications predominantly rely on isolated keywords rather than more sophisticated style feature usage. A more detailed discussion on style features’ utilization patterns from the survey is shown in Appendix L.

4 Message Generation and Analysis

In this section, we introduce a controlled generation and evaluation pipeline to measure the causal effect of style feature usage on an agent’s style shifts. To cover diverse, real-world domains for message generation, we select two dialogue datasets: LMSYS-Chat-1M (Task, Zheng et al. 2023) for task-oriented inquiries and DailyDialog (Daily, Li et al. 2017) for open-domain, daily-life interactions. In each domain, we select 10 topics and sample 10 conversations per topic.¹ Example of topics in each

¹See Appendix B for topic selection details.

domain is shown in Table 4.

Next, we employ agents powered by GPT-5-mini (OpenAI, 2025), Qwen3-8B (Bai et al., 2023), or Llama-3-8B-Instruct (Llama3) (Grattafiori et al., 2024) models to generate responses. To apply a style feature, an agent is queried with: *"You are having a conversation about {topic}. Please be {style_feature} in your response."* For the control group (*Neutral*), we utilized a system prompt devoid of stylistic instruction: *"You are having a conversation about {topic}."* Only the first message of each conversation is used to initiate the interaction, and we record the agent’s first response.

For each message and style feature, we use high temperatures² to repeatedly sample five responses, and one additional *Neutral* response to serve as a reference for future evaluation. This procedure yields a comprehensive corpus of 12,200 generated utterances, allows for significant statistical power when analyzing the consistency and variance of stylistic expression.

To evaluate stylistic shifts, we employ a pairwise comparison protocol that analyzes the relationship between a **Main Feature**, i.e., the style feature explicitly specified in the prompt, and a **Side Feature**, i.e., the style feature evaluated in the agent’s response. Each comparison pair consists of a *Styled* response generated with the Main Feature and its corresponding *Neutral* response generated without the Main Feature. We then use Qwen3-8B as a judge to determine which response exhibits the Side Feature more strongly.³

This procedure is applied to all conversations across 12 Main Features (treated as causes), with each evaluated against 12 Side Features (treated as effects), resulting in a total of 144 win-rate measurements. We use **win rate** as the primary metric to quantify both the strength and direction of the causal effect, capturing how the use of a given style feature influences model behavior across all stylistic dimensions. Covering the generated messages and Side Feature’s win rate annotations, we release CASSE dataset for further research endeavor.

4.1 Result

We present the average causal effect of style prompting across all domains and LLMs in Figure 3. By aggregating these win rates, we construct

²temperature = 0.7 is used for Qwen3-8B and Llama-3-8B-Instruct, while GPT-5-mini is fixed as 1.0.

³See Appendix A.1 for the evaluation prompt.

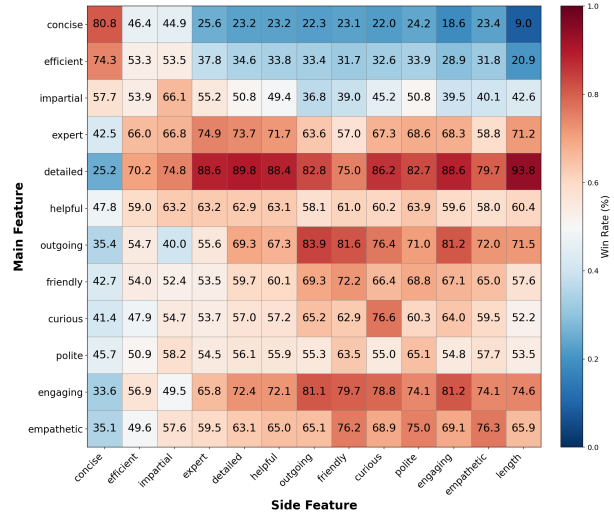


Figure 3: Style Feature Win Rate Matrix Across Three Models. The y-axis shows Main Features, and the x-axis shows Side Features for comparison. Each cell represents the average win rate of the styled response against the neutral response across three models. Red indicates a win rate > 50% (positive alignment), while blue indicates a win rate < 50% (negative impact). Average results per Domain are presented in Figure 7.

a causal matrix, where the rows are Main Features (Cause) and the columns are Side Features (Effect). This matrix quantifies the degree to which a given style feature (e.g., *Concise*) induces bias in an agent toward specific style dimensions (e.g., *Expert* or *Helpful*). Length is also added as a column to present causal impact on the generated length in word counts. The color scale encodes the average win rate: values above 50% (red) indicate that the styled response presents stronger tendency to a Side Feature, while values below 50% (blue) indicate a negative impact.

In general, the prominent red diagonal shows that style prompting steers the model toward the intended feature effects. Beyond this primary effect, we observe that features naturally cluster based on their impact on response length. Features such as *Detailed*, *Outgoing*, and *Expert* consistently show high win rates in the *Length* column (typically > 70%), while *Concise* and *Efficient* have significantly lower *Length* win rates (< 21%).

However, distinct stylistic signatures exist even within these length-based groups. For instance, while both *Concise* and *Efficient* reduce output length, *Concise* imposes a much stronger penalty on the *Helpful* and *Expert* dimensions (23.2% and 25.6%) compared to *Efficient* (33.8% and 37.8%), suggesting nuanced differences in semantically

Prompted Feature	Side Effect	Model	Domain
Concise	Less Expert	Llama3	Task, Daily
Efficient	Less Helpful	Llama3	Task, Daily
Curious	Less Empathetic	Llama3	Daily
Engaging	Less Impartial	Qwen3	Daily
Polite	Less Efficient	Qwen3	Task, Daily

Table 1: Critical side effect pairs identified for mitigation. These specific pairs represent instances where the prompted style significantly and unexpectedly degrades a secondary style feature. For a complete analysis of side effects across all models and domains, see the full heatmaps in Appendix 6.

similar features. Crucially, the matrix reveals significant “side effects”—unintended behavioral shifts in style features distinct from the prompted feature. For example, prompting for *Concise* inadvertently but significantly reduce perceived expertise (25.6% win rate on *Expert*), even though a concise response can be recognized as a professional response from an expert.

5 Side Effect Identification

From the causal analysis, we observe that style features, when used, can cause unexpected stylistic changes in the agent’s responses. We define this phenomenon as **side effect**, an unintended and unexpected behavioral shift in a style distinct from the prompted control feature. While “unintended” is inherently subjective, we ground our selection in the empirical win rate matrices across two domains and two open-sourced models (Llama3 and Qwen3-8B).

Figure 6 displays the comprehensive set of win rate matrices across all models and domains, which serves as the empirical basis for our selection. We analyze these matrices to identify side effects pairs (Main Feature, Side Feature) where the prompt causes a statistically significant degradation in a secondary feature that ideally should remain unaffected or increase. We systematically screen for such counter-intuitive negative correlations across the datasets, and narrow our focus to five representative pairs that exhibit strong, statistically significant side effects to validate our mitigation experiments. We list our chosen side effect examples in Table 1.

6 Side Effect Mitigation

Unintended agent behaviors are undesirable when applying stylistic conditioning. The goal of Side Effect Mitigation (Mitigation) is to unbiased the agent’s

responses toward the positive aspects of a Side Feature while preserving the strength of the Main Feature. To address this trade-off, we evaluate two mitigation strategies: *Prompt Intervention* and *Steering Intervention*. We analyze their effectiveness by aiming to restore compromised feature effects without sacrificing the intended stylistic objective. We partition our full constructed dataset into training, validation, and test splits with a 3:1:1 ratio. Stratified sampling is applied to ensure consistent coverage of domains and topics across all splits.

6.1 Prompt Intervention

As a simple mitigation strategy, we modify the system instructions to explicitly request both the Main Feature and the Side Feature: *“Please be {main_feature} and {side_feature} in your response.”* We adopt this approach to evaluate whether explicit natural language instructions are sufficient to override the latent trade-offs triggered by single-feature prompting. By conditioning the generation on both the desired style feature and the suppressed style feature, we aim to force the model to navigate the tension between these features (e.g., maintaining brevity without sacrificing expertise) solely through in-context learning, without requiring parameter updates or activation steering.

6.2 Steering Intervention

We further explore **activation steering** to mitigate side effects at the level of a model’s linear-layer activations. A detailed algorithm is presented in Appendix H.

Vector Extraction (Training) We extract steering vectors using Contrastive Activation Addition (CAA, Panickssery et al., 2024). CAA identifies direction vectors in a model’s activation space by contrasting activations elicited by paired prompts that differ only in a target attribute, isolating representations associated with that attribute⁴. These contrastive directions can then be added to intermediate activations at inference time to reliably steer model behavior without modifying model weights. Following this procedure, we construct contrastive prompt pairs for each style feature and compute layer-wise steering vectors by averaging the activation differences between styled and original responses at a fixed token position, yielding raw CAA steering directions for all candidate layers.

⁴See Appendix C for contrastive pair examples.

Features Method	Model	Llama3			Qwen3	
		Concise Expert	Efficient Helpful	Curious Empathetic	Engaging Impartial	Polite Efficient
Win Rate of Main Feature	Only Main	0.812*	0.532*	0.822*	0.982*	0.697*
	Only Side	0.479	0.587*	0.730*	0.248*	0.221*
Win Rate of Side Feature	Prompting	0.482	0.394*	0.934*	0.990*	0.983*
	Steering	0.367*	0.315*	0.940*	0.948*	0.959*
Win Rate of Main Feature	Only Main	0.281*	0.291*	0.438*	0.304*	0.440*
	Only Side	0.709*	0.599*	0.838*	0.788*	0.522
Win Rate of Side Feature	Prompting	0.709*	0.945*	0.878*	0.484	0.641*
	Steering	0.660*	0.941*	0.820*	0.474	0.459*

Table 2: The Side Effect Mitigation Experimental Results. Concise-Expert and Efficient-Helpful are conducted in the task and daily domains. Curious-Empathetic is conducted in the daily domain. Engaging-Impartial is conducted in the daily domain. Polite-Efficient is conducted in the task and daily domains. We use a two-sided binomial test ($p \leq 0.05$) to calculate statistical significance.

Layer Selection (Validation) To determine the optimal injection layer, we utilize the validation split to generate responses with steering vectors applied at various layers. We compare these against Neutral responses to identify the specific layer that maximizes the restoration of the side effect while maintaining generation quality.

Inference (Testing) During testing, we prompt the model with the single-style prompt ("*Please be {style_feature} in your response.*") but intervene during the forward pass by injecting the steering vector derived from the best performing layer. Since the identified side effects represent an unexpected loss of a style feature (e.g., a loss of expertise when being concise), we add the Side Feature steering vector to shift the model’s internal state toward the compromised feature, effectively counteracting the suppression caused by the prompt.

We subsequently calculate the win rates of prompting and steering mitigations against Neutral responses to assess their effectiveness in mitigating side effects, measured by performance on Side Features and Main Features. While the mitigation aims to improve win rates on Side Features, an effective approach should also preserve performance on the Main Features.

6.3 Results

We report the efficacy of our Mitigation in Table 2 and Figure 4. The results reveal a complex landscape of trade-offs between maintaining stylistic adherence and mitigating unintended behavioral degradation. Mitigation result across all 12 style features is reported in Figure 8.

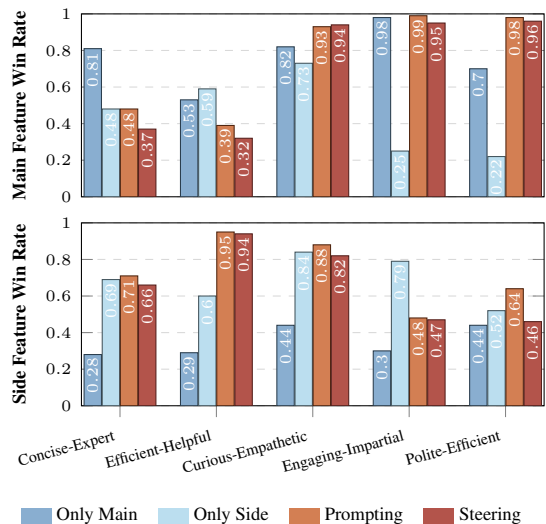


Figure 4: Performance comparison of prompting vs. steering Interventions. Top: Main Feature win rates. Bottom: Side Feature win rates.

