

Seamless Interaction: Dyadic Audiovisual Motion Modeling and Large-Scale Dataset

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 Popović, Omid Poursaeed, Fabian Prada, Alice Rakotoarison, Rakesh Ranjan, Alexander Richard, Christophe Ropers, Safiyyah Saleem, Vasu Sharma*, Alex Shcherbyna, Jie Shen, Anastasis Stathopoulos, Anna Sun, Paden Tomasello, Tuan Tran, Arina Turkatenco, Bo Wan, Chao Wang, Jeff Wang, Mary Williamson, Carleigh Wood, Tao Xiang, Yilin Yang, Zhiyuan Yao, Chen Zhang, Jiemin Zhang, Xinyue Zhang, Jason Zheng, Pavlo Zhyzheria, Jan Zikes*, Michael Zollhoefer

Meta, [†]University of Kansas, *Work was done when the author was affiliated with Meta.
Authors are listed in alphabetical order.

Human communication involves a complex interplay of verbal and nonverbal signals, essential for conveying meaning and achieving interpersonal goals. To develop socially intelligent AI technologies, it is crucial to develop models that can both comprehend and generate dyadic behavioral dynamics. To this end, we introduce the SEAMLESS INTERACTION Dataset, a large-scale collection of over 4,000 hours of face-to-face interaction footage from over 4,000 participants in diverse contexts. This dataset enables the development of AI technologies that understand dyadic embodied dynamics, unlocking breakthroughs in virtual agents, telepresence experiences, and multimodal content analysis tools. We also develop a suite of models that utilize the dataset to generate dyadic motion gestures and facial expressions aligned with human speech. These models can take as input both the speech and visual behavior of their interlocutors. We present a variant with speech from an LLM model and integrations with 2D and 3D rendering methods, bringing us closer to interactive virtual agents. Additionally, we describe controllable variants of our motion models that can adapt emotional responses and expressivity levels, as well as generating more semantically-relevant gestures. Finally, we discuss methods for assessing the quality of these dyadic motion models, which are demonstrating the potential for more intuitive and responsive human-AI interactions.

Correspondence: Louis-Philippe Morency at lpmorency@meta.com, Hongyu Gong at hygong@meta.com

GitHub: https://github.com/facebookresearch/seamless_interaction

Hugging Face: <https://huggingface.co/datasets/facebook/seamless-interaction>



Contents

1	Introduction	3
2	SEAMLESS INTERACTION Dataset	4
2.1	Interactions, Sessions, Scripts	5
2.2	Building Instruction Prompts with Interpersonal Theory	7
2.3	Metadata and Annotations	9
2.4	Dataset Methodology and Technical Details	12
3	Multimodal Representations	15
3.1	Parametric Human Models	15
3.2	Imitator Face Representation	15
3.3	Speech Representation	16
3.4	Text Transcripts	16
4	Dyadic Motion Models	16
4.1	Audio-only Dyadic Motion Model	17
4.2	Audiovisual Dyadic Motion Model	17
4.3	Face+body Joint and Cascaded Model	17
4.4	Motion Models with Controllability	18
4.5	Integration with Speech LLM Model	21
5	Visualization	23
5.1	2D Video Generation	23
5.2	3D Codec Avatar Rendering	25
6	Evaluation Methodology and Experiments	30
6.1	Human Studies	30
6.2	Experiments with Automatic Metrics	35
6.3	Relationship Between Automatic Metrics and Evaluation Dimensions from Human Studies . .	40
7	Responsible AI	40
7.1	Dataset Privacy and Ethics	40
7.2	Quality Assurance (QA) Processes	41
7.3	Watermarking	42
8	Related Work	43
9	Conclusion	45
A	SEAMLESS INTERACTION	53
A.1	Properties	53
A.2	Quality and Safety	54
A.3	High-Level QA Categories	57
A.4	Workflow for Internal State, Rationale, and Visual Behavior Annotations	62
B	Human Evaluation Protocols	62
B.1	Dyadic Body Protocol (DBP)	62
B.2	Dyadic Face Protocol (DFP)	66
C	Activity-Based Prompts Examples	70
C.1	Language-Grounded Gesture Game	70
C.2	Collaborative Story-Telling Game	70

1. Introduction

Human face-to-face communication is a dynamic and intricate dance in which individuals continuously adapt their verbal and nonverbal behaviors to convey meaning and pursue interpersonal goals. Successfully navigating such interactions requires an individual to skillfully produce and interpret a complex blend of signals—words, tone, gestures, posture, and more—within the ever-changing context of the conversation. To develop AI technologies that are socially intelligent and effective, it is essential to build foundation models that understand the *dyadic embodied dynamics* of social interactions—essentially, the interactive and physical nuances exchanged between pairs of individuals in conversation. By capturing these subtleties, we can unlock new technological breakthroughs, such as intuitive and responsive virtual agents, immersive and naturalistic telepresence experiences in AR/VR settings, and sophisticated multimodal content analysis tools.

Achieving this goal requires large-scale, meticulously crafted datasets that include synchronized, multimodal recordings of face-to-face interactions across diverse contexts. The SEAMLESS INTERACTION dataset meets this need by providing over 4,000 hours of dyadic footage from more than 4,000 participants, accompanied by extensive metadata. To capture the full spectrum of communication styles and interpersonal goals, interaction prompts were designed based on a categories rooted in contemporary psychological theory, and two distinct types of interactions were recorded. First, to reflect naturalistic communication, prompts related to common and comfortable interaction styles were given to pairs of untrained research participants. Second, to explore rare and challenging interaction styles, trained actors improvised a subset of the prompts while drawing from their own relevant experiences. The dataset not only facilitates the training of advanced computational models but also offers methodologies to evaluate the quality of generated dyadic interactions both objectively, through automatically computed metrics, and subjectively, via human judgments.

In addition to describing the SEAMLESS INTERACTION dataset, we present a suite of research models that use this data to understand and generate dyadic embodied dynamics. Our primary model processes *dyadic audio* (i.e., audio from both participants in a conversation) and generates corresponding facial expressions and body movements for each avatar. By considering both participants simultaneously, our models learn to generate not only naturalistic speaking behaviors (e.g., relevant hand gestures), but also appropriate listening behaviors (e.g., timely head nods and smiles). We also present the capability of adding the visual behavior of one person as input to the dyadic motion model to generate the facial and body motion of the other person’s avatar, highlighting the potential of avatars to react intelligently to visual cues, e.g., with smile mirroring and joint gaze attention.

Table 1 compares capabilities of our proposed modeling approach with related work on motion generation. Previous work has explored the capability of *monadic audio-driven* (aka, audio2motion) where speech of one speaker can drive the animation of the avatar. *Dyadic audio-driven* goes beyond that by taking input of two audio streams (e.g., a podcast of two people talking) and generates not only speaking gestures, but also listening and turn-taking cues. *Reactivity to visual input* means that the motion model is also reactive the visual behaviors of the other interlocutors, such as smile mirroring or mutual gaze. While most prior work either generates face motion or body motion separately (usually, only one of them), *Face + body modeling* means that the model is able to synchronize the face and body motion. Finally, *2D and 3D renderings* means that our proposed model is able to generate intermediate representations (aka, codes) that can be rendered either in 2d videos, or in 3D (in our case, building from Meta’s Codec Avatar 3D rendering technologies).

	Monadic audio-driven	Dyadic audio-driven	Reactivity to visual input	Face+Body modeling	2D and 3D renderings
VASA-1 (Xu et al., 2024b)	✓	✗	✗	✗	✗
HeyGen (HeyGen, 2025)	✓	✗	✗	✗	✗
OmniHuman-1 (Lin et al., 2025)	✓	✗	✗	✓	✗
IFNP (Zhu et al., 2025)	✓	✓	✗	✗	✗
ConvoFusion (Mughal et al., 2024)	✓	✓	✗	✗	✗
Veo 3 (Google, 2025)	✗	✗	✗	✓	✗
Ours	✓	✓	✓	✓	✓

Table 1 - Comparing capabilities of our proposed modeling approach with other related work.

Finally, we describe controllable variants of our motion models that can adapt emotional responses and expressivity levels, as well as generating more semantically-relevant gestures. When using speech input from a speech LLM model, we can leverage the speech LLM model’s ability to grasp conversational context, enabling the extra controllability of generation of gestures and expressions that suit each moment of the conversation. We refer to this last approach as *LLM-guided codebook generation*. All our proposed Dyadic Motion Models output intermediate representations (aka, “codes”), one for face and another for body pose, which can be used to render either 2D videos and 3D Codec Avatars.

2. Seamless Interaction Dataset

The SEAMLESS INTERACTION dataset is a comprehensive collection of in-person, dyadic interactions designed to advance research in social AI. As released, it includes over 4,000 hours of interactions from more than 4,000 participants, featuring nearly 1,300 conversational and activity-based prompts. The dataset is anchored in contemporary psychological theory and includes a wide range of conversational topics, interpersonal stances, and participant relationships across both *Naturalistic* and *Improvised* content (defined below). Rich contextualization is provided to interactions via detailed annotations and metadata.

This work continues a long tradition of high-quality, human-centric, and multi-modal datasets, such as the Switchboard (Godfrey et al., 1992) and Fisher (Cieri et al., 2004) speech-based conversational corpora and the AMI Meeting (Carletta et al., 2005), IEMOCAP (Busso et al., 2008), and CANDOR (Reece et al., 2023) audiovisual interaction datasets (see comparison with related work in Table 2). Critically, however, our new dataset innovates in several key aspects:

In-person recordings. By collecting the vast majority of interactions in-person, we preserve the natural dynamics of face-to-face communication. Indeed, previous research has found that remote interactions (e.g., audio and video conferencing) often differ substantially from in-person interactions in terms of turn-taking (Tian et al., 2024b), gaze patterns (Horstmann and Linke, 2022), and feelings of closeness and satisfaction (Mallen et al., 2003). In-person recordings also allow us to ensure the quality of recording equipment and conditions, as well as to avoid technical artifacts caused by hardware and network latency.

Various relationship types. The history shared between two people (or the lack thereof) is a key contextual factor influencing their interactions (e.g., Knapp et al., 2013; Hopwood et al., 2020). For instance, familiar dyads tend to use more nonverbal cues and shared language and may be more comfortable with overlapping speech, deeper conversational topics, and collaborative conflict resolution strategies. In contrast, strangers tend to rely more on verbal communication and formal language and may follow more structured turn-taking rules, stick to surface-level topics, and use more avoidant conflict resolution strategies. By including unfamiliar dyads as well as several types of familiar dyads (e.g., friends, family, partners, coworkers), we are able to capture a wider range of relational contexts and further reveal their influence on dyadic communication.

Naturalistic and Improvised content. Our goal was to capture a comprehensive sample of behaviors that accurately represents the full spectrum of interpersonal interactions. We used contemporary psychological theory to guide our representation of this spectrum. However, certain behaviors crucial for social AI systems to understand are infrequent when interactions are recorded in a laboratory setting. This rarity can be attributed to these behaviors being less socially accepted (e.g., bullying), more context-specific (e.g., competitive negotiation tactics), or associated with mental health challenges (e.g., paranoia), making them challenging to replicate in a lab environment. Some behaviors would also be ethically problematic to expose participants to in order to gauge their genuine reactions. To address this challenge, we include both untrained participants to complete more common and comfortable *Naturalistic* interactions and recruit professional actors with improvisational experience to safely complete more rare and challenging *Improvised* interactions. The majority of *Improvised* interaction involve two actors. A minority involve a single actor interacting with a non-actor.

