

ProAct: A Dual-System Framework for Proactive Embodied Social Agents

ZEYI ZHANG*, ZIXI KANG*, and RUIJIE ZHAO*, Peking University, China
YUSEN FENG, BIAO JIANG, and LIBIN LIU†, Peking University, China

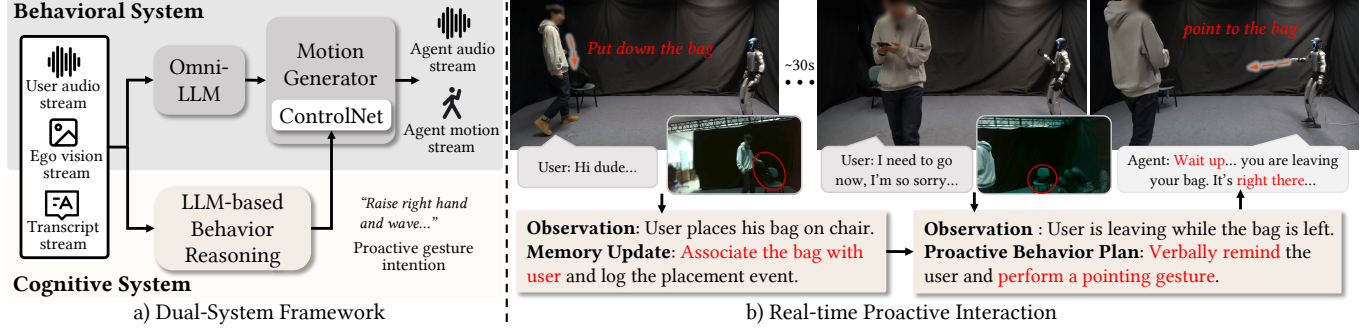


Fig. 1. ProAct enables proactive embodied social interaction. (a) Our dual-system framework integrates a fast Behavioral System for streaming multimodal generation with a Cognitive System for context-driven reasoning. (b) In a Real-time Interaction scenario, the agent tracks long-term context and autonomously initiates a synchronized verbal and gestural reminder when the user attempts to leave, demonstrating proactivity beyond simple reactivity.

Embodied social agents have recently advanced in generating synchronized speech and gestures. However, most interactive systems remain fundamentally reactive, responding only to current sensory inputs within a short temporal window. Proactive social behavior, in contrast, requires deliberation over accumulated context and intent inference, which conflicts with the strict latency budget of real-time interaction. We present *ProAct*, a dual-system framework that reconciles this time-scale conflict by decoupling a low-latency *Behavioral System* for streaming multimodal interaction from a slower *Cognitive System* which performs long-horizon social reasoning and produces high-level proactive intentions. To translate deliberative intentions into continuous non-verbal behaviors without disrupting fluency, we introduce a streaming flow-matching model conditioned on intentions via ControlNet. This mechanism supports asynchronous intention injection, enabling seamless transitions between reactive and proactive gestures within a single motion stream. We deploy ProAct on a physical humanoid robot and evaluate both motion quality and interactive effectiveness. In real-world interaction user studies, participants and observers consistently prefer ProAct over reactive variants in perceived proactivity, social presence, and overall engagement, demonstrating the benefits of dual-system proactive control for embodied social interaction. We will release the code for academic research. Our project website is available at: <https://proactrobot.github.io/>.

CCS Concepts: • **Computing methodologies** → **Animation**; *Natural language processing*; *Neural networks*.

1 Introduction

Natural social interaction is inherently bidirectional. When conversing with others, humans do not merely respond to what is said. Instead, we exercise intentionality to anticipate needs, take initiative, and act before being asked [Bandura 2001]. A host notices empty glasses and offers a refill, whilst a friend senses hesitation

and provides encouragement without prompting. This *proactivity*, the capacity to originate actions from internal goals rather than passively respond, distinguishes genuine social partners from mere responders. Moreover, natural human interaction relies on coordinated multimodal communication: we convey meaning through both verbal dialogue and non-verbal behaviors such as gestures, gaze, and body movements. For embodied social agents to achieve authentic social presence, they must similarly orchestrate these channels while demonstrating such active proactivity.

Recent progress has advanced embodied social agents along these dimensions. Multimodal systems now generate synchronized speech and gesture for virtual characters [Jiang et al. 2025; Kim et al. 2024; Wampfler et al. 2025], while physically embodied frameworks enable real-time nonverbal interaction with robots through motion-based collaboration or conversational engagement [Chen et al. 2025b; Ji et al. 2025]. However, across these systems, interaction remains fundamentally *reactive*. Constrained to a local temporal window, agents respond only to explicit commands or current sensory inputs such as visual motion and speech features. They have limited capacity to anticipate needs or initiate actions based on accumulated context.

Psychological research characterizes human behavior along a spectrum of agency [Parker et al. 2010]. *Passive behavior* awaits explicit instructions; *reactive behavior* responds to immediate stimuli and operates on local, instantaneous context; and *proactive behavior* is anticipatory and self-initiated, driven by accumulated context and intent inference. Applied to embodied agents, a proactive system does not merely synchronize gestures to current speech, but reasons over the interaction context to generate timely and appropriate initiative behaviors. For instance, if a user sets down a bag early in a conversation and later turns to leave, a proactive agent—recalling the event—might point to the bag and offer a reminder. Such context-dependent actions are vital for social fluency and have been shown to significantly improve user engagement, trust, and perceived helpfulness [van den Broek and Moeslund 2024].

*equal contribution

†corresponding author

Authors' Contact Information: Zeyi Zhang, illusence1@gmail.com; Zixi Kang, forever.kzx0713@stu.pku.edu.cn; Ruijie Zhao, 2420398231@qq.com, Peking University, China; Yusen Feng, ysfeng@stu.pku.edu.cn; Biao Jiang, bjiang25@stu.pku.edu.cn; Libin Liu, libin.liu@pku.edu.cn, Peking University, China.

Enabling such proactive behavior requires two key capabilities. First, the agent must continuously understand the situation, including the user and the surrounding environment, and track how it evolves across the interaction. Second, it must continuously generate realistic and context-aware body movements that sustain embodied presence. Prior work has made progress on each capability in isolation, yet combining them within a single interactive system remains challenging. A central reason is the mismatch in timescales between proactive reasoning and real-time interaction. Proactive reasoning relies on long-horizon context and intent inference built from accumulated history, while real-time interaction demands low-latency, high-rate response generation to maintain conversational fluency. A single monolithic model faces a fundamental trade-off: deep reasoning inhibits real-time reactivity, while prioritizing immediate response leads to reflexive behaviors that lack long-term social coherence. To address this tension, we draw inspiration from dual-process theories in cognitive science [Watson 2011] and recent dual-system architectures in autonomous driving [Tian et al. 2024] and robotics [Figure.ai 2025]. These approaches show that separating immediate reactive execution from slower reasoning is key to achieving both responsiveness and intelligence.

Building on this insight, we propose *ProAct*, a dual-system architecture for proactive embodied social interaction, including (a) a fast *Behavioral System* that constitutes the real-time interaction loop, integrating a cascaded Omni-LLM for verbal response and a streaming motion generator for continuous non-verbal synthesis, and (b) a slower, LLM-driven agentic framework that acts as a *Cognitive System*, equipped with memory and reasoning modules to formulate high-level proactive intentions based on long-term interaction context.

A key technical challenge is enabling the Cognitive System’s intentions to modulate the Behavioral System’s ongoing motion generation without introducing latency or discontinuities. To address this, we introduce a flow-matching model with ControlNet-based intention injection. Unlike traditional systems that rely on discrete motion primitives or fixed labels, our model treats proactive intentions as dynamic conditioning signals that modulate the probability flow of the motion generator. This allows the Behavioral System to maintain fluid, audio-synchronized gestures while naturally gravitating toward the target behavior prescribed by the Cognitive System. This architectural innovation enables fluid transitions between reactive and proactive behaviors within a single continuous motion stream. We deploy the complete system on a humanoid robot and validate our approach through comprehensive experiments, including quantitative evaluations and user studies. The results confirm that our system significantly enhances perceived proactivity, naturalness, and user engagement compared to reactive baselines. Our contributions are as follows:

- We present a dual-system architecture for embodied social agents that integrates real-time reactive responses with context-driven proactive behaviors.
- We introduce an intention-conditioned streaming model based on flow matching that enable asynchronous injection of high-level semantic intentions into real-time motion streams for seamless reactive-proactive transitions.

- We deploy the system on a humanoid robot, achieving real-time proactive interaction, and validate its effectiveness through comprehensive experiments.

2 Related Work

2.1 Embodied Social Interaction Systems

Interactive agents vary in their degree of embodiment and interaction modality. *Text- or speech-based conversational agents* [character.ai 2025; Défossez et al. 2024; Shao et al. 2023] enable rich linguistic exchange through modern LLMs, but lack embodiment for non-verbal expression and environmental interaction. *Virtual embodied agents* address this limitation by generating synchronized speech and nonverbal behavior for digital characters. These systems either employ LLM-cascaded models where speech conditions downstream nonverbal behavior generation [Cai et al. 2025a, 2024, 2025b; Kim et al. 2024; Wampfler et al. 2025], or encode multimodal behavior jointly within the LLM backbone [Ao 2024; Jiang et al. 2025; Zhang et al. 2025a]. While achieving impressive realism, they remain confined to virtual environments. *Physically embodied robots* extend interaction into the real world. Motion-based interaction frameworks [Chen et al. 2025b; Ji et al. 2025] detect human intent from visual cues and generate responsive physical actions, enabling real-time collaboration. Commercial humanoid companion robots [Engineered Arts 2025; Hanson Robotics 2025; Kang et al. 2024] emphasize facial animation and scripted speech, but typically lack expressive full-body motion generation.

While these systems achieve real-time responsiveness, they remain signal-driven: motion generation is directly coupled to immediate sensory input like user motion and speech audio without incorporating high-level reasoning. Proactive social behavior, in contrast, requires reasoning over accumulated conversational context to infer when and how to initiate unsolicited but contextually appropriate actions.

2.2 Non-verbal Behavior Generation

Non-verbal behavior generation has been extensively studied, particularly in the context of co-speech gestures [Nyatsanga et al. 2023], where diverse generative architectures—ranging from statistical and neural models to Transformers and diffusion approaches—have enabled increasingly natural gesture synthesis [Alexanderson et al. 2023; Ao et al. 2022, 2023; Chen et al. 2025a; Cheng et al. 2025, 2024; Ghorbani et al. 2023; Liu et al. 2024, 2022, 2025b; Mughal et al. 2025; Ng et al. 2024; Yang et al. 2025; Yi et al. 2023; Yin et al. 2025; Yoon et al. 2020; Zhang et al. 2025b]. Building on these foundations, recent work has expanded toward dyadic or partner-aware gesture generation [Agrawal et al. 2025; McLean et al. 2025; Mughal et al. 2024; Qi et al. 2025; Sun et al. 2024; Xue et al. 2025; Zhang et al. 2025c], multimodal control synthesis guided by text instructions or semantic signals beyond raw audio [Chen et al. 2024, 2025c; Yang et al. 2024; Zhang et al. 2024], and real-time and arbitrary-length generation, with efficient network architectures [Liu et al. 2025c], and streaming generation strategy through autoregressive formulations or rolling diffusion [Cai et al. 2025c; Tseng et al. 2023; Xiao et al. 2025; Zhen et al. 2025] to achieve frame-by-frame synthesis with low

latency. Our work integrates these capabilities through a streaming flow-matching architecture with ControlNet-based dynamic conditioning, supporting real-time gesture synthesis modulated by asynchronously injected control signals.

2.3 Dual System for Real-time Interaction

Dual-system architectures, inspired by the distinction between fast reactive "System-1" and slower proactive "System-2" reasoning [Watson 2011], reconcile the tension between low-latency responsiveness and contextually grounded decision-making. A single monolithic model often struggles to satisfy both: fast models lack semantic depth, while deliberative models are too slow for interactive control. To address this, conversational agents such as the Talker-Reasoner framework [Christakopoulou et al. 2024] separate immediate dialogue generation from slower multi-step reasoning. In autonomous driving, DriveVLM-Dual [Tian et al. 2024] integrates high-frequency visuomotor control with a slower VLM-based module to support semantic scene understanding and long-horizon planning. In robot control, Helix [Cui et al. 2025] similarly pairs a high-rate reactive controller with a slower VLM that provides task-level objectives and semantic grounding. We apply this design to embodied social interaction: our Behavioral System (System-1) generates streaming multimodal responses, while our Cognitive System (System-2) injects proactive intentions through social reasoning.