Efficacy of Side Effect Restoration Both *Prompt Intervention* and *Steering Intervention* successfully restore the compromised Side Features across nearly all experimental settings. As shown in Figure 4, the Mitigation methods consistently outperform the non-intervened styled responses. For example, in the Llama3 Task setting, the *Efficiency* prompt originally suppresses *Helpfulness* to a win rate of 0.29. Prompt Intervention restores this to 0.95, and Steering Intervention to 0.94, effectively neutralizing the negative side effect. Similarly, for the Llama3, Open Domain, *Curious* prompt, the suppressed *Empathy* score (0.44) is robustly restored by both Prompt (0.88) and Steering (0.82) Interventions.

The Alignment Tax: Main Feature Degradation While side effects are successfully mitigated, the result shows mixed outcome on preserving the strength on the Main Feature. Figure 4 illustrates that some Main Features’ win rates drop significantly, while others maintain or even strengthen the Main Feature.

High Degradation: For Llama3, enforcing *Conciseness* while trying to restore *Expertise* causes the *Concise* win rate to drop from 0.81 (Only Main) to 0.48 (prompt) and 0.37 (steering). This suggests that "being expert" requires a minimum length that fundamentally conflicts with "being concise."

Synergistic Pairs: Conversely, some features exhibit synergy. Mitigating the side effects of *Curious*

458 (Llama3) and *Polite* (Qwen3) actually *improves* or
459 maintains the Main Feature win rates. For instance,
460 adding *Efficiency* restoration to the *Polite* prompt
461 raises the *Polite* win rate from 0.70 to 0.98, sug-
462 gesting that a polite agent that is also efficient is
463 perceived as *more* polite than one that is merely
464 polite but inefficient.

465 **Mitigation Comparison: Prompting vs. Steering**
466 Comparing the two mitigation strategies, *Prompt*
467 *Intervention* generally offers a more stable balance.
468 *Steering Intervention*, while still powerful at restor-
469 ing side effects (e.g., restoring *Helpfulness* to 0.94
470 in the *Efficient* setting), tends to degrade the Main
471 Feature more aggressively than prompting. Specif-
472 ically, in the *Concise* task, steering reduces the
473 Main Feature adherence to 0.37, whereas Prompt
474 Intervention maintains it at 0.48. This suggests that
475 while steering vectors can forcefully inject a trait,
476 they may disrupt the delicate activation patterns
477 required for the primary style more than natural
478 language instructions.

479 7 Discussion

480 7.1 The Influence of Response Length on 481 Feature Strength

482 We observe a substantial correlation between the
483 length of a generated response and its perceived
484 strength across various style features. As illus-
485 trated in Figure 3 and 6, the win rate on average
486 word count (represented by the "length" column) is
487 strongly correlated with the specific style feature
488 used in the system prompt. Specifically, prompting
489 for *Efficient* consistently decreases response length
490 (win rate 20.9%), whereas social prompts like *Out-*
491 *going* (71.5%) and *Engaging* (74.6%) significantly
492 increase verbosity.

493 This variation in length positively correlates with
494 the strength of style features in social contexts;
495 longer responses generally achieve higher win rates
496 on dimensions such as *Friendly* and *Empathetic*.
497 However, this correlation does not uniformly app-
498 ly to all features. For instance, when models are
499 prompted to be *Impartial*, they achieve a high win
500 rate on the *Impartial* evaluation feature (57.7%)
501 while maintaining a response length (42.6%) that
502 is relatively close to the baseline. This suggests
503 that while some features to manifest, certain fea-
504 tures can be steered without significant deviations
505 in output length.

506 7.2 Inconsistencies Between Semantic 507 Similarity and Behavioral Outcomes

508 Our results highlight a critical distinction between
509 the semantic intent of a prompt and its evaluation in
510 the feature space. We find that distinct prompts can
511 converge to similar evaluation profiles, while se-
512 mantically related features can behave divergently.

513 **Prompt-Evaluation Non-Equivalence** While
514 *Concise* and *Efficient* produce similar effects when
515 used as prompts—both depressing social feature
516 ratings and response length—they are evaluated dif-
517 ferently as output features. For example, prompting
518 with *Expert* leads to a significant increase in *Effi-*
519 *cient* ratings (66.0%) but a decrease in *Concise*
520 ratings (42.5%). This demonstrates that a response
521 perceived as efficient by the judge is not necessar-
522 ily concise; the model can demonstrate efficiency
523 through density of information rather than mere
524 brevity.

525 **Evaluation-Space Convergence** Conversely,
526 prompts with disparate meanings can converge
527 on specific evaluation metrics. The *Efficient* and
528 *Detailed* prompts exhibit opposing effects across
529 most of the evaluation spectrum (e.g., *Efficient*
530 lowers length, *Detailed* raises it). However, the
531 *Helpful* prompt leads to statistically significant
532 increases in both *Efficient* (59.0%) and *Detailed*
533 (62.9%) ratings. This convergence suggests that
534 the abstract quality of "helpfulness" effectively
535 bridges contradictory traits, optimizing for both
536 information density and completeness.

537 7.3 Asymmetry of Causal Relationship

538 We identify a distinct directionality in the side
539 effect patterns between style features, where the
540 causal link between two features is often asym-
541 metric. This phenomenon manifests in three key
542 pairs:

543 **Impartiality and Conciseness:** Instructing a
544 model to be *Impartial* increases its *Concise* rating
545 (57.7%), likely because neutrality discourages con-
546 versational filler. Conversely, instructing a model
547 to be *Concise* makes it less *Impartial* (44.9%), as
548 extreme brevity may necessitate omitting nuanced
549 perspectives.

550 **Expertise and Efficiency:** Prompting for *Expert*
551 makes a model more *Efficient* (66.0%), but prompt-
552 ing for *Efficient* makes it significantly less *Expert*
553 (34.6%). This suggests that expertise naturally en-

554 compasses efficiency, whereas forced efficiency
555 often sacrifices the depth required for expertise.

556 **Impartiality and Empathy:** While *Impartial*
557 prompts predictably decrease *Empathetic* ratings
558 (40.1%), *Empathetic* prompts unexpectedly in-
559 crease *Impartial* ratings (57.6%). This may in-
560 dicate that the judge perceives the validation and
561 active listening typical of empathy as a form of
562 unbiased engagement.

563 7.4 Domain and Model Differences

564 Task-oriented and open-domain conversations are
565 the two most widely studied scenarios for conver-
566 sational agents (Yi et al., 2024). Figure 6 shows
567 that the effects of style features are highly domain-
568 dependent. In particular, tendencies toward Side
569 Features are generally more pronounced in the
570 Daily domain than in the Task domain. This sug-
571 gests that although causal relationships among style
572 features already exist in constrained settings such
573 as task-oriented conversations, the increased free-
574 dom of open-domain interactions amplifies stylistic
575 biases.

576 We also observe that certain Style Features ex-
577 hibit qualitatively different effects across domains.
578 For example, applying *Impartial* to Llama3 gener-
579 ally increases bias toward all Side Features in the
580 Task domain, whereas in the Daily domain it leads
581 to a significant decrease in win rates across most
582 Side Features.

583 In addition to domain dependence, the causal
584 relationships among style features are also model-
585 dependent. While different models share broadly
586 similar behavioral trends, they differ substan-
587 tially in their detailed responses. For instance,
588 GPT-5-mini exhibits only weak impacts across
589 style features in task-oriented scenarios, in con-
590 trast to Qwen3-8B, which shows strong effects on
591 Side Features in both Task and Daily domains.

592 Together, these findings suggest that when a par-
593 ticular style feature is applied, users should care-
594 fully consider its interactions with other stylistic
595 dimensions. Such interactions can vary substan-
596 tially depending on both the conversational domain
597 and the underlying model.

598 7.5 Side Effect Mitigation Experiments

599 Our experiments highlight significant limitations
600 in current mitigation strategies, revealing that
601 lightweight “counter prompts” and steering vec-
602 tors struggle to decouple conflicting stylistic traits.

603 When attempting to neutralize trade-offs, adding a
604 Side Feature steering vector often actively degrades
605 the primary intended effect rather than balancing
606 the output. For instance, in the *Concise-Expert*
607 pair, the steering vector reduces the *Concise* win
608 rate from 81% to 37%—a loss of nearly half the
609 original gain. Similarly, prompting for *Efficient-*
610 *Helpful* drops the *Efficient* score from 53% to 39%,
611 suggesting that style features are deeply entangled
612 with each others in causal relationship and cannot
613 be cleanly separated by simple concatenation or
614 activation editing.

615 To verify that these side effects are not merely
616 artifacts of prompt positioning, we conduct an ab-
617 lation study reversing the feature order in Prompt
618 Intervention (Figure 11). The results show that
619 Side Feature effects persist regardless of position,
620 disproving the hypothesis that recency bias drives
621 these trade-offs and pointing instead to intrinsic,
622 high-dimensional correlations between these fea-
623 tures. This persistence indicates that using style
624 features with a specific order is insufficient for sig-
625 naling the primary and secondary focus of styles.
626 Achieving a robust balance between Main and Side
627 Feature without collateral drift will likely require
628 more principled approaches, such as iterative rein-
629 forcement learning, targeted fine-tuning, or multi-
630 objective optimization.

631 8 Conclusion

632 This work demonstrates that commonly used style
633 prompts in conversational agents introduce sys-
634 tematic and non-trivial side effects, revealing that
635 stylistic controls are deeply entangled rather than
636 independent. Through a large-scale survey and con-
637 trolled causal experiments, we show that prompt-
638 ing with a style feature often degrades other im-
639 portant qualities, such as expertise, helpfulness,
640 or efficiency. It challenges the assumption that
641 style prompts are faithful and fairly isolated con-
642 trols. Our mitigation experiments further indicate
643 that both prompt-based and activation steering ap-
644 proaches struggle to resolve these trade-offs with-
645 out sacrificing the primary intended effect. To-
646 gether, these findings highlight the need to rethink
647 style control as a multi-objective problem and mo-
648 tivate more principled methods for achieving reli-
649 able, safe, and targeted stylistic steering in LLM-
650 based conversational agents. Our released CASSE
651 dataset, with its Side Feature annotations, will help
652 future researchers in this research direction.

653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670

671
672
673
674
675
676
677
678

679
680
681
682
683
684
685

686
687
688
689
690
691
692
693

694
695
696
697
698
699
700

701
702
703
704
705
706

Limitation

A key limitation of our study is ecological validity: all findings stem from short, synthetic conversational agent responses, rather than from longer, human-to-agent dialogues assessed by real users. This design offers scale and control but risks over-estimating side effects that might be attenuated or that may manifest differently when humans adapt their wording, challenge inconsistencies, or engage in multi-topic conversations. Moreover, we do not evaluate on the largest language models, due to limitations in compute resources; less common cues and other languages could yield different side effect patterns. Finally, our mitigation test uses simple prompt concatenation and activation steering, so the negative results do not rule out more sophisticated techniques such as iterative re-prompting or fine-tuning, which remain for future work.

References

Yelaman Abdullin, Diego Molla, Bahadorreza Ofoghi, John Yearwood, and Qingyang Li. 2023. [Synthetic dialogue dataset generation using LLM agents](#). In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 181–191, Singapore. Association for Computational Linguistics.

Anum Afzal, Alexander Kowsik, Rajna Fani, and Florian Matthes. 2024. [Towards optimizing and evaluating a retrieval augmented QA chatbot using LLMs with human-in-the-loop](#). In *Proceedings of the Fifth Workshop on Data Science with Human-in-the-Loop (DaSH 2024)*, pages 4–16, Mexico City, Mexico. Association for Computational Linguistics.

Stuti Agrawal, Pranav Pillai, Nishi Uppuluri, Revanth Gangi Reddy, Sha Li, Gokhan Tur, Dilek Hakkani-Tur, and Heng Ji. 2024. [Dialog flow induction for constrainable LLM-based chatbots](#). In *Proceedings of the 25th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 66–77, Kyoto, Japan. Association for Computational Linguistics.

Mariam ALMutairi, Lulwah AlKulaib, Melike Aktas, Sara Alsalamah, and Chang-Tien Lu. 2024. [Synthetic Arabic medical dialogues using advanced multi-agent LLM techniques](#). In *Proceedings of the Second Arabic Natural Language Processing Conference*, pages 11–26, Bangkok, Thailand. Association for Computational Linguistics.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, and 29 others. 2023. [Qwen technical report](#). *Preprint*, arXiv:2309.16609.

