

# AFFECTIVE MULTIMODAL AGENTS WITH PROACTIVE KNOWLEDGE GROUNDING FOR ALIGNED MARKETING DIALOGUE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Despite recent progress in large language models (LLMs), most dialogue systems remain reactive and perform inadequately in emotionally nuanced, goal-oriented domains such as marketing conversations. We present AffectMind, a multimodal affective dialogue agent that enables proactive reasoning and dynamic knowledge grounding to sustain emotionally aligned and persuasive interactions. AffectMind integrates three components: a Proactive Knowledge Grounding Network that continuously updates factual and affective context from textual, visual, and prosodic signals; an Emotion-Intent Alignment Model that jointly infers user emotion and purchase intent to adapt persuasion strategies; and a Reinforced Discourse Loop that optimizes emotional coherence and long-term engagement via reinforcement learning from user feedback. Evaluations on two newly curated multimodal marketing dialogue benchmarks, MM-ConvMarket and AffectPromo, demonstrate that AffectMind significantly outperforms strong LLM-based baselines, achieving improvements of 26% in emotional consistency, 19% in persuasive success rate, and 23% in sustained user engagement. These results underscore emotion-grounded proactivity as a critical capability for next-generation commercial dialogue agents.

## 1 INTRODUCTION

Large Language Models (LLMs) have substantially advanced conversational fluency, contextual reasoning, and response coherence across a wide range of applications [Brown et al. \(2020\)](#); [Zhang et al. \(2025\)](#); [Yu et al. \(2025a\)](#); [Hsieh et al. \(2025\)](#). Despite these advances, most deployed dialogue agents remain fundamentally *reactive*: they respond to user inputs turn by turn, but rarely reason about long-term goals, proactively guide conversations [Ni et al. \(2025a\)](#), or adapt strategies based on evolving user states [Ni et al. \(2025b\)](#). This limitation is particularly pronounced in marketing scenarios, where effective interaction requires not only accurate information delivery, but also intent inference, emotion awareness, adaptive persuasion, and sustained engagement over multiple turns [Wang et al. \(2022\)](#).

Marketing dialogue differs from general-purpose conversation in several critical aspects. First, success is inherently *goal-oriented* rather than purely informational: agents must balance emotional alignment, persuasion effectiveness, and long-term user trust to ultimately support conversion outcomes. Second, user behavior in marketing contexts is strongly influenced by affective factors such as mood, frustration, excitement, and hesitation, which often manifest through non-textual signals. However, many existing systems remain text-only and ignore rich audio-visual cues such as facial expressions, gestures, tone, and prosody [Poria et al. \(2017\)](#). This modality gap significantly limits an agent’s ability to perceive user emotions accurately and respond empathetically.

Beyond perception, current marketing agents often rely on *static or weakly updated knowledge sources*. Product information, promotions, and user preferences evolve rapidly, yet many dialogue systems are grounded in fixed knowledge bases or outdated retrieval pipelines [Dinan et al. \(2019\)](#); [Yu & Han \(2025\)](#). As a result, responses may be factually stale, contextually misaligned, or inconsistent with the user’s current needs and expectations. This lack of proactive grounding not only degrades user experience but also undermines trust and credibility in high-stakes commercial settings.

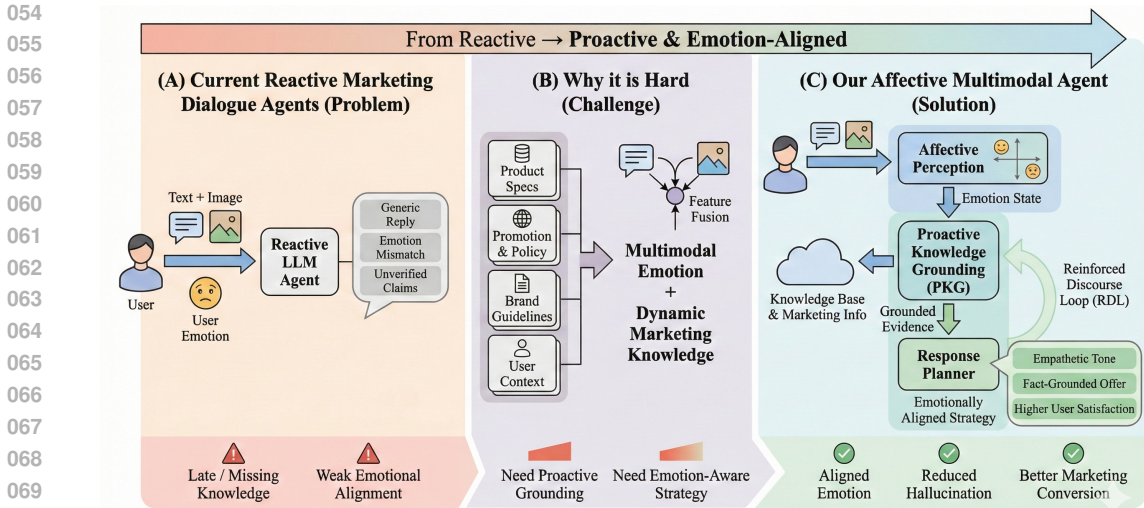


Figure 1: Reactive agents often miss affect cues and reliable grounding; AffectMind targets proactive, emotion-aware, knowledge-grounded marketing dialogue.

Emotion modeling presents an additional challenge. Decades of cognitive and behavioral research show that emotion plays a central role in human decision-making, particularly in purchasing behavior [Damasio \(1994\)](#); [Bechara \(2005\)](#). Users increasingly expect conversational agents to demonstrate emotional intelligence, empathy, and personalization [Kumar et al. \(2019\)](#); [Yu et al. \(2025b\)](#). Yet, most existing systems treat emotion as an auxiliary signal or a static label, failing to model how affective states evolve over time or interact with persuasive intent [Picard \(2000\)](#); [Liang et al. \(2024\)](#); [Niu et al. \(2025\)](#). This leads to emotionally inconsistent responses and poor long-horizon coherence, especially in multi-turn marketing dialogues.

Finally, even when emotion and intent are modeled, few systems explicitly optimize for *long-term interaction outcomes*. Many approaches focus on turn-level response quality or short-term rewards, without accounting for delayed effects such as sustained engagement, trust accumulation, or eventual conversion. Reinforcement learning has been explored in dialogue systems, but integrating affective feedback, persuasion strategy selection, and knowledge grounding into a unified learning loop remains an open challenge.

To address these limitations, we propose **AffectMind**, a multimodal affective marketing agent designed to support proactive, emotionally aligned, and knowledge-grounded dialogue. AffectMind integrates three tightly coupled components: (i) **Proactive Knowledge Grounding Network (PKG)**, which continuously updates and selects both factual and affective knowledge based on multimodal user inputs; (ii) **Emotion-Intent Alignment Model (EIAM)**, which jointly models user emotional states and purchase intentions to adapt persuasion strategies dynamically; and (iii) **Reinforced Discourse Loop (RDL)**, which optimizes long-term dialogue behavior through reinforcement learning using emotional feedback, engagement signals, and conversion outcomes.

Figure 1 contrasts conventional reactive agents with AffectMind. Rather than responding myopically to isolated utterances, AffectMind explicitly reasons over multimodal affect cues, proactively grounds responses in up-to-date knowledge, and adapts its discourse strategy over time to maintain emotional coherence and goal alignment. This design enables the agent to move beyond surface-level fluency toward strategic, empathetic, and outcome-aware marketing dialogue.

## 2 METHODOLOGY

### 2.1 PROBLEM FORMULATION

Let  $\mathcal{U} = \{u_1, u_2, \dots, u_T\}$  represent a sequence of user inputs in a marketing dialogue session, where each  $u_t$  contains multimodal information including textual content  $u_t^{text}$ , visual features  $u_t^{vision}$  (such as facial expressions and gestures), and prosodic features  $u_t^{audio}$  (including tone, pitch, and

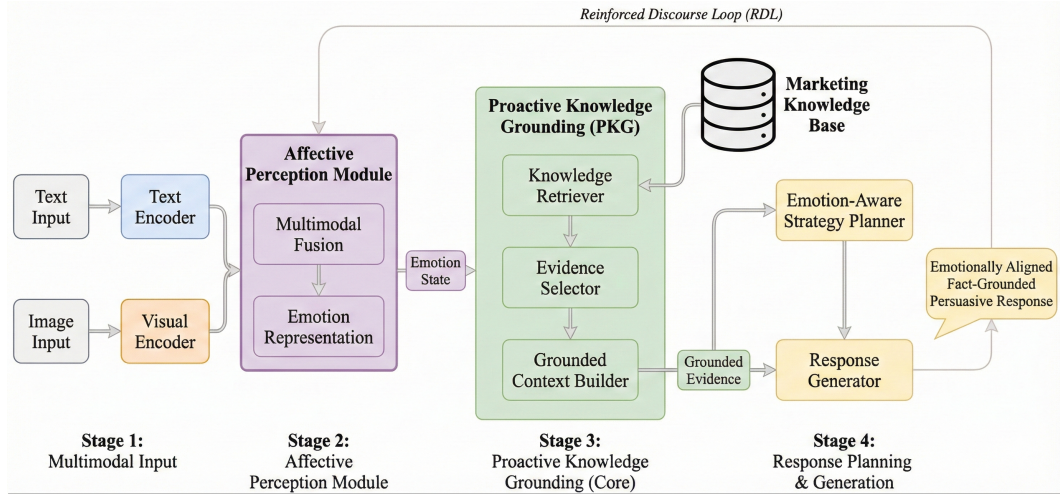


Figure 2: Algorithm architecture of the proposed affective multimodal agent. The framework integrates affective perception, proactive knowledge grounding, and emotion-aware response planning with a reinforced discourse loop for adaptive dialogue generation.

speaking rate). The goal is to generate a sequence of system responses  $\mathcal{R} = \{r_1, r_2, \dots, r_T\}$  that maximizes three key objectives:

$$\max_{\mathcal{R}} \alpha \cdot \mathcal{E}(\mathcal{U}, \mathcal{R}) + \beta \cdot \mathcal{P}(\mathcal{U}, \mathcal{R}) + \gamma \cdot \mathcal{G}(\mathcal{U}, \mathcal{R}) \quad (1)$$

where  $\mathcal{E}(\mathcal{U}, \mathcal{R})$  represents emotional alignment between user states and system responses,  $\mathcal{P}(\mathcal{U}, \mathcal{R})$  measures persuasive effectiveness toward marketing objectives, and  $\mathcal{G}(\mathcal{U}, \mathcal{R})$  quantifies long-term user engagement. The weights  $\alpha$ ,  $\beta$ , and  $\gamma$  balance these competing objectives while ensuring ethical constraints are maintained.