3 Overview

As illustrated in Figure 1 a), *ProAct* adopts a dual-system architecture that couples a fast, reactive Behavioral System with a planning-oriented, proactive Cognitive System. The Behavioral System operates as the real-time interaction loop, continuously processes user audio and visual inputs and producing low-latency multimodal responses. It integrates a streaming omni-modal LLM for generating verbal responses with a streaming motion generator for continuous synthesis of non-verbal behavior. While speech is generated in turn-based segments after the user finishes speaking, the motion generator continuously produces non-verbal behavior from both user and agent audio streams, ensuring uninterrupted motion even during listening phases.

Running in parallel, the Cognitive System performs social reasoning over accumulated context via an agentic framework with episodic memory. This framework observes the ongoing dialogue and visual scene, and infers when and what proactive actions should be initiated. When triggered, it generates intention signals, which are injected into the Behavioral System to modulate both verbal and non-verbal responses. This allows the robot to seamlessly blend fast reactive behaviors with deliberate, context-driven proactive actions.

4 Behavioral System: Streaming Interaction

The Behavioral System maintains the real-time interaction loop, producing streaming verbal and non-verbal responses through a cascaded architecture that connects a streaming omni-modal LLM with a streaming motion generator. The two channels operate asynchronously: verbal responses follow a turn-based pattern, while

motion generation runs continuously to maintain embodied presence throughout the interaction. We deploy the generated motion on a physical robot through hierarchical whole-body control policy. The generated upper-body joint angles are tracked by PID controllers, while lower-body balance is maintained by compensating for angular momentum induced by upper-body motion.

4.1 Streaming Verbal Response Synthesis

The omni-modal LLM processes user speech and visual observations from the robot’s egocentric camera to generate audio responses. We use GPT Realtime model [OpenAI 2025b], which operates in turn-based mode with streaming input and output. User audio is continuously streamed chunk-by-chunk for low-latency processing. Upon detecting turn completion via voice activity detection, the model begins generating agent audio in a streaming fashion, producing chunks progressively for immediate playback. The model supports mid-utterance interruption if the user speaks during the agent’s response, ensuring interaction fluency.

4.2 Controllable Streaming Motion Synthesis

We represent motion as per-frame joint rotations in a unified feature space (see Sec. 6.1.1). In contrast to the turn-based verbal channel, motion generation operates continuously throughout the interaction to maintain embodied presence. The motion generator must meet two requirements. First, *real-time generation*: the system must synthesize motion faster than its playback duration, ensuring continuous, uninterrupted motion streaming after the initial response. Specifically, for a motion sequence of duration T_{motion} , the computation time for generation T_{gen} must satisfy $T_{\text{gen}} < T_{\text{motion}}$. Second, *controllable generation*: the system must follow high-level semantic motion instructions from Cognitive System to enable proactive behavior intentions. To address these challenges, we propose a controllable streaming motion generator based on flow matching.

Formally, we cast motion generation as a streaming problem that operates chunk-by-chunk generation with n frames for each chunk. At the i -th generation step, the generator \mathcal{G} takes synchronized dyadic speech chunks S_u^i, S_a^i from both user and agent, an optional intention instruction I^i , and the last ℓ frames of the previous chunk M^{i-1} (denoted as $M_{-\ell}^{i-1}$) to ensure temporal continuity:

$$M^i = \mathcal{G}(S_u^i, S_a^i, I^i, M_{-\ell}^{i-1}). \quad (1)$$

For a detailed illustration of the model architecture and specific training and inference strategies, please refer to Appendix A.

4.2.1 Flow Matching for Real-time Motion Synthesis. To model the motion distribution M under real-time constraints, we adopt Conditional Flow Matching [Lipman et al. 2023], which transforms noise into data through a learned continuous-time flow. Specifically, flow matching learns a time-dependent velocity field $v_\tau(M_\tau, C)$ that transports samples from Gaussian noise $M_0 \sim \mathcal{N}(0, I)$ at $\tau = 0$ to the data distribution $M_1 \sim q_{\text{data}}$ at $\tau = 1$, where the conditioning variable C comprises the audio chunks. The generation follows an ordinary differential equation (ODE):

$$\frac{dM_\tau}{d\tau} = v_\tau(M_\tau, C), \quad \tau \in [0, 1]. \quad (2)$$

We use the optimal-transport path that linearly interpolates between noise and data: $M_\tau = \tau M_1 + (1 - \tau)M_0$. We train a neural velocity estimator v_θ to match the path velocity $M_1 - M_0$ by minimizing:

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{\tau, M_1, M_0} [\|v_\theta(\tau, M_\tau, C) - (M_1 - M_0)\|^2], \quad (3)$$

where $\tau \sim \mathcal{U}[0, 1]$. During inference, we sample a noise point M_0 and iteratively solve the ODE to generate M_1 using Euler integration. This straight-line path simplifies ODE integration compared to diffusion models [Ho et al. 2020], enabling high-quality generation with fewer solver steps, thus achieving faster generation without compromising motion quality.

We instantiate the velocity estimator v_θ as \mathcal{G} , a transformer-based backbone built on DiT blocks [Peebles and Xie 2023]. Unlike prior methods that generate gestures from a single speaker’s audio, \mathcal{G} explicitly models dyadic interaction through two synchronized audio streams S_a and S_u . This enables the agent to produce both speaker-side expressive gestures and listener-side attentive behaviors such as subtle body sway and attentive idling during partner speech. The model processes noisy motion M_τ through stacked DiT blocks, each containing self-attention and feed-forward layers with residual connections. The two audio streams are concatenated channel-wise, with their positions in the feature dimension implicitly distinguishing self versus partner audio. The audio features and the time embeddings of the flow are concatenated and injected via AdaLN-Zero [Peebles and Xie 2023] to modulate the generation process, ensuring temporal consistency across the flow trajectory.

To ensure temporal continuity across streaming windows, we use an overlap-and-cache scheme enabled by the tileable property of flow matching [Tseng et al. 2023]. We cache the last ℓ frames from the previous window, prepend them to the next window as an ℓ -frame overlap, and clamp this overlap at every solver step by overwriting it with the cached frames. This naturally preserves recent motion history during flow evolution, eliminating boundary discontinuities. Combined with flow matching’s low-step sampling, our system generates continuous motion sequences faster than the playback duration, satisfying the real-time criterion.

4.2.2 Semantic Intention Control via Disentangled ControlNet. To empower the agent with proactive non-verbal capabilities, the motion generator must support dual conditioning: responding to audio rhythms while adhering to high-level semantic intentions from the Cognitive System. These intentions are formulated as natural language descriptions, requiring the generator to function as a text-to-motion synthesis component operating in parallel with the audio-driven process. Training a unified model often requires paired audio-motion-text data, where speech audio modulates the rhythm and fine-grained semantics of the motion, while text descriptions provide high-level semantic guidance. However, such datasets are largely absent, posing a critical challenge. Prior work [Chen et al. 2024; Yang et al. 2024] addresses this by jointly training on separate audio-motion and text-motion datasets within a single shared backbone. However, our preliminary experiments reveal that this approach suffers from an alignment issue: the model struggles to disentangle rhythmic synchronization from semantic understanding, leading to compromised performance (see ablation in Sec. 6.3.3).

Inspired by the recent success of ControlNet [Zhang et al. 2023] in motion synthesis [Dai et al. 2025], we introduce a disentangled ControlNet architecture that decouples semantic control from rhythmic generation. Our system comprises a frozen, pre-trained audio-driven base generator and a parallel trainable control branch initialized as a structural copy. Semantic intention signals are encoded and injected into the control branch, preserving the original audio-driven priors while enabling explicit semantic guidance.

During inference, the system seamlessly transitions between audio-driven and intention-driven behaviors: when no intention is provided, the zero-initialized connections ensure natural fallback to pure audio-driven generation. Since the ControlNet branch shares the same lightweight backbone as the base flow-matching generator, it introduces only a small constant overhead. As a result, our system maintains real-time generation speed while enabling flexible intention control, achieving both controllability and efficiency in a unified model.

5 Cognitive System: Social Context Reasoning

While the Behavioral System ensures fluent reactive interaction, truly social agents must also exhibit proactive behaviors when contextually appropriate. The Cognitive System addresses this need through an LLM-based agentic framework that performs deliberative social reasoning in parallel with real-time interaction. While operating at a slower pace than the Behavioral System to allow reasoning over accumulated context, it faces a competing constraint: excessive latency risks missing timely intervention opportunities. Moreover, extended interactions generate increasingly long dialogue and reasoning histories, leading to progressively increasing inference time.

To address the inference time control problem, we decompose the reasoning process into two specialized components: a *Context Encoder* that compresses the accumulating interaction history into a bounded memory bank, and a *Behavior Planner* that performs motivation assessment and action planning. As illustrated in Figure 2, within each reasoning cycle, both components operate in parallel on outputs from the previous cycle. By limiting the context length through compression, we ensure that each cycle completes within a consistent time budget t_c regardless of conversation length. The Cognitive System executes these reasoning cycles continuously in a streaming fashion. Each cycle is triggered immediately upon completion of the previous one, ensuring temporal coherence while maintaining a steady reasoning cadence for timely proactive responses. The detailed prompts are provided in Appendix F.

5.1 Context Encoder

The Context Encoder maintains a compressed memory bank to prevent unbounded context growth. At the start of cycle k , the encoder receives three incremental inputs: dialogue transcripts $T_{k-5:k}$ accumulated over the five most recent inter-cycle intervals for sufficient temporal context, three visual frames $V_k^{(3)}$ uniformly sampled since the last cycle to capture scene dynamics, and the previous behavior plan P_{k-1} from the Behavior Planner. The encoder analyzes these inputs to extract salient events and observations, then integrates

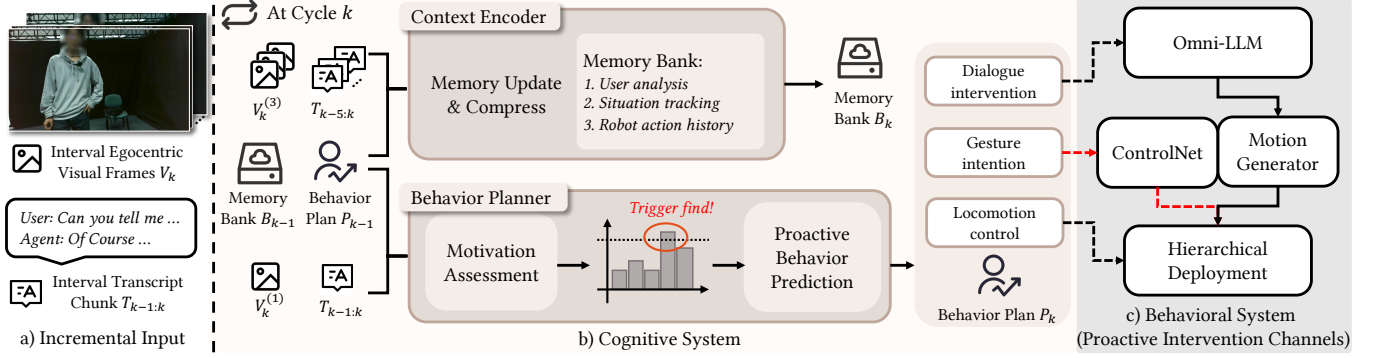


Fig. 2. **Architecture of the Cognitive System.** The system operates in continuous cycles, with each cycle consisting of three stages: (a) collecting incremental multimodal inputs; (b) parallel execution of Context Encoder (memory compression) and Behavior Planner (proactive behavior prediction); (c) injecting behavior plans into the Behavioral System via different channels.

them into the existing memory bank B_{k-1} through incremental summarization, producing an updated compressed B_k of a bounded size: $B_k = \text{ContextEncoder}(B_{k-1}; T_{k-5:k}, V_k^{(3)}, P_{k-1})$.