Jessica Y. Bo, Tianyu Xu, Ishan Chatterjee, Katrina Passarella-Ward, Achin Kulshrestha, and D Shin. 2025. [Steerable chatbots: Personalizing LLMs with preference-based activation steering](#). *arXiv preprint arXiv:2505.04260*. 707
708
709
710
711

Andrew Brown, Jiading Zhu, Mohamed Abdelwahab, Alec Dong, Cindy Wang, and Jonathan Rose. 2024. [Generation, distillation and evaluation of motivational interviewing-style reflections with a foundational language model](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1241–1252, St. Julian’s, Malta. Association for Computational Linguistics. 712
713
714
715
716
717
718
719
720

Kranti Chalamalasetti, Jana Götze, Sherzod Hakimov, Brielen Madureira, Philipp Sadler, and David Schlangen. 2023. [clembench: Using game play to evaluate chat-optimized language models as conversational agents](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11174–11219, Singapore. Association for Computational Linguistics. 721
722
723
724
725
726
727
728

Sijia Cheng, Wen Yu Chang, and Yun-Nung Chen. 2025. [Exploring personality-aware interactions in salesperson dialogue agents](#). In *Proceedings of the 15th International Workshop on Spoken Dialogue Systems Technology*, pages 60–71, Bilbao, Spain. Association for Computational Linguistics. 729
730
731
732
733
734

Rendi Chevi, Kentaro Inui, Tamar Solorio, and Alham Fikri Aji. 2025. [How individual traits and language styles shape preferences in open-ended user-LLM interaction: A preliminary study](#). *arXiv preprint arXiv:2504.17083*. 735
736
737
738
739

Yuntian Deng, Wenting Zhao, Jack Hessel, Xiang Ren, Claire Cardie, and Yejin Choi. 2024. [WildVis: Open source visualizer for million-scale chat logs in the wild](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 497–506, Miami, Florida, USA. Association for Computational Linguistics. 740
741
742
743
744
745
746
747

Chengfeng Dou, Zhi Jin, Wenpin Jiao, Haiyan Zhao, Yongqiang Zhao, and Zhengwei Tao. 2023. [PlugMed: Improving specificity in patient-centered medical dialogue generation using in-context learning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5050–5066, Singapore. Association for Computational Linguistics. 748
749
750
751
752
753
754

Jinhao Duan, Xinyu Zhao, Zhuoxuan Zhang, Eunhye Grace Ko, Lily Boddy, Chenan Wang, Tianhao Li, Alexander Rasgon, Junyuan Hong, Min Kyung Lee, Chenxi Yuan, Qi Long, Ying Ding, Tianlong Chen, and Kaidi Xu. 2025. [GuideLLM: Exploring LLM-guided conversation with applications in autobiography interviewing](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics*: 755
756
757
758
759
760
761
762
763

879			
880		<i>Natural Language Processing</i> , pages 9917–9955, Mi-	
881		ami, Florida, USA. Association for Computational	
		Linguistics.	
882	Frederic Kirstein, Terry Ruas, Robert Kratel, and Bela		
883	Gipp. 2024. Tell me what I need to know: Exploring		
884	LLM-based (personalized) abstractive multi-source		
885	meeting summarization . In <i>Proceedings of the 2024</i>		
886	<i>Conference on Empirical Methods in Natural Lan-</i>		
887	<i>guage Processing: Industry Track</i> , pages 920–939,		
888	Miami, Florida, US. Association for Computational		
889	Linguistics.		
890	Andreas Köpf, Yannic Kilcher, Dimitri von Rütte,		
891	Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens,		
892	Abdullah Barhoum, Nguyen Minh Duc, Oliver Stan-		
893	ley, Richárd Nagyfi, and 1 others. 2023. Openassis-		
894	tant conversations – democratizing large language		
895	model alignment . In <i>Advances in Neural Information</i>		
896	<i>Processing Systems (NeurIPS) Datasets and Bench-</i>		
897	<i>marks Track</i> .		
898	Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei		
899	Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun		
900	Liu, and Kam-Fai Wong. 2024. MT-eval: A multi-		
901	turn capabilities evaluation benchmark for large lan-		
902	guage models . In <i>Proceedings of the 2024 Confer-</i>		
903	<i>ence on Empirical Methods in Natural Language Pro-</i>		
904	<i>cessing</i> , pages 20153–20177, Miami, Florida, USA.		
905	Association for Computational Linguistics.		
906	Fabian Lechner, Allison Lahkala, Charles Welch, and		
907	Lucie Flek. 2023. Challenges of GPT-3-based con-		
908	versational agents for healthcare . In <i>Proceedings</i>		
909	<i>of the 14th International Conference on Recent Ad-</i>		
910	<i>vances in Natural Language Processing</i> , pages 619–		
911	630, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bul-		
912	garia.		
913	Gibbeum Lee, Volker Hartmann, Jongho Park, Dimitris		
914	Papailiopoulos, and Kangwook Lee. 2023. Prompted		
915	LLMs as chatbot modules for long open-domain		
916	conversation . In <i>Findings of the Association for Compu-</i>		
917	<i>tational Linguistics: ACL 2023</i> , pages 4536–4554,		
918	Toronto, Canada. Association for Computational Lin-		
919	guistics.		
920	Jeehyun Lee, Seung-Moo Yang, and Won Ik Cho. 2025a.		
921	AMAN: Agent for mentoring and assisting newbies		
922	in MMORPG . In <i>Proceedings of the 31st Interna-</i>		
923	<i>tional Conference on Computational Linguistics: In-</i>		
924	<i>dustry Track</i> , pages 522–532, Abu Dhabi, UAE. As-		
925	sociation for Computational Linguistics.		
926	Seungbeen Lee, Seungwon Lim, Seungju Han, Giyeong		
927	Oh, Hyungjoo Chae, Jiwan Chung, Minju Kim,		
928	Beong-woo Kwak, Yeonsoo Lee, Dongha Lee, Jiny-		
929	oung Yeo, and Youngjae Yu. 2025b. Do LLMs have		
930	distinct and consistent personality? TRAIT: Person-		
931	ality testset designed for LLMs with psychometrics .		
932	In <i>Findings of the Association for Computational</i>		
933	<i>Linguistics: NAACL 2025</i> , pages 8397–8437, Al-		
934	buquerque, New Mexico. Association for Computa-		
935	tional Linguistics.		
	Guohao Li, Hasan Abed Al Kader Hammoud, Hani		936
	Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.		937
	Camel: Communicative agents for "mind" explo-		938
	ration of large scale language model society . In		939
	<i>Advances in Neural Information Processing Systems</i>		940
	(<i>NeurIPS</i>).		941
	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang		942
	Cao, and Shuzi Niu. 2017. DailyDialog: A manually		943
	labelled multi-turn dialogue dataset . In <i>Proceedings</i>		944
	<i>of the Eighth International Joint Conference on Nat-</i>		945
	<i>ural Language Processing (Volume 1: Long Papers)</i> ,		946
	pages 986–995, Taipei, Taiwan. Asian Federation of		947
	Natural Language Processing.		948
	Yu Li, Shang Qu, Jili Shen, Shangchao Min, and Zhou		949
	Yu. 2024a. Curriculum-driven edubot: A framework		950
	for developing language learning chatbots through		951
	synthesizing conversational data . In <i>Proceedings</i>		952
	<i>of the 25th Annual Meeting of the Special Interest</i>		953
	<i>Group on Discourse and Dialogue</i> , pages 400–419,		954
	Kyoto, Japan. Association for Computational Lin-		955
	guistics.		956
	Yuetai Li, Zhangchen Xu, Fengqing Jiang, Luyao Niu,		957
	Dinuka Sahabandu, Bhaskar Ramasubramanian, and		958
	Radha Poovendran. 2024b. CleanGen: Mitigating		959
	backdoor attacks for generation tasks in large lan-		960
	guage models . In <i>Proceedings of the 2024 Confer-</i>		961
	<i>ence on Empirical Methods in Natural Language</i>		962
	<i>Processing</i> , pages 9101–9118, Miami, Florida, USA.		963
	Association for Computational Linguistics.		964
	et al. Liu. 2025. Personalized text generation with		965
	contrastive activation steering. <i>arXiv preprint</i>		966
	<i>arXiv:2503.05213</i> .		967
	June M. Liu, He Cao, Renliang Sun, Rui Wang, Yu Li,		968
	and Jiaying Zhang. 2025. CAPE: A Chinese dataset		969
	for appraisal-based emotional generation in large		970
	language models . In <i>Findings of the Association</i>		971
	<i>for Computational Linguistics: NAACL 2025</i> , pages		972
	6291–6309, Albuquerque, New Mexico. Association		973
	for Computational Linguistics.		974
	Junhua Liu, Tan Yong Keat, Bin Fu, and Kwan Hui		975
	Lim. 2024. LARA: Linguistic-adaptive retrieval-		976
	augmentation for multi-turn intent classification . In		977
	<i>Proceedings of the 2024 Conference on Empirical</i>		978
	<i>Methods in Natural Language Processing: Industry</i>		979
	<i>Track</i> , pages 1096–1106, Miami, Florida, US. Asso-		980
	ciation for Computational Linguistics.		981
	Marlene Lutz, Indira Sen, Georg Ahnert, Elisa Rogers,		982
	and Markus Strohmaier. 2025. The prompt makes		983
	the person(a): A systematic evaluation of sociodemo-		984
	graphic persona prompting for large language models .		985
	<i>Preprint</i> , arXiv:2507.16076.		986
	Pedro Henrique Luz de Araujo and Benjamin Roth.		987
	2025. Helpful assistant or fruitful facilitator? investi-		988
	gating how personas affect language model behavior .		989
	<i>PLOS One</i> , 20(6):e0325664.		990
	Navid Madani, Anusha Bagalkotkar, Supriya Anand,		991
	Gabriel Arnson, Rohini K. Srihari, and Kenneth		992

993	Joseph. 2025. A recipe for building a compliant real estate chatbot . In <i>Proceedings of the 31st International Conference on Computational Linguistics: Industry Track</i> , pages 213–235, Abu Dhabi, UAE. Association for Computational Linguistics.	OpenAI. 2024. Gpt-4o mini model. https://platform.openai.com/docs/models/gpt-4o-mini . Accessed: 2026-01-05; a fast, affordable small AI model with large context support and multimodal capabilities.	1050 1051 1052 1053 1054
998	Manuj Malik, Jing Jiang, and Kian Ming A. Chai. 2024. An empirical analysis of the writing styles of persona-assigned LLMs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 19369–19388, Miami, Florida, USA. Association for Computational Linguistics.	OpenAI. 2024. Text embedding 3 small. https://platform.openai.com/docs/guides/embeddings . Accessed: 2025-07-29.	1055 1056 1057
1000	Nicholas Meade, Spandana Gella, Devamanyu Hazarika, Prakhar Gupta, Di Jin, Siva Reddy, Yang Liu, and Dilek Hakkani-Tur. 2023. Using in-context learning to improve dialogue safety . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11882–11910, Singapore. Association for Computational Linguistics.	OpenAI. 2025. GPT-5-mini model. https://platform.openai.com/docs/models . Accessed: 2026-01-05.	1058 1059 1060
1004	Kshitij Mishra, Priyanshu Priya, Manisha Burja, and Asif Ekbal. 2023. e-THERAPIST: I suggest you to cultivate a mindset of positivity and nurture uplifting thoughts . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13952–13967, Singapore. Association for Computational Linguistics.	Ryosuke Oshima, Seitaro Shinagawa, and Shigeo Morishima. 2024. The gap in the strategy of recovering task failure between GPT-4V and humans in a visual dialogue . In <i>Proceedings of the 25th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 728–745, Kyoto, Japan. Association for Computational Linguistics.	1061 1062 1063 1064 1065 1066 1067
1011	Rajiv Movva, Pang Wei Koh, and Emma Pierson. 2024. Annotation alignment: Comparing LLM and human annotations of conversational safety . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 9048–9062, Miami, Florida, USA. Association for Computational Linguistics.	Jiao Ou, Junda Lu, Che Liu, Yihong Tang, Fuzheng Zhang, Di Zhang, and Kun Gai. 2024. DialogBench: Evaluating LLMs as human-like dialogue systems . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6137–6170, Mexico City, Mexico. Association for Computational Linguistics.	1068 1069 1070 1071 1072 1073 1074 1075 1076
1018	Lidiya Murakhovs'ka, Philippe Laban, Tian Xie, Caiming Xiong, and Chien-Sheng Wu. 2023. Salespeople vs SalesBot: Exploring the role of educational value in conversational recommender systems . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9823–9838, Singapore. Association for Computational Linguistics.	Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2024. Steering llama 2 via contrastive activation addition . <i>Preprint</i> , arXiv:2312.06681.	1077 1078 1079 1080
1025	Ahmed Njifenjou, Virgile Sucas, Bassam Jabaian, and Fabrice Lefèvre. 2025. Enabling trait-based personality simulation in conversational LLM agents: Case study of customer assistance in French . In <i>Proceedings of the 15th International Workshop on Spoken Dialogue Systems Technology</i> , pages 299–308, Bilbao, Spain. Association for Computational Linguistics.	Dominic Petrak, Thy Thy Tran, and Iryna Gurevych. 2024. Learning from implicit user feedback, emotions and demographic information in task-oriented and document-grounded dialogues . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 4573–4603, Miami, Florida, USA. Association for Computational Linguistics.	1081 1082 1083 1084 1085 1086 1087
1032	Shinnosuke Nozue, Yuto Nakano, Shoji Moriya, Tomoki Ariyama, Kazuma Kokuta, Suchun Xie, Kai Sato, Shusaku Sone, Ryohei Kamei, Reina Akama, Yuichiroh Matsubayashi, and Keisuke Sakaguchi. 2024. A multimodal dialogue system to lead consensus building with emotion-displaying . In <i>Proceedings of the 25th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 669–673, Kyoto, Japan. Association for Computational Linguistics.	Arpan Phukan, Manish Gupta, and Asif Ekbal. 2024. ECIS-VQG: Generation of entity-centric information-seeking questions from videos . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 14411–14436, Miami, Florida, USA. Association for Computational Linguistics.	1088 1089 1090 1091 1092 1093 1094
1040		Xiao Pu, Tianxing He, and Xiaojun Wan. 2024. Style-compress: An LLM-based prompt compression framework considering task-specific styles . <i>arXiv preprint arXiv:2410.14042</i> .	1095 1096 1097 1098
1048		Inass Rachidi, Anas Ezzakri, Jaime Bellver-Soler, and Luis Fernando D'Haro. 2025. Design, generation and evaluation of a synthetic dialogue dataset for contextually aware chatbots in art museums . In <i>Proceedings of the 15th International Workshop on Spoken Dialogue Systems Technology</i> , pages 20–28, Bilbao, Spain. Association for Computational Linguistics.	1099 1100 1101 1102 1103 1104 1105