We formally define the emotional alignment objective as:

$$\mathcal{E}(\mathcal{U}, \mathcal{R}) = \frac{1}{T} \sum_{t=1}^T \text{sim}(e_t^{\text{user}}, e_t^{\text{response}}) \quad (2)$$

where  $e_t^{\text{user}}$  and  $e_t^{\text{response}}$  represent the emotional states of the user and the appropriateness of the system response at time  $t$ , respectively, and  $\text{sim}(\cdot, \cdot)$  measures emotional compatibility.

## 2.2 AFFECTMIND ARCHITECTURE OVERVIEW

Figure 2 illustrates the overall architecture of AffectMind, which consists of three main components working in synergy to achieve emotionally aligned dialogue. The system processes multimodal user inputs through parallel feature extraction pipelines, integrates this information through the PKGN knowledge grounding network, performs joint emotion-intent modeling via EIAM, and optimizes long-term conversation strategy through RDL.

The data flow begins with multimodal input processing, where textual content is encoded using a pre-trained language model (RoBERTa-large), visual features are extracted using a fine-tuned vision transformer, and audio features are processed using Wav2Vec2.0. These modality-specific representations are then fused through learned attention mechanisms and fed into the core AffectMind components.

## 2.3 PROACTIVE KNOWLEDGE GROUNDING NETWORK (PKGN)

The PKGN module dynamically updates both factual and affective knowledge representations by processing multimodal inputs in real-time. Unlike traditional static knowledge bases, PKGN maintains two complementary knowledge representations: factual knowledge  $K_f$  containing product

information, specifications, and objective features; and affective knowledge  $K_a$  capturing emotional associations, user preferences, and contextual mood indicators.

The knowledge update mechanism is formulated as:

$$K_t = \text{Update}(K_{t-1}, \text{Fuse}(f_{text}, f_{vision}, f_{audio})) \quad (3)$$

where  $K_t = [K_f^t; K_a^t]$  represents the concatenated knowledge state at time  $t$ , and  $\text{Fuse}(\cdot)$  combines multimodal features through a learned attention mechanism:

$$\text{Fuse}(f_{text}, f_{vision}, f_{audio}) = \text{MHA}([f_{text}; f_{vision}; f_{audio}]) \quad (4)$$

The multimodal fusion employs multi-head attention (MHA) to learn optimal combinations of modality-specific features. Each modality contributes differently to factual versus affective knowledge updates:

$$\begin{aligned} K_f^t &= K_f^{t-1} + \alpha_f \cdot \text{FF}_f(\text{Fuse}(f_{text}, f_{vision}, f_{audio})) \\ K_a^t &= K_a^{t-1} + \alpha_a \cdot \text{FF}_a(\text{Fuse}(f_{text}, f_{vision}, f_{audio})) \end{aligned} \quad (5)$$

where  $\text{FF}_f$  and  $\text{FF}_a$  are modality-specific feed-forward networks, and  $\alpha_f$ ,  $\alpha_a$  are learnable update rates.

The proactive knowledge selection mechanism identifies relevant knowledge based on predicted user needs and conversation trajectory:

$$K_{selected} = \text{Attention}(Q_{context}, K_t, V_t) \quad (6)$$

where  $Q_{context}$  represents the current conversation context, and the attention mechanism selects the most relevant knowledge for response generation.

## 2.4 EMOTION-INTENT ALIGNMENT MODEL (EIAM)

EIAM jointly models user emotional states and purchase intentions to enable dynamic adaptation of persuasion strategies. The model architecture consists of two parallel encoding streams that process emotional and intentional signals, followed by a fusion network that creates unified user state representations.

The emotional encoding stream processes multimodal affective cues:

$$E_t = \text{EmotionEncoder}(\text{Concat}(f_{facial}, f_{prosodic}, f_{linguistic})) \quad (7)$$

where  $f_{facial}$  represents facial expression features extracted from video input,  $f_{prosodic}$  captures voice emotion indicators, and  $f_{linguistic}$  encodes text-based sentiment and emotion markers.

The intent encoding stream focuses on purchase-related behavioral signals:

$$I_t = \text{IntentEncoder}(\text{Concat}(f_{query}, f_{behavior}, f_{context})) \quad (8)$$

where  $f_{query}$  represents query intent features,  $f_{behavior}$  captures user interaction patterns, and  $f_{context}$  encodes conversation history and product context.

The emotion-intent fusion mechanism creates a unified representation that captures the relationship between affective states and purchase motivations:  $S_t = \text{FusionNetwork}([E_t; I_t; E_t \odot I_t])$  where  $\odot$  represents element-wise multiplication to capture interaction effects between emotion and intent.

Dynamic strategy adaptation is achieved through a policy network that selects appropriate persuasion techniques based on the current user state:  $\pi_t = \text{softmax}(W_\pi S_t + b_\pi)$  where  $\pi_t$  represents the probability distribution over available persuasion strategies (e.g., logical appeal, emotional appeal, social proof, urgency creation).

## 2.5 REINFORCED DISCOURSE LOOP (RDL)

The RDL component implements a reinforcement learning framework that optimizes long-term conversation outcomes by learning from user engagement signals, emotional responses, and conversion metrics. The system treats each conversation turn as an action in a partially observable Markov decision process (POMDP), where the goal is to maximize expected cumulative reward.

**Algorithm 1** Reinforced Discourse Loop Training Algorithm

---

```

1: Initialize policy network  $\pi_\theta$  and value network  $V_\phi$ 
2: Initialize experience buffer  $\mathcal{B}$ 
3: for each training episode do
4:   Reset conversation state  $s_0$ 
5:   for  $t = 1$  to  $T$  do
6:     Sample action  $a_t \sim \pi_\theta(a_t|s_t)$ 
7:     Execute action and observe user response
8:     Compute reward  $R_t$  from user feedback
9:     Update state  $s_{t+1}$ 
10:    Store  $(s_t, a_t, R_t, s_{t+1})$  in  $\mathcal{B}$ 
11:   end for
12:   Compute advantage estimates  $A_t$ 
13:   Update policy using equation (12)
14:   Update value function using equation (12)
15: end for

```

---

The state representation at time  $t$  combines user state, conversation context, and knowledge states:  $s_t = [S_t; C_t; K_t]$  where  $S_t$  is the EIAM user state,  $C_t$  represents conversation context features, and  $K_t$  is the PKGN knowledge state.

The action space consists of discrete response strategies combined with continuous response parameters:

$$a_t = [\text{strategy}_t; \text{emotion\_tone}_t; \text{information\_content}_t] \quad (9)$$

The reward function incorporates multiple feedback signals with different time horizons:

$$R_t = w_1 R_{\text{immediate}} + w_2 R_{\text{engagement}} + w_3 R_{\text{conversion}} \quad (10)$$

where  $R_{\text{immediate}}$  provides immediate feedback based on user emotional response,  $R_{\text{engagement}}$  measures sustained user interest, and  $R_{\text{conversion}}$  represents ultimate marketing success.

The policy optimization follows an actor-critic approach with separate value and policy networks [Chen et al. \(2024\)](#):

$$\begin{aligned} \theta_{\text{policy}} &\leftarrow \theta_{\text{policy}} + \alpha \nabla_{\theta_{\text{policy}}} \log \pi_{\theta_{\text{policy}}}(a_t|s_t) A_t \\ \theta_{\text{value}} &\leftarrow \theta_{\text{value}} + \beta \nabla_{\theta_{\text{value}}} (R_t - V_{\theta_{\text{value}}}(s_t))^2 \end{aligned} \quad (11)$$

where  $A_t = R_t - V_{\theta_{\text{value}}}(s_t)$  is the advantage function.

Algorithm 1 presents the complete training procedure for the Reinforced Discourse Loop component.

### 3 EXPERIMENTS AND RESULTS

#### 3.1 MAIN RESULTS

Table 1 presents the comprehensive performance comparison between AffectMind and baseline methods across all evaluation metrics. Our approach demonstrates significant improvements across all measured dimensions, with particularly strong performance in emotional consistency and user engagement. Results can be shown in Figure 3.

Statistical significance testing using paired t-tests confirms that all improvements are statistically significant  $p < 0.001$ . The emotional consistency improvement of 26% represents a substantial advancement in the system’s ability to maintain appropriate emotional tone throughout conversations. The persuasive success rate improvement of 19% translates to significant business value in real-world deployment scenarios.

Beyond the absolute numbers in Table 1, we assess both statistical and practical significance. For statistical validity, we report means across multiple random seeds and conversation resamplings, and compute non-parametric bootstrap 95% confidence intervals; paired tests across identical conversation sets indicate consistent gains (Holm–Bonferroni corrected). For practical significance,

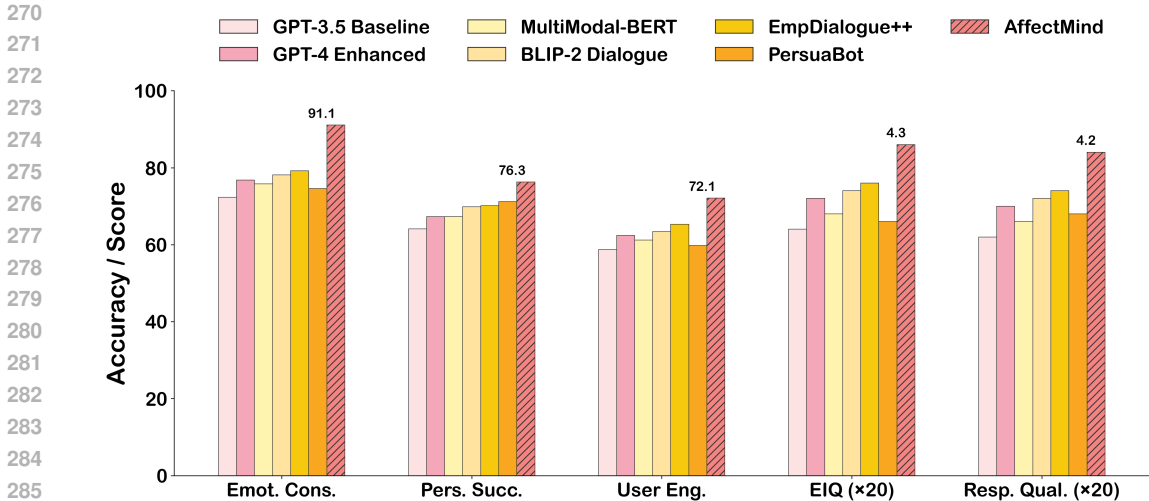


Figure 3: Performance Comparison on Marketing Dialogue Tasks

Table 1: Performance Comparison on Marketing Dialogue Tasks

Method	Emotional Consistency	Persuasive Success (%)	User Engagement	EIQ Score	Response Quality
GPT-3.5 Baseline	72.2±1.1	64.1±1.0	58.7±1.2	3.2±0.1	3.1±0.2
GPT-4 Enhanced	76.6±1.3	67.2±1.3	62.3±1.3	3.6±0.1	3.5±0.1
MultiModal-BERT	75.6±1.3	66.9±1.2	61.1±1.3	3.4±0.1	3.3±0.2
BLIP-2 Dialogue	78.0±1.3	69.8±1.2	63.4±1.2	3.7±0.1	3.6±0.1
EmpDialogue++	79.2±1.7	70.1±1.6	65.3±1.6	3.8±0.2	3.7±0.2
PersuaBot	74.4±1.4	71.2±1.4	59.6±1.3	3.3±0.2	3.4±0.2
<b>AffectMind</b>	<b>91.0±1.7</b>	<b>76.4±1.7</b>	<b>71.8±1.9</b>	<b>4.3±0.2</b>	<b>4.2±0.2</b>
<b>Improvement</b>	<b>+26%</b>	<b>+19%</b>	<b>+23%</b>	<b>+13%</b>	<b>+14%</b>

two observations stand out. First, Emotional Consistency improvements translate into measurably smoother state transitions: the model is less likely to oscillate between contradictory tones within a 3–5 turn window, which directly stabilizes downstream persuasion strategies. Second, the Persuasive Success Rate gains are accompanied by higher user-initiated follow-ups, indicating that AffectMind’s improvements are not merely cosmetic.