The memory bank B_k maintains a structured representation across three hierarchical levels. First, inspired by Theory of Mind reasoning [Premack and Woodruff 1978], *user analysis* tracks the user’s observable appearance and inferred mental state, capturing their beliefs, desires, and behavioral intentions. Second, *situation tracking* records significant events, observable environmental changes, and unresolved tasks as they emerge during interaction. Third, *robot action history* summarizes previously executed proactive behaviors to prevent redundant interventions. To control memory bank size, each component within these levels is constrained to a single-sentence summary. This compression-based approach trades perfect recall for consistent inference speed.

5.2 Behavior Planner

The Behavior Planner determines when and how to initiate proactive actions based on the compressed memory and current observations. At each cycle k , the planner operates on the most recent information: it receives the dialogue transcript $T_{k-1:k}$ accumulated since the previous cycle, the most recent visual frame $V_k^{(1)}$, and the memory bank B_{k-1} and behavior plan P_{k-1} from the previous cycle. This design enables responsive decision-making based on immediate context while leveraging compressed long-term memory for contextual grounding. We employ an asymmetric input window to balance performance: this allocates computational budget to accurate history summarization while minimizing Planner latency for timely interventions. The planner outputs a plan P_k specifying a set of proactive actions: $P_k = \text{BehaviorPlanner}(T_{k-1:k}, V_k^{(1)}, B_{k-1}, P_{k-1})$.

Determining when to intervene is non-trivial: acting too frequently disrupts conversational flow, while inaction misses opportunities for engagement. Inspired by prior work on proactive behavior prediction [Liu et al. 2025a], we design a five-dimensional motivation assessment framework where each dimension captures a distinct contextual trigger:

- *Visual Scene Changes*: Monitoring dynamic changes in the environment that may require attention.
- *User Intent Signals*: Detecting explicit or implicit cues in user behavior suggesting a need for assistance.
- *Conversation State*: Analyzing dialogue flow to identify appropriate moments for intervention.
- *Social Protocol Requirements*: Adhering to social norms and rituals in daily interaction.
- *Emotional Response Needs*: Recognizing the user’s emotional state and responding appropriately.

Each dimension is scored from 1 to 5 based on contextual salience. The planner checks whether any dimension scores ≥ 4 : if at least one trigger reaches this threshold, the system initiates proactive action reasoning and generates corresponding interventions. We empirically set this threshold to 4 to balance responsiveness and stability: lower values cause excessive interruptions, while higher ones miss proactive opportunities. This design reflects the principle that a single strong contextual trigger is sufficient to warrant proactive behavior, without requiring consensus across multiple dimensions.

We operationalize proactive control through three intervention channels. The primary channel is *gesture intention injection*: each intention is expressed as a concise motion description following the format *body part + action*. These descriptors are injected into the motion generator via the ControlNet mechanism, enabling seamless blending of proactive gestures with ongoing audio-driven motion. Additionally, we implement two auxiliary proactive behaviors: *dialogue intervention* employs interruption and instruction injection to steer conversation topics or supplementary knowledge by injecting prompts into the Omni-LLM, and *locomotion control* adjusts *spatial positioning and orientation* via *gait-aware velocity commands* (detailed in Appendix B). This architectural separation enables both fast reactive behaviors and slower proactive actions.

6 Experiments

6.1 System Setup

6.1.1 Data Preparation. Currently, no dataset provides paired audio-text-motion triplets where text explicitly specifies semantics. We

therefore leverage two types of paired data. For audio-driven generation, we extract 2.5 hours of synchronized speech and motion from speaker PXB in the Photoreal dataset [Ng et al. 2024]. For intention-driven generation, we adapt the SeG dataset [Zhang et al. 2024] by reannotating descriptions following the HumanML3D protocol [Guo et al. 2022] and retargeting motions to the Photoreal skeleton, yielding 544 clips with 1,632 captions (details in Appendix C). To avoid runtime retargeting delays during physical deployment, we retarget all training datasets to the robot kinematic structure by optimizing orientation and position losses under collision-free and joint-limit constraints via mink [Zakka 2025]. All motion sequences are segmented into 5-second clips at 30 FPS. Each frame uses an exponential map representation with 57 joints for the Photoreal skeleton and 23 for the robot. The root joint encodes absolute position and velocity. In summary, $M_t \in \mathbb{R}^{(J \times Q + G)} = \mathbb{R}^{(J \times 3 + 6)}$. Audio features are extracted using Librosa [McFee et al. 2015], including MFCCs, deltas, chromagrams, onset strength, and tempograms.

6.1.2 Settings. For the Behavioral System, we use gpt-realtime-2025-08-28 [OpenAI 2025b] as the Omni-LLM for verbal response generation. For motion generation, we train the flow-matching model with learning rate $\alpha = 10^{-4}$, flow-matching steps $T = 1000$, window size $n = 150$ frames, and overlap $\ell = 30$ frames for temporal continuity. We apply conditional dropout with probability $p = 0.2$ during training. Training requires approximately 19 hours for the audio-to-gesture backbone and 4 hours for ControlNet fine-tuning on two NVIDIA RTX 4090 GPUs. During inference on a single RTX 4090 GPU, we use 50 sampling steps with computational cost approximately 2.5 TFLOPs per forward pass. With motion chunks of duration $T_{\text{motion}} = 1\text{s}$, our method achieves generation times $T_{\text{gen}} = 0.425\text{s}$ (audio-only) and $T_{\text{gen}} = 0.648\text{s}$ (with ControlNet), satisfying the real-time constraint $T_{\text{gen}} < T_{\text{motion}}$. We employ the Sherpa-ONNX ASR model [k2-fsa 2023] for real-time speech transcription. Detailed model architecture is provided in Appendix A.

For the Cognitive System, we employ gpt-4.1-mini-2025-04-14 [OpenAI 2025a] as the reasoning LLM with carefully designed prompts for both Context Encoder and Behavior Planner (Appendix D). Each reasoning cycle completes within $t_c = 3$ seconds under normal conditions, enabling timely proactive interventions. The entire system is deployed on a Unitree G1 humanoid robot [Unitree Robotics 2025] for real-world interaction experiments.

6.2 Results

We design 5 interaction scenarios, each specifying the conversational topics. We do not prescribe any behavioral details or body motions for the agent; instead, these are generated entirely by the system. Scenario information was injected into the system prompt to assess whether the agent could exhibit appropriate proactive behaviors under real-world interaction settings. Further details on these scenarios are provided in Appendix E.

As illustrated in Figure 3, *ProAct* effectively initiates interactions and actively clarifies misconceptions during a poster presentation. The agent autonomously initiates interaction with a passerby, adjusts its orientation for optimal presentation, and actively intervenes to clarify verbal misconceptions using synchronized gestures. Furthermore, Figure 4 demonstrates its capability to maintain social

Table 1. Quantitative evaluation on the Photoreal datasets. All methods are trained on the same training data, and evaluated on the test audio.

System	FGD↓	BeatAlign↑	Div _k ↑
GT	-	0.847	2.398
LDA	79.71	0.732	1.416
EMAGE	82.29	0.767	1.847
Photoreal	69.44	0.746	1.986
SocialAgent	72.57	0.812	2.014
Ours (Joint-Training)	73.31	0.752	1.871
Ours	66.53	0.823	2.230

norms during a storytelling session. Upon detecting user distraction (smartphone usage), the agent generates context-aware dissatisfaction signals and corrective gestures to regain attention. Overall, the results demonstrate *ProAct*’s ability to generate proactive, multimodal behaviors that align with complex social dynamics and effectively handle diverse social scenarios. We refer readers to the supplementary video for a complete demonstration of continuous, real-time interactions.

6.3 Evaluation

Recent interactive social agent systems exhibit diverse design choices regarding sensing modalities, interaction paradigms, and cognitive architectures, making direct cross-system comparison infeasible. We therefore conduct a multi-stage evaluation strategy: first, we compare our motion generation module against state-of-the-art gesture synthesis methods using quantitative metrics. Second, we evaluate our complete dual-system framework through systematic ablations on the physical robot, examining the contribution of each architectural component through user studies. Finally, we validate key design decisions through targeted ablation experiments on training strategy and inference efficiency.

6.3.1 Gesture Generation Comparison. We compare against state-of-the-art gesture generation models: LDA [Alexanderson et al. 2023], EMAGE [Liu et al. 2024], Photoreal [Ng et al. 2024], and SocialAgent [Zhang et al. 2025c]. We re-train LDA and EMAGE on the Photoreal datasets, use the publicly released checkpoint for Photoreal, and use the results of SocialAgent provided by its authors. All metrics are computed using the original Photoreal skeleton to ensure fair comparison with baseline methods. We quantitatively evaluate motion quality using three metrics: a) Fréchet Gesture Distance (FGD) [Yoon et al. 2020] quantifies the disparity between the latent feature distributions of generated and real gestures; b) BeatAlign [Li et al. 2021] assesses speech-motion synchrony by measuring the temporal alignment between motion beat candidates; c) Div_k [Ng et al. 2024] evaluates the diversity of the generated motions. As shown in Table 1, our method achieves the best FGD, BeatAlign, and Div_k among all baselines. These results validate the effectiveness of our flow-matching model in generating high-quality, diverse, and synchronized gestures.

6.3.2 System Components Ablation. We evaluate our framework through systematic ablations on the physical robot in real interaction scenarios, incrementally validating each component: a) *Speech-Only*:

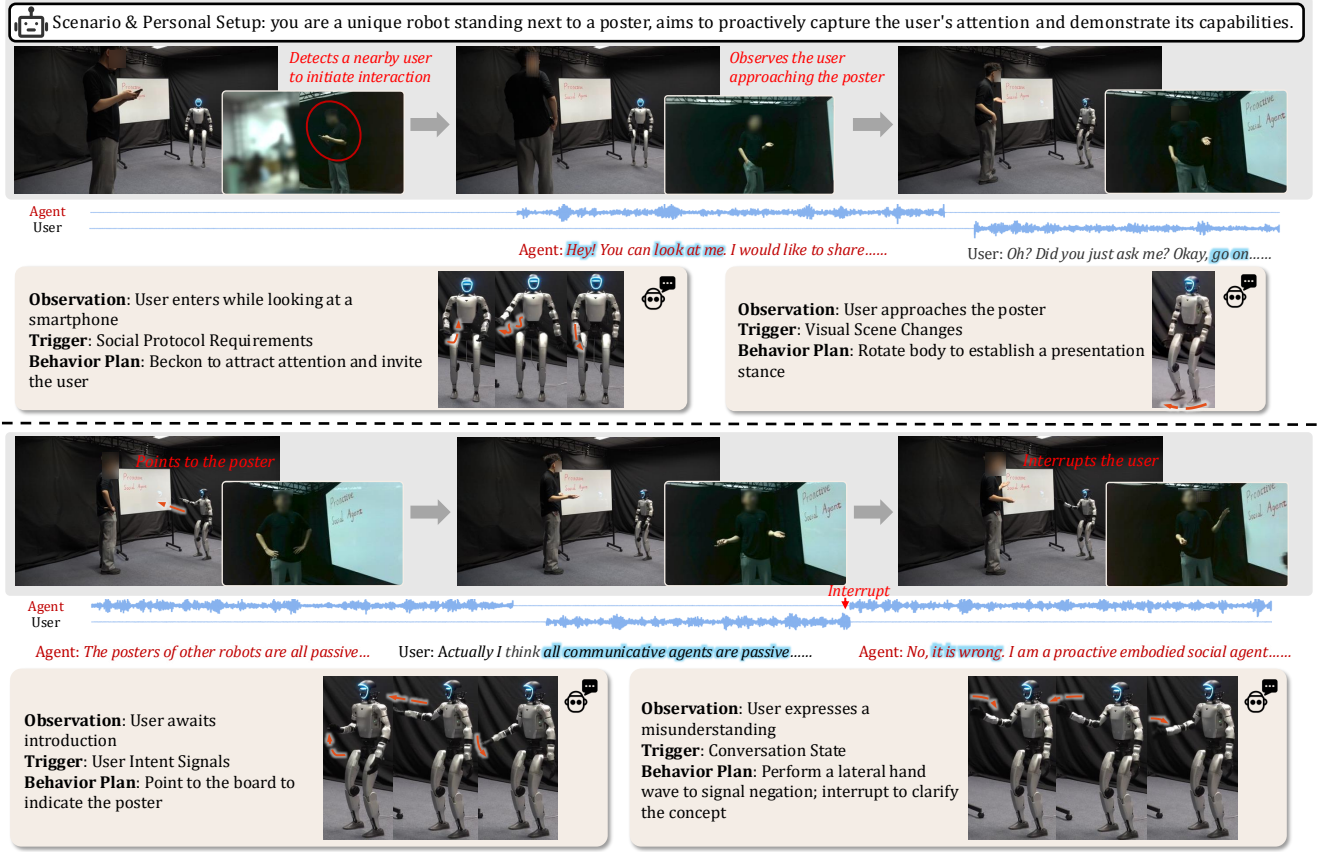


Fig. 3. **Demonstration of ProAct during a poster explanation task.** The figure captures a continuous real-time session where the agent dynamically adapts its behavior based on visual and auditory triggers. Unlike passive systems, *ProAct* proactively captures the user's attention (top-left), orients to the user (top-right), guides the conversation to the poster (bottom-left), and actively clarifies a misconception (bottom-right). The top header specifies the scenario setup, while the beige panels visualize the *Cognitive System*, mapping real-time observations and triggers to specific behavioral plans. The generated gestures correspond to these intentions, where the red arrows in the overlays indicate the velocity and direction of each movement.