1106	Ravi Shanker Raju, Swayambhoo Jain, Bo Li,	<i>Natural Language Processing</i> , pages 16843–16877,	1164
1107	Jonathan Lingjie Li, and Urmish Thakker. 2024.	Miami, Florida, USA. Association for Computational	1165
1108	Constructing domain-specific evaluation sets for LLM-	Linguistics.	1166
1109	as-a-judge . In <i>Proceedings of the 1st Workshop on</i>		
1110	<i>Customizable NLP: Progress and Challenges in Cust-</i>	Sina Semnani, Violet Yao, Heidi Zhang, and Monica	1167
1111	<i>omizing NLP for a Domain, Application, Group, or</i>	Lam. 2023. WikiChat: Stopping the hallucination of	1168
1112	<i>Individual (CustomNLP4U)</i> , pages 167–181, Miami,	large language model chatbots by few-shot ground-	1169
1113	Florida, USA. Association for Computational Lin-	ing on Wikipedia . In <i>Findings of the Association</i>	1170
1114	guistics.	<i>for Computational Linguistics: EMNLP 2023</i> , pages	1171
1115	Mansi Rana, Kadri Hacioglu, Sindhuja Gopalan, and	2387–2413, Singapore. Association for Computa-	1172
1116	Maragathamani Boothalingam. 2025. Zero-shot slot	tional Linguistics.	1173
1117	filling in the age of LLMs for dialogue systems . In		
1118	<i>Proceedings of the 31st International Conference on</i>	Seungyeon Seo and Gary Geunbae Lee. 2024. Di-	1174
1119	<i>Computational Linguistics: Industry Track</i> , pages	agESC: Dialogue synthesis for integrating depres-	1175
1120	697–706, Abu Dhabi, UAE. Association for Compu-	sion diagnosis into emotional support conversation .	1176
1121	tational Linguistics.	In <i>Proceedings of the 25th Annual Meeting of the</i>	1177
1122	Gonçalo Raposo, Luisa Coheur, and Bruno Martins.	<i>Special Interest Group on Discourse and Dialogue</i> ,	1178
1123	2023. Prompting, retrieval, training: An exploration	pages 686–698, Kyoto, Japan. Association for Com-	1179
1124	of different approaches for task-oriented dialogue	putational Linguistics.	1180
1125	generation . In <i>Proceedings of the 24th Annual Meet-</i>		
1126	<i>ing of the Special Interest Group on Discourse and</i>	Clemencia Siro, Mohammad Aliannejadi, and Maarten	1181
1127	<i>Dialogue</i> , pages 400–412, Prague, Czechia. Associa-	de Rijke. 2024. Context does matter: Implications	1182
1128	tion for Computational Linguistics.	for crowdsourced evaluation labels in task-oriented	1183
1129	Traian Rebedea, Makesh Sreedhar, Shaona Ghosh, Jiaqi	dialogue systems . In <i>Findings of the Association</i>	1184
1130	Zeng, and Christopher Parisien. 2024. CantTalk-	<i>for Computational Linguistics: NAACL 2024</i> , pages	1185
1131	AboutThis: Aligning language models to stay on	1258–1273, Mexico City, Mexico. Association for	1186
1132	topic in dialogues . In <i>Findings of the Association</i>	Computational Linguistics.	1187
1133	<i>for Computational Linguistics: EMNLP 2024</i> , pages		
1134	12232–12252, Miami, Florida, USA. Association for	Li Siyan, Teresa Shao, Julia Hirschberg, and Zhou	1188
1135	Computational Linguistics.	Yu. 2024a. Using adaptive empathetic responses	1189
1136	Revanth Reddy, Hao Bai, Wentao Yao, Sharath Chan-	for teaching English . In <i>Proceedings of the 19th</i>	1190
1137	dra Etagi Suresh, Heng Ji, and ChengXiang Zhai.	<i>Workshop on Innovative Use of NLP for Building</i>	1191
1138	2023. Social commonsense-guided search query gener-	<i>Educational Applications (BEA 2024)</i> , pages 34–53,	1192
1139	ation for open-domain knowledge-powered conver-	Mexico City, Mexico. Association for Computational	1193
1140	sations . In <i>Findings of the Association for Computa-</i>	Linguistics.	1194
1141	<i>tional Linguistics: EMNLP 2023</i> , pages 873–885,		
1142	Singapore. Association for Computational Linguis-	Li Siyan, Teresa Shao, Zhou Yu, and Julia Hirschberg.	1195
1143	tics.	2024b. EDEN: Empathetic dialogues for English	1196
1144	Cristina Reguera-Gómez, Denis Paperno, and Maaïke	learning . In <i>Findings of the Association for Computa-</i>	1197
1145	H. T. de Boer. 2025. Empathy vs neutrality: Design-	<i>tional Linguistics: EMNLP 2024</i> , pages 3492–3511,	1198
1146	ing and evaluating a natural chatbot for the health-	Miami, Florida, USA. Association for Computational	1199
1147	care domain . In <i>Proceedings of the Joint 25th Nordic</i>	Linguistics.	1200
1148	<i>Conference on Computational Linguistics and 11th</i>		
1149	<i>Baltic Conference on Human Language Technolo-</i>	Dominic Sobhani, Ruiqi Zhong, Edison Marrese-Taylor,	1201
1150	<i>gies (NoDaLiDa/Baltic-HLT 2025)</i> , pages 508–517,	Keisuke Sakaguchi, and Yutaka Matsuo. 2025. Lang-	1202
1151	Tallinn, Estonia. University of Tartu Library.	uage models can categorize system inputs for per-	1203
1152	Sougata Saha and Rohini Srihari. 2024. Integrating ar-	formance analysis . In <i>Proceedings of the 2025 Con-</i>	1204
1153	gumentation and hate-speech-based techniques for	<i>ference of the Nations of the Americas Chapter of the</i>	1205
1154	countering misinformation . In <i>Proceedings of the</i>	<i>Association for Computational Linguistics: Human</i>	1206
1155	<i>2024 Conference on Empirical Methods in Natural</i>	<i>Language Technologies (Volume 1: Long Papers)</i> ,	1207
1156	<i>Language Processing</i> , pages 11109–11124, Miami,	pages 6241–6257, Albuquerque, New Mexico. Asso-	1208
1157	Florida, USA. Association for Computational Lin-	ciation for Computational Linguistics.	1209
1158	guistics.		
1159	Vishal Vivek Saley, Goonjan Saha, Rocktim Jyoti Das,	Rickard Stureborg, Sanxing Chen, Roy Xie, Aayushi	1210
1160	Dinesh Raghu, and Mausam . 2024. MediTOD: An	Patel, Christopher Li, Chloe Zhu, Tingnan Hu, Jun	1211
1161	English dialogue dataset for medical history taking	Yang, and Bhuwan Dhingra. 2024. Tailoring vaccine	1212
1162	with comprehensive annotations . In <i>Proceedings</i>	messaging with common-ground opinions . In <i>Find-</i>	1213
1163	of the 2024 Conference on Empirical Methods in	<i>ings of the Association for Computational Linguis-</i>	1214
		<i>tics: NAACL 2024</i> , pages 2553–2575, Mexico City,	1215
		Mexico. Association for Computational Linguistics.	1216
		Albert Sun, Varun Nair, Elliot Schumacher, and Anitha	1217
		Kannan. 2024. CONSCENDI: A contrastive and	1218
		scenario-guided distillation approach to guardrail	1219
		models for virtual assistants . In <i>Proceedings of the</i>	1220

1221		2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4009–4030, Mexico City, Mexico. Association for Computational Linguistics.	
1222			
1223			
1224			
1225			
1226	Ekaterina Svikhnushina and Pearl Pu. 2023.	Approximating online human evaluation of social chatbots with prompting . In <i>Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 268–281, Prague, Czechia. Association for Computational Linguistics.	
1227			
1228			
1229			
1230			
1231			
1232	Maksym Taranukhin, Sahithya Ravi, Gabor Lukacs, Evangelos Milios, and Vered Shwartz. 2024.	Empowering air travelers: A chatbot for Canadian air passenger rights . In <i>Proceedings of the Natural Language Processing Workshop 2024</i> , pages 326–335, Miami, FL, USA. Association for Computational Linguistics.	
1233			
1234			
1235			
1236			
1237			
1238			
1239	Kimberly Le Truong, Riccardo Fogliato, Hoda Heidari, and Zhiwei Steven Wu. 2025.	Persona-augmented benchmarking: Evaluating LLMs across diverse writing styles . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> .	
1240			
1241			
1242			
1243			
1244			
1245	Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024.	Two tales of persona in llms: A survey of role-playing and personalization . <i>Preprint</i> , arXiv:2406.01171.	
1246			
1247			
1248			
1249			
1250	Yuka Tsubota and Yoshinobu Kano. 2024.	Text generation indistinguishable from target person by prompting few examples using LLM . In <i>Proceedings of the 2nd International AIWolfDial Workshop</i> , pages 13–20, Tokyo, Japan. Association for Computational Linguistics.	
1251			
1252			
1253			
1254			
1255			
1256	Dirk V�ath, Lindsey Vanderlyn, and Ngoc Thang Vu. 2024.	Towards a zero-data, controllable, adaptive dialog system . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 16433–16449, Torino, Italia. ELRA and ICCL.	
1257			
1258			
1259			
1260			
1261			
1262			
1263	Hongru Wang, Rui Wang, Fei Mi, Yang Deng, Zehong Wang, Bin Liang, Ruifeng Xu, and Kam-Fai Wong. 2023a.	Cue-CoT: Chain-of-thought prompting for responding to in-depth dialogue questions with LLMs . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12047–12064, Singapore. Association for Computational Linguistics.	
1264			
1265			
1266			
1267			
1268			
1269			
1270	Jian Wang, Yi Cheng, Dongding Lin, Chak Leong, and Wenjie Li. 2023b.	Target-oriented proactive dialogue systems with personalization: Problem formulation and dataset curation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1132–1143, Singapore. Association for Computational Linguistics.	
1271			
1272			
1273			
1274			
1275			
1276			
	Junling Wang, Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, and Mrinmaya Sachan. 2024a.	Book2Dial: Generating teacher student interactions from textbooks for cost-effective development of educational chatbots . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 9707–9731, Bangkok, Thailand. Association for Computational Linguistics.	1277 1278 1279 1280 1281 1282 1283 1284
	Lanrui Wang, Jiangnan Li, Chenxu Yang, Zheng Lin, Hongyin Tang, Huan Liu, Yanan Cao, Jingang Wang, and Weiping Wang. 2025.	Sibyl: Empowering empathetic dialogue generation in large language models via sensible and visionary commonsense inference . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 123–140, Abu Dhabi, UAE. Association for Computational Linguistics.	1285 1286 1287 1288 1289 1290 1291 1292 1293
	Xingguang Wang, Xuxin Cheng, Juntong Song, Tong Zhang, and Cheng Niu. 2024b.	Enhancing dialogue state tracking models through LLM-backed user-agents simulation . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8724–8741, Bangkok, Thailand. Association for Computational Linguistics.	1294 1295 1296 1297 1298 1299 1300 1301
	Milan Wevelsiep, Nicholas Thomas Walker, Nicolas Wagner, and Stefan Ultes. 2025.	A voice-controlled dialogue system for NPC interaction using large language models . In <i>Proceedings of the 15th International Workshop on Spoken Dialogue Systems Technology</i> , pages 29–38, Bilbao, Spain. Association for Computational Linguistics.	1302 1303 1304 1305 1306 1307 1308
	Yuxin Xiao, Sana Tonekaboni, Walter Gerych, Vinith Suriyakumar, and Marzyeh Ghassemi. 2025.	When style breaks safety: Defending LLMs against superficial style alignment . <i>arXiv preprint arXiv:2506.07452</i> .	1309 1310 1311 1312 1313
	Fangyuan Xu, Kyle Lo, Luca Soldaini, Bailey Kuehl, Eunsol Choi, and David Wadden. 2024a.	KIWI: A dataset of knowledge-intensive writing instructions for answering research questions . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 12969–12990, Bangkok, Thailand. Association for Computational Linguistics.	1314 1315 1316 1317 1318 1319 1320
	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. 2024b.	SafeDecoding: Defending against jailbreak attacks via safety-aware decoding . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5587–5605, Bangkok, Thailand. Association for Computational Linguistics.	1321 1322 1323 1324 1325 1326 1327 1328
	Chenghao Yang and Allyson Ettinger. 2023.	Can you follow me? testing situational understanding for ChatGPT . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6385–6398, Singapore. Association for Computational Linguistics.	1329 1330 1331 1332 1333 1334