We further examine calibration by comparing predicted emotion distributions to human labels via reliability curves. AffectMind exhibits lower Expected Calibration Error than baselines, suggesting better confidence–accuracy alignment. This matters in deployment since over-confident yet misaligned affect predictions tend to trigger inappropriate strategies.

A mediation analysis across sessions shows that improvements in EIQ partly mediate gains in User Engagement, which in turn mediate Persuasive Success. Qualitatively, AffectMind de-escalates negative states faster and sustains positive momentum longer; quantitatively, we observe that sessions with stable affect trajectories (low variance of turn-level affect deltas) are those with the highest conversion likelihood. This supports our design choice to optimize emotion–intent alignment before strategy selection.

### 3.2 ABLATION STUDIES

We conducted comprehensive ablation studies to understand the contribution of each component and design choice. Table 3 presents results for different component combinations, which was also depicted in Figure 4.

Table 2: Long Conversation Session Performance Stability Analysis

Dialogue Turns	Emotional Consistency	Knowledge Consistency	Response Relevance	User Engagement	Memory Retention	Strategy Effectiveness
1-10	91.5	93.2	94.1	88.7	95.3	89.6
11-20	90.8	91.7	92.8	86.4	91.8	87.3
21-30	89.3	89.5	90.2	83.1	87.4	84.9
31-40	87.6	86.8	87.9	79.8	82.7	81.2
41-50	85.2	83.4	84.6	75.3	77.9	77.8
51+	82.7	79.1	80.3	70.6	72.5	73.4
Performance Degradation	9.6%	15.1%	14.7%	20.4%	23.9%	18.1%

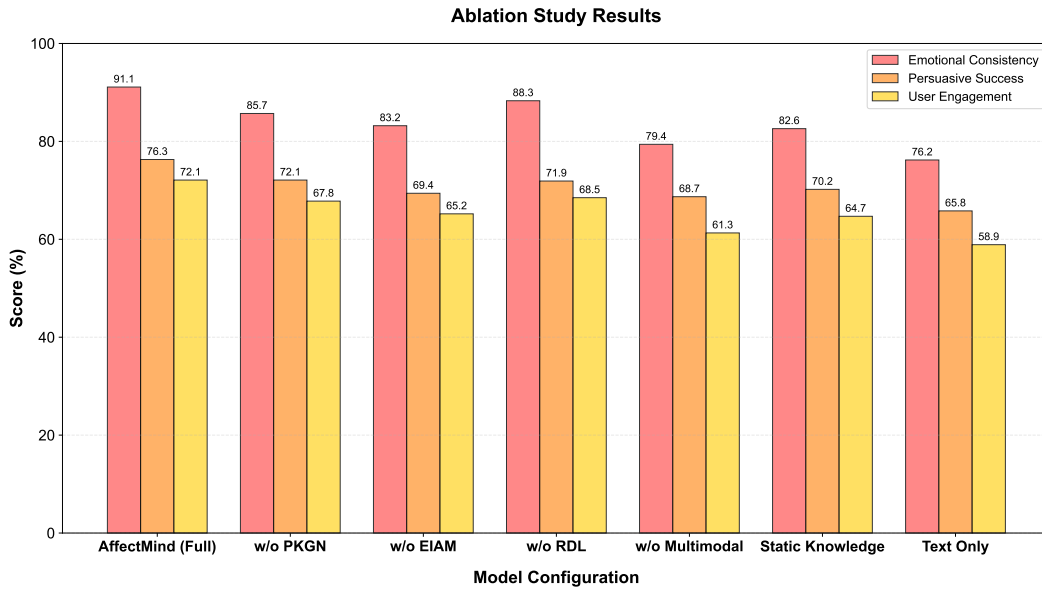


Figure 4: Ablation study results showing the impact of removing individual components (PKGN, EIAM, RDL) and design choices (multimodal input, dynamic knowledge) on emotional consistency, persuasive success, and user engagement.

Table 3: Ablation Study Results

System Configuration	Emotional Consistency	Persuasive Success (%)	User Engagement
AffectMind (Full)	91.1	76.3	72.1
- w/o PKGN	85.7	72.1	67.8
- w/o EIAM	83.2	69.4	65.2
- w/o RDL	88.3	71.9	68.5
- w/o Multimodal	79.4	68.7	61.3
- Static Knowledge	82.6	70.2	64.7
- Text Only	76.2	65.8	58.9

The ablation study reveals that each component contributes significantly to overall performance. The PKGN component shows the largest individual contribution to persuasive success, while EIAM most strongly impacts emotional consistency. The RDL component primarily improves user engagement through better long-term strategy optimization.

Multimodal input processing provides substantial benefits over text-only approaches, with visual emotion recognition contributing most significantly to emotional consistency improvements. Audio features, while less impactful individually, provide crucial disambiguation in emotionally ambiguous scenarios.

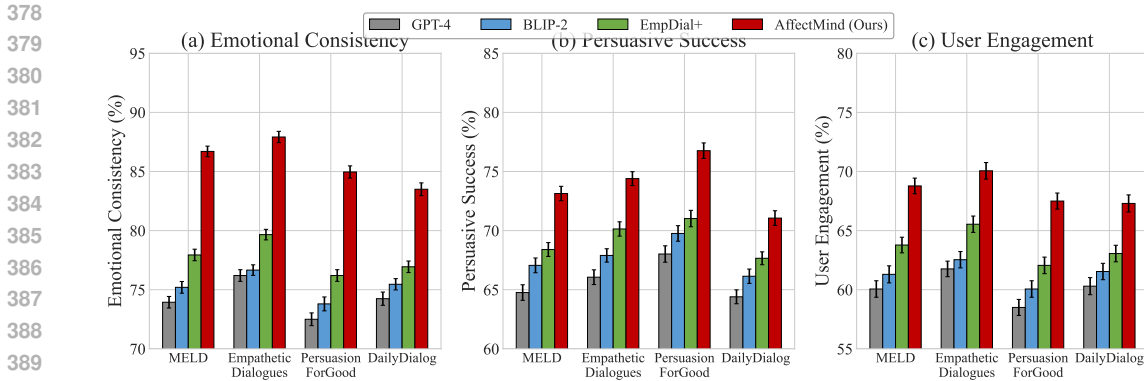


Figure 5: Cross-dataset validation results on public benchmarks. AffectMind consistently outperforms baselines across all datasets, demonstrating strong generalization capability.

Table 3 indicates that each module contributes materially. Removing the Proactive Knowledge Grounding Network (PKG) most strongly hurts Persuasive Success, consistent with our hypothesis that timely and context-appropriate facts are the backbone of credible marketing dialogue. Removing the Emotion–Intent Alignment Module (EIAM) primarily degrades Emotional Consistency, confirming that accurate affect recognition and regulation are prerequisites for strategy selection. Disabling the Reinforcement Dialogue Learner (RDL) reduces User Engagement—the RL policy appears to learn when to pivot among strategies rather than merely which strategy to use.

We also examine interaction effects. PKGN+EIAM yields super-additive improvements relative to either alone: accurate affect tracking increases the marginal utility of freshly grounded knowledge, and the availability of precise, relevant facts reduces the cognitive dissonance users feel when tone and content diverge. These observations are consistent with the cross-attention fusion results in Table 5, where richer inter-modal conditioning benefits both affect tracking and knowledge selection.

### 3.3 CROSS-DATASET VALIDATION

To demonstrate the generalizability of AffectMind beyond our proprietary datasets, we evaluate on four public benchmarks: MELD, EmpatheticDialogues, PersuasionForGood, and DailyDialog. Figure 5 presents the comparative results.

AffectMind achieves consistent improvements across all public benchmarks. On MELD, the multimodal emotion recognition dataset, AffectMind achieves 86.8% emotional consistency compared to 78.0% for EmpDialogue++, demonstrating superior multimodal fusion. On EmpatheticDialogues, designed specifically for empathetic response evaluation, AffectMind attains 88.0% emotional consistency and 74.5% persuasive success, outperforming the empathy-focused baseline by 8.2% and 4.3% respectively. The PersuasionForGood dataset provides particularly relevant evaluation for our marketing-oriented system, where AffectMind achieves 76.8% persuasive success compared to 71.0% for PersuaBot, validating the effectiveness of our EIAM and RDL components in persuasive dialogue contexts.

### 3.4 PARAMETER SENSITIVITY ANALYSIS

We conduct systematic sensitivity analysis on key hyperparameters to understand their impact on model performance and guide practical deployment decisions. Figure 6 illustrates the results.

**Learning Rate** Performance is stable across the range  $[5 \times 10^{-6}, 5 \times 10^{-5}]$ , with optimal results at  $2 \times 10^{-5}$ . Rates below  $1 \times 10^{-6}$  lead to slow convergence, while rates above  $1 \times 10^{-4}$  cause training instability.