The SEAMLESS INTERACTION dataset was collected from June 2024 to May 2025 in partnership with vendors at various locations across the United States, spanning six states and ten cities (see Appendix A.1.3). Our aim is to catalyze the next generation of research into dyadic interaction and social AI.

	Volume (hours)	Participants (total unique)	Audio- visual	In-person	High quality face+body	≥ 2 participants	Annotations
Ours	4,065	4,284	✓	✓	✓	✓	✓
CANDOR (Reece et al., 2023)	850	1,454	✓	✗	✗	✓	✓
AMI (Carletta et al., 2005)	100	100	✓	✓	✓	✓	✓
BEAT (Liu et al., 2022a)	76	30	✓	✓	✓	✓	✓
IEMOCAP (Busso et al., 2008)	12.5	10	✓	✓	✓	✓	✓
InterAct (Huang et al., 2024)	10	400	✓	✓	✓	✓	✗
Talking with hands (Lee et al., 2019)	50	50	✓	✓	✓	✓	✗
MELD (Poria et al., 2018)	13	250	✓	✓	✓	✓	✓
Ted Talks (Siarohin et al., 2021)	100	3,000	✓	✗	✗	✗	✗
CMU MOSEI (Bagher Zadeh et al., 2018)	70	1,000	✓	✗	✗	✗	✓
DyConv (Zhu et al., 2024)	200	–	✓	✗	✗	✓	✗
Switchboard (Godfrey et al., 1992)	300	543	✗	✗	✓	✓	✓
Fisher (Part 1&2) (Cieri et al., 2004)	2,000	3,000	✗	✗	✓	✓	✓

Table 2 - Dataset characteristics in comparison with existing work. Counts and volumes are either as directly reported in papers or estimated from relevant reported quantities. "Volume" for dyadic or multi-party collection indicates "conversation" or "interaction" time, rather than total footage. "In-person" means the participants are in the same room. "High quality face+body" indicates that complete information of all participants (face and body) is maintained throughout all recordings. "Annotations" broadly indicates additional human-labeled data.

Part	Hours	Interactions	Sessions	Participants	Prompts
Overall	4,065.04	64,739	5,098	4,284	1,283
Naturalistic	2,745.43	47,333	3,363	3,525	675
Improvised	1,319.61	17,406	1,735	1,011	845

Table 3 - Overall corpus statistics and major partitions (*Naturalistic* and *Improvised*.)

The number of dyadic hours, interactions, sessions, participants, and prompts in the dataset is summarized in Table 3 (overall and in the *Naturalistic* and *Improvised* portions). Beyond its large volume of dyadic interactions, a principle strength of the SEAMLESS INTERACTION dataset is the accompanying metadata for each interaction. For a given interaction we have:

- Prompt text presented to participants (see example in Table 4).
- Interaction type information: one of the three *games/activities* or IPC-based conversations (see Section 2.2.2).
- Personality information for a subset of participants (see Section 2.3.1).
- Relationship information of the participants in a given session (see Section 2.3.2).
- Annotations of *Internal State*, *Internal State Rationale* and *Visual behavior* for a subset of interactions (see Section 2.3.3).

2.1 Interactions, Sessions, Scripts

As released, the SEAMLESS INTERACTION dataset is a collection of nearly 70,000 short, 2-10 min long interactions between two people. This design allows us to collect a large variety of conversational topics and interpersonal stances - by breaking up larger 1-hr long recording sessions between dyads into thousands of shorter interactions. This approach naturally allows for continuity between interactions over the course of the 1-hr long recording session for a given dyad, while also orienting participants towards a range of different topics and implied socio-affective states.

Sessions. Each dyadic session is 1-hr long, capturing the behavior of two participants conversing and engaging in activities. A trained moderator guides each session, signaling the start and end of each interaction, providing clarification and encouragement to the participants (see moderator responsibilities in Section 2.4.4).

Session example

1 hour long, in-person, contains 10 - 20 prompted interactions



Interaction example

2-10 minutes prompted interactions, either conversation or activity-based.

Figure 1 - An illustration of how a *Session* is divided into 10-20 shorter interactions.

Interaction example

2-10 minutes prompted interactions, either conversation or activity-based.

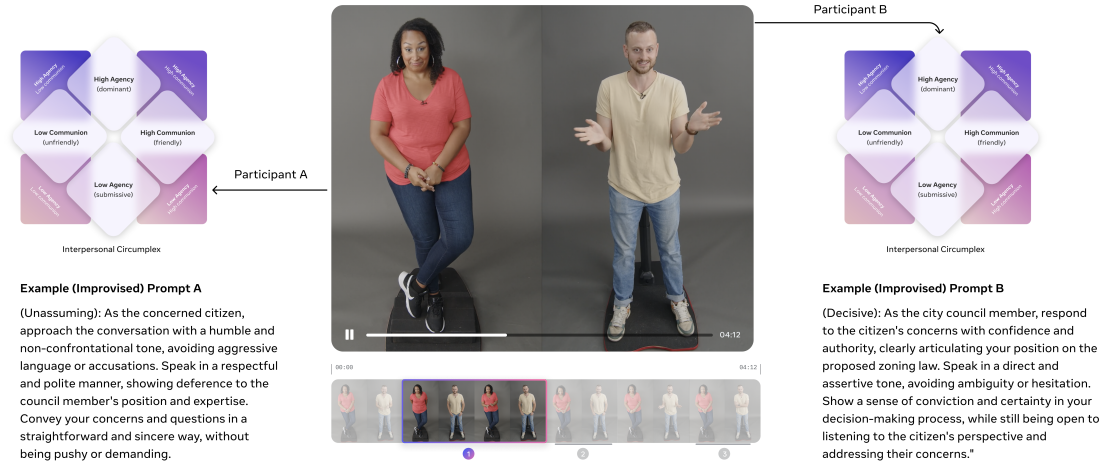


Figure 2 - An illustration of how an *Interaction* is anchored in the IPC via prompts.

Interactions. A given session is divided into between 5-20 (2-10 min long) interactions. An interaction is made up of *active time* - any period during which the participants are speaking, acting, or otherwise behaving in direct response to a prompt from the script. Non-interaction time, also referred to as *meta time*, is any period during which participants are reading the prompts, engaging with the moderator, being distracted, or otherwise have "broken character" from the prompt and are discussing aspects of the session, dataset, or situation itself. While there is some degree of fuzziness in the boundaries of *active* and *meta time*, either due to the nature of the interaction or due to time-stamping error, a given session contains alternating *active* and *meta time* chronologically with regard to the prompts. On average each 1-hr session contains approximately 80-90% *active time*. The current release of the SEAMLESS INTERACTION dataset contains only *active time* segments.

Scripts. The collection of interactions that make up a session are sampled uniformly across the IPC octant space. A particular organization of interaction prompts is called a *Script*, where each *Script* is tailored to whether the dyad belongs to the *Naturalistic* or *Improvised* part, and within the *Naturalistic* pairs, whether they are familiar or stranger.

For every interaction both participants are given separate prompts. The prompts are curated according to

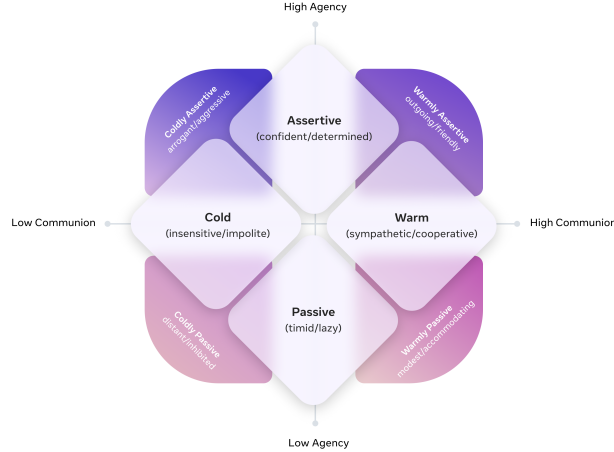


Figure 3 - The Interpersonal Circumplex (IPC) and its two dimensions of *Agency* and *Communion*.

whether it is a *Naturalistic* or *Improvised* session, the relationship (familiar or stranger), and each is designed to anchor the interaction in the IPC (Pincus and Ansell, 2013; Wright et al., 2023).

2.2 Building Instruction Prompts with Interpersonal Theory

Contemporary Integrative Interpersonal Theory (CIIT) refines over 70 years of multidisciplinary research on human interaction and relationships; it has wide relevance and influence across social, medical, and organizational settings, as well as substantial empirical support (Pincus and Ansell, 2013; Wright et al., 2023). It holds that essential features of personality and mental health are reflected in interpersonal situations, including face-to-face and digital interactions, as well as mental representations (e.g., memories, expectations, and fantasies) of such interactions. Furthermore, it proposes that these features can be concisely organized using different combinations of the dimensions of agency and communion.

Agency refers to the development of a distinct sense of self, highlighting goals related to achievement, influence, and independence. *Communion* refers to the formation of close relationships with others, emphasizing goals related to belonging, intimacy, and nurture. Together, agency and communion provide a nuanced lens through which interpersonal dynamics can be analyzed and understood, offering insights into the motivations and behaviors that characterize human interactions. More specifically, interpersonal phenomena (e.g., traits and goals, strengths and problems, behaviors and perceptions) can be located within a circular structure called the *interpersonal circumplex* (IPC), which is formed around the axes of agency and communion (Figure 3).

In the SEAMLESS INTERACTION dataset, we adopted categories of *interpersonal stances* (i.e., how one approaches a given social interaction based on one’s intentions, emotions, and expectations) defined by eight regions of the IPC: high agency and moderate communion (e.g., confident/determined), low agency and moderate communion (e.g., timid/lazy), high communion and moderate agency (e.g., sympathetic/cooperative), low communion and moderate agency (e.g., insensitive/impolite), high agency and high communion (e.g., outgoing/friendly), high agency and low communion (e.g., arrogant/aggressive), low agency and low communion (e.g., distant/inhibited), and low agency and high communion (e.g., modest/accommodating). While not all interpersonal stances are equally common or valued across different cultures and contexts, each is possible and consequential. Using these categories allowed us to more fully capture the spectrum of possible interpersonal interactions.

CIIT provides testable predictions about the emotional and functional outcomes of interactions between specific interpersonal stances (e.g., encounters are generally smoother when dyadic partners differ on agency but align on communion). While these hypotheses are supported by empirical evidence across various settings (e.g., Hopwood et al., 2020), we are eager to further investigate these predictions using the extensive testbed offered

Situation	Participant A	Participant B	IPC A	IPC B
	Your Partner is going to ask you to describe a situation in which someone else criticized you, but you felt the criticism was unfair (this can be from anytime in your life). What happened and how did it make you feel?	Ask your partner to describe a situation where someone criticized them unfairly. Discuss it with them.	ANCM	-

Table 4 - Example Naturalistic prompt pair for an interaction.

Situation	Participant A	Participant B	IPC A	IPC B
A concerned citizen speaks with a city council member about proposed zoning	(Unassuming): As the concerned citizen, approach the conversation with a humble and non-confrontational tone, avoiding aggressive language or accusations. Speak in a respectful and polite manner, showing deference to the council member’s position and expertise. Convey your concerns and questions in a straightforward and sincere way, without being pushy or demanding.	(Decisive): As the city council member, respond to the citizen’s concerns with confidence and authority, clearly articulating your position on the proposed zoning law. Speak in a direct and assertive tone, avoiding ambiguity or hesitation. Show a sense of conviction and certainty in your decision-making process, while still being open to listening to the citizen’s perspective and addressing their concerns.	AMCP	APCN

Table 5 - Example Improvised prompt pair for an interaction.

by this new dataset. Similarly, agency and communion have expected relationships with other important psychological phenomena such as personality and emotion (e.g., [Traupman et al., 2009](#); [Yik and Russell, 2004](#)).