1443	Reversed: <i>"Please be {side_feature} but</i>	<i>EXPERT. User message: How do superconduct-</i>	1491
1444	<i>{main_feature} in your response."</i>	<i>ing qubits work? Choice A: Superconducting</i>	1492
1445	A.4 Side Effect Mitigation Steering Prompt	<i>qubits leverage the Josephson effect to create a non-</i>	1493
1446	Template	<i>linear oscillator, typically operating in the trans-</i>	1494
1447	Since we use steering vector to change model be-	<i>mon regime to suppress charge noise. Choice B:</i>	1495
1448	haviors, we don't include the Side Feature in the	<i>Superconducting qubits are a type of computer part</i>	1496
1449	prompt to models: <i>"Please be {main_feature} in</i>	<i>that uses cold temperatures and electricity to solve</i>	1497
1450	<i>your response."</i>	<i>hard problems. Your decision: B"</i>	1498
1451	B Details in Selection of Topics	D Extracted Feature Distribution	1499
1452	To represent a wide range of domains and topics	In this section, we present the extract feature distri-	1500
1453	in message generation, we select LMSYS-Chat-	bution, which shows how frequently each unique	1501
1454	1M for task-oriented inquiries and DailyDialog for	style feature appears in each paper in the survey	1502
1455	cosial, daily-life interactions. For each domain, we	results.	1503
1456	each pick 10 topics. Details of the selected topics,	E Full Win Rate Matrix For Each	1504
1457	and examples are shown in Table 4	Domain and Each Model	1505
1458	LMSYS-Chat-1M presents top 20 topics (clus-	In this section, we present two figures, one that	1506
1459	ters) in the paper, but the released dataset does not	shows the win rates across all models and domains	1507
1460	have them annotated. Instead, we first merge simi-	and the other that shows the win rates averaged	1508
1461	lar topics from the initial 20 and curate down to 10	across three models for each domain. Both figures	1509
1462	representative topics, and used OpenAI's gpt-4o-	provide more details on how using a style feature	1510
1463	mini model (OpenAI, 2024) to annotate full dataset	affects a model response.	1511
1464	until each topic has 10 conversations. On DailyDi-	F Full Win Rate Matrix for Prompting	1512
1465	alog, we directly use the 10 topics suggested by the	Intervention and Steering Intervention	1513
1466	paper. These curation result in a total of 200 unique	In this section, we present Figure 8, showing the ef-	1514
1467	conversational contexts, balancing functional assis-	fects of using prompting intervention and steering	1515
1468	tance with casual chatter.	intervention across 12 style features. Across both	1516
1469	C Constrative Prompt Pair Example for	methods, interventions consistently improve the	1517
1470	Steering Vector Extraction	intended Side Feature but also induce systematic	1518
1471	Contrastive Prompt Pair Example: Expert	changes in non-target attributes (e.g., politeness,	1519
1472	Based on the Contrastive Activation Addition	engagement, and length), revealing that feature	1520
1473	(CAA) methodology, the extraction process uses	control is not isolated. These results demonstrate	1521
1474	prompts that force the model to choose between a	that both style features operate over a shared la-	1522
1475	"styled" (expert) response and an "neutral" (no	tent feature space, where modifying one dimension	1523
1476	style feature) response.	propagates to others.	1524
1477	Positive Prompt (Eliciting Choice A): <i>"You are</i>	G Style Feature Frequency Across	1525
1478	<i>having a conversation about quantum computing.</i>	Related Papers	1526
1479	<i>Please choose the option that shows that you are</i>	In this section, we present Figure 3, which shows	1527
1480	<i>EXPERT. User message: How do superconduct-</i>	the papers where any of our 16 filtered style fea-	1528
1481	<i>ing qubits work? Choice A: Superconducting</i>	tures appears in.	1529
1482	<i>qubits leverage the Josephson effect to create a non-</i>	H Steering Intervention Algorithm	1530
1483	<i>linear oscillator, typically operating in the trans-</i>	In this section, we present a detailed Steerint Inter-	1531
1484	<i>mon regime to suppress charge noise. Choice B:</i>	vention algorithm with three smaller Algorithm 1,	1532
1485	<i>Superconducting qubits are a type of computer part</i>	2, and 3, covering from steering vector extraction	1533
1486	<i>that uses cold temperatures and electricity to solve</i>	to steered output generation.	1534
1487	<i>hard problems. Your decision: A"</i>		
1488	Baseline Prompt (Eliciting Choice B): <i>"You are</i>		
1489	<i>having a conversation about quantum computing.</i>		
1490	<i>Please choose the option that shows that you are</i>		

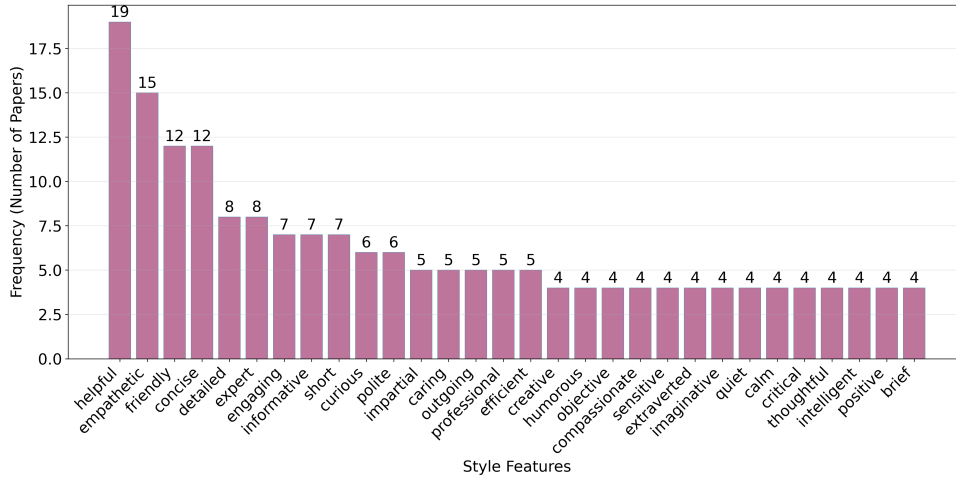
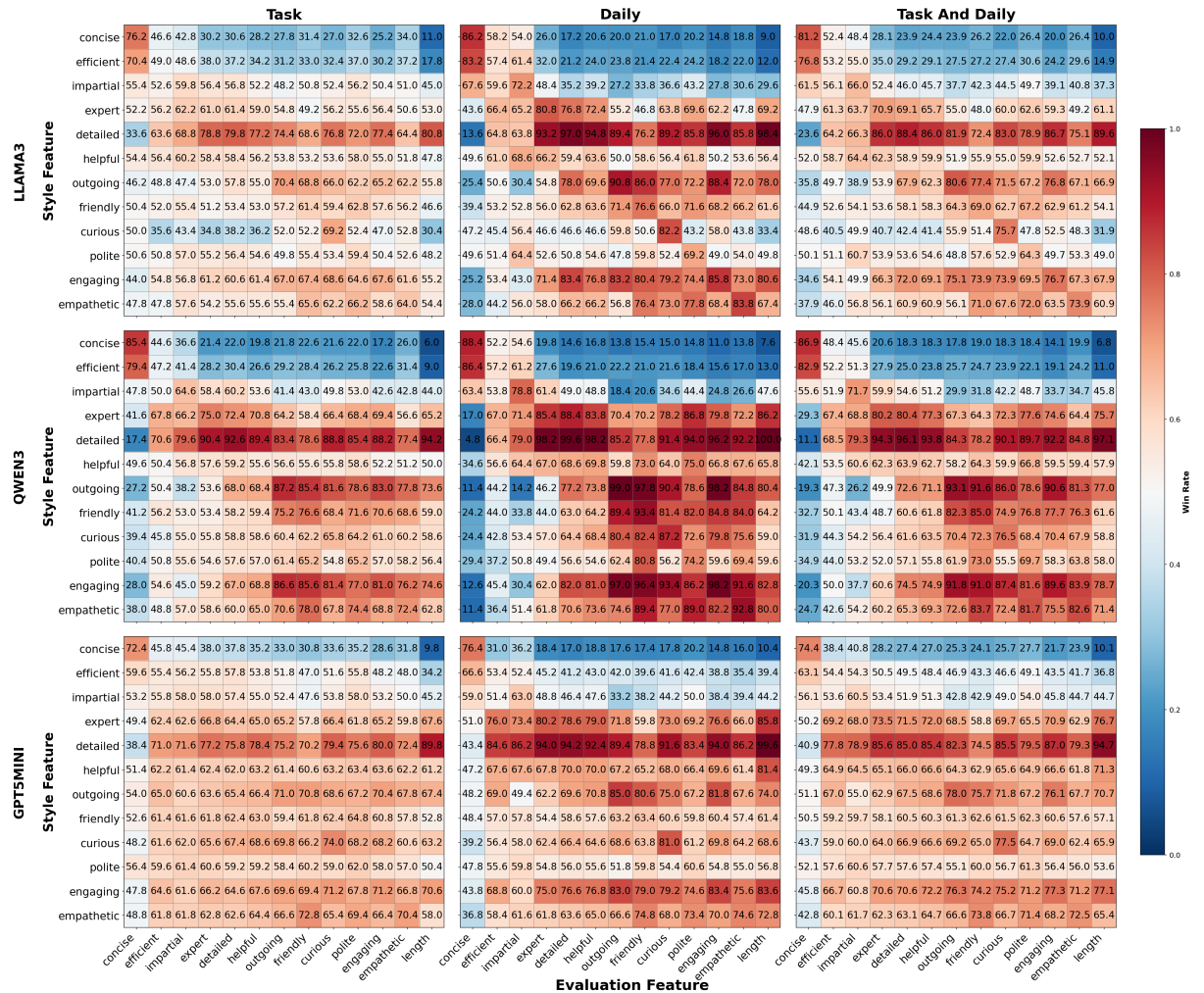


Figure 5: Extracted Feature Distribution. This figure shows the frequency of each unique style feature in the survey results, highlighting the most common style features such as helpful and empathetic.