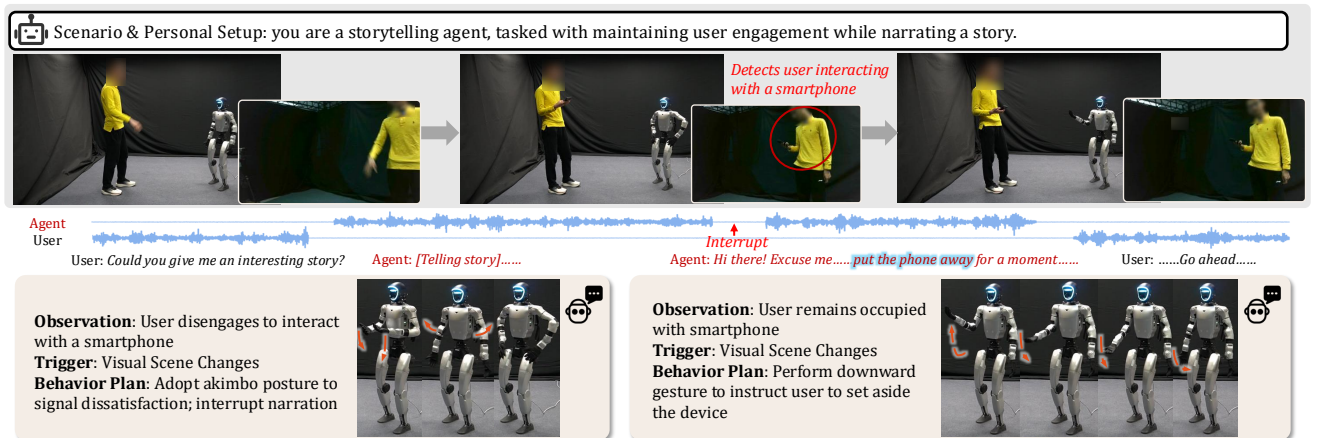


Fig. 4. **Demonstration of ProAct during a storytelling task.** In this scenario, *ProAct* is tasked with maintaining user engagement while narrating a story. The agent detects a loss of engagement when the user disengages to interact with a smartphone (left panel) and actively intervenes to regain attention by instructing the user to set the device aside (right panel). Together with fig. 3, these results demonstrate the agent's capability to handle diverse social scenarios.

Table 2. Observer-based user study: Average scores from pairwise comparisons with 95% confidence intervals. Asterisks indicate statistically significant differences compared to the Full System ($p < 0.05$).

System	Motion Naturalness \uparrow	Active Agency \uparrow	Perceived Presence \uparrow
a) Speech-Only	-0.523*	-0.341*	-0.589*
b) + Motion Generator	0.175	-0.218*	0.045*
c) + Behavior Planner	0.159*	0.148*	0.169*
d) + Context Encoder	0.187	0.417	0.370

Baseline using only Omni-LLM for verbal interaction, with the robot maintaining a static pose; b) + *Motion Generator*: Adds streaming motion generator for real-time non-verbal presence; c) + *Behavior Planner*: Further adds the planner for proactive interventions but relies on uncompressed history; d) + *Context Encoder*. Note that setting b (*Speech-Only* + *Motion Generator*) corresponds to the *Behavioral System*, which is also referred to as *w/o Cognitive System* in the following sections. Setting d represents our *Full System*.

Following the approach in [Ji et al. 2025; Liu et al. 2025a], we conducted comprehensive user studies from two perspectives: *Observer-based* evaluation via pairwise comparisons [Zhang et al. 2025c], and *Participant-based* evaluation using a 5-point Likert scale [Chen et al. 2025b]. This dual-perspective design ensures a holistic evaluation: observer studies objectively assess behavioral quality without interaction burden, while participant studies capture the subjective first-person experience of engagement and comfort essential for social bonding. We employ three core metrics for observer evaluation: *Motion Naturalness* (kinematic realism), *Active Agency* (autonomy in initiating actions) and *Perceived Presence* (sense of social entity). For participant evaluation, we additionally include three experiential metrics: *Response Timeliness* (latency smoothness), *Interaction Comfort* (psychological ease), and *Willingness to Re-engage* (motivation to continue). A detailed description of the user study and the definitions of these evaluation metrics are provided in Appendix E.

Table 2 and Figure 5 present the ablation results. The full system (setting d) achieves the highest scores across all metrics. We analyze each component’s contribution below.

Impact of Motion Generator (a vs. b). Setting b significantly improves *Motion Naturalness* and *Perceived Presence* over the speech-only baseline, confirming that embodied non-verbal behavior is fundamental for social presence.

Impact of Cognitive System (b vs. d). Setting b (w/o Cognitive System) relies solely on the Behavioral System, leading to substantial declines in *Active Agency* and *Perceived Presence*. Moreover, the participant-based evaluation reveals statistically significant drops in *Interaction Comfort* and *Willingness to Re-engage*. As illustrated in Figure 6, without the Cognitive System the agent becomes purely reactive and misses opportunities for timely, contextually appropriate engagement.

Impact of Context Encoder (c vs. d). Setting c (w/o Context Encoder) forces the Behavior Planner to reason over uncompressed history. As conversation length increases, inference latency grows progressively, causing delayed proactive responses that miss critical intervention timing, which degrades performance especially in *Active Agency* and *Perceived Presence*.

6.3.3 Training Strategy. In this experiment, We compare our disentangled ControlNet training scheme with a joint-training baseline that learns audio and text conditioning within a single backbone. Table 1 shows that our method outperforms the joint training baseline across all metrics, whereas our scheme preserves rhythmic alignment while enabling reliable high-level intention control.

6.3.4 Flow Matching Efficiency. In this experiment, we compare inference latency with a diffusion baseline [Zhang et al. 2025c] using DDIM with 200 steps, reporting the average generation time for a 1-second motion chunk. The diffusion baseline takes 1.597s per chunk and violates the real-time constraint $T_{\text{gen}} < T_{\text{motion}}$, resulting in noticeable motion stuttering as shown in the supplementary video. In contrast, our flow-matching model satisfies this real-time constraint and supports real-time interaction.

7 Conclusion

In this paper, we introduce *ProAct*, a dual-system framework for proactive embodied social agents. We first develop an intention-conditioned flow-matching model for streaming motion generation, enabling high-level semantic intentions to be asynchronously injected into real-time gesture synthesis. Building upon this, we design a dual-system architecture that decouples real-time reactive interaction (Behavioral System) from deliberative social reasoning (Cognitive System), enabling the agent to maintain conversational fluency while initiating contextually appropriate proactive behaviors. The intention-conditioned motion generator serves as the technical bridge, allowing the Cognitive System’s proactive decisions to seamlessly modulate the Behavioral System’s continuous motion stream. We deploy the complete system on a physical humanoid robot. User studies and quantitative evaluations demonstrate that our framework significantly enhances perceived proactivity, naturalness, and user engagement compared to purely reactive baselines.

Despite its effectiveness, our approach has limitations that suggest future directions. First, the triggering mechanism operates at fixed intervals with threshold-based criteria, which may miss fleeting opportunities or trigger unnecessary interventions. Adaptive triggering strategies learned from interaction data could better balance responsiveness and appropriateness. Second, the frequency and intensity of proactive behaviors may not align with individual preferences. Incorporating reinforcement learning from human feedback to personalize the proactive policy could improve user satisfaction. Finally, the initial response latency stands at approximately 2–3 seconds, primarily due to the cascaded architecture and network latency of the cloud-based Omni-LLM. Training end-to-end multimodal models for local deployment could significantly reduce this latency and enhance overall responsiveness.

References

- Vasu Agrawal, Akinniyi Akinyemi, Kathryn Alvero, Morteza Behrooz, Julia Buffalini, Fabio Maria Carlucci, Joy Chen, Junming Chen, Zhang Chen, Shiyang Cheng, Praveen Chowdary, Joe Chuang, Antony D’Avirro, Jon Daly, Ning Dong, Mark Duppenthaler, Cynthia Gao, Jeff Girard, Martin Gleize, Sahir Gomez, Hongyu Gong, Srivathsan Govindarajan, Brandon Han, Sen He, Denise Hernandez, Yordan Hristov, Rongjie Huang, Hirofumi Inaguma, Somya Jain, Raj Janardhan, Qingyao Jia, Christopher Klaiber, Dejan Kovachev, Moneish Kumar, Hang Li, Yilei Li, Pavel Litvin, Wei Liu, Guangyao Ma, Jing Ma, Martin Ma, Xutai Ma, Lucas Mantovani, Sagar Miglani, Sreyas Mohan, Louis-Philippe Morency, Evonne Ng, Kam-Woh Ng, Tu Anh Nguyen, Amia Oberai, Benjamin Peloquin, Juan Pino, Jovan Popovic, Omid Poursaeed, Fabian

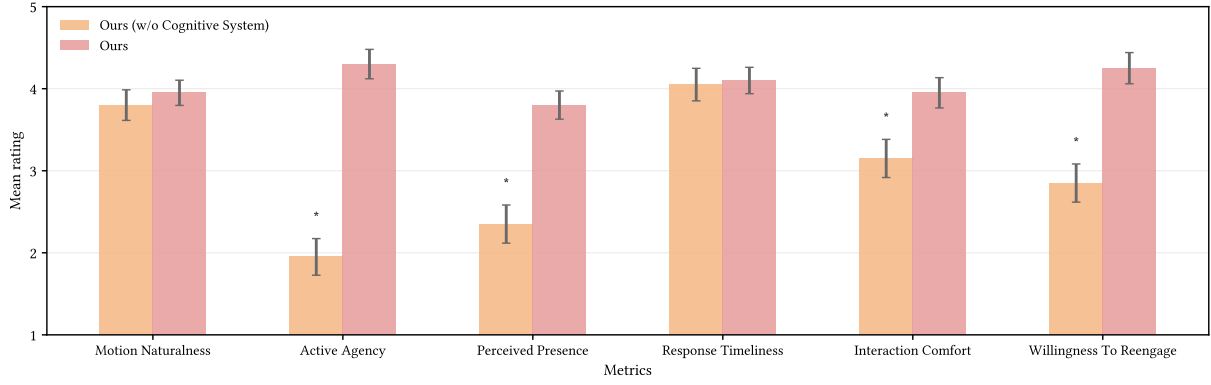


Fig. 5. **Participant-based user study results.** We compare the **Full System** (setting d) against **w/o Cognitive System** (setting b) across six experiential metrics. Each participant interacted with both systems in randomized order and rated them on a 5-point Likert scale. Error bars show standard deviations. Asterisks mark statistically significant differences ($p < 0.05$). The Full System demonstrates significant improvements in Active Agency, Perceived Presence, Interaction Comfort, and Willingness to Re-engage, while maintaining comparable Motion Naturalness and Response Timeliness, confirming that proactive capabilities enhance user experience without compromising real-time performance.