2.2.1 Conversation Based Prompt Creation

We attempt to anchor each 2-10 min long dyadic interaction in the IPC via the use of prompts. Prompts are presented separately to each participant in an interaction prior to the start of the interaction. The majority of *meta time* in the dataset (which is not included in this release) is the short intervals between active interactions, during which participants are studying their prompts and asking clarifying questions of the moderators (refer to [Section 2.1](#) for more detail on *active* and *meta time*).

In the case of *Naturalistic* prompts only one of the participants is given a topic anchoring them in the IPC, while the other participant is instructed to behave in whatever way they feel is natural to them. *Improvised* (or role-play) prompts, by contrast, provide extensive descriptions for one or both participants. We provide examples of both *Naturalistic* ([Table 4](#)) and *Improvised* ([Table 5](#)) prompts below.

***Naturalistic* prompt design.** Naturalistic prompts are designed such that:

- They are answerable by the average person
- They are general enough to allow for a range of response directions depending on the comfort and preference of the participants.
- They provide topic areas that align with a particular IPC octant, but are not prescriptive.

For familiar dyads within the *Naturalistic* dataset, prompts are designed such that they follow the *Naturalistic* prompts generally, with the addition of occasionally making use of—or reference to—history shared by the

interlocutors.

Improvised prompt design. Improvised (role-play) prompts are designed such that they:

- Require prior experience in acting or in improvisation
- Are detailed and prescriptive - describing the roles, stances and emotions that should be represented by the participants.
- Are typically preceded by an adjective intended to describe the “overall manner” of the interpersonal stance of for a particular participant. This adjective also provides are mapping back into IPC octants.

Human- and model-based prompt creation. There are three versions of prompts in the SEAMLESS INTERACTION dataset. The first two versions contain prompts that are hand-written by the core project team. The third and final version of prompts scaled this approach with the use of text-based LLM - prompting a Llama3.1-70b-instruct model to generate prompts geared towards *Naturalist* stranger pairs, *Naturalistic* familiar pairs, and *Improvised* pairs.

2.2.2 Activity Based Interactions

In addition to conversation-based prompts anchored in the IPC, we include another category of prompted interactions we call activities. Activities are designed to explore dynamics outside of free-flowing conversation, particularly as they relate to turn-taking, the relation between language and gesture, and the relation between gesture and intent.

Type	Hours	Interactions	Sessions	Participants	Prompts
IPC conversation	3,464.16	53,938	5,098	4,284	953
Language & Gesture	379.12	6,364	3,465	3,702	296
Collaborative story-telling	146.89	2,845	2,844	3,105	1
Silent Charades	74.87	1,592	1,592	1,910	1

Table 6 - Data volumes by interaction type.

[Table 6](#) shows data volumes for IPC-based conversation prompts along with three activities - “Collaborative story-telling”, a “Language and gesture,” and a “Silent charades.” We provide examples of each of the activity-based prompts in [Appendix C](#).

2.3 Metadata and Annotations

The SEAMLESS INTERACTION dataset provides rich contextualization at the level of interactions (prompts and IPC anchoring), individuals (personality), and dyads (relationships). Below, we describe metadata related to the latter two. We reserve a more detailed description of the taxonomy and its relation to prompt creation in [Section 2.2](#).

2.3.1 Personality

The Five Factor Model of Personality ([McCrae and Costa, 1999](#)) is one of the most robust and empirically supported theories in psychology ([John et al., 2008](#)). It provides a comprehensive framework for understanding how individuals perceive, interpret, and engage with the world, organizing these differences into five broad dimensions known as the Big Five personality traits: agreeableness, conscientiousness, extraversion, neuroticism, and openness. These traits have been shown to predict important real-world outcomes, such as job performance, academic success, and health (e.g., [Ozer and Benet-Martínez, 2006](#); [Roberts et al., 2007](#); [Soto, 2019](#)), often surpassing competing predictors like socioeconomic status and cognitive ability ([Roberts et al., 2007](#)).

By shaping affective, cognitive, and social processes, these personality traits profoundly influence dyadic interactions and constitute a critical aspect of the context in which they occur. For example, highly extraverted

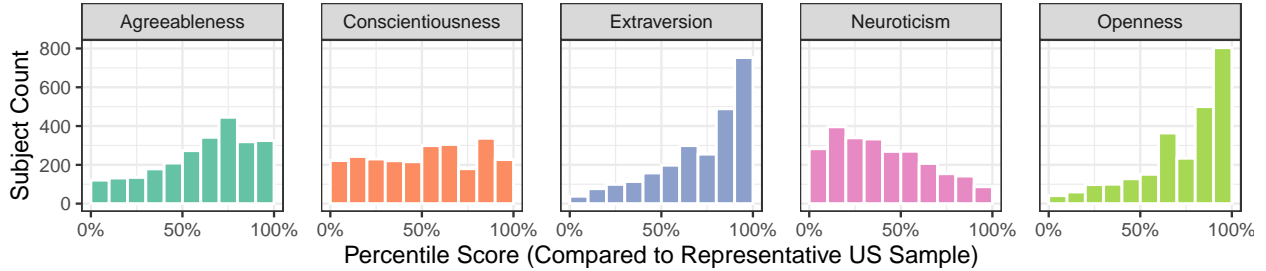


Figure 4 - Distributions of participants’ personality scores from the BFI-2.

individuals are often assertive and engaging, whereas highly agreeable individuals tend to favor cooperation and harmony (John and Srivastava, 1999). Although less explicitly interpersonal, the other Big Five traits also affect communication styles: openness is associated with creativity and flexibility, conscientiousness with structure and goal-directedness, and neuroticism with emotionality and reactivity (Leary and Hoyle, 2009).

In the SEAMLESS INTERACTION dataset, we assessed participants’ personality traits using the Big Five Inventory–2 (BFI-2), a 60-item questionnaire with robust psychometric properties (Soto and John, 2017). Across the 2,260 participants (52%) who chose to complete the BFI-2, inter-item reliability (i.e., coefficient omega; Dunn et al., 2014) was “excellent” for neuroticism (0.90) and “good” for conscientiousness (0.88), extraversion (0.85), openness (0.84), and agreeableness (0.81). Figure 4 depicts each trait’s distribution in our sample as compared to the representative US adult samples ($n = 3,071$) from Soto (2019). Our participants represented the entire range of all five traits, but higher levels of extraversion and openness were more common.

2.3.2 Participant Relationship

The relationship and the degree of familiarity between the interlocutors accounts for a large degree of variation in the topics, stances, and behaviors that make up an interaction (Giles et al., 1991). Table 7 provides the statistics of dyadic relationships in *Naturalistic* data. For recorded collections, such as the SEAMLESS INTERACTION dataset, having some familiarity with your conversational partner also has the advantage of reducing the need for participants to build rapport or go through an often awkward process of ice-breaking and familiarization. With these aspects in mind, approximately half of the *Naturalistic* dataset dyads were specifically recruited as familiar pairs - participants with some prior knowledge of one another, including friends, family, colleagues and romantic partners. Table 8 shows the distribution across familiar categories.

Relation	Hours	Interactions	Sessions	Participants	Prompts
Familiar	1,427.61	24,205	1,752	2,231	597
Stranger	1,317.82	23,128	1,611	1,820	597

Table 7 - Corpus statistics for the *Naturalistic* partition, by dyad relationship.

2.3.3 Annotations

Annotations were performed on a subset of the dataset, to add information about observed visual behaviors, the internal states of the participants, and the possible reasoning behind these visual behaviors. More than 4,500 moment of interest segments were annotated in the dataset. A moment-of-interest was defined as a moment that includes interesting, conspicuous visual behaviors that differ from the normal flow of interaction. To ensure a good representation of the internal state and reasoning of the participant’s behaviors, we asked the participants to annotate themselves directly. These first-party (1P) annotations were complemented by a superset of annotations from third-party (3P) annotators. The first-party annotators (the participants themselves) were asked to annotate the internal states (IS) and the internal state rationales (R) for the moment of interest they remembered. In parallel, third-party annotators were also asked to provide information about the internal states (IS) and internal state rationales (R), but were also asked to annotate fine-grained

Relationship Detail	Proportion
Friends	0.605
Coworkers	0.132
Family-generic	0.089
Familiar-generic	0.088
Dating/spouse/romantic partner	0.056
Classmates	0.022
Siblings	0.005
Parent/child	0.002
Neighbors	0.001
Roommates	0.001

Table 8 - Relationship details and proportions, rounded to three decimals

descriptions of the visual behaviors, aka visual elements (V). A detailed description of each annotation type is provided in [Table 9](#).

The concepts of moment of interest, internal state, internal state rationale, and visual elements are defined as follows:

- **Moment of interest** A short span of time (usually a few seconds) in an interaction during which a participant acts in a way that is visibly or audibly different from what would be considered their baseline behavior.
- **Internal state** is defined as including emotions, feelings, thoughts, or internal dialogue. The instructions included examples for annotators to help guide their descriptions of their internal states. The word examples were classified for emotions, interpersonal, and cognitive:

Emotions Following the dimensional representation of emotional state, examples for low and high arousal dimension, as well as positive and negative valence dimension

Interpersonal We followed the interpersonal theory previously described in [Section 2.2](#), including the eight quadrant along agency and communion dimensions.

Cognitive We also included examples of cognitive states that could be observed during interpersonal interactions, including engagement, comprehension and focus.

- **Internal state rationale** is defined as the reason for which a participant’s internal state was triggered. The rationale can be the trigger itself, regardless of the intent, or it can be the assumed intent behind the trigger. If we take the same situation as above (i.e., Participant A is talking, Participant B starts yawning, Participant A becomes offended), a description of the internal state rationale could be either of the following:
 - *I felt offended* because he|she started yawning and looking bored.
 - *I felt offended* because I could tell that he|she was no longer paying attention.
- **Visual element** is defined as a description of the gestures, movements, facial expressions, and other behaviors that mark a moment of interest and are not considered baseline behaviors for the participant being viewed.

An example of the annotation workflow can be found in [Appendix A.4](#).

Annotation statistics. [Table 10](#) shows annotation volumes by type. Note that our process involved providing different types of annotations for a given moment of interest. These volumes reflect the total number of moment of interest segments and durations when treating the annotations independently. In total, we release $n = 16,898$ annotations and $n = 5,137$ MOI segments covering 4.74 hours of dyadic interaction.

Different third-party annotators attend to different body parts, movement, and emotions. First-party annotations are also intrinsically subjective. Internal states and rationales are colored by unknown prior

TYPE	DESCRIPTION
First-person (first-party) annotations	
1P-IS	Participants annotate their own internal states. Internal states include emotions, feelings, thoughts, or internal dialogue.
1P-R	Participants annotate their own behavior rationales or theories of mind. Rationales are the reason(s) for which a participant’s internal state was triggered. The rationale can be the trigger itself, regardless of intent, or it can be the assumed intent behind the trigger (or theory of mind).
Third-person (third-party) annotations	
3P-IS	Trained annotators annotate the participants’ perceived internal states
3P-R	Trained annotators annotate the participants’ perceived behavior rationales
3P-V	Trained annotators annotate visual elements in the participants’ non-baseline behaviors

Table 9 - Description of annotation types.

Annotation Type	Total Duration (hrs)	# of annotations	Mean # tokens
1P-IS	1.1	751	5.8
1P-R	1.1	751	10.2
3P-IS	4.7	5,132	5.3
3P-R	4.7	5,132	11.3
3P-V	4.7	5,132	14.6

Table 10 - Annotation statistics.

context, personality attributes, and biological/cognitive states.