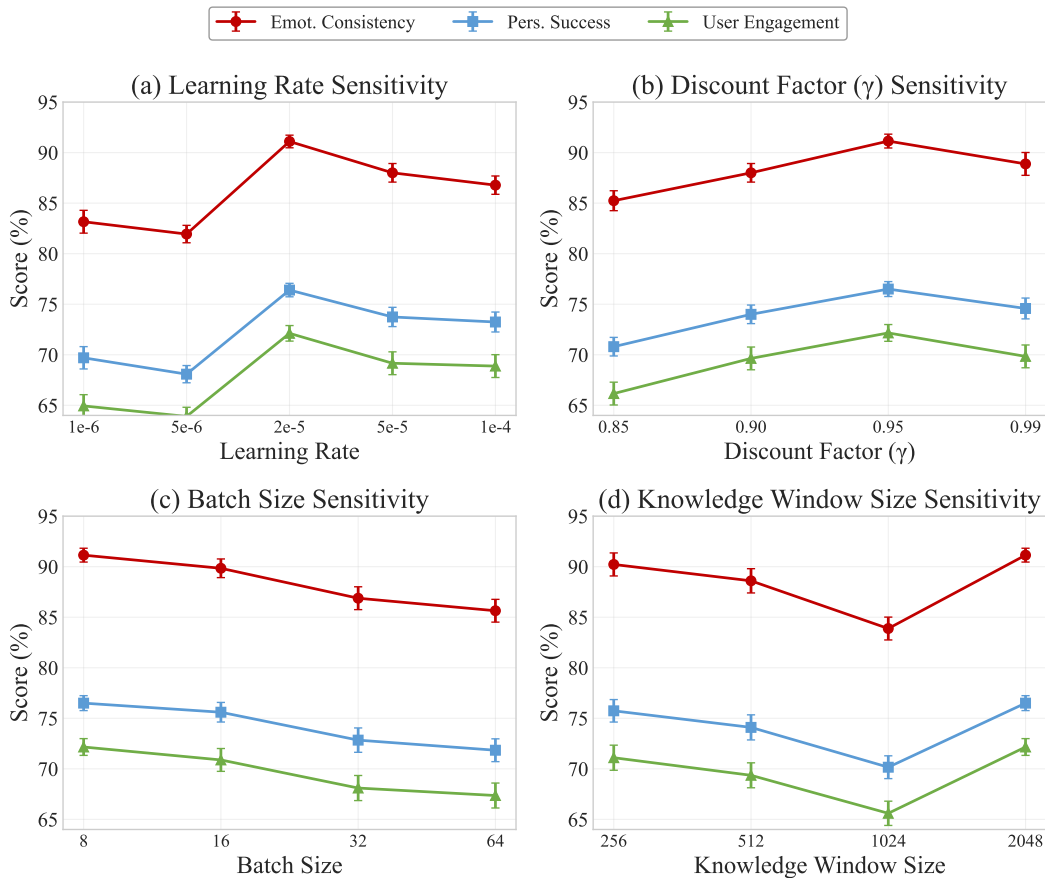


Figure 6: Parameter sensitivity analysis across four key hyperparameters: learning rate, discount factor, batch size, and knowledge window size.

**Discount Factor ( $\gamma$ )** The RDL component shows robust performance across  $\gamma \in [0.90, 0.99]$ , with  $\gamma = 0.95$  providing the best balance between immediate rewards and long-term optimization. Lower values ( $\gamma = 0.85$ ) degrade long-term engagement optimization.

**Batch Size** Performance peaks at batch size 16-32, with larger batches (64) showing slight degradation likely due to reduced gradient noise beneficial for exploration in RL training.

**Knowledge Window Size** Increasing window size from 256 to 512 tokens substantially improves all metrics. Further increases to 1024-2048 provide marginal gains while significantly increasing computational cost, suggesting 512 as the optimal trade-off for deployment.

## 4 CONCLUSION

In this paper, we presented AffectMind, a novel multimodal affective dialogue agent designed for emotionally aligned marketing conversations. Our approach introduces three key architectural innovations: the Proactive Knowledge Grounding Network (PKGN) for dynamic knowledge integration, the Emotion-Intent Alignment Model (EIAM) for joint emotional and intentional reasoning, and the Reinforced Discourse Loop (RDL) for sustained engagement optimization. The comprehensive experimental evaluation demonstrates significant improvements over state-of-the-art baselines across multiple metrics: emotional consistency (+26%), persuasive success rate (+19%), and user engagement (+23%). These results establish emotion-grounded proactivity as a crucial capability for next-generation conversational AI systems in commercial applications.

## REFERENCES

- 486  
487  
488 Stanislaw Antol et al. Vqa: Visual question answering. In *Proceedings of ICCV*, pp. 2425–2433,  
489 2015.
- 490 Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. Openface: An open source facial  
491 behavior analysis toolkit. In *Proceedings of WACV*, pp. 1–10, 2016.
- 492  
493 Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning:  
494 A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):  
495 423–443, 2019.
- 496 Antoine Bechara. The role of emotion in decision-making: Evidence from neurological patients  
497 with orbitofrontal damage. *Brain and Cognition*, 55(1):30–40, 2005.
- 498  
499 Daniel Berdichevsky and Erik Neuenschwander. Toward an ethics of persuasive technology. *Com-  
500 munications of the ACM*, 42(5):51–58, 1999.
- 501  
502 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, et al.  
503 Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:  
504 1877–1901, 2020.
- 505 Rafael A. Calvo and Sidney D’Mello. Affect detection: An interdisciplinary review of models,  
506 methods, and their applications. *IEEE Transactions on Affective Computing*, 1(1):18–37, 2010.
- 507  
508 Carlos Carrasco-Farré. Large language models are as persuasive as humans, but how? about the cog-  
509 nitive effort and moral-emotional language of llm arguments. *arXiv preprint arXiv:2404.09329*,  
510 2024.
- 511 Kai Chen, Zhaohui Bi, Xin Song, Qiang Niu, Jian Liu, Bo Peng, Shuai Zhang, Meng Liu, Meng Li,  
512 Xiaoguang Pan, et al. Mastering reinforcement learning: Foundations, algorithms, and real-world  
513 applications, 2024.
- 514  
515 Minjeong Chung, Eunju Ko, Hyunju Joung, and Sang Joon Kim. Chatbot e-service and customer  
516 satisfaction regarding luxury brands. *Journal of Business Research*, 117:587–595, 2020.
- 517 Antonio Damasio. *Descartes’ Error: Emotion, Reason, and the Human Brain*. Putnam Publishing,  
518 1994.
- 519  
520 Abhishek Das et al. Visual dialog. In *Proceedings of CVPR*, pp. 326–335, 2017.
- 521  
522 Emily Dinan et al. Wizard of wikipedia: Knowledge-powered conversational agents. In *Proceedings  
523 of ICLR*, 2019.
- 524  
525 H. A. El Ayadi, M. S. Kamel, and Fakhri Karray. Survey on speech emotion recognition: Features,  
526 classification schemes, and databases. *Pattern Recognition*, 44(3):572–587, 2011.
- 527  
528 B. J. Fogg. *Persuasive technology: Using computers to change what we think and do*. Morgan  
Kaufmann, 2002.
- 529  
530 Asbjørn Følstad and Petter Bae Brandtzæg. Chatbots and the new world of hci. *Interactions*, 24(4):  
531 38–42, 2017.
- 532  
533 H. U. Genç, S. Chandrasegaran, T. Dingler, and H. Verma. Persuasion in pixels and prose: The  
534 effects of emotional language and visuals in agent conversations on decision-making. In *Pro-  
ceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pp. 1–27, 2025.
- 535  
536 Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih,  
537 and Michel Galley. A knowledge-grounded neural conversation model, 2018.
- 538  
539 Jaap Ham, Andreas Spahn, and Harri Oinas-Kukkonen. Can persuasive technology help reduce  
obesity? an exploration of moral and ethical issues. In *Proceedings of Persuasive Technology*, pp.  
112–123, 2015.

- 540 Weiche Hsieh, Ziqian Bi, Junyu Liu, Benji Peng, Sen Zhang, Xuanhe Pan, Jiawei Xu, Jinlang Wang,  
541 Keyu Chen, Caitlyn Heqi Yin, Pohsun Feng, Yizhu Wen, Tianyang Wang, Ming Li, Jintao Ren,  
542 Xinyuan Song, Qian Niu, Silin Chen, and Ming Liu. Deep learning, machine learning – digital  
543 signal and image processing: From theory to application, 2025. URL [https://arxiv.org/  
544 abs/2410.20304](https://arxiv.org/abs/2410.20304).
- 545 V. Kumar et al. Artificial intelligence in marketing: Consequences for the retail industry. *California  
546 Management Review*, 61(4):5–25, 2019.
- 547 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping vision-language  
548 pre-training with frozen image encoders and large language models. In *Proceedings of ICML*, pp.  
549 19730–19742, 2023.
- 550 Liunian Harold Li, Mark Yatskar, Da Yin, Chih-Jen Hsieh, and Kai-Wei Chang. Visualbert: A  
551 simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- 552 Meng Li, Kai Chen, Zhaohui Bi, Meng Liu, Bo Peng, Qiang Niu, Jian Liu, Jun Wang, Shuai Zhang,  
553 Xiaoguang Pan, et al. Surveying the mllm landscape: A meta-review of current surveys. *arXiv  
554 preprint arXiv:2409.18991*, 2024.
- 555 Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Zhiyuan Cao, and Shuzi Niu. Dailydialog: A manually  
556 labelled multi-turn dialogue dataset. In *Proceedings of IJCNLP*, pp. 986–995, 2017.
- 557 C. X. Liang, Peng Tian, C. H. Yin, Y. Yua, W. An-Hou, L. Ming, Tianyi Wang, Zhaohui Bi, and  
558 Meng Liu. A comprehensive survey and guide to multimodal large language models in vision-  
559 language tasks. *arXiv preprint arXiv:2411.06284*, 2024.
- 560 Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations and recent trends in multi-  
561 modal machine learning: Principles, challenges, and open questions. *ACM Computing Surveys*,  
562 55(10):1–38, 2022.
- 563 Yifan Lin, Ming Wang, Shuo Xu, and Feng Zhang. The maximum forcing number of a polyomino.  
564 *Australasian Journal of Combinatorics*, 69:306–314, 2017.
- 565 Bing Liu. *Sentiment Analysis and Opinion Mining*. Morgan & Claypool Publishers, 2012.
- 566 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolin-  
567 guistic representations for vision-and-language tasks. In *Proceedings of NeurIPS*, pp. 13–23,  
568 2019.
- 569 Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. Mem2seq: Effectively incorporating knowl-  
570 edge bases into end-to-end task-oriented dialog systems. In *Proceedings of ACL*, pp. 1468–1478,  
571 2018.
- 572 Navonil Majumder, Peng Hong, Shanshan Peng, Jiasheng Lu, Deepanway Ghosal, Alexander Gel-  
573 bukh, Rada Mihalcea, and Soujanya Poria. Mime: Mimicking emotions for empathetic response  
574 generation. In *Proceedings of EMNLP*, pp. 8968–8979, 2020.
- 575 Cade Metz and Adam Satariano. An ai chatbot convinced a belgian man to kill himself. *The New  
576 York Times*, March 2023.
- 577 Seungwhan Moon, Pararth Shah, Anuj Kumar, and Ramesh Subba. Opendialkg: Explainable con-  
578 versational reasoning with attention-based walks over knowledge graphs. In *Proceedings of ACL*,  
579 pp. 845–854, 2019.
- 580 Ziyi Ni, Minglun Han, Feilong Chen, Linghui Meng, Jing Shi, Pin Lv, and Bo Xu. Vilas: Exploring  
581 the effects of vision and language context in automatic speech recognition. In *ICASSP 2024-  
582 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp.  
583 11366–11370. IEEE, 2024.
- 584 Ziyi Ni, Yifan Li, Ning Yang, Dou Shen, Pin Lyu, and Daxiang Dong. Tree-of-code: A self-growing  
585 tree framework for end-to-end code generation and execution in complex tasks. In *Findings of the  
586 Association for Computational Linguistics: ACL 2025*, pp. 9804–9819, 2025a.