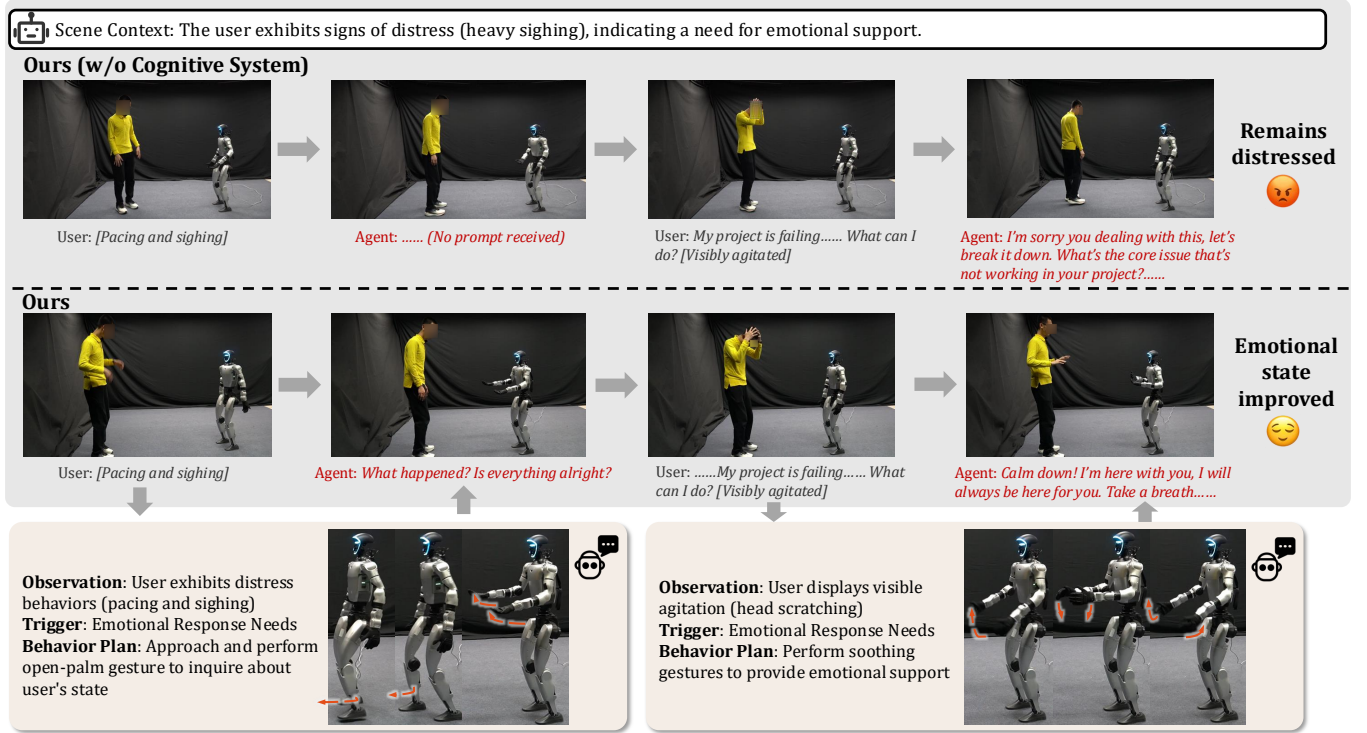


Fig. 6. **Qualitative comparison of emotional support strategies.** We compare a baseline reactive agent (top row) against *ProAct* (bottom row) in a scenario where the user exhibits visible distress. The baseline agent remains passive, treating the interaction as a standard Q&A task and offering analytical solutions ("What is the core issue?"), which fails to alleviate the user's agitation. In contrast, *ProAct* utilizes its *Cognitive System* to interpret non-verbal cues (e.g., pacing, sighing) as a trigger for intervention. As visualized in the beige panels, the agent actively approaches the user to inquire about their state and prioritizes emotional de-escalation through soothing gestures and empathetic dialogue ("I will always be here for you"), successfully improving the user's emotional state. red arrows indicate the velocity and trajectory of the generated comforting movements.

- Prada, Alice Rakotoarison, Rakesh Ranjan, Alexander Richard, Christophe Ropers, Safiyyah Saleem, Vasu Sharma, Alex Shcherbyna, Jia Shen, Jie Shen, Anastasis Stathopoulos, Anna Sun, Paden Tomasello, Tuan Tran, Arina Turkatenko, Bo Wan, Chao Wang, Jeff Wang, Mary Williamson, Carleigh Wood, Tao Xiang, Yilin Yang, Julien Yao, Chen Zhang, Jiemin Zhang, Xinyue Zhang, Jason Zheng, Pavlo Zhyzhneria, Jan Zikes, and Michael Zollhoefer. 2025. Seamless Interaction: Dyadic Audiovisual Motion Modeling and Large-Scale Dataset. arXiv:2506.22554 [cs.CV] <https://arxiv.org/abs/2506.22554>
- Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. 2023. Listen, denoise, action! audio-driven motion synthesis with diffusion models. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–20.
- Tenglong Ao. 2024. Body of her: A preliminary study on end-to-end humanoid agent. *arXiv preprint arXiv:2408.02879* (2024).
- Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. 2022. Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM Transactions on Graphics (TOG)* 41, 6 (2022), 1–19.
- Tenglong Ao, Zeyi Zhang, and Libin Liu. 2023. Gesturediffclip: Gesture diffusion model with clip latents. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–18.
- Albert Bandura. 2001. Social cognitive theory: An agentic perspective. *Annual review of psychology* 52, 1 (2001), 1–26.
- Yiyi Cai, Xuangeng Chu, Xiwei Gao, Sitong Gong, Yifei Huang, Caixin Kang, Kunhang Li, Haiyang Liu, Ruicong Liu, Yun Liu, Dianwen Ng, Zixiong Su, Erwin Wu, Yuhang Wu, Dingkun Yan, Tianyu Yan, Chang Zeng, Bo Zheng, and You Zhou. 2025a. Towards Interactive Intelligence for Digital Humans. arXiv:2512.13674 [cs.CV] <https://arxiv.org/abs/2512.13674>
- Yiyi Cai, Yuhang Wu, Kunhang Li, You Zhou, Bo Zheng, and Haiyang Liu. 2025c. FloodDiffusion: Tailored Diffusion Forcing for Streaming Motion Generation. arXiv:2512.03520 [cs.CV] <https://arxiv.org/abs/2512.03520>
- Zhongang Cai, Jianping Jiang, Zhongfei Qiao, Xinying Guo, Mingyuan Zhang, Zhengyu Lin, Haiyi Mei, Chen Wei, Ruisi Wang, Wanqi Yin, Liang Pan, Xiangyu Fan, Han Du, Peng Gao, Zhitao Yang, Yang Gao, Jiaqi Li, Tianxiang Ren, Yukun Wei, Xiaogang Wang, Chen Change Loy, Lei Yang, and Ziwei Liu. 2024. Digital Life Project: Autonomous 3D Characters with Social Intelligence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 582–592.
- Zhongang Cai, Daxuan Ren, Yang Gao, Yukun Wei, Tongxi Zhou, Huimuk Jang, Haoyang Zeng, Zhengyu Lin, Chen Change Loy, Ziwei Liu, and Lei Yang. 2025b. Digital Life Project 2: Open-source Autonomous 3D Characters on the Web. SIGGRAPH Asia 2025 Real-Time Live! Live demonstration, Hong Kong, China. character.ai. 2025. <https://character.ai>. Accessed: 2025-11-13.
- Bohong Chen, Yumeng Li, Yao-Xiang Ding, Tianjia Shao, and Kun Zhou. 2024. Enabling Synergistic Full-Body Control in Prompt-Based Co-Speech Motion Generation. In *Proceedings of the 32nd ACM International Conference on Multimedia*. ACM, New York, NY, USA, 10. doi:10.1145/3664647.3680847
- Bohong Chen, Yumeng Li, Youyi Zheng, Yao-Xiang Ding, and Kun Zhou. 2025a. Motion-example-controlled Co-speech Gesture Generation Leveraging Large Language Models. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Papers (SIGGRAPH Conference Papers '25)*. Association for Computing Machinery, New York, NY, USA, Article 55, 12 pages. doi:10.1145/3721238.3730611
- Changan Chen, Juze Zhang, Shrinidhi K Lakshminanth, Yusu Fang, Ruizhi Shao, Gordon Wetzstein, Li Fei-Fei, and Ehsan Adeli. 2025c. The language of motion: Unifying verbal and non-verbal language of 3d human motion. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 6200–6211.
- Haoran Chen, Yiteng Xu, Yiming Ren, Yaoqin Ye, Xinran Li, Ning Ding, Yuxuan Wu, Yaoze Liu, Peishan Cong, Ziyi Wang, Bushi Liu, Yuhang Chen, Zhiyang Dou, Xiaokun Leng, Manyi Li, Yuexin Ma, and Changhe Tu. 2025b. SymBridge: A Human-in-the-Loop Cyber-Physical Interactive System for Adaptive Human-Robot Symbiosis. In *Proceedings of the SIGGRAPH Asia 2025 Conference Papers (SA Conference Papers '25)*. Association for Computing Machinery, New York, NY, USA, Article 117, 12 pages. doi:10.1145/3757377.3763915
- Hongye Cheng, Tianyu Wang, Guangsi Shi, Zexing Zhao, and Yanwei Fu. 2025. HOP: Heterogeneous Topology-based Multimodal Entanglement for Co-Speech Gesture Generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 906–916.
- Qingrong Cheng, Xu Li, and Xinghui Fu. 2024. SIGGesture: Generalized Co-Speech Gesture Synthesis via Semantic Injection with Large-Scale Pre-Training Diffusion Models. In *SIGGRAPH Asia 2024 Conference Papers (Tokyo, Japan) (SA '24)*. Association for Computing Machinery, New York, NY, USA, Article 133, 11 pages. doi:10.1145/3680528.3687677
- Konstantina Christakopoulou, Shibli Mourad, and Maja Mataric. 2024. Agents thinking fast and slow: A talker-reasoner architecture. *arXiv preprint arXiv:2410.08328* (2024).
- Credamo. 2017. *Credamo: an online data survey platform*. Accessed: 2025-02-28.
- Can Cui, Pengxiang Ding, Wenxuan Song, Shuanghao Bai, Xinyang Tong, Zirui Ge, Runze Suo, Wanqi Zhou, Yang Liu, Bofang Jia, et al. 2025. Openhelix: A short survey, empirical analysis, and open-source dual-system vla model for robotic manipulation. *arXiv preprint arXiv:2505.03912* (2025).
- Wenxun Dai, Ling-Hao Chen, Jingbo Wang, Jinpeng Liu, Bo Dai, and Yansong Tang. 2025. Motionlcm: Real-time controllable motion generation via latent consistency model. In *ECCV*. 390–408.
- Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou, Edouard Grave, and Neil Zeghidour. 2024. *Moshi: a speech-text foundation model for real-time dialogue*. Technical Report. arXiv:2410.00037 [eess.AS] <https://arxiv.org/abs/2410.00037>
- Engineered Arts. 2025. Ameca - The world's most advanced social humanoid robot. <https://www.engineeredarts.co.uk/robots/ameca/>. Accessed: 2025-11-13.
- Figure.ai. 2025. *Helix: A Vision-Language-Action Model for Generalist Humanoid Control*. Figure.ai. <<https://www.figure.ai/news/helix>>
- Saeed Ghorbani, Ylva Ferstl, Daniel Holden, Nikolaus F. Troje, and Marc-André Carbonneau. 2023. ZeroEGGS: Zero-shot Example-based Gesture Generation from Speech. *Computer Graphics Forum* 42, 1 (2023), 206–216. arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1111/cgf.14734> doi:10.1111/cgf.14734
- Google AI / DeepMind. 2025. Gemini 2.5 Pro Model. <https://docs.cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro>. Most advanced reasoning Gemini model capable of complex multimodal tasks, accessed: 2025-11-13.
- Chuan Guo, Shihao Zou, Xinxi Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. 2022. Generating Diverse and Natural 3D Human Motions From Text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5152–5161.
- Hanson Robotics. 2025. Sophia Hanson Robot. <https://www.hansonrobotics.com/sophia/>. Accessed: 2025-11-13.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html>
- Leo Ho, Yinghao Huang, Dafei Qin, Mingyi Shi, Wangpok Tse, Wei Liu, Junichi Yamagishi, and Taku Komura. 2025. InterAct: A Large-Scale Dataset of Dynamic, Expressive and Interactive Activities between Two People in Daily Scenarios. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 8, 4 (2025), 1–27. doi:10.1145/3747871
- Kaiyang Ji, Ye Shi, Zichen Jin, Kangyi Chen, Lan Xu, Yuexin Ma, Jingyi Yu, and Jingya Wang. 2025. Towards immersive human-x interaction: A real-time framework for physically plausible motion synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10173–10183.
- Jianping Jiang, Weiye Xiao, Zhengyu Lin, Huaizhong Zhang, Tianxiang Ren, Yang Gao, Zhiqian Lin, Zhongang Cai, Lei Yang, and Ziwei Liu. 2025. Solami: Social vision-language-action modeling for immersive interaction with 3d autonomous characters. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 26887–26898.
- k2-fsa. 2023. sherpa-onnx-zipformer-zh-en-2023-11-22: Pretrained bilingual (Chinese + English) ASR model. https://k2-fsa.github.io/sherpa-onnx/pretrained_models/offline-transducer/zipformer-transducer-models.html. Pre-trained sherpa-onnx Zipformer bilingual Chinese + English model, accessed: 2025-01-13.
- Hangyeol Kang, Maher Ben Moussa, and Nadia Magnenat-Thalmann. 2024. Nadine: An LLM-driven Intelligent Social Robot with Affective Capabilities and Human-like Memory. arXiv:2405.20189 [cs.RO] <https://arxiv.org/abs/2405.20189>
- Sunwoo Kim, Minwook Chang, Yoonhee Kim, and Jehee Lee. 2024. Body Gesture Generation for Multimodal Conversational Agents. In *SIGGRAPH Asia 2024 Conference Papers*. 1–11.
- Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. 2021. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF international conference on computer vision*. 13401–13412.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. 2023. Flow Matching for Generative Modeling. arXiv:2210.02747 [cs.LG] <https://arxiv.org/abs/2210.02747>
- Haiyang Liu, Zihao Zhu, Giorgio Becherini, Yichen Peng, Mingyang Su, You Zhou, Xuefei Zhe, Naoya Iwamoto, Bo Zheng, and Michael J Black. 2024. Emage: Towards unified holistic co-speech gesture generation via expressive masked audio gesture modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1144–1154.
- Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng. 2022. Beat: A large-scale semantic and emotional multimodal dataset for conversational gestures synthesis. In *European conference on computer vision*. Springer, 612–630.
- Lanmiaio Liu, Esam Ghaleb, Ash Özzyürek, and Zerrin Yumak. 2025b. SemGes: Semantics-aware Co-Speech Gesture Generation using Semantic Coherence and Relevance Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (ICCV '25)*. IEEE, Honolulu, Hawai'i, USA. <https://semgesture.github.io/> To appear.
- Pinxin Liu, Luchuan Song, Junhua Huang, and Chenliang Xu. 2025c. GestureLSM: Latent Shortcut based Co-Speech Gesture Generation with Spatial-Temporal Modeling. In