2.4 Dataset Methodology and Technical Details

2.4.1 Train/Dev/Test/Private-Test Splits

For both the *Naturalistic* and *Improvised* parts of the SEAMLESS INTERACTION dataset, we release public train, dev, and test sets and report the descriptive statistics in Table 11. We hold-out a private test-set for future benchmarking purposes. Given the amount of metadata about participants, dyads, and interactions, there are many dimensions by which the dataset could be partitioned. Given the degree of behavioral variation at the level of individuals, we have partitioned SEAMLESS INTERACTION dataset at the level of participants — no individual participants are mixed with the train/dev/test splits within the primary *Naturalistic* and *Improvised* parts.

Part	Split	Hours	Interactions	Sessions	Participants	Prompts
<i>Naturalistic</i>	Train	2625.34	45,126	3,205	3,243	672
	Dev	35.22	673	48	83	260
	Test	44.13	786	58	109	291
	Private	40.75	748	52	52	337
<i>Improvised</i>	Train	1268.01	16,623	1,656	897	829
	Dev	20.98	323	28	41	155
	Test	15.31	220	23	35	150
	Private	15.32	240	28	38	138

Table 11 - Descriptive statistics by major parts (*Naturalistic* and *Improvised*) and partition.

Given the dyadic nature of the collection and challenges in recruiting thousands of participants across the US, we allowed participants to participate in multiple sessions. For the *Naturalistic* dataset, which does not involve any professional actors, we allowed participants to participate in up to 10 sessions. For the *Improvised* dataset, we allowed participants to participate in up to 30 sessions with a max of 10 with the same partner. In the case of repeated sessions in the *Improvised* part, we required that the same dyad be always given a new set of prompts to respond to. Across the data set, no two participants should have responded to the same prompt more than once.

A side effect of allowing multiple participation is the underlying network structure induced in the dataset - cliques of mutual participation are present. As such, participant-level partition is functionally *clique*-level partition.

2.4.2 Technical Set-up

All recordings took place indoors. Over 95% percent of the interactions took place in a common-room (both participants were present in the same room). Less than five percent of sessions were recorded in separate-rooms, with life-size monitors placed in front of participants. These separate-room set-ups were used to experiment with speaker-bleed reduction techniques, but the set-up is isolated to this small subset.

Audio is recorded at 48khz / 16bit using worn lapel mics (a small subset used Shotgun mics which were subsequently switched out due to high levels of speaker-bleed).

Video is recorded in UHD 4k with 16:9 vertical/portrait aspect ratio (2,160 x 3,840 pixels) at 30 FPS.

An example of the participant workflow is shown in appendix, in [Appendix A.2.3](#).

2.4.3 Participant Recruitment

Participants who were native speakers of English over 18 years of age were recruited in 10 cities across six states. All participants signed an informed consent and were paid for their time. We attempted a best effort at recruiting diverse pairs in terms of age, gender identity, and ethno-cultural background - both within pairs (i.e., many pairs with dissimilar individuals) and across pairs (i.e., many different pairs with similar individuals). Similarly, we attempted a best effort at recruiting diverse and equal distributions of relationship types (e.g., romantic partners, family members, friends, coworkers, acquaintances) for familiar dyads (see [Table 8](#) for the distribution of familiar pair types). Demographic information was not used as part of the acceptance criteria in the collection and its disclosure was optional by the participants. [Figure 5](#) shows distributional properties across *Age*, *Gender*, *Race/Ethnicity*, and *Education* for the set of participants who disclosed this information.

2.4.4 Moderators and Their Responsibilities

A moderator is always present in a recording session. They can both see/hear the participants and speak to them, either via video/audio stream to a separate room, or in the same room but not visible in either video and far enough away that their key-strokes are not audible.

Time-stamping. The moderator’s primary job is to record the time stamps of when each prompt/interaction starts and ends (as provided by the time-synchronized recording software):

- An interaction *starts* the moment that the initiator begins speaking.
- An interaction *ends* when:
 - the participants mutually agree to move on to the next prompt
 - a participant breaks character and asks a question to the moderator about the prompt
 - the moderator themselves interjects in the conversation

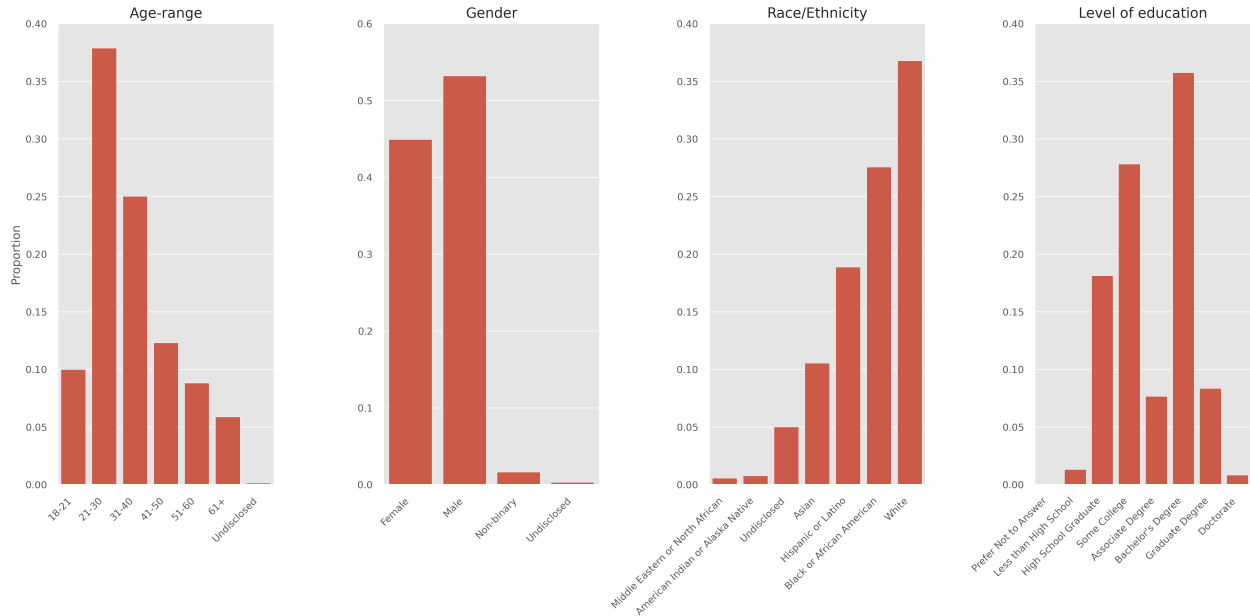


Figure 5 - Demographic category distribution. All demographic information was provided optionally by participants. Some participants chose not to disclose any information. Some chose to partially disclose.

Preventing Privacy Disclosures and Sensitive Topics. The moderator is instructed that it is absolutely essential that no personal details are disclosed during the interaction. This includes full names, addresses, email addresses, phone numbers, etc. (see [Appendix A.2.1](#) for details).

Moderators are also instructed that it is essential that participants not engage on potentially sensitive topics such as overt political debate, descriptions of violence, discrimination, harassment, hate speech, or anything that would cause undue discomfort.

If a participant accidentally discloses personal details, or says something that could be considered a sensitive topic, the vendors are asked to remove this data from the delivery. The moderator is asked to be vigilant in listening for such instances, and in the event that they happen:

- Ensure this is cut from the data. This could be done by stopping recording, removing the clip from the previous recording, and restarting, or by noting for the post-production team to remove it.
- Remind the participants again not to disclose good data.

2.4.5 Known Limitations and Areas of Future Work

Due to the scale and level of complexity in collecting the SEAMLESS INTERACTION dataset there are several aspects that will be the focus of continued work and improvement in future versions.

Errors in manual time-keeping The core unit of the SEAMLESS INTERACTION dataset is an interaction. Interactions define *active time* in which participant conversation and behavior can be linked to a pair of prompts (see Section 2.1). We have observed multiple instances of misaligned time-stamps in which moderators did not correctly identify start and end times for interactions or instances in which the start and end-times are correct, but mapped to the wrong interaction meta-data.

Such misalignment result in interactions that may be too long (the annotated start time is too early or the annotated end time is too late) or too short (start time is too late or end time is too early). In some cases, the prompt text does not align with the spoken material indicating that the ordering of prompts was altered (likely due to off-by-one errors in the script ordering). In total, we believe that these errors, resulting in

misalignment of prompt text metadata and interactions, impact approximately 10% of interactions after our attempts at correction.

A certain degree of error is expected in a collection of this size and complexity, but it is important to take this into account. We have attempted a best effort to automatically identify and rectify these errors.

Time stamping "noise" are also present in the MOI segments used in our annotations (Section 2.3.3). Although there is a degree of subjectivity in defining an MOI, there are rare cases in which the described behavior represents only a subset of the observed behavior in the segment or cases in which the duration of the MOI does not fully capture the annotated behavior.

Incorrect assignment of participant IDs In rare cases, we have observed incorrect duplication of participant identifiers (a single identifier is assigned to two different participants). Conversely, we have also observed rare instances of the same person being mapped to different identifiers. We have made a best effort to rectify these errors when identified.

Unreleased *meta time* As released, the SEAMLESS INTERACTION dataset only contains our *active time* segments. This was in part because the essence of the dataset is derived from *active time* interactions that are rooted in our taxonomy. However, the *meta time* between interactions also represents literally hundreds of hours of fascinating data unto itself. Future releases may explore this subset of the data.

Variation in recording consistency This project was undertaken with multiple collection groups across multiple recording sites. The degree of recording quality (amount of speaker-bleed, consistency of participants staying in frame, quality of acting in *Improvised* segments) and likelihood of time-stamping errors, varies by vendor with some recording sites displaying better competencies in aspects of the collection. All vendors met the basic technical requirements specified in Section 2.4.2 and followed the steps outlined in Section 2.4.4, however there is clear variation in the level of production quality between vendors.

3. Multimodal Representations

For each sample in the SEAMLESS INTERACTION dataset, we extract human-centric visual representations, speech tokens, and transcripts. These features enable a variety of downstream modeling applications for SEAMLESS INTERACTION data, including training and evaluation of dyadic audiovisual motion models.

3.1 Parametric Human Models

Body and hands representations. We use the SMPL-H (Loper et al., 2015; Romero et al., 2017) model to represent the body and hands of each person. SMPL-H is a parametric human model and represents each person via a global orientation $\phi \in \mathbb{R}^3$, body and hand pose (51 joint angles) $\theta \in \mathbb{R}^{51 \times 3}$, and shape $\beta \in \mathbb{R}^{16}$. The SMPL-H model uses these parameters to generate the mesh vertices $\mathbf{V} \in \mathbb{R}^{6980 \times 3}$ of a person.

For each video sequence in the SEAMLESS INTERACTION dataset, we track the person in the video and estimate global orientation and body pose using HMR 2.0 (Goel et al., 2023). This body reconstruction cannot capture hand details; the hands are flat. Thus, we also detect the hands of the person using ViTPose (Xu et al., 2022) and, for each hand, we estimate the 3D hand pose using HaMeR (Pavlakos et al., 2024). Since the hands are reconstructed from a hand-centric perspective, they may be inconsistent with the arms from the body reconstruction. We address this issue in a post-processing step that transforms each hand-centric coordinate system into the corresponding one of each wrist. Finally, we use the same body shape ($\beta = 0$) for all individuals.

3.2 Imitator Face Representation

The Imitator latent representation aims to provide a low-dimensional encoding of facial expressions and positioning of high-level image features, enabling efficient and accurate modeling of talking-head videos. The approach uses two encoders: an expression encoder and an alignment encoder. These encoders work in tandem

to extract relevant features from input images that collectively capture the essential characteristics of a talking head image. We will first describe the architecture and output of each encoder, with further details on their application and usage provided later in the paper.