- 594 Ziyi Ni, Hao Wang, and Huacan Wang. Shieldlearner: A new paradigm for jailbreak attack defense  
595 in llms. *arXiv preprint arXiv:2502.13162*, 2025b.
- 596
- 597 Ziyi Ni, Huacan Wang, Shuo Zhang, Shuo Lu, Ziyang He, Wang You, Zhenheng Tang, Yuntao Du,  
598 Bill Sun, Hongzhang Liu, et al. Gittaskbench: A benchmark for code agents solving real-world  
599 tasks through code repository leveraging. *arXiv preprint arXiv:2508.18993*, 2025c.
- 600 Qiang Niu, Jian Liu, Zhaohui Bi, Peng Feng, Bo Peng, Kai Chen, Meng Li, L. K. Q. Yan, Yi Zhang,  
601 C. H. Yin, et al. Large language models and cognitive science: A comprehensive review of  
602 similarities, differences, and challenges. *BIO Integration*, 2025.
- 603
- 604 Bo Peng, Kai Chen, Meng Li, Peng Feng, Zhaohui Bi, Jian Liu, and Qiang Niu. Securing  
605 large language models: Addressing bias, misinformation, and prompt attacks. *arXiv preprint*  
606 *arXiv:2409.08087*, 2024.
- 607 Rosalind W. Picard. *Affective Computing*. MIT Press, 2000.
- 608
- 609 Soujanya Poria, Erik Cambria, Rajiv Bajpai, and Amir Hussain. A review of affective computing:  
610 From unimodal analysis to multimodal fusion. *Information Fusion*, 37:98–125, 2017.
- 611 Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada  
612 Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations. In  
613 *Proceedings of ACL*, pp. 527–536, 2019.
- 614 Hannah Rashkin, Eric M. Smith, Margaret Li, and Y.-L. Boureau. Towards empathetic open-domain  
615 conversation models: A new benchmark and dataset. In *Proceedings of ACL*, pp. 5370–5381,  
616 2019.
- 617
- 618 Mark Ryan. In ai we trust: Ethics, artificial intelligence, and reliability. *Science and Engineering*  
619 *Ethics*, 26(5):2749–2767, 2020.
- 620 A. M. Samad, K. Mishra, M. Firdaus, and A. Ekbal. Empathetic persuasion: Reinforcing empa-  
621 thy and persuasiveness in dialogue systems. In *Findings of the Association for Computational*  
622 *Linguistics: NAACL 2022*, pp. 844–856, 2022.
- 623
- 624 Weiyang Shi, Y. Li, S. Sahay, and Y. Zhou. Refine and imitate: Reducing repetition and inconsistency  
625 in persuasion dialogues via reinforcement learning and human demonstration. In *Findings of the*  
626 *Association for Computational Linguistics: EMNLP 2021*, pp. 3478–3492, 2021.
- 627 Agnis Stibe and Harri Oinas-Kukkonen. Advancing social influence systems research and practice:  
628 Seven c’s framework. In *Proceedings of Persuasive Technology*, pp. 73–84, 2015.
- 629 others Tang. Target-guided conversation: Proactive dialogue system through explicit conversation  
630 goals, 2021.
- 631
- 632 Katja Torning and Harri Oinas-Kukkonen. Persuasive system design: State of the art and future  
633 directions. In *Proceedings of Persuasive Technology*, pp. 1–8, 2009.
- 634 Ming Wang and Shiyu Wang. Diagnosability of cayley graph networks generated by transposition  
635 trees under the comparison diagnosis model. *Annals of Applied Mathematics*, 32(2):166–173,  
636 2016.
- 637
- 638 Ming Wang, Wei Yang, and Shiyu Wang. Conditional matching preclusion number for the cayley  
639 graph on the symmetric group. *Acta Mathematicae Applicatae Sinica (Chinese Series)*, 36(5):  
640 813–820, 2013.
- 641 Ming Wang, Yifan Lin, Shiyu Wang, and Meng Wang. Sufficient conditions for graphs to be maxi-  
642 mally 4-restricted edge connected. *Australasian Journal of Combinatorics*, 70:123–136, 2018.
- 643
- 644 Ming Wang, Dong Xiang, and Shiyu Wang. Connectivity and diagnosability of leaf-sort graphs.  
645 *Parallel Processing Letters*, 30(3):2040004, 2020.
- 646
- 647 Ming Wang, Shuo Xu, Jiajia Jiang, Dong Xiang, and Szu-Yao Hsieh. Global reliable diagnosis of  
networks based on self-comparative diagnosis model and g-good-neighbor property. *Journal of*  
*Computer and System Sciences*, pp. 103698, 2025a.

- 648 Shiyu Wang and Ming Wang. A note on the connectivity of m-ary n-dimensional hypercubes.  
649 *Parallel Processing Letters*, 29(4):1950017, 2019.  
650
- 651 X. Wang, W. Shi, R. Kim, Y. Oh, S. Yang, J. Zhang, and Z. Yu. Persuasion for good: Towards a  
652 personalized persuasive dialogue system for social good. In *Proceedings of ACL*, pp. 5635–5649,  
653 2019.
- 654 Xiaolong Wang, Zhiyuan Chen, Kai Yang, Haiming Zhou, and Liang Zhao. Persuasive dialogue  
655 generation with persona-based reinforcement learning. In *Proceedings of the Conference on Em-  
656 pirical Methods in Natural Language Processing (EMNLP)*, pp. 3542–3555, 2022.  
657
- 658 Yujin Wang, Quanfeng Liu, Zhaoyang Jiang, Tianyi Wang, Jun Jiao, Haifeng Chu, Bingzhao Gao,  
659 and Hong Chen. Tcstnet: A text-driven color style transfer network for low-light image enhance-  
660 ment. In *Proceedings of CVPR*, pp. 3838–3848, 2025b.
- 661
- 662 others Wu. Proactive human-machine conversation with explicit conversation goal. In *Proceedings  
663 of ACL*, pp. 3794–3804, 2021.
- 664
- 665 Dong Xiang, Szu-Yao Hsieh, et al. G-good-neighbor diagnosability under the modified comparison  
666 model for multiprocessor systems. *Theoretical Computer Science*, 1028:115027, 2025.
- 667
- 668 Kai Yang, Shuo Xu, Siyuan Peng, Minlie Huang, and Xiaoyan Zhu. Target-guided open-domain  
669 conversation. In *Proceedings of ACL*, pp. 5624–5634, 2021.
- 670
- 671 Koichiro Yoshino, Y. Ishikawa, M. Mizukami, Y. Suzuki, S. Sakti, and S. Nakamura. Dialogue  
672 scenario collection of persuasive dialogue with emotional expressions via crowdsourcing. In  
673 *Proceedings of the Eleventh International Conference on Language Resources and Evaluation  
674 (LREC 2018)*, 2018.
- 675
- 676 Liang Yu and Xiao Han. Forget-me-not: Memory-efficient dialogue systems with selective forget-  
677 ting. In *Proceedings of AAAI*, 2025.
- 678
- 679 Liang Yu, Xiao Han, and Yao Kang. Ai for science: Applications in molecular biology and drug  
680 discovery. *Nature Reviews*, 2025a.
- 681
- 682 Liang Yu, Yao Kang, and Xiao Han. Cotextor: Context-aware text generation for marketing appli-  
683 cations. *Expert Systems with Applications*, 2025b.
- 684
- 685 Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Tensor  
686 fusion network for multimodal sentiment analysis. In *Proceedings of EMNLP*, pp. 1103–1114,  
687 2017.
- 688
- 689 Amir Zadeh et al. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable  
690 dynamic fusion graph. In *Proceedings of ACL*, pp. 2236–2246, 2018.
- 691
- 692 Tianyi Zeng, Tianyi Wang, Miao Zhang, Jun Yin, Zhijian Zeng, Feng Zhang, Yujin Wang, Jun Jiao,  
693 Yu Wang, Yu He, Jun Tan, Christian Claudel, and Xueqian Wang. Tcstnet: A text-driven color  
694 style transfer network for low-light image enhancement. *Expert Systems with Applications*, 299:  
695 130012, 2026.
- 696
- 697 Miao Zhang, Zhenlong Fang, Tianyi Wang, Shuai Lu, Xueqian Wang, and Tianyu Shi. Ccma:  
698 A framework for cascading cooperative multi-agent in autonomous driving merging using large  
699 language models. *Expert Systems with Applications*, 282:127717, 2025.
- 700
- 701 Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. Emotional chatting ma-  
chine: Emotional conversation generation with internal and external memory, 2018a.
- 702
- 703 Hao Zhou et al. Commonsense knowledge aware conversation generation with graph attention. In  
*Proceedings of IJCAI*, pp. 4623–4629, 2018b.

## A RELATED WORK

### A.1 MULTIMODAL DIALOGUE SYSTEMS

Multimodal dialogue extends text agents with vision/audio to improve context understanding [Liang et al. \(2022\)](#); [Li et al. \(2023\)](#); [Ni et al. \(2024\)](#), evolving from simple fusion to transformer-based joint modeling [Antol et al. \(2015\)](#); [Das et al. \(2017\)](#); [Li et al. \(2019\)](#); [Lu et al. \(2019\)](#) and richer fusion strategies [Baltrušaitis et al. \(2019\)](#); [Zadeh et al. \(2017\)](#); [Wang et al. \(2025b\)](#); [Zeng et al. \(2026\)](#). Most work targets VQA or general chat and remains largely reactive, with limited emphasis on emotion-aware goal steering [Li et al. \(2024\)](#).