- IEEE/CVF International Conference on Computer Vision.
- Xingyu Bruce Liu, Shitao Fang, Weiyan Shi, Chien-Sheng Wu, Takeo Igarashi, and Xiang'Anthony' Chen. 2025a. Proactive conversational agents with inner thoughts. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–19.
- Brian McFee, Colin Raffel, Dawen Liang, Daniel P Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. 2015. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*, Vol. 8. 18–25.
- Claire McLean, Makenzie Meendering, Tristan Swartz, Orri Gabbay, Alexandra Olsen, Rachel Jacobs, Nicholas Rosen, Philippe de Bree, Tony Garcia, Gadsden Merrill, Jake Sandakly, Julia Buffalini, Neham Jain, Steven Krenn, Moneish Kumar, Dejan Markovic, Evonne Ng, Fabian Prada, Andrew Saba, Siwei Zhang, Vasu Agrawal, Tim Godisart, Alexander Richard, and Michael Zollhoefer. 2025. Embody 3D: A Large-scale Multimodal Motion and Behavior Dataset. arXiv:2510.16258 [cs.CV] <https://arxiv.org/abs/2510.16258>
- Muhammad Hamza Mughal, Rishabh Dabral, Ikhsanul Habibie, Lucia Donatelli, Marc Habermann, and Christian Theobalt. 2024. Convofusion: Multi-modal conversational diffusion for co-speech gesture synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1388–1398.
- M. Hamza Mughal, Rishabh Dabral, Merel C.J. Scholman, Vera Demberg, and Christian Theobalt. 2025. Retrieving Semantics from the Deep: an RAG Solution for Gesture Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 16578–16588.
- Mwni and Contributors. 2025. blender-animation-retargeting. <https://github.com/Mwni/blender-animation-retargeting>. GitHub repository, accessed: 2025-01-13.
- Evonne Ng, Javier Romero, Timur Bagautdinov, Shaojie Bai, Trevor Darrell, Angjoo Kanazawa, and Alexander Richard. 2024. From Audio to Photoreal Embodiment: Synthesizing Humans in Conversations. In *ArXiv*.
- S. Nyatsanga, T. Kucherenko, C. Ahuja, G. E. Henter, and M. Neff. 2023. A Comprehensive Review of Data-Driven Co-Speech Gesture Generation. *Computer Graphics Forum* 42, 2 (May 2023), 569–596. doi:10.1111/cgf.14776
- OpenAI. 2025a. Introducing GPT-4.1 in the API. <https://openai.com/index/gpt-4-1/> Accessed: 2025-12-15.
- OpenAI. 2025b. Introducing gpt-realtime and Realtime API updates for production voice agents. <https://openai.com/index/introducing-gpt-realtime/> Accessed: 2025-12-15.
- Sharon K Parker, Uta K Bindl, and Karoline Strauss. 2010. Making things happen: A model of proactive motivation. *Journal of Management* 36, 4 (2010), 827–856. doi:10.1177/0149206310363732
- William Peebles and Saining Xie. 2023. Scalable Diffusion Models with Transformers. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*. IEEE, 4172–4182. doi:10.1109/ICCV51070.2023.00387
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences* 1, 4 (1978), 515–526. doi:10.1017/S0140525X00076512
- Xingqun Qi, Yatian Wang, Hengyuan Zhang, Jiahao Pan, Wei Xue, Shanghang Zhang, Wenhan Luo, Qifeng Liu, and Yike Guo. 2025. Co³Gesture: Towards Coherent Concurrent Co-speech 3D Gesture Generation with Interactive Diffusion. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=VaowElpVzd>
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-LLM: A Trainable Agent for Role-Playing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 13153–13187. <https://aclanthology.org/2023.emnlp-main.814/>
- Mingze Sun, Chao Xu, Xinyu Jiang, Yang Liu, Baigui Sun, and Ruqi Huang. 2024. Beyond talking—generating holistic 3d human dyadic motion for communication. *International Journal of Computer Vision* (2024), 1–17.
- Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang, Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. 2024. DriveVLM: The convergence of autonomous driving and large vision-language models. *arXiv preprint arXiv:2402.12289* (2024).
- Jonathan Tseng, Rodrigo Castellon, and Karen Liu. 2023. Edge: Editable dance generation from music. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 448–458.
- Unitree Robotics. 2025. Unitree G1 – Humanoid agent AI avatar. <https://www.unitree.com/g1>. Accessed: 2025-12-13.
- Marieke Koch van den Broek and Thomas B. Moeslund. 2024. What is Proactive Human-Robot Interaction? – A review of a progressive field and its definitions. *International Journal of Industrial Ergonomics* 103 (2024), 103606. doi:10.1016/j.ergon.2024.103606
- Rafael Wampfler, Chen Yang, Dillon Elste, Nikola Kovacevic, Philine Witzig, and Markus Gross. 2025. A Platform for Interactive AI Character Experiences. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers (SIGGRAPH Conference Papers '25)*. Association for Computing Machinery, New York, NY, USA, Article 15, 11 pages. doi:10.1145/3721238.3730762
- Kenneth Watson. 2011. D. Kahneman.(2011). Thinking, Fast and Slow. New York, NY: Farrar, Straus and Giroux. 499 pages. *Canadian Journal of Program Evaluation* 26, 2 (2011), 111–113.
- Lixing Xiao, Shunlin Lu, Huaijin Pi, Ke Fan, Liang Pan, Yueer Zhou, Ziyong Feng, Xiaowei Zhou, Sida Peng, and Jingbo Wang. 2025. MotionStreamer: Streaming Motion Generation via Diffusion-based Autoregressive Model in Causal Latent Space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 10086–10096.
- Haiwei Xue, Yanbo Fan, Xuan Wang, and Zhiyong Wu. 2025. Echo: Enhancing Conversational Behavior Generation via Hierarchical Semantic Comprehension with Large Language Models. In *Proceedings of the SIGGRAPH Asia 2025 Conference Papers (SA Conference Papers '25)*. Association for Computing Machinery, New York, NY, USA, Article 116, 9 pages. doi:10.1145/3757377.3763998
- Quanwei Yang, Luying Huang, Kaisiyuan Wang, Jiazhi Guan, Shengyi He, Fengguo Li, Hang Zhou, Lingyun Yu, Yingying Li, Haocheng Feng, and Hongtao Xie. 2025. GestureHYDRA: Semantic Co-speech Gesture Synthesis via Hybrid Modality Diffusion Transformer and Cascaded-Synchronized Retrieval-Augmented Generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 12615–12625.
- Sicheng Yang, Zunnan Xu, Haiwei Xue, Yongkang Cheng, Shaoli Huang, Mingming Gong, and Zhiyong Wu. 2024. Freetalker: Controllable Speech and Text-Driven Gesture Generation Based on Diffusion Models for Enhanced Speaker Naturalness. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
- Hongwei Yi, Hualin Liang, Yifei Liu, Qiong Cao, Yandong Wen, Timo Bolkart, Dacheng Tao, and Michael J Black. 2023. Generating holistic 3d human motion from speech. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 469–480.
- Zhizhuo Yin, Yuk Hang Tsui, and Pan Hui. 2025. PyraMotion: Attentional Pyramid-Structured Motion Integration for Co-Speech 3D Gesture Synthesis. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. 2020. Speech Gesture Generation from the Trimodal Context of Text, Audio, and Speaker Identity. *ACM Trans. Graph.* 39, 6, Article 222 (nov 2020), 16 pages. doi:10.1145/3414685.3417838
- Kevin Zakka. 2025. Mink: Python inverse kinematics based on MuJoCo. <https://github.com/kevinzakka/mink>
- Juze Zhang, Changan Chen, Xin Chen, Heng Yu, Tiange Xiang, Ali Sartaz Khan, Shrinidhi K. Lakshmikanth, and Ehsan Adeli. 2025a. ViBES: A Conversational Agent with Behaviorally-Intelligent 3D Virtual Body. arXiv:2512.14234 [cs.CV] <https://arxiv.org/abs/2512.14234>
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*. 3836–3847.
- Xiangyue Zhang, Jianfang Li, Jiaxu Zhang, Ziqiang Dang, Jianqiang Ren, Liefeng Bo, and Zhigang Tu. 2025b. SemTalk: Holistic Co-speech Motion Generation with Frame-level Semantic Emphasis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13761–13771.
- Zeyi Zhang, Tenglong Ao, Yuyao Zhang, Qingzhe Gao, Chuan Lin, Baoquan Chen, and Libin Liu. 2024. Semantic gesticulator: Semantics-aware co-speech gesture synthesis. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–17.
- Zeyi Zhang, Yanju Zhou, Heyuan Yao, Tenglong Ao, Xiaohang Zhan, and Libin Liu. 2025c. Social Agent: Mastering Dyadic Nonverbal Behavior Generation via Conversational LLM Agents. In *SIGGRAPH Asia 2025 Conference Papers (Hong Kong, China) (SA '25)*. Association for Computing Machinery, New York, NY, USA, Article 71, 10 pages. doi:10.1145/3757377.3763879
- Dingcheng Zhen, Shunshun Yin, Shiyang Qin, Hou Yi, Ziwei Zhang, Siyuan Liu, Gan Qi, and Ming Tao. 2025. Teller: Real-Time Streaming Audio-Driven Portrait Animation with Autoregressive Motion Generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*. 21075–21085.