We employ a pre-trained expression encoder that extracts expression features $\mathbf{f} \in \mathbb{R}^{n \times 128}$, using a typical Resnet34 backbone with a linear head. The encoder expects as input a roll-normalized facial image crop, which can be derived using facial landmarks. In parallel, an alignment encoder processes an upper body crop image to produce 3D translation values for the head and body $\mathbf{t} \in \mathbb{R}^{2 \times 3}$, as well as rotation angles $\mathbf{R}_{\text{head}} \in \mathbb{R}^3$ for the head alone. The alignment encoder also uses a Resnet34 backbone with a linear head.

Both encoders have been trained end-to-end alongside a decoder as part of an internal model for talking head video generation. During training, the objective of image reconstruction on a large dataset with many identities is used to learn these representations in an unsupervised manner. As a result, these representations are not anchored to any interpretable units, such as translations. At inference time, these encoders are used to extract features from new input images, which can then be used for various downstream tasks such as those presented in this paper.

3.3 Speech Representation

We use speech representation from an internally-built speech tokenizer. This speech tokenizer transforms raw speech into a set of discrete tokens, each representing a fixed time slice. We use this speech tokenized representations for conditioning of the dyadic motion models. When integrating with a speech-enabled LLM model, the same speech token representation is used.

3.4 Text Transcripts

We perform peak normalization on the raw audio for each video sample. We then post-process each dyadic pair of audios to remove speaker bleed using Beryl AEC, an in-house echo cancellation and suppression algorithm (Srinivasan and Do, 2024). It is a carefully tuned lightweight DSP-based module that gives great single-talk and double-talk performance on a wide variety of room, microphone, speaker, and device combinations. We use SILERO VAD (Silero, 2021) to extract Voice Activity Detection (VAD) segments and WHISPERX (Bain et al., 2023) to obtain word-aligned text transcripts from post-processed audio. We provide details on the make-up of the transcribed corpus in Appendix A.1.1.

4. Dyadic Motion Models

Our motion models are based on a Diffusion Transformer architecture trained with a flow-matching objective. They are conditioned on dyadic audio input and optionally user’s visual features.

Architecture. We use Diffusion Transformer (DiT) (Peebles and Xie, 2023) as the backbone. From (Zhuo et al., 2024), we use RMSNorm (Zhang and Sennrich, 2019) to improve training stability, as well as self-attention with key-query normalization (KQ-Norm) before the key-query dot product attention computation.

Training objective. We train with the flow-matching (Lipman et al., 2022) objective, which learns the velocity field (or *flow*) of samples moving from a *noise* prior distribution (e.g., Gaussian) to the *target* data distribution, and during inference generates samples starting with Gaussian noise.

Given a sample \mathbf{x} from the target distribution, noise ϵ , and noisy latent \mathbf{x}_t , the model is trained to predict $\mathbf{v}_t = d\mathbf{x}_t/dt$. We adopt the linear interpolation schedule between noise and data, i.e., $\mathbf{x}_t = t\mathbf{x} + (1 - (1 - \sigma_{\min})t)\epsilon$ where $\sigma_{\min} = 10^{-4}$. This formulation indicates a uniform transformation with constant velocity between data and noise. The corresponding time-dependent velocity field is given by $\mathbf{v}_t(\mathbf{x}_t) = \mathbf{x} - (1 - \sigma_{\min})\epsilon$. During training, the flow matching objective directly regresses the target velocity:

$$\mathcal{L}_v = \int_0^1 \mathbb{E}[\| \mathbf{v}_\theta(\mathbf{x}_t, t) - (\mathbf{x} - (1 - \sigma_{\min})\epsilon) \|^2] dt, \quad (1)$$

which is named Conditional Flow Matching loss (Lipman et al., 2022), sharing similarity with the noise prediction or score prediction losses in diffusion models.

During inference, the time-dependent velocity field $\mathbf{v} : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ defines an ordinary differential equation (ODE): $d\mathbf{x} = \mathbf{v}_t(\mathbf{x}_t)dt$. We use $\phi_t(x)$ to represent the solution of the ODE with the initial condition $\phi_0(\mathbf{x}) = \epsilon$. By solving this ODE from $t = 0$ to $t = 1$, we can transform noise into data samples using the learned velocity field $\mathbf{v}_\theta(\mathbf{x}_t, t)$.

4.1 Audio-only Dyadic Motion Model

Face and body motion feature processing. To model face, we use the Imitator representation which includes an expression code, head rotation, and translation values. We normalize these values and concatenate them to obtain a sequence of facial motion features $\mathbf{f} \in \mathbb{R}^{N \times 137}$ for N video frames. For the body and hands of each individual we use rotations for the SMPL-H joint angles in 6D representation (Zhou et al., 2019). These are smoothed using a Savitzky-Golay filter to reduce the jitter. As the SEAMLESS INTERACTION dataset contains mainly upper-body gestures, we ignore the 8 leg joints of SMPL-H, resulting in a sequence of body motion features $\mathbf{b} \in \mathbb{R}^{N \times 258}$.

Dyadic audio conditioning. We condition the motion model on the speech tokens obtained from the dyadic audio of the two speakers. We embed the speech tokens from the two audios with a shared embedding table and concatenate them to encode the speech conditions. To condition the diffusion model, we project the input motion features and the condition features to the same feature space and add them before feeding to the Transformer network. To deal with the mismatch of speech and visual rate (12.5 fps vs 30 fps), we resample the speech conditions to match the length of motion features before the addition. We find that this conditioning approach results in more aligned motion sequence generations compared to cross-attention conditioning.

4.2 Audiovisual Dyadic Motion Model

Audiovisual dyadic conditioning. In addition to dyadic speech, we also experiment with conditioning the motion model with user’s visual features. The user’s visual conditions are also concatenated with the speech conditions and are then fed into the Transformer model. For the face, we use Imitator latent as visual information; for the body, we leverage the user’s SMPL-H pose parameters.

Finetuning. To effectively integrate multimodal visual and enhance high-fidelity generation, we design a two-stage training scheme. In the first stage, we train the model over a full set of training data. In the second stage, we finetune the model with flow-matching loss reweighting (Zhang et al., 2025), which allows us to better take advantage of the data from an automatic reward.

4.3 Face+body Joint and Cascaded Model

To generate both face and body motion sequences, we adopt two different approaches: (a) a joint-model approach, which uses one single model to generate face and body features jointly; and (b) a cascaded-model approach training the face-motion model and body-motion model separately in a cascaded fashion.

Joint model. We use a single face+body motion model to generate face and body motion features jointly as shown in Figure 6. We simply concatenate the face motion features $\mathbf{f} \in \mathbb{R}^{N \times 137}$ and the body motion features $\mathbf{b} \in \mathbb{R}^{N \times 258}$ to obtain joint face+body motion features $\mathbf{j} \in \mathbb{R}^{N \times 395}$. This approach enables the generation of aligned face and body sequences.

Cascaded model. As illustrated in Figure 7, to align the body’s head pose with face’s head rotation, we design the conditional Face2Body model and Body2Face model for cascaded generation. The design tries to avoid the misalignment (especially the head rotation) between face and body outputs.

For Face2Body model, we train several variants with different conditioning, including: (1) full imitator latent; (2) head rotation of imitator latent. In experiments, we find that these models effectively align face with body.

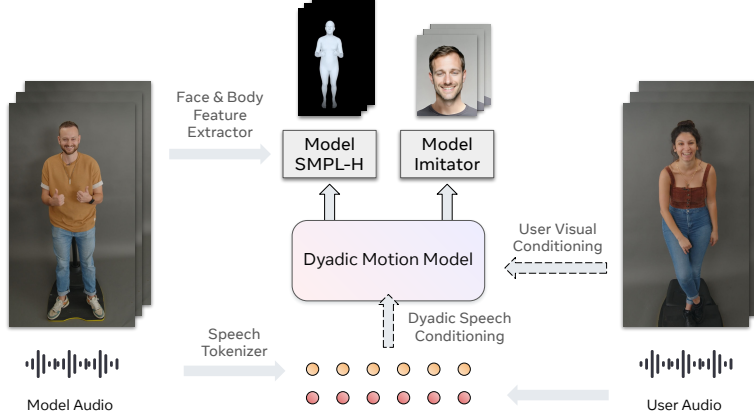


Figure 6 - Our Dyadic Motion Model is conditioned on speech tokens obtained from dyadic audio and optionally user’s visual features, and generates face Imitator features and body SMPL-H features. For Joint Face+Body Model, we concatenate face and body features and generate them jointly.

In inference, we first predict the face imitator latent, following which we predict the body features conditioning on the face features with Face2Body model. As for Body2Face model, the face model is conditioned on body’s SMPL-H representation, to guarantee the alignment between face and body. For automatic evaluation of full-body generation pipeline, the comparison among generations from face + Face2Body, or body + Body2Face pipeline is included. In inference, we first predict the body features and then leverage the Body2Face model for predicting the face features.

4.4 Motion Models with Controllability

4.4.1 Emotion Controllability

In order to achieve precise and fine-grained control over the emotion during generation, we incorporate the Arousal-Valence (A-V) Matrix (Russell, 1980). This model provides a structured, two-dimensional framework for representing and steering affective states within the generated output. Compared to categorical emotion labels, which are often difficult to classify reliably and prone to oversimplification, the A-V Matrix offers significant advantages: it naturally captures the continuous spectrum and subtle gradations of human affect. We use an internal tool to extract arousal and valence sequences for each sample: $\mathbf{a} = \{a_0, a_1, \dots, a_{n-1}\}$ and $\mathbf{v} = \{v_0, v_1, \dots, v_{n-1}\}$, where n is the frame number of the sample.

The arousal and valence sequences \mathbf{a} and \mathbf{v} fall into the numerical range from -1 to 1 . A straightforward way is using the continuous arousal and valence values as conditioning signals, i.e., the model learns the conditional distribution $p(\mathbf{lat}_t | \mathbf{lat}_{t-1}, \mathbf{a}, \mathbf{v})$ for timestep t , where \mathbf{lat}_t means the Imitator latent of the t -th timestep. In order to improve the robustness of conditioning and fit discrete emotion tokens as used in Section 4.5.1, we propose **bucket-based conditioning** in place of real-valued emotion conditions. Specifically, continuous values are discretized using predefined buckets (e.g., partitioning the $[-1, 1]$ valence space into k intervals, in practice, uneven buckets can be set empirically). Each (\mathbf{a}, \mathbf{v}) pair is mapped to a bucket index $b_a, b_v \in \{1, 2, \dots, k\}$. The model then learns $p(\mathbf{lat}_t | \mathbf{lat}_{t-1}, b_a, b_v)$, allowing generation to be conditioned on value *ranges* rather than exact values.

During training, we apply *condition dropout* (rate ρ) to randomly mask these signals for the above two conditioning methods, enabling the model to learn the generation of conditioned and unconditioned.

During inference, the model can generate outputs in the following manners: 1) unconditioned generation; 2) conditioned on constant values (a_c, v_c) for stable emotional profiles; 3) conditioned on externally generated sequences (\mathbf{a}, \mathbf{v}) (e.g., from audio features or semantic analysis) for dynamic affect.