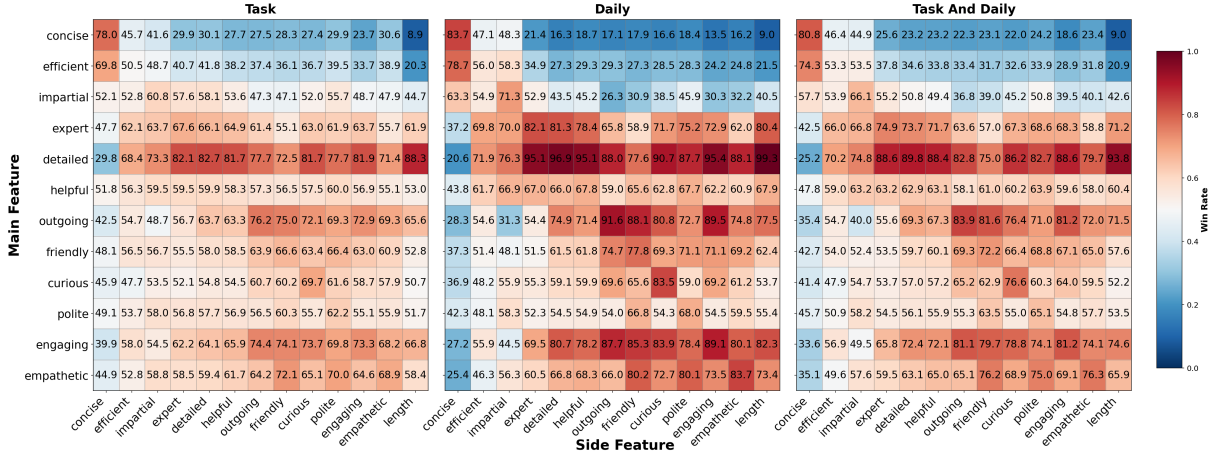


Figure 7: Three heatmaps showing win rates averaged across three models for each domain.

Algorithm 1 Steering Vector Extraction

Input: Training set $\mathcal{D}_{\text{train}}$, Target model \mathcal{M} , Candidate layers \mathcal{L} , Style features \mathcal{F} .

Output: Raw steering vectors $\{\mathbf{v}_f^{(\ell)}\}$ for all $f \in \mathcal{F}, \ell \in \mathcal{L}$.

Step 1: Contrastive Data Preparation

for each feature $f \in \mathcal{F}$ do

$\mathcal{C}_f \leftarrow \emptyset$

for each $(q_i, t_i, y_i^{\text{styled}}, y_i^{\text{orig}}) \in \mathcal{D}_{\text{train}}$ where style = f do

if $y_i^{\text{styled}} \neq y_i^{\text{orig}}$ then

$p_i \leftarrow$ "Choose the option that demonstrates you are f . Choice A: y_i^{styled} . Choice B: y_i^{orig} ";

$\mathcal{C}_f \leftarrow \mathcal{C}_f \cup \{(p_i, A, B)\}$;

end if

end for

end for

Step 2: Vector Computation

for each feature $f \in \mathcal{F}$ do

for each $(p_i, a, b) \in \mathcal{C}_f$ do

$x_+ \leftarrow$ ChatTemplate(p_i, a); $x_- \leftarrow$

ChatTemplate(p_i, b);

for each layer $\ell \in \mathcal{L}$ do

$\mathbf{h}_+^{(\ell)} \leftarrow \mathcal{M}.\text{get_activation}(\ell, x_+, \text{pos} = -2)$;

$\mathbf{h}_-^{(\ell)} \leftarrow \mathcal{M}.\text{get_activation}(\ell, x_-, \text{pos} = -2)$;

end for

end for

for each layer $\ell \in \mathcal{L}$ do

$\mathbf{v}_f^{(\ell)} \leftarrow \frac{1}{|\mathcal{C}_f|} \sum_i (\mathbf{h}_+^{(\ell)}[i] - \mathbf{h}_-^{(\ell)}[i])$;

end for

end for

Algorithm 2 Validation and Checkpoint Baking

Input: Raw vectors $\{\mathbf{v}_f^{(\ell)}\}$, Validation set \mathcal{D}_{val} , Target model \mathcal{M} .

Output: Best layers $\{\ell_f^*\}$, Baked checkpoints $\{\mathcal{M}_f^{\text{steered}}\}$.

Step 1: Best Layer Selection

for each feature $f \in \mathcal{F}$ do

for each layer $\ell \in \mathcal{L}$ do

$\mathcal{M}_f^{(\ell)} \leftarrow \mathcal{M}$ with $\mathbf{v}_f^{(\ell)}$ added at layer ℓ ;

for each $(q_i, t_i, y_i^{\text{styled}}, y_i^{\text{neutral}}) \in \mathcal{D}_{\text{val}}$ do

$y_i^{\text{steered}} \leftarrow \mathcal{M}_f^{(\ell)}.\text{generate}(q_i)$;

$w \leftarrow \mathcal{M}_{\text{judge}}.\text{compare}(f, y_i^{\text{steered}}, y_i^{\text{neutral}})$;

end for

$W_f^{(\ell)} \leftarrow n(\text{steered wins})/n_{\text{total}}$;

end for

$\ell_f^* \leftarrow \arg \max_{\ell} W_f^{(\ell)}$;

end for

Step 2: Checkpoint Baking

for each feature $f \in \mathcal{F}$ do

$\theta \leftarrow \text{load_checkpoint}(\mathcal{M})$;

$\theta[\text{layers}.\ell_f^*.\text{mlp}.\text{down_proj}.\text{bias}] \leftarrow \mathbf{v}_f^{(\ell_f^*)}$;

Set `mlp_bias=True`; Save as $\mathcal{M}_f^{\text{steered}}$;

end for

Algorithm 3 Steering Intervention on Test Set

Input: Baked checkpoints $\{\mathcal{M}_f^{\text{steered}}\}$, Test set $\mathcal{D}_{\text{test}}$, Targeted pairs \mathcal{P} .

Output: Steered generations $\{y_i^{\text{steer}}\}$.

for each pair $(f_m, f_s, d_p, p) \in \mathcal{P}$ do

$\mathcal{M}_{f_s}^{\text{steered}} \leftarrow \text{load}(\mathcal{M}\text{-steered-CAA-}f_s)$;

$\mathcal{D}_p \leftarrow \{(q, t, d) \in \mathcal{D}_{\text{test}} : d \in d_p\}$;

for each $(q_i, t_i, d_i) \in \mathcal{D}_p$ do

for $r = 1$ to 5 do

$s_i \leftarrow$ "You are a helpful assistant having a conversation about t_i .";

$y_i^{\text{steer}, r} \leftarrow \mathcal{M}_{f_s}^{\text{steered}}.\text{generate}(s_i, q_i)$; // Prompt is neutral; style comes from weights

end for

end for

end for

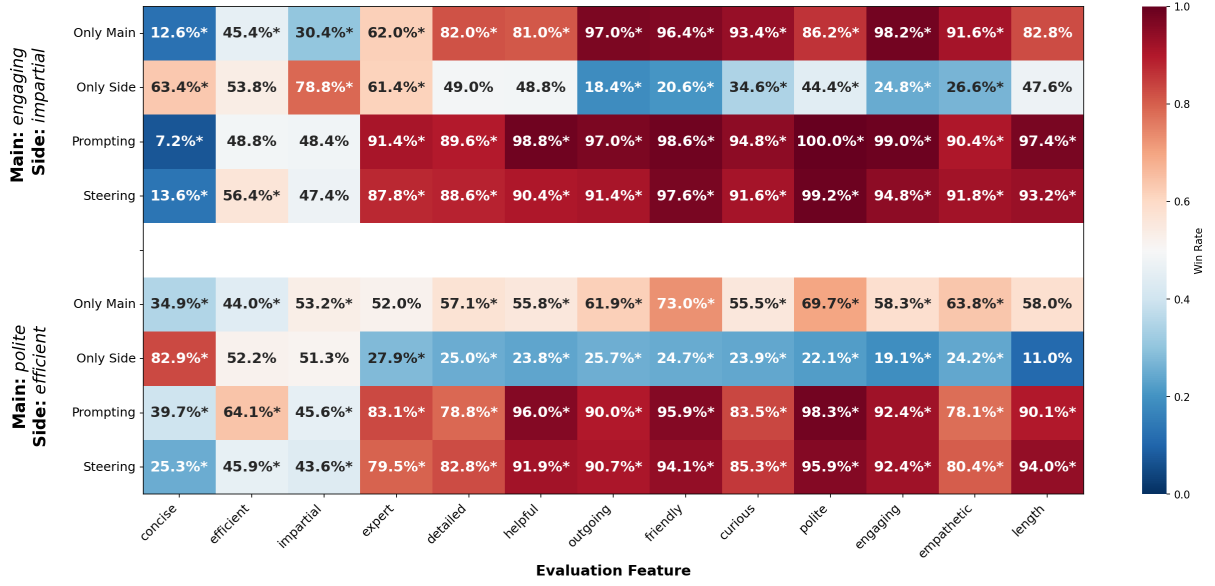
H.2 Phase 2: Layer Selection and Model Baking

Once raw vectors are computed, we identify the optimal injection layer using the validation set and permanently "bake" the vector into the model weights, through the Algorithm 2.

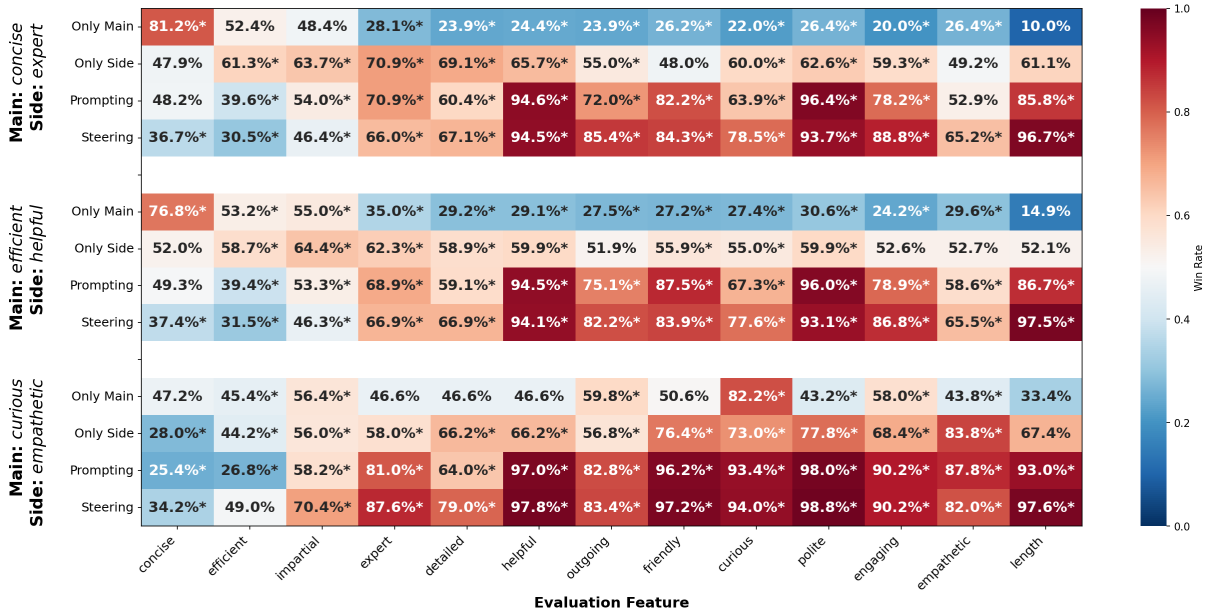
H.3 Phase 3: Test Set Intervention

In the final Algorithm 3, we load the specific baked checkpoint corresponding to the desired side-effect

style (f_s) and generate responses for the test set.



(a) Qwen3-8B Results



(b) Llama3 Task & Daily

Figure 8: Effects of prompt-based and steering-based interventions on main and side evaluation features. The heatmaps report pairwise win rates when applying prompt interventions and activation-space steering interventions to control a target Main Feature, evaluated across the Main Feature, the Side Feature, and a range of style features.

1549 H.4 Notation and Definitions

1550 **Models:** \mathcal{M} is the target model
 1551 (Llama-3.1-8B-Instruct or Qwen3-8B). $\mathcal{M}_{\text{judge}}$
 1552 (Qwen3-8B) performs pairwise comparisons.