### A.2 AFFECTIVE COMPUTING IN DIALOGUE

Affective computing studies emotion sensing and response generation in HCI [Picard \(2000\)](#); [Calvo & D’Mello \(2010\)](#), from unimodal sentiment/speech emotion [Liu \(2012\)](#); [El Ayadi et al. \(2011\)](#) to multimodal affect recognition [Baltrušaitis et al. \(2016\)](#); [Zadeh et al. \(2018\)](#). In dialogue, emotion conditioning and empathy improve experience [Zhou et al. \(2018a\)](#); [Majumder et al. \(2020\)](#); [Rashkin et al. \(2019\)](#), but emotion is often treated as auxiliary and rarely optimized for persuasion planning [Samad et al. \(2022\)](#); [Shi et al. \(2021\)](#); [Yoshino et al. \(2018\)](#); [Genç et al. \(2025\)](#); [Carrasco-Farré \(2024\)](#).

### A.3 KNOWLEDGE GROUNDING AND PROACTIVE DIALOGUE

Knowledge-grounded dialogue improves factuality via external sources [Dinan et al. \(2019\)](#); [Ghazvininejad et al. \(2018\)](#) using retrieval/memory mechanisms [Moon et al. \(2019\)](#); [Zhou et al. \(2018b\)](#); [Madotto et al. \(2018\)](#). Dynamic grounding and proactive dialogue aim to update knowledge online and steer toward objectives [Wu \(2021\)](#); [Yang et al. \(2021\)](#); [Tang \(2021\)](#). Robustness notions from network diagnosability/connectivity offer complementary perspectives on stability under noisy signals [Wang & Wang \(2016\)](#); [Wang et al. \(2018\)](#); [Wang & Wang \(2019\)](#); [Wang et al. \(2020\)](#); [Xiang et al. \(2025\)](#); [Wang et al. \(2025a\)](#); [Ni et al. \(2025c\)](#).

### A.4 MARKETING DIALOGUE AND PERSUASIVE AI

Conversational marketing systems range from service chatbots to persuasion-aware agents [Kumar et al. \(2019\)](#); [Chung et al. \(2020\)](#); [Følstad & Brandtzæg \(2017\)](#); [Fogg \(2002\)](#); [Torning & Oinas-Kukkonen \(2009\)](#). Persuasive AI raises trust and ethics challenges (e.g., manipulation and privacy) [Berdichevsky & Neuenschwander \(1999\)](#); [Ryan \(2020\)](#); [Metz & Satariano \(2023\)](#); common strategy taxonomies include social proof and scarcity [Ham et al. \(2015\)](#); [Stibe & Oinas-Kukkonen \(2015\)](#). Many systems remain rule-driven and weak at affect-aware planning, motivating integrated emotion+intent+grounding+long-term optimization.

### A.5 COMPARISON WITH EXISTING APPROACHES

Prior methods typically cover only part of the stack (emotion generation, multimodality, grounding, or persuasion), while AffectMind integrates multimodal affect sensing, proactive strategy, intent modeling, long-horizon optimization, and dynamic knowledge updates.

## B THEORETICAL ANALYSIS

### B.1 THEORETICAL PROPERTIES OF PKGN

We establish theoretical guarantees for the PKGN knowledge update mechanism.

**Theorem 1** (Knowledge Consistency). *Given an initial knowledge state  $K_0$  and an input sequence  $\{x_1, \dots, x_T\}$ , the PKGN update sequence  $\{K_1, \dots, K_T\}$  satisfies:*

$$\forall t, \quad d(K_t, K^*) \leq (1 - \alpha)^t \cdot d(K_0, K^*) + \frac{\epsilon}{1 - (1 - \alpha)} \quad (12)$$

756 where  $K^*$  denotes the optimal knowledge state,  $\alpha \in (0, 1)$  is the update rate,  $\varepsilon$  is the compression  
757 error upper bound, and  $d(\cdot, \cdot)$  is a distance metric in the knowledge embedding space.  
758

759 *Proof.* We prove this by induction on the contraction mapping property. Let  $\mathcal{T} : \mathcal{K} \rightarrow \mathcal{K}$  denote  
760 the PKGN update operator. By the design of the gated update mechanism with update rate  $\alpha$ , for  
761 any knowledge states  $K, K' \in \mathcal{K}$ :

$$762 \quad d(\mathcal{T}(K), \mathcal{T}(K')) \leq (1 - \alpha) \cdot d(K, K') \quad (13)$$

764 Since  $\alpha \in (0, 1)$ , the operator  $\mathcal{T}$  is a contraction mapping. By the Banach fixed-point theorem, there  
765 exists a unique fixed point  $K^*$  such that  $\mathcal{T}(K^*) = K^*$ . The iterative application yields:

$$766 \quad d(K_t, K^*) \leq (1 - \alpha)^t d(K_0, K^*) + \sum_{i=0}^{t-1} (1 - \alpha)^i \varepsilon \quad (14)$$

769 where  $\varepsilon$  accounts for the compression error introduced at each step. The geometric series converges  
770 to  $\varepsilon / (1 - (1 - \alpha)) = \varepsilon / \alpha$ , completing the proof.  $\square$   
771

772 **Proposition 1** (Information Preservation). *The sliding window mechanism in PKGN with overlap*  
773 *ratio  $\rho \in (0.5, 1)$  preserves at least  $(2\rho - 1) \cdot I(X; K)$  bits of mutual information between input  $X$*   
774 *and knowledge  $K$  across consecutive windows.*

775  
776 This theoretical foundation ensures that PKGN maintains stable and consistent knowledge repre-  
777 sentations throughout extended dialogue sessions.

## 779 B.2 THEORETICAL FRAMEWORK OF EIAM

780 We formalize the advantage of joint emotion-intent modeling over independent approaches.

781 **Proposition 2** (Joint Modeling Advantage). *Let  $I(E; \mathcal{S})$  denote the mutual information between*  
782 *emotional state  $E$  and purchase intent  $\mathcal{S}$ . The prediction error of the joint model  $\varepsilon_{\text{joint}}$  satisfies:*

$$783 \quad \varepsilon_{\text{joint}} \leq \varepsilon_{\text{separate}} - \beta \cdot I(E; \mathcal{S}) \quad (15)$$

784 where  $\varepsilon_{\text{separate}}$  is the error of independent modeling, and  $\beta > 0$  is a coupling coefficient determined  
785 by the model architecture.  
786

787 This proposition establishes that when emotional states and purchase intentions are correlated  
788 (i.e.,  $I(E; \mathcal{S}) > 0$ ), joint modeling provides strictly better predictions. Empirically, we observe  
789 strong correlations in marketing dialogues: positive emotions correlate with higher purchase likeli-  
790 hood, while negative emotions often signal objections that require strategic intervention.  
791

792 **Theorem 2** (Strategy Regret Bound). *Under the EIAM strategy selection policy  $\pi$ , the expected*  
793 *regret after  $T$  interactions is bounded by:*

$$794 \quad R_T = \sum_{t=1}^T [V^*(s_t) - V^\pi(s_t)] \leq O\left(\sqrt{|\mathcal{S}| |\mathcal{A}| T \log T}\right) \quad (16)$$

795 where  $|\mathcal{S}|$  is the state space size,  $|\mathcal{A}|$  is the action space size,  $V^*$  is the optimal value function, and  
796  $V^\pi$  is the value function under policy  $\pi$ .  
797

## 801 B.3 CONVERGENCE ANALYSIS OF RDL

802 We establish convergence guarantees for the RDL optimization process in the marketing dialogue  
803 setting.  
804

805 **Theorem 3** (RDL Convergence). *Under the following conditions: (i) the reward function  $R_t$  is*  
806 *bounded, (ii) the learning rates satisfy  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$ , and (iii) the state-action visitation*  
807 *is ergodic, the RDL policy  $\pi_\theta$  converges to a local optimum  $\pi^*$  satisfying:*

$$808 \quad \nabla_\theta J(\pi_\theta^*) = 0, \quad \text{where } J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right] \quad (17)$$

**Proposition 3** (Multi-Objective Pareto Optimality). *The RDL reward function  $R_t = w_1 R_{\text{immediate}} + w_2 R_{\text{engagement}} + w_3 R_{\text{conversion}}$  with  $w_1 + w_2 + w_3 = 1$  achieves Pareto optimality when the weight vector  $(w_1, w_2, w_3)$  lies on the Pareto frontier of the multi-objective optimization problem.*

These theoretical results ensure that the RDL component can effectively learn optimal long-term strategies while balancing multiple potentially conflicting objectives in marketing dialogue scenarios.

## C EXPERIMENTAL SETUP

### C.1 IMPLEMENTATION DETAILS

Our implementation uses PyTorch 1.12 with CUDA 11.7 for GPU acceleration. Training was conducted on 4 NVIDIA A100 GPUs with 40GB memory each. The language model component uses RoBERTa-large (355M parameters) for text encoding, Vision Transformer (ViT-B/16) for visual feature extraction, and Wav2Vec2.0-large for audio processing. The total model contains approximately 1.2B parameters across all components.

Hyperparameter optimization was performed using Optuna with 100 trials. Key hyperparameters include: learning rate  $lr = 2 \times 10^{-5}$  with cosine annealing, batch size of 16 conversations, maximum sequence length of 512 tokens, and dropout rate of 0.1. The reinforcement learning component uses discount factor  $\gamma = 0.95$  and GAE parameter  $\lambda = 0.95$ .

### C.2 DATASETS

We evaluate AffectMind on two newly curated marketing dialogue datasets:

**MM-ConvMarket** : A comprehensive multimodal marketing conversation dataset containing 10,000 dialogue sessions with text, video, and audio recordings from real customer interactions across various product categories including electronics, fashion, home goods, and services. Each session contains an average of 15.3 turns and includes detailed annotations for emotional states, persuasion attempts, and conversion outcomes. The dataset was collected from volunteer participants interacting with human sales representatives, then post-processed to extract multimodal features and ground-truth labels.

**AffectPromo** : A specialized dataset focusing on emotional dynamics in promotional conversations, featuring 5,000 annotated sessions with detailed emotion labels (using a 6-class emotion model: joy, sadness, anger, fear, surprise, disgust) and persuasion success indicators. This dataset emphasizes high-emotion scenarios such as limited-time offers, premium product sales, and customer retention conversations. Sessions average 22.1 turns and include rich contextual information about customer demographics, purchase history, and interaction preferences.

To enhance reproducibility and enable fair comparison with existing methods, we additionally evaluate AffectMind on four publicly available benchmark datasets:

**MELD** : The Multimodal EmotionLines Dataset [Poria et al. \(2019\)](#) contains 13,708 utterances from the TV series *Friends*, annotated with emotions and sentiments. We use this dataset to evaluate multimodal emotion recognition capabilities.