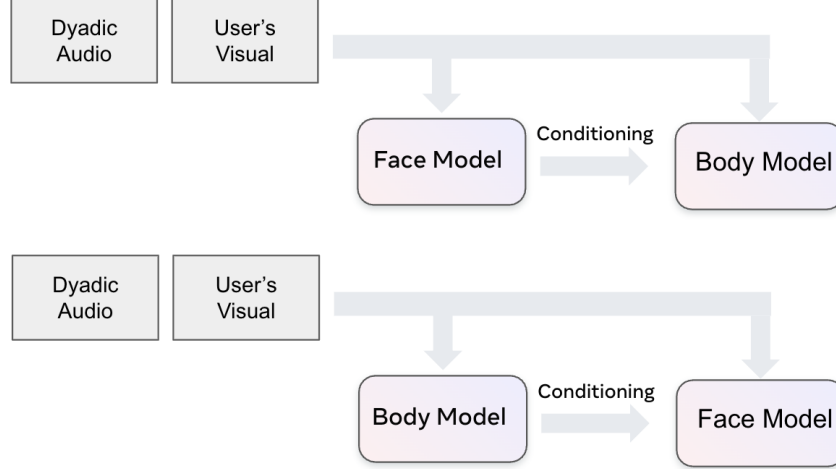


Figure 7 - Cascaded model (Face2Body or Body2Face generation).

4.4.2 Expressivity Controllability

Control signals. The facial expressivity is related to the movement of multiple correlated components (i.e., head pose, eye brows, mouth region, etc.). Hence, we design our architecture to jointly conditioned on these control signals. Specifically, we consider the following control signals:

- **Head rotation:** For each video frame sequence of length n , we extract a 3-dim (pitch, yaw, roll) head rotation sequence $\mathbf{R}_{\text{head}} = \{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{n-1}\}$, where each $\mathbf{r}_i \in \mathbb{R}^3$ is the instantaneous Tait-Bryan rotation vector:

$$r_i = \begin{pmatrix} r_{\text{pitch}} \\ r_{\text{yaw}} \\ r_{\text{roll}} \end{pmatrix} (\text{radians}). \quad (2)$$

The values of r_{pitch} , r_{yaw} , r_{roll} represent the radian angle along pitch, yaw, roll axis respectively. In concrete terms, r_{pitch} embodies the rotation of the head in the sagittal plane, with positive values denoting chin downward and negative values denoting chin upward, thus we can regard there is a head nod motion if $\Delta(r_{\text{pitch}})$ is larger than a threshold. Similarly, r_{yaw} means horizontal plane rotation, i.e., head turns to left or right. r_{roll} represents coronal plane rotation, describing the head tilt.

- **Eye brows:** To encode the person-independent facial expressions, we use the Facial Action Coding System (FACS) (Ekman and Rosenberg, 1997) which describes a taxonomy of facial action units (FAU) to capture the movements of different muscles or muscle groups. For each video sequence with length n , we first use an internal tool to extract the FAU sequence $\mathbf{F}_{\text{faul}} = \{\mathbf{f}_0, \mathbf{f}_1, \dots, \mathbf{f}_{n-1}\}$, where each $\mathbf{f}_i \in \mathbb{R}^{46}$ is the FAU value vector with 46 facial motion unit values. Next, we select the 3 FAU types only related to eye brows and use their average value vector \mathbf{m}_{eb} as the control signal for eye brows movement.
- **Mouth:** Similar to the eye brows control signal, we select the 20 FAU types related to mouth and use their average value vector \mathbf{m}_{m} as the control signal for the mouth movement.
- **Eye Gaze:** We introduce gaze values to control the angle and direction of eye gaze. For each video sequence of length n , we extract the 2-dim pitch and yaw eye gaze sequence $\mathbf{g}_{\text{gaze}} = \{\mathbf{g}_0, \mathbf{g}_1, \dots, \mathbf{g}_{n-1}\}$, where $\mathbf{g}_i \in \mathbb{R}^2$ is the rotation vector of pitch and yaw values for both eyes.

Conditioning methods. Given the continuity of each control signal, we apply a moving average to it to mitigate the interference of outliers in the original data. After that, for each aforementioned control signal, we propose several ways to introduce them into the model training process (along with *condition dropout* of ρ for each signal to enable both conditioned and unconditioned generation setting). Similar to the conditioning

ways mentioned in emotion controllability (Sec 4.4.1), the easiest way to involve those control signal is using the original sequences (\mathbf{R}_{head} , \mathbf{m}_{eb} , \mathbf{m}_{m} , \mathbf{g}_{gaze}) directly.

In order to focus more on movement patterns rather than static poses, we calculate and condition the model on motion dynamism $\dot{\mathbf{s}}$ for each $\mathbf{s} \in \{\mathbf{R}_{\text{head}}, \mathbf{m}_{\text{eb}}, \mathbf{m}_{\text{m}}, \mathbf{g}_{\text{gaze}}\}$:

$$\begin{aligned}\dot{\mathbf{s}} &= \{\dot{s}_0, \dot{s}_1, \dots, \dot{s}_{n-1}\}, \\ \dot{s}_t &= \text{abs}(s_t - s_{t-1}).\end{aligned}$$

For dynamism ($\dot{\mathbf{s}}$) of each control signal, we run over a random selected subset including 300-hour videos to obtain the statistics, including the maximum value, the minimum value, and different quartiles. We then build up a set of threshold $\tau = \{\tau_0, \tau_1, \dots, \tau_{k-1}\}$, where k is the total number of buckets in total. Thus, it is easy to convert the signal into bucket index. Take the head rotation as an example:

$$\mathbf{b}^j = \text{bucket}(\mathbf{r}_j) = \sum_{i=1}^{k-1} \mathbb{I}(\mathbf{r}_j > \tau_i^j), \quad j \in \{\text{pitch, yaw, roll}\}.$$

Then the model takes the bucket index b_j as the condition:

$$p(\mathbf{lat}_t \mid \mathbf{lat}_{t-1}, \mathbf{b}^{\text{head}}, \mathbf{b}^{\text{eb}}, \mathbf{b}^{\text{m}}, \mathbf{b}^{\text{gaze}}),$$

where $\mathbf{b}^{\text{head}}, \mathbf{b}^{\text{eb}}, \mathbf{b}^{\text{m}}, \mathbf{b}^{\text{gaze}}$ are the bucket index sequences of $\mathbf{R}_{\text{head}}, \mathbf{m}_{\text{eb}}, \mathbf{m}_{\text{m}}$, and \mathbf{g}_{gaze} respectively. Taking advantage of the statistics, the bucket-based conditioning approach is robust to perceptually equivalent classes.

Expressiveness level. While the model can achieve fine-grained single-signal control, different combinations of multiple signals can present different levels of expressiveness, allowing the overall modulation of expressiveness. For example, when we want the avatar to perform more nodding motions, we can set a higher pitch-axis head rotation dynamism $\dot{\mathbf{r}}_{\text{pitch}}$ during inference. With more head nodding, the avatar tends to be regarded as more expressive. In a similar way, setting large eyebrow motion dynamism mouth motion can also help with improving the avatar expressiveness.

4.4.3 Semantic Gesture Controllability

Semantic gesture generation plays a crucial role in enhancing the naturalness and expressiveness of virtual agents in speech-driven gesture generation tasks. However, with only speech as a condition, generative models face significant challenges in producing semantic gestures due to their rarity and long-tail distribution. To tackle this challenge, we propose semantic controllability.

Methodology. To realize zero-shot arbitrary gesture conditioning in the co-speech diffusion model, we introduce a random **temporal gesture condition dropping** strategy. This would allow generating any unseen gestures with a smooth transition as long as we have the gesture condition. Let $\mathcal{G} = \{\mathbf{g}_i\}_{i=1}^T$ denote the ground truth gesture sequence of length T , where \mathbf{g}_i represents the gesture at time step i . Let $\mathbf{s} = \{s_i\}_{i=1}^T$ denote the corresponding speech condition sequence, where s_i is the speech feature at time step i .

During training, we apply a random temporal gesture condition drop to simulate partial gesture absence. For each frame index i , the gesture condition \mathbf{g}_i is retained with probability $1 - p_{\text{drop}}$ or replaced with a null condition (e.g., a zero vector or mask token) with probability p_{drop} . The speech condition \mathbf{s} is always retained to ensure consistent contextual grounding. The modified gesture sequence $\mathcal{G}_{\text{drop}}$ is defined as:

$$\mathcal{G}_{\text{drop}} = \{\tilde{\mathbf{g}}_i\}_{i=1}^T, \quad \tilde{\mathbf{g}}_i = \begin{cases} \mathbf{g}_i & \text{with probability } 1 - p_{\text{drop}}, \\ \mathbf{0} & \text{with probability } p_{\text{drop}}, \end{cases}$$

where $\mathbf{0}$ denotes the null condition. The diffusion model is trained to predict the original gesture sequence $\mathcal{G}_{\text{pred}}$ given the dropped gesture sequence $\mathcal{G}_{\text{drop}}$ and the full speech condition \mathbf{s} at each diffusion step. This

approach encourages the model to generate rhythm-aware gestures when only speech is provided as a condition, and to follow semantic gestures when gesture conditions are available.

In addition to using SMPL-H as a gesture condition, we also experimented with using VQ-VAE codes of the SMPL parameters as the gesture condition. A VQ-VAE is trained to encode SMPL-H parameters into a discrete codebook \mathcal{C} of size $|\mathcal{C}|$, with codebook indices ranging from 0 to $|\mathcal{C}| - 1$. For each SMPL gesture \mathbf{g}_i , the VQ-VAE encoder produces a latent embedding, quantized to the nearest codebook entry, yielding a discrete code $\mathbf{c}_i \in \mathcal{C}$. The gesture sequence is represented as $\mathcal{G}_{\text{VQ}} = \{\mathbf{c}_i\}_{i=1}^T$. Similarly, we have the temporarily randomly dropped VQ ID sequences as a gesture condition during training:

$$\mathcal{G}_{\text{VQ,drop}} = \{\tilde{\mathbf{c}}_i\}_{i=1}^T, \quad \tilde{\mathbf{c}}_i = \begin{cases} \mathbf{c}_i & \text{with probability } 1 - p_{\text{drop}}, \\ |\mathcal{C}| & \text{with probability } p_{\text{drop}}, \end{cases}$$

The diffusion model is trained to predict the original gesture sequence \mathcal{G} (in SMPL-H parameters) given the dropped VQ code sequence $\mathcal{G}_{\text{VQ,drop}}$ and the full speech condition \mathbf{s} at each diffusion step.

Temporal Condition Integration. For discrete temporal conditions, such as gesture VQ IDs and audio token IDs, we map the IDs to randomly initialized embeddings, which are tuned during diffusion training, and encode them to the target dimension. For continuous temporal conditions, such as the SMPL-H condition, we employ an MLP to get encoded features. These conditions are then concatenated with the noisy diffusion input along the channel dimension, and the concatenated output is linearly projected back to the original dimension of the noisy input.

4.5 Integration with Speech LLM Model

We integrated our dyadic motion models with a speech LLM model trained on transcribed dialogs.

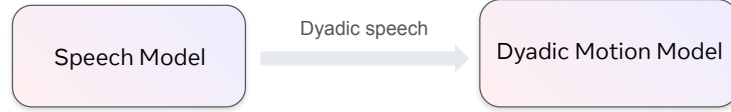


Figure 8 - Cascaded integration of dyadic diffusion with speech LLM model.

For this work, we always freeze the speech model to ground our diffusion generation. During inference, we initiate two speech models as two agents to talk to each other, each treating the other agent as a user. Two agents together can generate coherent and expressive dialogues following our dyadic speech prompts, showing great in-context learning capability. As shown in Figure 8, we feed the generated dialogue (in the form of speech tokens) to ground our dyadic motion models.