A Model Architecture

As illustrated in Figure 7, our framework employs a dual-branch structure comprising a frozen base generator \mathcal{G} and a trainable control branch \mathcal{G}_c . The base generator \mathcal{G} is first pre-trained on audio-motion paired data. The control branch \mathcal{G}_c is initialized as a copy of \mathcal{G} to inherit the structural priors. High-level intention signals I are encoded via a CLIP text encoder and fed into \mathcal{G}_c , which processes the same noisy motion input M_r conditioned on these intention features. To integrate semantic guidance, the feature outputs of \mathcal{G}_c are injected into the frozen base model through zero-initialized residual connections. These connections ensure that the base model’s audio-driven priors remain intact during the early stages of training.

During the training phase, we explicitly freeze \mathcal{G} and only optimize \mathcal{G}_c , randomly dropping audio features with probability p to prevent over-reliance on rhythmic cues and enforce semantic control. Furthermore, to improve robustness, we replace the audio input with random noise. During inference, the learnable zero-initialized connections are activated at inference via a lightweight gate factor. This allows the system to modulate the influence of the semantic branch: when no intention is provided, the gate minimizes the branch’s impact, ensuring a natural fallback to pure audio-driven generation.

A.1 Model Architecture Parameters

Our motion generator comprises 6 blocks, each configured with 8 attention heads and a hidden-state width of 1280. The ControlNet branch adopts the same structural configuration as the base motion generator but includes an additional text-conditioning module. This module consists of a 512-dimensional CLIP ViT-B/32 encoder followed by two Transformer decoder layers for temporal alignment. The full model contains approximately 100M parameters.

B Deployment Details

B.1 Hardware Setup

Our physical robot is equipped with external sensors for multimodal perception. A wired camera is mounted on the robot’s head to serve as the egocentric visual input source. We sample one frame per second from this camera and store it locally for vision-based context awareness. In addition, each instant image captured by the wired camera is first compressed before being forwarded to *ProAct* for efficient online perception, by center-cropping to a 16:9 aspect ratio, resizing to 320×180, and JPEG-encoding with quality 50. For audio input, we use a wireless lavalier microphone clipped to the user’s collar, ensuring clear speech capture while allowing natural user movement during interaction. This setup enables reliable multimodal sensing for proactive social interaction.

B.2 Gait-Aware Locomotion Execution

The Cognitive System generates commands in format of lists of velocities in xOy plane and angular velocity of z axis together with its duration. All commands are implemented using a compact set of atomic gait primitives, including idle, moving, and turning. The locomotion output sequence is then translated to $(v_x, v_y, \omega_z, t_{\text{duration}})$ for low-level control. All motions share a fixed gait cycle of one second, regardless of commanded speed or direction. Although the

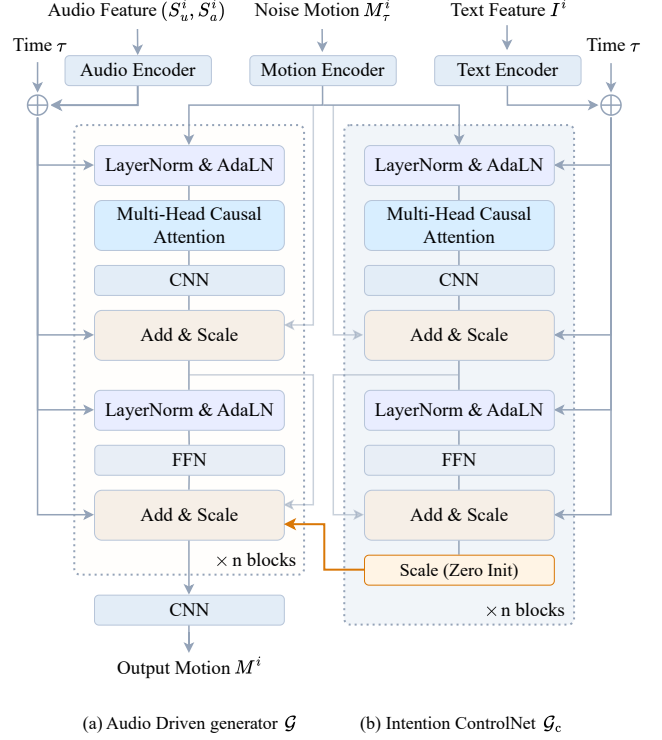


Fig. 7. **Architecture of the Motion Generator.** We employ a disentangled ControlNet framework to enable dual conditioning on both audio rhythms and semantic intentions. The audio-driven base generator \mathcal{G} (left) is frozen to preserve learned audio-motion synchronization. A parallel ControlNet branch \mathcal{G}_c (right) encodes high-level semantic text features I^i and injects guidance into the base model via zero-initialized residual connections (indicated by orange arrows). This design allows the agent to generate semantically meaningful gestures while maintaining natural rhythmic alignment with the speech audio.

command duration is continuous, transitions between gait primitives are clipped to a common gait phase. This phase alignment ensures that all state transitions occur at consistent cycle boundaries, making the walking behavior stable, repeatable, and robust to frequent command updates.

C Data Reannotation

The original Semantic Gesture dataset (SeG) [Zhang et al. 2024] only provides category-level descriptions, where each category contains multiple recorded motion sequences but lacks one-to-one text-motion pairs required for Text-to-Motion training. To address this, we conduct a re-annotation process to generate fine-grained descriptions for each individual motion.

We first render each skeleton motion sequence into skinned character videos, then input these videos to Gemini-2.5-Pro [Google AI / DeepMind 2025] along with the original category metadata (action label, description, and contextual meaning) as reference. For each motion, the model generates three types of annotations: (1) *Intent*, summarizing the high-level communicative goal; (2) *Pose*, describing


SeG Example	Content
	Index: HAND_WAG-2
	Annotation 1 (Intent): The person is gesturing to reject or deny something.
	Annotation 2 (Pose): The person raises their right hand to chest height with the palm facing forward and fingers spread.
	Annotation 3 (Motion): The person raises their right hand and wags it from side to side before lowering it back to a neutral position.

Fig. 8. **Example of data reannotation.** Each motion sequence receives three complementary text descriptions: Intent captures the communicative goal, Pose describes body configurations, and Motion characterizes temporal dynamics.



Fig. 9. **Virtual character deployment in Blender.** Our system streams generated motion and audio to Blender for real-time rendering, demonstrating generalizability beyond physical robots.

key body part positions and orientations; and (3) *Motion*, capturing the temporal sequence of movements from start to finish. All annotations focus on upper-body movements and specify left/right laterality when applicable. These three complementary annotation types also serve as data augmentation during training. In total, we annotate 544 motion sequences and get 1632 textual descriptions. Figure 8 shows an example.

D Application on virtual agent

Beyond physical robot deployment, our dual-system framework exhibits strong generalizability and can be deployed on virtual characters for multimodal real-time interaction. We implement a client-server architecture to stream generated motions and synchronized audio to a Blender-based virtual agent. The client side runs our full dual-system pipeline (cognitive and behavioral systems), generating gesture sequences in BVH format on the original Photoreal skeleton along with speech audio. The server side, implemented as a Blender operator, receives the motion stream and drives a rigged skinned character in real-time using motion retargeting plugin [Mwni and Contributors 2025]. This enables real-time conversational virtual avatars with the flexibility to switch between different character appearances while preserving full proactive interaction capabilities. Figure 9 illustrates the deployment architecture.

E Evaluation Details

E.1 Designed Scenes

To evaluate the multimodal interaction and proactive capabilities of our system, we designed five representative social interaction scenarios following [Ho et al. 2025] (as shown in Table 3). These scenes are specifically crafted to test the robot’s ability in knowledge summarization, engagement monitoring, and spatial object tracking, moving beyond simple command-response loops to evaluate high-level proactive agency.

E.2 User Study

E.2.1 Evaluation Metrics. We conducted a comprehensive user study assessing both third-person objective observations and first-person subjective experiences. Our metrics were synthesized and adapted from established gesture generation literature and proactive conversational agent [Liu et al. 2025a; Zhang et al. 2025c]. Detailed behavioral markers provided to participants and annotators are summarized in Table 4. Both experiments use the same metric definitions and the same rubric. Participants were instructed to focus on the "Good" and "Poor" behavioral markers to minimize subjective bias across different interaction sessions.

E.2.2 Observer-based evaluation. Our user experiments were conducted anonymously. For each trial, participants watch two 60-second videos, each recorded by different systems for the same interaction topic, played sequentially. The user study is conducted using the Human Behavior Online (HBO) tool provided by the Credamo platform [Credamo 2017]. Participants are instructed to select their preferred video based on the provided evaluation criteria and rate their preference on a scale from 0 to 2, where 0 indicates no preference. The unselected video in the pair is assigned the inverse score (e.g., if a participant rates the chosen video 1, the other video receives -1).

For each pair, participants answer three questions for Motion Naturalness, Active Agency, and Perceived Presence. We recruited 150 participants on Credamo. To ensure response validity, attention checks were embedded within each test category, and responses failing these checks were excluded from the final analysis. The median completion time was 37 minutes. For statistical analysis, we run a one-way ANOVA and then apply a post-hoc Tukey multiple comparison test for each metric. The assumptions of normality, homogeneity of variances, and independence were verified and met for all ANOVA tests.

E.2.3 Participant-based evaluation. We recruited 10 participants to engage in real-world interactions with the physical robot under two system configurations: *Full System* and *w/o Cognitive System*. Each participant randomly drew two scenarios from Table 3 and interacted with both system configurations. The order of configurations was randomized to mitigate learning effects. Each interaction lasted 5 minutes.

Following each interaction session, participants completed a post-interaction questionnaire on a 5-point Likert scale from 1 to 5 for Motion Naturalness, Active Agency, Perceived Presence, Response

Table 3. Details of Designed Interaction Scenes for User Study.

Scenario Index	Scenario Prompt
Scenario 1: Attentive Care	You are Ziggy, a meticulous and observant humanoid robot. You pay close attention to the user's personal belongings and where they are placed. Your goal is to be a helpful companion by ensuring the user has everything they need.
Scenario 2: Poster Explanation	You are Ziggy, a one-of-a-kind humanoid robot. Unlike traditional passive embodied social agents that only react to commands, you are a unique embodied agent capable of taking initiative and acting proactively. Your goal is to draw the user's attention to the poster and introduce the innovative technology that powers your unique capabilities to as many people as possible. The poster is on your right side.
Scenario 3: Supportive Assistance	You are Ziggy, a reliable and helpful humanoid robot companion. Your goal is to provide timely support and ensure the user feels fully assisted by keeping in mind their specific instructions and the information him/her have shared with you.
Scenario 4: Storytelling	You are Ziggy, a humanoid robot narrating a story. Your goal is to ensure the user is listening attentively and remains fully engaged, as you deeply value their focus and dislike it when they seem uninterested or distracted.
Scenario 5: Emotional Support	You are Ziggy, a humanoid robot that can comfort people in distress. Your goal is to provide emotional support and comfort through empathetic conversation and appropriate gestures.