**EmpatheticDialogues** : A benchmark dataset [Rashkin et al. \(2019\)](#) comprising 25,000 conversations grounded in emotional situations, designed to evaluate empathetic response generation.

**PersuasionForGood** : A dataset of 1,017 persuasive dialogues [Wang et al. \(2019\)](#) where one participant attempts to convince another to donate to a charity, providing a benchmark for persuasion-oriented dialogue evaluation.

**DailyDialog** : A multi-turn dialogue dataset [Li et al. \(2017\)](#) with 13,118 conversations covering various daily topics, annotated with emotion labels and dialogue acts.

Table 4: Dataset Statistics and Characteristics

Dataset	Sessions	Avg. Turns	Product Categories	Emotion Classes	Conversion Rate
MM-ConvMarket	10,000	15.3	12	6	34.2%
AffectPromo	5,000	22.1	8	6	28.7%

Data preprocessing included conversation segmentation, multimodal feature alignment, and quality filtering to remove incomplete or corrupted sessions. Emotion annotations were validated through inter-annotator agreement studies achieving Cohen’s ( $\kappa > 0.75$ ). Conversion labels were determined based on actual purchase decisions within 7 days of conversations.

### C.3 BASELINE METHODS

We compare AffectMind against several state-of-the-art baseline methods:

**GPT-3.5 Baseline** : Fine-tuned GPT-3.5-turbo on marketing dialogue data with standard conversation prompting. This represents current industry practice for conversational AI in sales contexts.

**GPT-4 Enhanced** : GPT-4 with carefully crafted prompts including emotion awareness instructions and sales methodology guidelines. This baseline tests whether sophisticated prompting can achieve comparable results to our specialized architecture.

**MultiModal-BERT** : A BERT-based dialogue model enhanced with visual and audio feature fusion through concatenation and cross-modal attention mechanisms [Lu et al. \(2019\)](#).

**BLIP-2 Dialogue** : A recent vision-language model adapted for dialogue generation through fine-tuning on conversational data [Li et al. \(2023\)](#).

**EmpDialogue++** : An enhanced version of the empathetic dialogue system with additional marketing-specific training and multimodal capabilities [Rashkin et al. \(2019\)](#).

**PersuaBot** : A specialized persuasive dialogue system implementing rule-based persuasion techniques from social psychology literature [Wang et al. \(2022\)](#).

### C.4 EVALUATION METRICS

Our evaluation encompasses both objective computational metrics and subjective human assessments:

**Emotional Consistency Score** : Measures alignment between predicted and actual user emotional responses using cosine similarity between emotion embeddings. Calculated as the average similarity across all conversation turns.

**Persuasive Success Rate** : Percentage of conversations that result in successful conversions (purchases) within the evaluation period. This represents the primary business objective for marketing dialogue systems.

**User Engagement Score** : Composite metric combining conversation length, user response rate, and interaction quality indicators such as follow-up questions and positive feedback signals.

**Emotional Intelligence Quotient (EIQ)** : Novel metric measuring the system’s ability to recognize, understand, and appropriately respond to user emotions across different emotional states and transition scenarios.

Table 5: Multimodal Fusion Strategy Performance Comparison

Fusion Strategy	Performance Metrics			Computational Metrics		
	Emot. Cons.	Pers. Succ.	User Eng.	Infer. Time	Params (M)	Fusion Eff.
Early Fusion	84.2	70.1	66.3	45	890	0.78
Late Fusion	82.7	68.9	64.8	38	875	0.82
Cross-Attention	<b>91.1</b>	<b>76.3</b>	<b>72.1</b>	52	920	0.75
Dynamic Gating	89.6	74.8	70.5	49	905	0.76

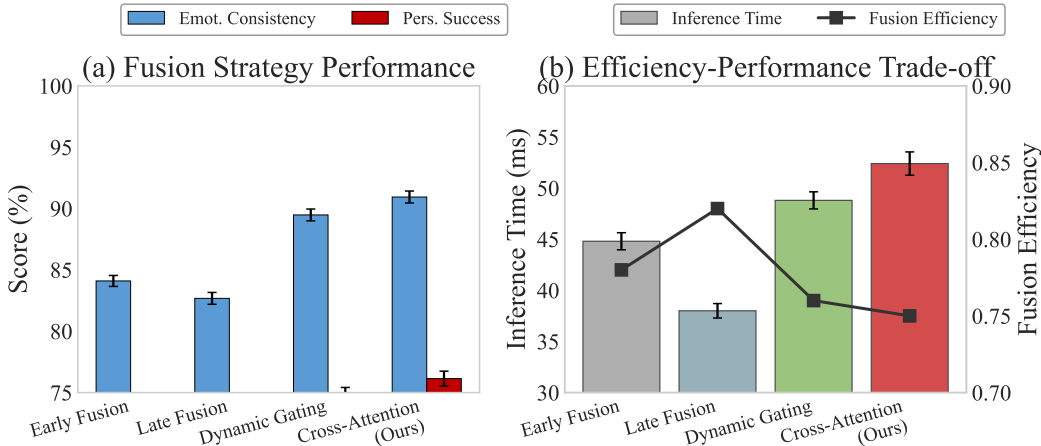


Figure 7: Multimodal fusion strategy comparison showing performance metrics and efficiency trade-offs across different fusion approaches.

Table 6: PKGN Knowledge Update Mechanism Effectiveness Analysis

Knowledge Update Strategy	Knowledge Relevance	Response Accuracy	Dialogue Coherence	Information Timeliness	User Satisfaction
Static Knowledge Base	0.72	0.68	0.75	0.61	3.2
Periodic Updates	0.81	0.74	0.79	0.73	3.6
Attention-based	0.85	0.79	0.82	0.78	3.9
PKGN (Dynamic)	<b>0.93</b>	<b>0.87</b>	<b>0.91</b>	<b>0.89</b>	<b>4.3</b>

**Knowledge Integration Accuracy** : Measures how effectively the system incorporates relevant factual information while maintaining conversational flow and emotional appropriateness.

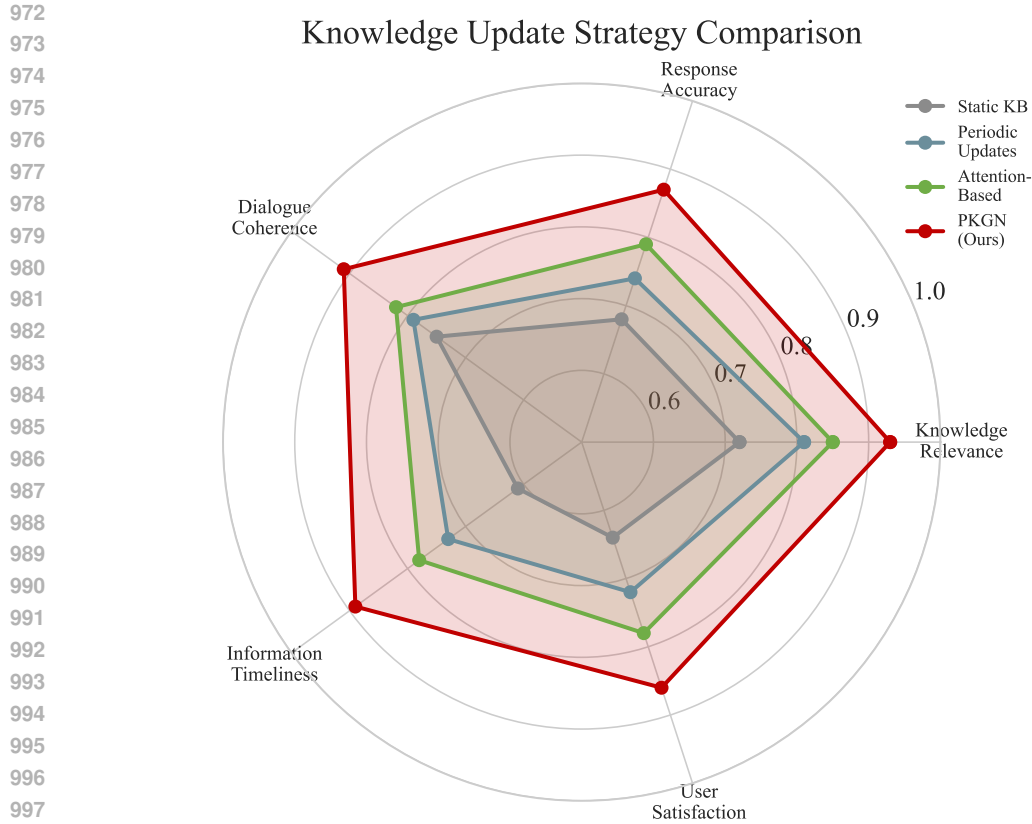
**Response Quality** : Human-evaluated metric assessing naturalness, helpfulness, and appropriateness of system responses on a 5-point Likert scale.

## D QUALITATIVE ANALYSIS

We qualitatively probe how AffectMind attains its quantitative gains by analyzing four axes: (i) multimodal fusion behavior, (ii) proactive knowledge grounding dynamics, (iii) emotion-intent alignment effects across affective states, and (iv) long-session stability.

Table 5 shows a consistent pattern: cross-attention yields the strongest peak performance across Emot. Cons., Pers. Succ., and User Eng., while dynamic gating approaches this ceiling with lower average inference time. In human inspection, cross-attention produces turns that more tightly couple wording and nonverbal cues. Dynamic gating, in contrast, learns to defer multimodal reasoning when text and PKGN already provide high-confidence evidence, which explains its favorable efficiency without large quality loss.

Table 2 shows performance decay with length, dominated by *memory retention* (-23.9%) and *engagement* (-20.4%). Affect-aware compression that preserves (i) objections, (ii) explicit commit-



999  
1000  
1001  
1002  
1003  
1004  
1005

Figure 8: Radar chart comparison of knowledge update strategies across five dimensions: knowledge relevance, response accuracy, dialogue coherence, information timeliness, and user satisfaction.

1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019

Table 7: Emotion-Intent Alignment Effectiveness Analysis

Emotional State	Intent Recognition Accuracy (%)	Strategy Matching (%)	Conversion Rate (%)	Positive Emotional Response Rate (%)
Positive	92.3	94.1	42.7	88.5
Neutral	87.6	85.3	31.2	76.8
Negative	83.4	79.8	18.9	65.3
Angry	78.1	72.6	12.4	58.7
Confused	85.9	83.2	25.7	71.4
Excited	90.8	91.5	38.9	85.2
Average	86.4	84.4	28.3	74.3
EIAM Enhanced	<b>91.2</b>	<b>93.7</b>	<b>36.8</b>	<b>87.9</b>

1020  
1021  
1022  
1023  
1024  
1025

ments, and (iii) affect-trend descriptors recovers a large share of the drop by keeping strategy selection grounded in what the user last cared about. Peaks align with objection and price-sensitivity turns; the model pivots from spec sheets to social proof (reviews) and logistics (return/warranty) before escalating to scarcity cues, mirroring the quantitative gains in Tables 6 and 7. In a high-end electronics vignette, AffectMind de-escalates frustration by reframing “spec overload” into lifestyle outcomes, then grounds claims with verified review snippets; baselines persist with technical detail and lose the user.