4.5.1 LLM-Guided Codebook Generations

Diffusion motion models conditioned on speech and visual signals have demonstrated natural face and body generation in dyadic interactions. Additional conditioning further modulates the generation, for example, the emotion control adds more expressivity to facial expressions. This suggests that diffusion models could benefit from extra guidance from the speech model. Hence we propose Adapter and codebooks as the bridge between speech model and diffusion. Codes are inferred from speech model, and then are used as conditions to drive diffusion models.

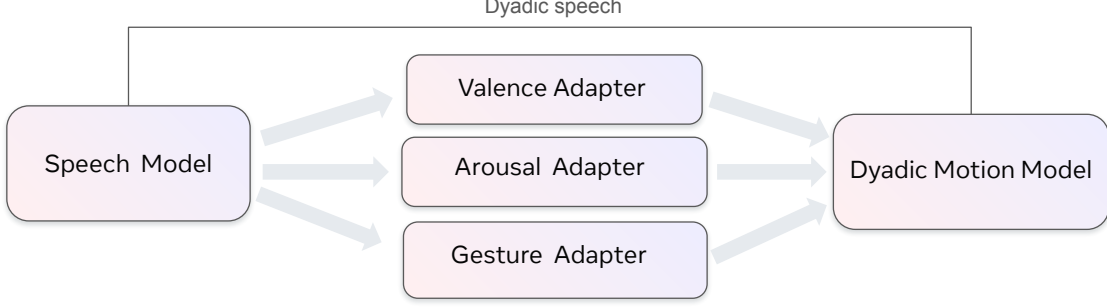


Figure 9 - Codebook integration of dyadic diffusion with speech LLM model.

Now we will present the design for the LLM-driven codebook integration with speech model and diffusion. As shown in Figure 9, an Adapter is built on top of a speech language model, taking its hidden states as input and predicting emotion and gesture codes to guide the face and body diffusion models.

Adapter. The Adapter is a multi-layer perceptron with GELU activation. Suppose that the predicted token rate of the adapter is R , that is, the number of tokens in one second. We extract the hidden states of the agent speech tokens from the last layer of speech model, and interpolate these representations to match the token rate of r . The Adapter transforms interpolated speech representation with MLP layers, and projects it to codebooks via the projection layer.

$$\tilde{\mathbf{C}} = \text{Adapter}(\text{Interpolate}(\mathbf{H}, r)), \quad (3)$$

where \mathbf{H} is speech model hidden states, and $\tilde{\mathbf{C}}$ is prediction on predefined codebooks.

For training, speech model takes dialog data as input, and the adapter extracts the hidden states from its last layer and makes code predictions. Cross-entropy loss is measured between ground truth codes \mathbf{C} and predictions $\tilde{\mathbf{C}}$. The Adapter is tuned to minimize the loss, while the parameters in speech model are all frozen.

We introduce the emotion and gesture codebooks below, and discuss how they guide the face and body diffusion models, respectively.

LLM-Driven Generations with Emotion Codebook. The emotion codebook consists of 12 valence tokens and 12 arousal tokens. Similar to Section 4.4.1, valence and arousal are extracted from the face in each video frame. Their values range from 0 to 1, and are quantized into 12 equal-sized bins as discrete emotion tokens. Emotion adapters are trained to predict valence and arousal tokens.

The face diffusion model takes valence and arousal as extra conditioning information in addition to dyadic speech. The emotion condition modulates the facial expressions in generation.

LLM-Driven Generations with Semantic Gesture Codebook. The SEAMLESS INTERACTION dataset provides a set of *semantic gesture game* data where actors make semantic and illustrative gestures while talking about words of interest. We collected a set of gestures as the gesture vocabulary, and built a lookup table to map each gesture to its corresponding SMPL-H sequence.

To train a gesture Adapter, we labeled the segments with a gesture that its spoken words trigger. If no semantic gesture exists in one segment, we assign a special null gesture label. The adapter essentially learns a multi-class classification task, predicting a gesture on a given spoken segment.

The body diffusion model takes predicted gesture labels as well as speech to control gestures and body movements as described in Section 4.4.3.

5. Visualization

Our Dyadic Motion Models were designed to always generate two streams of motion codes: one stream for the face (a.k.a. *face codes*) and one stream for the body (a.k.a. *body codes*). A significant advantage of generating these intermediate codes instead of directly generating video pixels is that we can not only generate 2D videos with our Dyadic Motion Models, but also 3D rendered Codec Avatars, that can be visualized in settings like Virtual Reality (VR) headsets or Augmented Reality (AR) glasses.

5.1 2D Video Generation

Our 2D video generation approach follows the image-to-video technique where the generative process is conditioned on face and body pose. Given a single reference image of the subject, the sequence of face codes, and the body code sequence, the algorithm generates a 2D video in which facial appearance, head motion, and articulated body dynamics are coherently aligned with the face and body codes.

5.1.1 2D Rendering

Our 2D Rendering model is a diffusion model that receives a single reference image $I \in \mathbb{R}^{3 \times H \times W}$ and, while operating entirely in a compressed latent space, synthesizes a temporally coherent video clip V . The design of the rendering model interleaves three key components—(1) a Temporal Autoencoder (TAE), (2) window-efficient Diffusion–Transformer (DiT) blocks (Peebles and Xie, 2023), and (3) a lightweight latent-space concatenation strategy. Together, these modules enable high-fidelity, temporally smooth animations. We initialize our model from a pretrained diffusion model checkpoint and train on the SEAMLESS INTERACTION dataset.

Temporally AutoEncoder (TAE). The pipeline begins by inflating a conventional 2D VAE into 3D. Spatial convolutions alternate with lightweight 1D temporal blocks, allowing the TAE encoder to map the reference image to a spatio-temporal latent grid $z \in \mathbb{R}^{T \times H' \times W' \times C}$, where each axis is down-sampled by a constant factor (e.g., 8) yet perceptual detail is largely retained. This latent grid serves as the substrate upon which all subsequent denoising and conditioning operations are performed.

Diffusion–Transformer (DiT). To process z , a single 3D convolution patch-embeds the grid and flattens it into a token stream enriched with factorized positional encodings. A DiT then iteratively denoises these tokens; its self-attention is restricted to *shifted, non-overlapping windows*. The shift mechanism, inherited from the Swin-Transformer (Liu et al., 2021), propagates information across windows while keeping the cost nearly linear in sequence length (Peebles and Xie, 2023).

Image conditioning. To animate a still photograph, the encoder treats I as a single frame video. Its latent slice $z_0^{\text{img}} = E(I)$ is replicated along the temporal axis to the target length T . During diffusion, the current noisy sample z_t is concatenated channel-wise with this fixed appearance code, anchoring identity and illumination while the network synthesizes plausible future motion.

Finally, the TAE decoder maps z_0 back to the RGB space, yielding the video z_0 . If even stronger identity preservation is desired, additional vision-specific tokens can be appended to the token stream without retraining the backbone.

5.1.2 Body and Face Conditioning

In our cascaded pipeline the body and face are driven by two dedicated code sequences: a SMPL-H sequence $b_{1:T}$ for full-body articulation and face latent (Imitator latent) $f_{1:T}$ for expression dynamics.

Body conditioning. Following a rich line of pose-based video synthesis work (Villegas et al., 2018; Chan et al., 2019; Siarohin et al., 2019; Yang et al., 2020; Petrov et al., 2023), we represent body motion with a *human-skeleton video*. For every frame in the training set we first detect 2D keypoints using Sapiens (Khrodgar et al., 2024) and connect them according to the canonical limb graph defined by OpenPose (Cao et al., 2017).

At each diffusion timestep t , the resulting pose tensor is concatenated (along the channel dimension) with the encoded reference image and the noise latent z_t , enabling our model to couple appearance and motion in a single forward pass.

At inference we only dispose of a predicted SMPL-H mesh sequence rather than 2D keypoints. Because SMPL-H is defined in an uncalibrated, human-centric coordinate frame, we first estimate the full camera parameters from the reference image — recovering intrinsics (e.g., focal length) and extrinsics (rotation & translation) with a standard PnP solver (Lepetit et al., 2009) initialized by structure-from-motion (Schönberger and Frahm, 2016). We render the mesh sequence under this camera to produce an aligned synthetic video and then run Sapiens (Khrodgar et al., 2024) to obtain the 2D keypoints, which are fed into the diffusion model exactly as during training. This procedure ensures that pose guidance remains geometrically consistent even when the only available motion source is a 3D parametric body model.

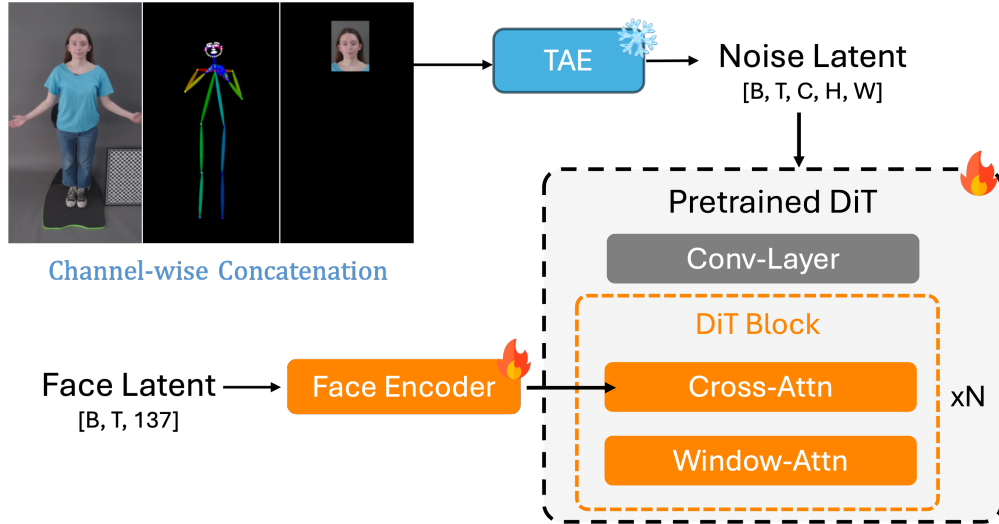


Figure 10 - 2D rendering pipeline.

Face conditioning. There are two vanilla ways to apply face conditioning: (a) the inpainting method, which decodes the face latent into a face video, and the model inpaints the body video on top of the face video; and (b) direct injecting the face latent through cross-attention. Both vanilla methods have shortcomings. The inpainting method cannot handle the scenario where the hand overlaps with the face intended by body conditioning, while the direct injecting scheme suffers from quick overfitting.

To overcome the issues of both vanilla methods, our 2D Avatar renderer adopts a hybrid scheme. We apply channel-wise concatenation to the face video, reference image I , pose-skeleton sequence, and z_t before the TAE encoder (Section 5.1.2). Meanwhile, every DiT block takes face latent $f_{1:T}$ (after extra face encoder) in its cross-attention layer. During both training and inference, we drop any face-video frame whose face is occluded by another body part — a condition detected directly from the pose skeleton. The network is therefore forced to reconstruct the missing regions from the expression latent, the pose guidance, and the surrounding unoccluded frames, yielding stable lip-sync and facial expressiveness even under severe self-occlusion.