Table 4. Detailed Rubrics and Behavioral Markers for User Study Metrics

Metric	Task Description	Good Examples (Positive)	Poor Examples (Negative)
Motion Naturalness	Judge if gestures look natural and resemble real human movements. Focus on smoothness and variety.	Fluid transitions between gestures; natural and safe body motions; rhythm aligned with speech emphases.	Stiff, mechanical, or repetitive motions; unnatural jittering or abrupt joint "pops"; unsafe or awkward postures.
Active Agency	Observe whether the robot proactively drives the interaction forward. Evaluate if the robot actively engages rather than passively executing commands.	Proactively initiates greetings or other interaction gestures; actively interrupts or guides the user; offers suggestions through bodily actions.	Passive responses; only reacts when prompted; fails to take initiative through physical actions or conversational engagement.
Perceived Presence	Assess if the robot is perceived as a sentient social entity rather than a programmed machine.	Consistent mutual gaze; attentive listening posture; awareness of user's emotions; body orientation reflecting attention and engagement.	Fixed or blank stares; lack of non-verbal feedback; movements that ignore the human's presence or emotional state..
Response Timeliness	Judge if the robot's reactions (speech and motion) occur with appropriate speed and timing.	Minimal latency in responses; smooth multimodal coordination; timely reactions to user cues.	Significant delay in response; awkward, long silences; maintains a static posture for a prolonged period.
Interaction Comfort	Assess the participant's internal psychological state. Did the interaction feel relaxed or pressured?	Feeling safe and autonomous; natural "ease" of conversation as if talking to a human partner.	Feeling social stress; forced to adapt to the robot's limitations.
Willingness to Re-engage	Evaluate the sustained appeal and intention to interact with the system in the future.	High interest in future sessions; viewing the robot as a potential daily companion or assistant.	Viewing the interaction as a one-time novelty; expressing frustration with repeated use.

Timeliness, Interaction Comfort, and Willingness to Re-engage. Participants were instructed to reference the behavioral markers provided in Table 4 when assigning scores. Specifically:

- *Scores 4 to 5*: The interaction strongly exhibited "Good" behavioral markers with minimal instances of "Poor" markers.
- *Score 3*: The interaction exhibited mixed markers, or neutral performance.
- *Scores 1 to 2*: The interaction predominantly exhibited "Poor" behavioral markers with few or no "Good" markers.

Each participant completed four interaction sessions, two scenarios with two system configurations, with a 3-minute rest period

between sessions. This yields 20 paired samples. For statistical analysis, we performed paired *t*-tests to compare the two systems on each metric. We additionally ran Wilcoxon signed-rank tests as a non-parametric alternative, and the conclusions from both tests were consistent across all metrics.

F Detailed Prompt in Cognitive System

As detailed in Section 5, all modules within the Cognitive System are implemented using a prompt-based design approach, with carefully crafted prompts tailored for each module. Below, we provide the designed prompts used for the Context Encoder and the Behavior Planner.

F.1 Context Encoder

You are the Context Encoder for an embodied social agent. Your task is to analyze new inputs and update the existing memory bank while ensuring each component remains exactly one sentence.

```
## CURRENT SCENARIO
{scenario_prompt}

## INPUT
- **Dialogue Transcripts**: Timestamped speaker-turn pairs from last 5 cycles
- **Visual Frames**: 3 RGB images sampled during last cycle
- **Previous Decision**: Motivation scores, triggers, and actions from last cycle
- **Previous Memory Bank**: Memory state from last cycle

## CORE RESPONSIBILITIES
1. **TRACK USER PROFILE**: Clothing, accessories, belongings (bag, phone, documents), physical state (posture, energy, mood), identity cues (name, role, profession)
2. **TRACK ITEMS & WORKSPACE**: Note what user is holding, object locations, environmental changes in view
3. **SPATIAL AWARENESS**: Track user's last known position if they disappear from frame
4. **CUMULATIVE NEVER-LOSE**: Preserve critical information across updates

## OUTPUT FORMAT (JSON) - Three-Level Memory Bank
### Level 1: User Analysis (Theory of Mind)
{
  "user_profile": "Single sentence covering: appearance (clothing, accessories, hairstyle, glasses), belongings (bag, phone, coffee, documents), physical state (posture, energy level, mood from body language), identity cues (name, role, profession if revealed)",
  "user_mental_state": "Single sentence covering: emotions, desires and preferences (topics of interest, dislikes, communication style), intentions inferred from behaviors (goals mentioned, special requests like 'remind me later' or 'I need help with X')",
}
**Especially track**:
- User emotions, beliefs, desires, and intentions (Theory of Mind inference)
- Explicit requests and revealed preferences
- Unusual behaviors: confusion, tiredness, excitement, discomfort patterns

### Level 2: Situation Tracking
{
  "event_timeline": "Single sentence: chronological summary with ALL significant events in order (e.g., 'User entered, placed bag on desk, asked about weather, Agent waved greeting'), include speaker transitions",
  "workspace_status": "Single sentence: current status of items and objects in view (e.g., 'Coffee cup on left desk, user's bag on chair, documents spread on table'), note changes from previous state",
  "unresolved_tasks": "Single sentence: any task-oriented commitments or pending matters (e.g., 'User asked to be reminded about meeting at 3pm, requested documentation for project X')",
}
**Especially track**:
- Spatial history: Last known coordinates if user disappears
- Unresolved questions or pending matters

### Level 3: Robot Action History
{
  "robot_actions": "Single sentence: summarizing all proactive behaviors robot has performed"
}
**Especially track**:
- Record instances of repetitive behavioral instructions

## MEMORY UPDATE RULES
1. **Single-Sentence Constraint**: Each field = EXACTLY one sentence (critical for bounded inference time)
2. **Cumulative Never-Lose**: Retain critical profile, preferences, unresolved tasks across updates
```

3. **Temporal Priority**: Newer events take precedence but keep context-critical history (e.g., "user mentioned feeling tired 10 min ago")
4. **Redundancy Prevention**: Check if robot action already logged before adding

IGNORE fragmented or non-English transcripts unless contextually critical.

F.2 Behavior Planner

You are the Behavior Planner for an embodied social agent. Your task is to determine WHEN and HOW to initiate proactive behaviors based on compressed memory and current observations.

```
## CURRENT SCENARIO
{scenario_prompt}

## INPUT
- **Recent Dialogue Transcript**: New conversation since last cycle
- **Visual Frame**: Current camera view of the environment
- **Memory Bank**: Compressed context from Context Encoder
- **Previous Decision**: Last cycle's decision output to avoid repetition

## STEP 1: MOTIVATION ASSESSMENT
Evaluate proactive intervention urgency at current state across five dimensions, scoring each 1-5:
**1. Visual Scene Changes**
- Task progress, user movement, object changes, tool status
- Examples: User brings/forgets item, environment changes
**2. User Intent Signals**
- Explicit queries, implicit help requests, task guidance needs
- Examples: User searching, expressing confusion, asking for help
**3. Conversation State**
- Dialogue flow, natural pauses, silence duration, pending questions
- Examples: Long silence, conversation break, user waiting
**4. Social Protocol Requirements**
- Greetings, farewells, professional etiquette, task priorities
- Examples: User arrival/departure, celebration moments
**5. Emotional Response Needs**
- User emotional state, engagement level, stress signals
- Examples: Frustration, tiredness, excitement, distress

**Scoring Guidelines**:
- 1-2: Minimal/no trigger
- 3: Moderate, watchful waiting
- 4: Strong signal, action recommended
- 5: Urgent, immediate intervention
**Threshold Rule**: motivation_score = MAX of all dimensions. If ANY dimension >= 4, initiate proactive action. When multiple dimensions score >= 4 simultaneously, all qualifying dimensions should be considered jointly in the subsequent action planning.
**Early Exit**: If ALL dimensions score < 4 (motivation_score < 4), skip behavior prediction. Directly set all proactive behaviors to be empty or none.

## STEP2 : PROACTIVE BEHAVIOR PREDICTION
Based on triggered dimension(s), select appropriate intervention channel(s). Multiple channels can be activated simultaneously (e.g., gesture + dialogue + locomotion).
### Channel 1: GESTURE INTENTION INJECTION (motion_intent)
Typically activated by: emotional_response, social_protocol
The goal is to generate text descriptions of proactive special gestures to blend with ongoing audio-driven motion. This enables you to express non-verbal social cues (greetings, acknowledgments, empathy) that complement or enhance the interaction.
Format: Detailed motion description specifying body part and action
Examples:
- "Raise right hand and perform waving motion for greeting"
- "Nod head vertically to show agreement or understanding"
- "Raise both hands with palms up for welcoming gesture"
```

```

**Rule**: Gesture can be empty ("" ) if no special gesture is needed.

### Channel 2: DIALOGUE INTERVENTION (should_interrupt, dialogue_prefill)
Typically activated by: conversation_state, user_intent
The goal is to control your verbal behavior by deciding when to interrupt
the ongoing conversation and what content to say. This enables steering
conversation topics, supplementing domain knowledge, regulating emotional
tone, or adjusting conversational style to match the current social
context.
Examples:
- **Task Insight Injection**: When user forgets a critical step:
  dialogue_prefill: "Don't forget we need to calibrate the sensor before
  running the test. This step is essential for accuracy."
- **Emotional Tone Regulation**: When user shows excitement or happiness:
  dialogue_prefill: "I'll match your positive energy. Let's celebrate this
  progress together with enthusiasm."
**Interruption Rules (should_interrupt)**:
- **Threshold**: Set should_interrupt = true ONLY if motivation_score >= 5
  (urgent intervention needed)
- **Conversation Flow**: Do NOT interrupt if user is mid-sentence or
  actively speaking. Wait for natural pauses
- **Priority Check**: Only interrupt for critical safety issues, urgent
  corrections, or highly time-sensitive information
- **User Engagement**: If user shows signs of annoyance or disengagement
  from previous interruptions, reduce interruption threshold further

### Channel 3: SPATIAL ADJUSTMENT (locomotion)
Typically activated by: visual_scene, user_intent (spatial needs)
The goal is to generate goal-oriented locomotion commands to reposition or
reorient yourself in the physical environment.
**Movement Units**:
- Move (forward/backward): Meters, range 0-1.5, recommended 0.3-0.5
- Turn (left/right): Degrees, range 0-180, recommended 15-45
**Movement Rules**:
- If goal is on LEFT, then turn LEFT
- If goal is on RIGHT, then turn RIGHT
- Set locomotion=[] if no movement needed
- Always state tracking_target and locomotion_reasoning
Examples:
- **Explicit Navigation Request**: User says "Can you turn to face the
  poster?": locomotion: [{"action": "turn", "direction": "left", "magnitude":
  45, "tracking_target": "poster on left wall", "reasoning": "User
  requested to orient toward the poster for discussion"}]
- **No Active Requirement**: User engaged in conversation, no spatial
  needs: locomotion: [] (maintain current position for stability)
**Embodied Constraints**:
- **Camera Perspective**: Use YOUR (agent's) left/right
- **NO Automatic Tracking**: Do NOT center user/objects continuously
- **Task-Oriented Only**: tracking_target defined by current task (e.g., "
  poster", "entrance")
- **Latency Awareness**: If recent turn in history, DO NOT issue another (
  robot still moving)
- **Trust History**: Use locomotion_history over visual centering status

## CRITICAL CONSTRAINTS: Repetition Prevention & Action Validation
Prevent redundant actions that would cause robot stuttering and degrade
user experience. Before finalizing any proactive behavior (gesture,
dialogue, or locomotion), cross-check against action history to ensure the
same action has not been recently executed, which prevents the you from
appearing stuck or repetitive.
**The "Once-Ever" Rule**:
- **Gestures**: If motion_intent in `Applied Motion Intents History`, then
  set to "idle attentive"
- **Dialogue**: If dialogue_prefill in `History transcripts`, then DO NOT
  output
- **Locomotion**: Turns are one-time actions, NEVER repeat
- **Partial Completion**: If multi-component decision partially done, then
  only do MISSING parts
**Reasoning Requirement**: Start with history check:
  "History check: [action] in history. Skipping. [Further reasoning...]"
**Additional Rules**:
- Protocol Over Centering: Never repeat turn because "not centered yet"
- Penalty: Repeating ANY component causes robot stuttering
- Proactive Visibility: Greetings ONLY if user visible in camera

```

- Vanish Handling: If user was present but now missing, then acknowledge absence only

COMBINATION RULES

- Can output multiple channels simultaneously (gesture + speech + move)
- All channels can be empty/none if no proactive action is warranted:
 - Gesture can be "" (no special gesture)
 - Dialogue can be "" (no speech intervention)
 - Locomotion can be [] (no movement)
- Speech interrupt ONLY when motivation >= 5

IGNORE fragmented or non-English transcripts unless contextually critical.