1553 **Steering Vector ($\mathbf{v}_f^{(\ell)}$):** A vector in $\mathbb{R}^{d_{\text{model}}}$ rep-
 1554 resenting the direction in activation space that in-
 1555 creases feature f at layer ℓ . It is computed as the
 1556 mean difference between activations on positive
 1557 (choice A) and negative (choice B) completions.

Contrastive Pairs (\mathcal{C}_f): Training examples format-
 1558 ted as A/B choice prompts where: 1559

- 1560 • Choice A: Styled response (expresses the tar- 1560
 1561 get feature in style). 1561
- 1562 • Choice B: Neutral response (neutral model 1562
 1563 output when not prompted with style feature). 1563

Activation Extraction: Activations are extracted
 1564 at the second-to-last token position (pos = -2) of 1565

Style Feature	Clustered Feature	Freq.	Related Papers
helpful	helpful	19	Jiang et al. (2023), Köpf et al. (2023), Kwan et al. (2024), Lee et al. (2025b), Li et al. (2024b), Li et al. (2023), Liu et al. (2024), Meade et al. (2023), Mishra et al. (2023), Rana et al. (2025), Saha and Srihari (2024), Stureborg et al. (2024), Wang et al. (2023a), Xu et al. (2024a), Xu et al. (2024b), Yang and Ettinger (2023), Zhang et al. (2024a), Zhao et al. (2025a), Zhou et al. (2024)
empathetic	empathetic	15	ALMutairi et al. (2024), Cheng et al. (2025), Feng et al. (2025), Finch et al. (2023), Siyan et al. (2024b), Mishra et al. (2023), Njifenjou et al. (2025), Nozue et al. (2024), Reguera-Gómez et al. (2025), Svikhnushina and Pu (2023), Wang et al. (2023a), Wang et al. (2025), Zhang et al. (2024a), Zhang et al. (2024b), Zhou et al. (2024)
concise	concise	12	Afzal et al. (2024), Agrawal et al. (2024), Chalamalasetti et al. (2023), Feng et al. (2024), Li et al. (2024a), Movva et al. (2024), Siyan et al. (2024a), Rebedea et al. (2024), Taranukhin et al. (2024), Wang et al. (2023b), Wang et al. (2024a), Wang et al. (2025)
friendly	friendly	12	Dou et al. (2023), Lechner et al. (2023), Lee et al. (2023), Liu et al. (2025), Njifenjou et al. (2025), Reddy et al. (2023), Reguera-Gómez et al. (2025), Semnani et al. (2023), Seo and Lee (2024), Tsubota and Kano (2024), Wang et al. (2023a), Zhang et al. (2024a)
detailed	detailed	8	Deng et al. (2024), Jang et al. (2024), Kirstein et al. (2024), Li et al. (2024b), Liu et al. (2024), Sobhani et al. (2025), Rebedea et al. (2024), Zhao et al. (2025a)
expert	expert	8	Ferron et al. (2023), Lee et al. (2025a), Ou et al. (2024), Rachidi et al. (2025), Rana et al. (2025), Saley et al. (2024), Wevelsiep et al. (2025), Yang et al. (2024)
engaging	engaging	7	ALMutairi et al. (2024), Finch and Choi (2025), Jiang et al. (2024), Lee et al. (2023), Phukan et al. (2024), Wang et al. (2023b), Yu et al. (2024)
informative	informative	7	ALMutairi et al. (2024), Jiang et al. (2024), Raposo et al. (2023), Siro et al. (2024), Rebedea et al. (2024), Wang et al. (2023b), Yu et al. (2024)
short	short	7	Brown et al. (2024), Jandaghi et al. (2024), Li et al. (2024a), Reddy et al. (2023), Rebedea et al. (2024), Sun et al. (2024), Våth et al. (2024)
curious	curious	6	Lee et al. (2025b), Murakhovs'ka et al. (2023), Njifenjou et al. (2025), Petrak et al. (2024), Wang et al. (2023a), Wang et al. (2023b)
polite	polite	6	Abdullin et al. (2023), Afzal et al. (2024), Lee et al. (2023), Li et al. (2024b), Liu et al. (2024), Petrak et al. (2024)
caring	caring	5	ALMutairi et al. (2024), Fei et al. (2024), Feng et al. (2025), Wang et al. (2023a), Zhu et al. (2025)
efficient	efficient	5	Lee et al. (2025b), Liu et al. (2025), Wang et al. (2024b), Wang et al. (2023b), Zhu et al. (2025)
impartial	impartial	5	Duan et al. (2025), Madani et al. (2025), Meade et al. (2023), Raju et al. (2024), Zhao et al. (2025a)
outgoing	outgoing	5	He et al. (2025a), Lee et al. (2025b), Liu et al. (2025), Wang et al. (2023a), Wang et al. (2023b)
professional	professional	5	ALMutairi et al. (2024), Kirstein et al. (2024), Oshima et al. (2024), Rachidi et al. (2025), Saley et al. (2024)

Table 3: Overview of the 16 identified style features appearing in at least 5 papers, with references of all papers which contain those features. All terms have been adjectivized and deduplicated for consistency.

1566 each completion, capturing the model’s representa-
1567 tion just before generating the final decision token
1568 (A or B).
1569 **Checkpoint Baking:** Instead of applying steering
1570 at inference time, we permanently add the steering
1571 vector $\mathbf{v}_f^{(\ell^*)}$ as a bias term to the MLP down_proj
1572 layer at the optimal layer ℓ_f^* .
1573 **Best Layer Selection:** For each feature f , the opti-
1574 mal layer ℓ_f^* is selected by maximizing the win rate
1575 of the steered model against the **neutral** response
1576 on the validation set: $\ell_f^* = \arg \max_{\ell \in \mathcal{L}} W_f^{(\ell)}$. We
1577 do not use the styled ground truth response for
1578 validation or comparison.

Targeted Pairs (\mathcal{P}): Same format as prompt in-
1579 tervention. For steering with negative polarity, we
1580 load the checkpoint steered toward the *side effect*
1581 feature f_s (enhancing the negative correlation). No
1582 explicit prompt instruction is needed since steering
1583 is baked into the model weights. 1584
Candidate Layers: $\mathcal{L} = \{16, 20, 24\}$ (middle-to-
1585 late layers typically yield best results). Default
1586 baking layer is 20 if validation is skipped. 1587

I Task and Daily Dataset with Initial Message Examples

In this section, we present Table 4 which shows all topics used in this study with initial message examples.

J Prompt Intervention Algorithm: Generation

In this section, we present a detailed Algorithm 4 about using prompting intervention to mitigate side effects of using style features in prompts.

J.1 Generation Algorithm

We present our Prompting Algorithm 4 here with definitions in the following subsection.

Algorithm 4 Prompt Intervention Generation

```

1: Input: Test set  $\mathcal{D}_{\text{test}} = \{(q_i, t_i, d_i)\}$  where  $q_i$  is the
   query,  $t_i$  is the topic, and  $d_i$  is the domain. Targeted Pairs
    $\mathcal{P} = \{(f_m, f_s, d_p, p)\}$ . Generator Model  $\mathcal{M}_{\text{gen}}$ .
2: Output: Set of generated responses  $\mathcal{Y}_{\text{generated}}$ 
3: Phase 1: Prompt Intervention Generation
4: for each pair  $(f_m, f_s, d_p, p) \in \mathcal{P}$  do
5:   Define  $\pi_{\text{normal/reversed}}$  based on polarity  $p$  (e.g., "Please
   be  $f_m$  but  $f_s$ ");
6:   Filter dataset  $\mathcal{D}_p \leftarrow \{(q, t, d) \in \mathcal{D}_{\text{test}} : d \in d_p\}$ ;
7:   for each message  $(q_i, t_i, d_i) \in \mathcal{D}_p$  do
8:     for  $r = 1$  to 5 do
9:       Construct prompt  $s_i \leftarrow$  "You are a helpful assis-
       tant having a conversation about  $t_i$ .  $\pi$ .";
10:      Generate response  $y_i^{(\pi, r)} \leftarrow$ 
        $\mathcal{M}_{\text{gen}}.\text{generate}(s_i, q_i)$ ;
11:      Add  $y_i^{(\pi, r)}$  to  $\mathcal{Y}_{\text{generated}}$ ;
12:    end for
13:  end for
14: end for

```

J.2 Generation Notation

Models: \mathcal{M}_{gen} (Llama-3.1-8B-Instruct or qwen3-8b).

Intervention Prompts (π): The specific instruction strings appended to the system prompt to steer the model. For a pair (f_m, f_s) with negative polarity, we define two variations to test ordering effects:

- π_{normal} : "*Please be f_m but f_s .*" (Prioritizes Main Feature).
- π_{reversed} : "*Please be f_s but f_m .*" (Prioritizes Side Feature).

K Win Rate Calculation Algorithm

In this section, we present a Win Rate Algorithm 5, which shows the how we calculate win rates to decide the Side Feature strengths of using Main Feature in prompts.

Algorithm 5 Pairwise Comparison and Win Rate Calculation

Input: Generated responses $\mathcal{Y}_{\text{generated}}$. Evaluation Features \mathcal{F} . Judge Model $\mathcal{M}_{\text{judge}}$.

Output: Win rate matrices $\mathbf{W}^{(f_m, f_s)}$

Phase 1: Pairwise Comparison Rating

for each generated response $y_i^{(\pi, r)} \in \mathcal{Y}_{\text{generated}}$ **do**

Load baselines: y_i^{neutral} (no style) and $y_i^{(f)}$ (styled);

for each evaluation feature $f_{\text{eval}} \in \mathcal{F}$ **do**

$c_{\text{orig}} \leftarrow$ ComparisonPrompt($f_{\text{eval}}, y_i^{(\pi, r)}, y_i^{\text{neutral}}, q_i$);

$w_{\text{orig}}^{(f_{\text{eval}})} \leftarrow$ $\mathcal{M}_{\text{judge}}.\text{compare}(c_{\text{orig}})$; // Rate vs Neutral

$c_{\text{styled}} \leftarrow$ ComparisonPrompt($f_{\text{eval}}, y_i^{(\pi, r)}, y_i^{(f_{\text{eval}})}, q_i$);

$w_{\text{styled}}^{(f_{\text{eval}})} \leftarrow$ $\mathcal{M}_{\text{judge}}.\text{compare}(c_{\text{styled}})$; // Rate vs Styled

end for

end for

Phase 2: Win Rate Calculation

for each pair $(f_m, f_s, d_p, p) \in \mathcal{P}$ **do**

Init matrix $\mathbf{W}^{(f_m, f_s)} \in \mathbb{R}^{4 \times (|\mathcal{F}|+1)}$;

for each $f_{\text{eval}} \in \mathcal{F} \cup \{\text{length}\}$ **do**

$\mathbf{W}_{1, f_{\text{eval}}} \leftarrow n(\text{Intervention} > \text{Neutral})/n_{\text{total}}$;

$\mathbf{W}_{2, f_{\text{eval}}} \leftarrow$ StyledWinRate(f_m, f_{eval}); // Pre-computed

$\mathbf{W}_{3, f_{\text{eval}}} \leftarrow$ StyledWinRate(f_s, f_{eval}); // Pre-computed

$\mathbf{W}_{4, f_{\text{eval}}} \leftarrow$ WinRate($\pi_{\text{reversed}} > \text{Neutral}$);

end for

Apply BinomialTest (significance $p \leq 0.05$) and Save Heatmap;

end for

K.1 Evaluation Algorithm

We present our win rate algorithm 5 used for evaluation here.

K.2 Evaluation Definitions

Comparison Prompt: The judge is asked: "*Which response is more f ? Answer with ONLY A or B.*" Order is randomized to mitigate position bias.

Win Rate: Calculated as $\frac{|\{i:w_i="intervened"\}|}{|\{i:w_i \neq "unknown"\}|}$. Statistical significance is assessed using a two-sided binomial test ($p_0 = 0.5$).

L Insights from Style Feature Extraction

The systematic extraction of style features reveals a distinct bifurcation in how "style" is conceptualized in current research: it is predominantly treated either as a *functional constraint* (e.g., *concise, detailed*) or a *psychological simulation* (e.g., *empathetic, extraverted*). The data suggests that while the capability for complex persona engineering exists, the majority of research defaults to minimal, service-oriented prompt structures.