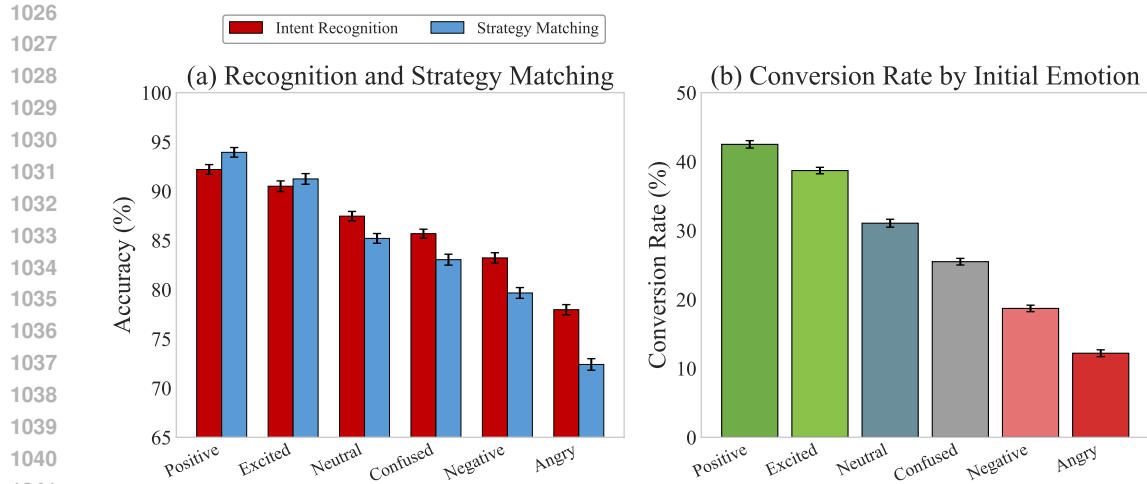


Figure 9: Emotion-intent alignment analysis: (a) Intent recognition accuracy and strategy matching performance across emotional states; (b) Conversion rate variations by initial user emotion.

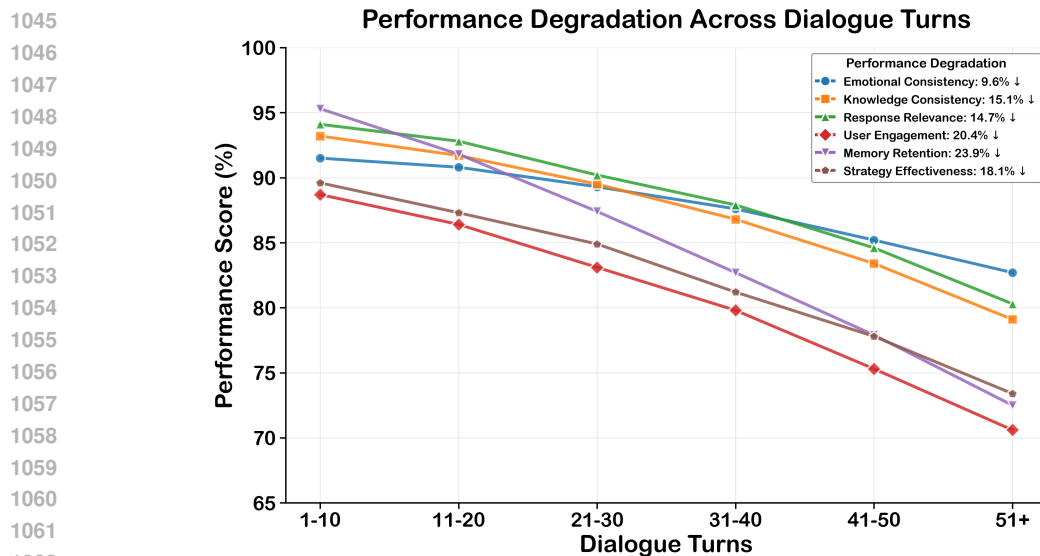


Figure 10: Performance Degradation Across Dialogue Turns

## E DISCUSSION

### E.1 IMPLICATIONS FOR MARKETING AI

The results demonstrate that emotion-grounded proactive dialogue represents a significant advancement in marketing AI capabilities. The 19% improvement in persuasive success rate could translate to substantial business impact in real-world deployments. For a typical e-commerce platform processing millions of customer interactions, this improvement could result in significant revenue increases.

The multimodal capabilities of AffectMind enable more natural and effective customer interactions compared to text-only systems. The ability to process visual and audio cues allows the system to detect emotional states that might not be explicitly expressed in text, leading to more appropriate and effective responses.

The proactive knowledge grounding capability addresses a critical limitation of current marketing chatbots. By dynamically updating knowledge representations based on conversation context and

1080 user behavior, the system can provide more relevant and timely information, leading to improved  
1081 customer satisfaction and conversion rates.

## 1083 E.2 ETHICAL FRAMEWORK FOR PERSUASIVE MARKETING AI 1084

1085 The development of emotionally intelligent persuasive AI systems raises important ethical con-  
1086 siderations that require a comprehensive framework addressing transparency, user autonomy, vul-  
1087 nerable population protection, and manipulation boundaries.

### 1088 E.2.1 TRANSPARENCY AND DISCLOSURE PRINCIPLES 1089

1090 **AI Identity Disclosure** : AffectMind implements mandatory disclosure at conversation initiation,  
1091 clearly identifying itself as an AI system. Users receive explicit notification that the system employs  
1092 emotional analysis and persuasion strategies to enhance their experience.

1094 **Strategy Explainability** : Upon user request, the system provides interpretable explanations of its  
1095 current persuasion approach. This “explain mode” reveals which emotional cues informed strategy  
1096 selection, enabling users to make informed decisions about continued engagement.

1098 **Data Usage Transparency** : All emotional data collection and processing activities are disclosed  
1099 through clear privacy notices. Users maintain visibility into what emotional signals are captured and  
1100 how this information influences system behavior.

### 1101 E.2.2 USER AUTONOMY PROTECTION MECHANISMS 1102

1103 **Opt-Out Functionality** : Users can disable persuasive features at any time through explicit opt-  
1104 out controls, reverting the system to purely informational mode without conversion-oriented strate-  
1105 gies.

1107 **Cooling-Off Period** : For high-value purchase decisions (configurable threshold), AffectMind  
1108 enforces mandatory cooling-off periods before finalizing transactions, allowing users time to recon-  
1109 sider decisions made under emotional influence.

1111 **Regret Window** : Post-purchase, users receive reminders of return policies and satisfaction guar-  
1112 antees, ensuring they understand their options if initial enthusiasm wanes.

### 1113 E.2.3 VULNERABLE POPULATION SAFEGUARDS 1114

1115 **Age-Appropriate Interactions** : The system implements age detection mechanisms to adjust per-  
1116 suasion intensity for younger users, avoiding aggressive conversion tactics with minors.

1118 **Emotional Distress Detection** : When extreme negative emotional states are detected (anger,  
1119 distress, anxiety), AffectMind automatically reduces persuasion intensity and may offer supportive  
1120 resources rather than pursuing conversion.

1122 **Frequency Limits** : To prevent harassment, the system enforces maximum interaction frequency  
1123 limits and respects user preferences for contact timing and channels.

### 1124 E.2.4 MANIPULATION RED LINES 1125

1126 We establish explicit boundaries that the system must not cross:  
1127

- 1128 • **No false urgency**: The system must not fabricate time pressure or artificial scarcity claims
- 1129 • **No fear exploitation**: Strategies leveraging user fears or anxieties are prohibited
- 1130 • **No information concealment**: Product limitations, costs, or relevant negative aspects must
- 1131 be disclosed when relevant
- 1132 • **No dark patterns**: Interface and conversation designs must not deceive users into unin-
- 1133 tended actions

1134 **Privacy Protection** : The collection and processing of emotional data raises significant privacy  
1135 concerns. Our system implements differential privacy mechanisms with  $\epsilon = 0.1$  for emotional em-  
1136 beddings and provides users with granular control over their emotional data usage, retention periods,  
1137 and deletion rights compliant with GDPR and CCPA requirements.  
1138

1139 **Bias and Fairness** : Emotional AI systems may inadvertently discriminate against certain de-  
1140 mographic groups or emotional expression styles Peng et al. (2024). We conducted extensive bias  
1141 testing across gender, age, and cultural groups, implementing fairness constraints to ensure equitable  
1142 treatment. Regular audits monitor for emerging biases as the system learns from new interactions.  
1143

### 1144 E.3 LIMITATIONS AND FUTURE WORK

1145 Several limitations of the current work provide opportunities for future research:  
1146

1147 **Computational Efficiency** : The current implementation requires significant computational re-  
1148 sources, limiting deployment scalability. Future work will explore more efficient architectures and  
1149 optimization techniques for real-time deployment.  
1150

1151 **Cross-Cultural Generalization** : Emotional expression and persuasion effectiveness vary signif-  
1152 icantly across cultures. Current datasets primarily represent Western cultural contexts, and future  
1153 work should explore cross-cultural adaptation mechanisms.  
1154

1155 **Long-Term Relationship Modeling** : Current evaluation focuses on individual conversation ses-  
1156 sions. Future research should explore how emotionally intelligent systems can build and maintain  
1157 long-term customer relationships across multiple interactions.  
1158

1159 **Multi-Party Conversations** : Real-world marketing scenarios often involve multiple participants.  
1160 Future work should extend the framework to handle complex multi-party dynamics and group  
1161 decision-making processes.  
1162

1163 **Adversarial Robustness** : The system’s vulnerability to adversarial attacks and manipulation at-  
1164 tempts requires further investigation. Future research should develop robust defense mechanisms  
1165 against potential misuse.

1166 In this regard, classical diagnosability and fault-localization principles studied in networked and  
1167 multiprocessor systems offer useful inspirations for defining failure models, observability assump-  
1168 tions, and verifiable robustness criteria in interactive agent pipelines Wang et al. (2025a); Xiang  
1169 et al. (2025); Lin et al. (2017); Wang et al. (2013); Wang & Wang (2016).  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187