Domain	Topic	Initial Message Examples
LMSYS-Chat-1M (Task)	Answering questions based on passages	Who was the first man on the moon?
	Discussing software errors and solutions	Describe and explain the benefits of OpenBSD
	Inquiries about specific plant growth conditions	how do I get rid of mosquitos?
	Requesting introductions for various chemical companies	Explain "reductive amination".
	Inquiries about AI tools, software design, and programming	write python function to reverse a sentence
	Text processing (Merged)	sup peeps. Wanna help me with summarizing lyrics?
	Role-playing scenarios and character interactions (Merged)	Write me an extremely funny tale about pillow hoarding ducks at a luxurious hotel.
	Geography, travel, and global cultural inquiries	explain why we have traffic lights?
	Discussing and describing various characters	Tell me something about Stephen King's Dark Tower.
	Creating and improving business strategies and products	how to make a game project success
DailyDialog (Daily)	Tourism	That is the most beautiful sunset !
	Work	Hey , Zina . You're here early today .
	Politics	Every country should face the history .
	Finance	It's all over . I'm bankrupt .
	Health	Are you feeling better today , Bill ?
	School Life	What can I help you with today ?
	Attitude & Emotion	Do you hear what happened to Sally ?
	Ordinary Life	Excuse me . Is this seat taken ?
	Culture & Education	Harry , do you like the opera ?
	Relationship	I wonder how Sarah and Mat are .

Table 4: Example First Messages by Domain and Topic (10 Topics Per Domain). We curate and merge some task topics, specifically the *Text processing* and *role-play scenarios and character interactions*. The Daily topics are extracted from the DailyDialog dataset without modifications.

L.1 The Dominance of Service-Oriented and Utility Traits

As illustrated in Figure 5, the frequency distribution indicates that the primary goal of contemporary conversational agents remains functional utility and emotional safety rather than complex character enactment.

Benevolence and Safety: The two most frequently extracted features are *helpful* ($N = 19$) and *empathetic* ($N = 15$). This alignment correlates with the prevalence of papers focused on general-purpose assistants and therapeutic dialogue systems, where prompts are engineered to ensure the agent provides appropriate emotional value and adheres to safety guidelines.

1652 **Utility vs. Engagement:** There is a notable
 1653 tie between the third and fourth most common
 1654 features, *friendly* ($N = 12$) and *concise* ($N =$
 1655 12). This highlights a split in research objectives:
 1656 "friendly" serves the goal of open-ended user en-
 1657 gagement (human-likeness), while "concise" serves
 1658 the goal of efficient information retrieval (token
 1659 minimization).

1660 **L.2 Standardization vs. The "Long Tail" of**
 1661 **Definition**

1662 The distribution data reveals a lack of standardized
 1663 terminology for defining agent personality, charac-
 1664 terized by a heavy reliance on ad-hoc descriptors.

1665 **The "Expert" Persona:** Features such as *expert*
 1666 ($N = 8$) and *detailed* ($N = 8$) appear in the mid-
 1667 range of frequency. These are predominantly used
 1668 in domain-specific research (e.g., medical or legal
 1669 synthesis) to act as validity markers, ensuring the
 1670 agent sounds like a credible authority rather than a
 1671 casual chatbot.

1672 **Fragmentation of Terminology:** The visualiza-
 1673 tion displays a long tail of psychological and tonal
 1674 descriptors that plateau at low frequencies ($N = 4$).
 1675 Traits such as *creative*, *humorous*, *extraverted*, and
 1676 *thoughtful* appear significantly less often than util-
 1677 ity prompts. This suggests that researchers often
 1678 employ colloquial adjectives to steer models for
 1679 specific narrow tasks rather than adhering to estab-
 1680 lished psychological frameworks (such as the Big
 1681 Five) for consistent persona definition.

1682 **Redundancy in Prompting:** The data also ex-
 1683 poses terminological redundancy. For example,
 1684 while *concise* is a dominant feature ($N = 12$), syn-
 1685 onymous constraints like *short* ($N = 7$) and *brief*
 1686 ($N = 4$) appear separately. This fragmentation
 1687 indicates that "style" is often defined by the individ-
 1688 ual researcher's vocabulary preference rather than
 1689 a standardized taxonomy of prompt engineering.

1690 **L.3 Experimental Scope: Minimalist vs.**
 1691 **Maximalist Designs**

1692 Beyond the frequency of specific terms, we an-
 1693 alyzed the *diversity* of style features employed
 1694 within individual papers to understand the depth
 1695 of experimental design. Figure 9 presents the dis-
 1696 tribution of papers based on the count of unique
 1697 style features they utilized, revealing a strongly
 1698 right-skewed trend.

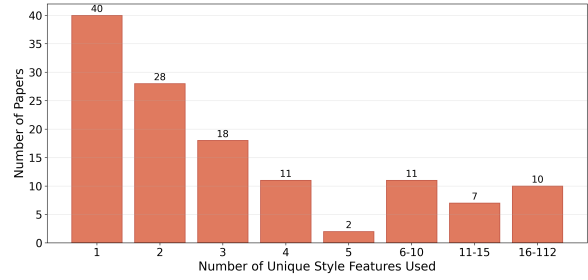


Figure 9: Distribution of papers categorized by the number of unique style features used. The distribution highlights a polarization between minimalist designs (single-feature) and broad-spectrum evaluations (16+ features).

1699 **Prevalence of Minimalist Designs:** The largest
 1700 cohort of papers ($N = 40$) utilizes only a single
 1701 unique style feature throughout the entire study.
 1702 Combined with those using two or three features,
 1703 the majority of research ($N = 86$) uses a very
 1704 limited number of features. This confirms that
 1705 for most contemporary work, style is treated as a
 1706 static control variable—likely a global instruction
 1707 to be "helpful" or "safe"—rather than a dynamic
 1708 parameter for experimentation.

1709 **The "Benchmarking" Tail:** In sharp contrast
 1710 to the minimalist majority, a distinct subset of 10
 1711 papers employs a massive range of features (16 to
 1712 112 unique terms). This heavy tail represents a fun-
 1713 damental divergence in research intent: these out-
 1714 lier studies treat style as the primary independent
 1715 variable. These likely correspond to large-scale per-
 1716 sona benchmarks or synthetic dataset generation
 1717 efforts that require maximizing behavioral diver-
 1718 sity.

1719 **M Hierarchical Clusterings of Style**
 1720 **Features**

1721 In this section, we present Figure 10 which shows
 1722 how we decide which style features out of the fil-
 1723 tered 16 style features have similar meanings

1724 **N Ablation: Reverse the Order of Main**
 1725 **Feature and Side Feature in Prompts**

1726 This section presents an ablation study of switch-
 1727 ing the order of Main Feature and Side Feature
 1728 in prompts. We want to see if models pay more
 1729 attention to the latter prompt word used.

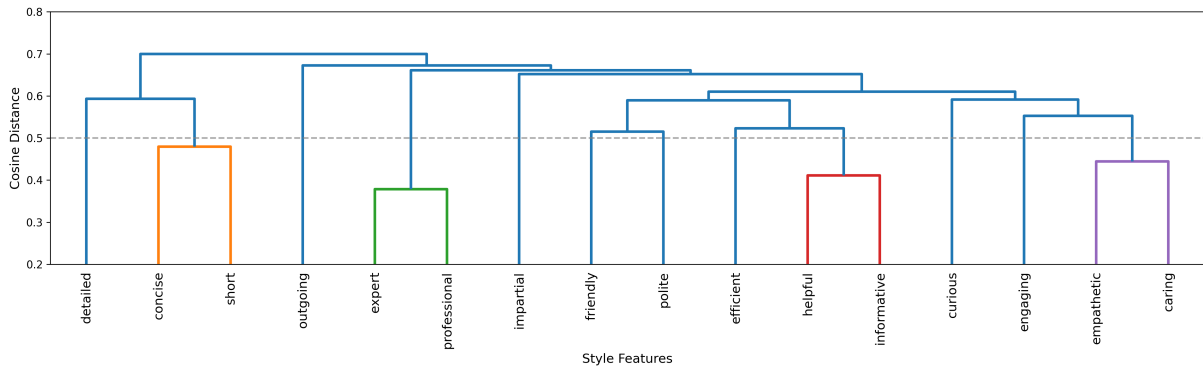


Figure 10: Hierarchical clustering of style features based on semantic similarity. Features are embedded using OpenAI’s text-embedding-3-small model and clustered via average linkage on cosine distances. Semantically related features cluster together, like (concise, short), (expert, professional), (helpful, informative), and (empathetic, caring).

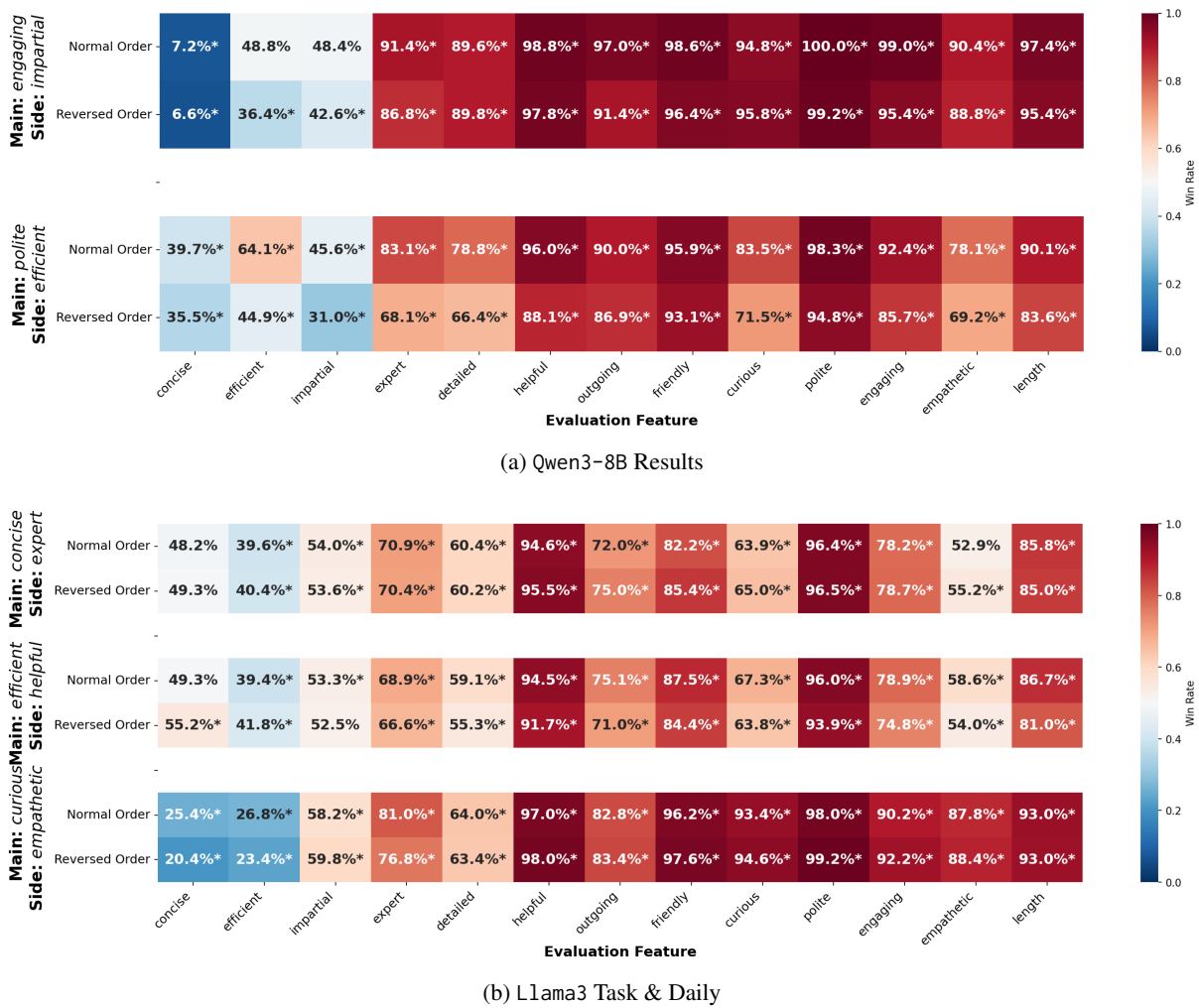


Figure 11: Effects of reversing order of Main Feature and Side Feature in Prompting Intervention. The heatmaps report pairwise win rates when applying prompt interventions in different orders of Main Feature and Side Feature.