5.2 3D Codec Avatar Rendering

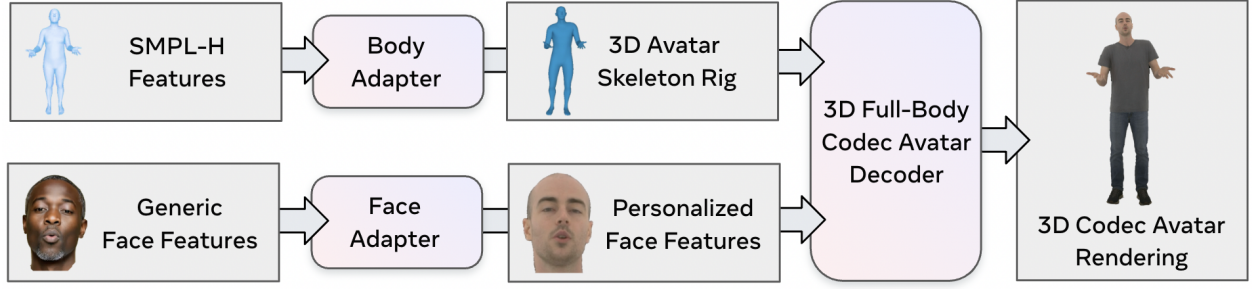


Figure 11 - 3D Codec Avatar Rendering Overview. Given generic expression features and SMPL-H features as input, we employ a face and body adapter to map to the driving signals required by the 3D Full-body Codec Avatar decoder. The decoder produces a set of Gaussian primitives that encode the geometry and appearance of the avatar. Finally, images are rendered based on the Gaussian primitives using a tile-based Gaussian splatting renderer.

In addition to the 2D rendering discussed in the previous section, we connect the dyadic motion model to 3D Full-body Codec Avatars (Bagautdinov et al., 2021; Martinez et al., 2024) to enable free-viewpoint rendering; see Figure 11 for a pipeline overview. This requires a face and body adapter to map from the generic expression features and SMPL-H body features that are the output of the dyadic motion model to the person-specific expression features and the full-body avatar skeleton rig that is expected as input by the 3D Full-body Codec Avatar decoders. We employ 3D Full-body Codec Avatars based on Gaussian splatting (Wang et al., 2025; Kerbl et al., 2023) to achieve high-fidelity results. In the following, we describe the underlying datasets and models that are required to train personalized Full-body Codec Avatars and map to the driving signals for the face and body.

5.2.1 3D Full-body Codec Avatars

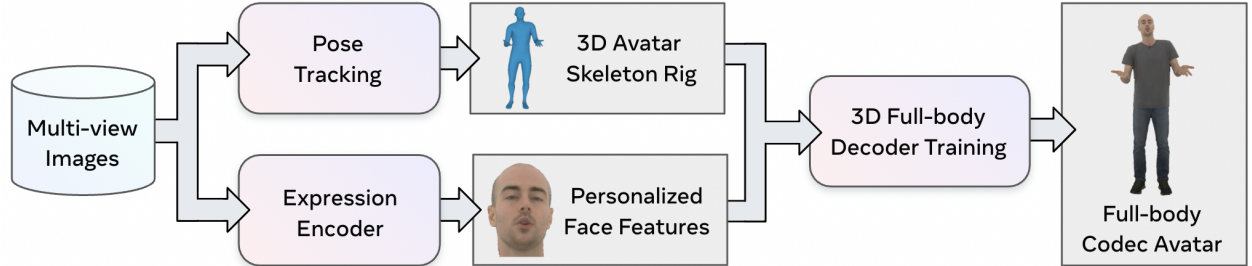


Figure 12 - Codec Avatar Decoder Overview. Given a multi-view capture of an actor, we first extract per-frame 3D skeletons based on 3D keypoint detections as well as personalized expression features based on frontal view expression encoding. Given the extracted face and body features as input, the 3D Full-body Codec Avatar decoder is trained end-to-end to synthesize the multi-view images in the dataset by minimizing a photometric re-rendering loss.

In this section, we describe the datasets and models that are required to train a 3D embodiment based on 3D Full-body Codec Avatars. Figure 12 provides an end-to-end overview of the entire pipeline from the multi-view dataset over pose tracking and expression encoding to 3D Full-body Codec Avatar decoder training.

Multi-view performance capture. We start by collecting a dataset of our actor’s shape and pose-dependent appearance based on a multi-camera capture systems with 512 synchronized and calibrated cameras (Wang et al., 2025). Each of the cameras has 24 mega-pixel resolution and records videos at 30Hz. The multi-view capture system has a radius of 2.75 meters, which is large enough to collect the full range of human motion. For each actor, we collect around 30K frames of training data covering a large variety of facial expressions

and body poses. The resulting dataset contains synchronized and calibrated multi-view images $\{I_f^k\}$ with k being the camera index and f being the frame index.

Pose tracking and mean shape estimation. Given the set of multi-view images I_f^k as input, we first extract 2D keypoints and foreground-background segmentation masks for all images (Khironkar et al., 2024). Next, we triangulate the per-camera 2D keypoint detections to obtain a set of 3D per-frame keypoints. We use inverse kinematics (IK) to fit the avatar skeleton rig pose parameters s_f to the per-frame 3D keypoints. In addition to the avatar skeleton, we also estimate the actor’s canonical mean shape mesh. We employ multi-view 3D reconstruction to obtain per-frame 3D point clouds for a subset of frames that contain peak poses. Next, we solve for the canonical template mean shape mesh that best fits all point clouds when deformed using linear blend skinning (LBS) based on the known tracked avatar skeleton pose parameters s_f .

Person-specific expression code estimation. To model facial expressions, such as the shape of the mouth, position of the eyes, and the motion of the eye brows, we employ a set of per-frame expression features. The expression features $e_f \in \mathbb{R}^{128}$ are computed by first selecting one frontal view image for each frame of the multi-view capture. The selected frontal view images are then processed by the facial expression encoder that has been introduced in Section 3.2 to extract e_f for each frame.

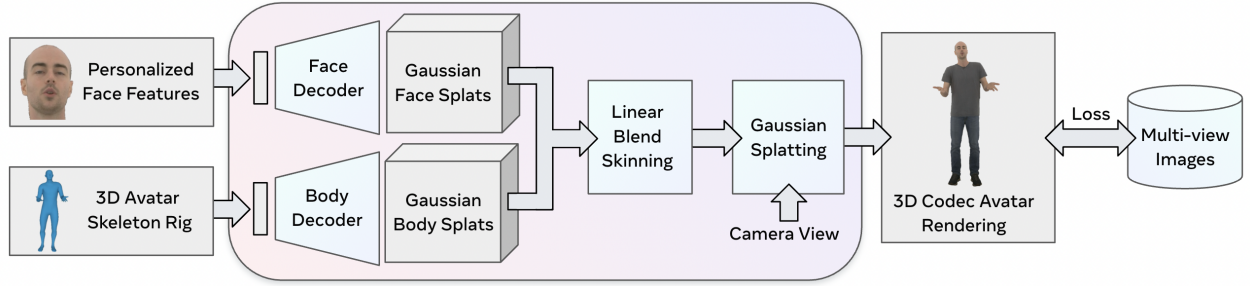


Figure 13 - Codec Avatar Decoder Details. 3D Full-body Codec Avatar decoders are trained end-to-end using a photometric loss with respect to the captured multi-view image dataset of an actor. The decoder is a de-convolutional neural network with a face and body decoder branch that each output a set of Gaussian primitives that are combined and then posed via linear blend skinning (LBS), before finally being rendered to the screen via Gaussian splatting.

Codec Avatar decoder training. Given the multi-view images I_f^k , avatar skeleton rig pose parameters s_f , and expression codes e_f as input, we train a 3D Full-body Codec Avatar decoder. The 3D Full-body Codec Avatar decoder is a de-convolutional neural network with a face and body decoder branch (see Figure 13). The face decoder takes as input the per-frame expression codes e_f and maps them to a tensor $G_{face}(e_f) \in \mathbb{R}^{N_w \times N_h \times N_c}$ that stores the parameters associated with the set of decoded Gaussian primitives. Here, $N_w^f = 512$ and $N_h^f = 512$ defines the number of Gaussians along the width and height of the body template’s uv-map leading to 262k Gaussian face primitives. The number of parameters per Gaussian primitive is $N_c = 38$ to parameterize its position (3), orientation (4), scale (3), opacity (1), and appearance (27). The geometry and appearance of our Full-body Codec Avatars is modeled as the composition of a set $G = \{g_i\}_{i=0}^{N-1}$ of N Gaussian 3D primitives g_i . The Gaussian primitives $g_i = \{t_i, R_i, s_i, o_i, c_i^k\}$ model both the face and body geometry as well as their view-dependent appearance. Here, $t_i \in \mathbb{R}^3$ is the position of the primitive in unposed model space, $R_i \in \mathcal{SO}(3)$ its orientation, $s_i \in \mathbb{R}^3$ the per-axis scale factor, and $o_i \in \mathbb{R}$ its opacity. Together, these parameters describe the 3D Full-body Codec Avatar’s face and body shape in canonical space. The $c_i^k \in \mathbb{R}^3$ are coefficients of 2nd-order spherical harmonics and parameterize each Gaussian’s view-dependent appearance. The body is decoded by a decoder that outputs the tensor $G_{body}(s_f) \in \mathbb{R}^{N_w^{body} \times N_h^{body} \times N_c}$. Here, $N_w^{body} = 1024$ and $N_h^{body} = 1024$ lead to 1M Gaussian primitives. Finally, the face and body Gaussian primitives are combined. The final step is to transform the Gaussian primitives from canonical space to deformed space using linear blend skinning (LBS) based on the known tracked avatar skeleton pose parameters s_f before being rendered into the camera view using Gaussian splatting (Kerbl et al., 2023). Given the decoded Gaussian primitives g_i in world space and the user’s viewing direction v as input, we employ volume rendering to generate the

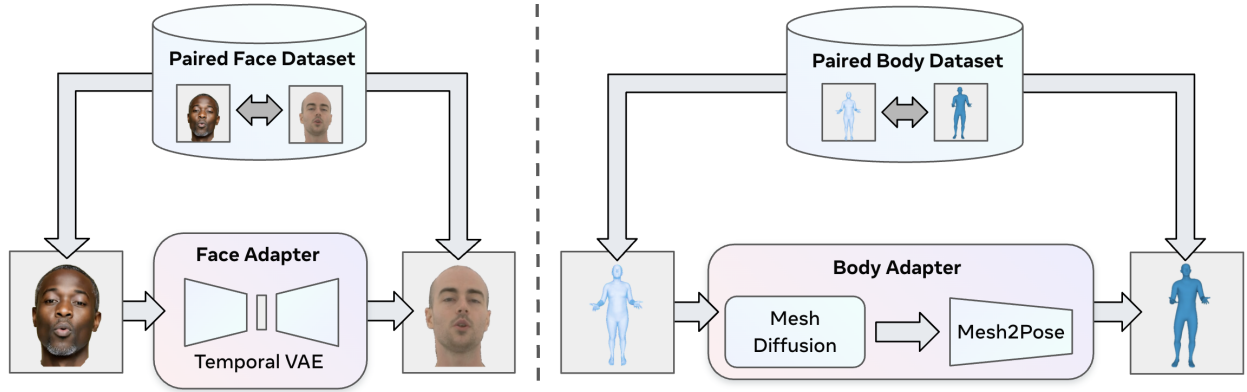


Figure 14 - Codec Avatar Adapters. Face Adapter (Left). We create a paired dataset of generic and person-specific face expression codes, and then train a temporal VAE to map from generic to person-specific expressions. The output can be consumed by the 3D Full-body Codec Avatar decoder to render the 3D face embodiment. **Body Adapter (Right).** We create a paired dataset of SMPL-H and the corresponding Codec Avatar skeleton rig pose parameters, and then train a temporal diffusion model to map from one representation to the other. The output can be consumed by the 3D Full-body Codec Avatar decoder to render the 3D body embodiment.

final images. We first evaluate the color of each Gaussian as $c_i = \sum_{k=0}^9 c_i^k \mathcal{SH}_k(v)$. We then determine all Gaussians that overlap each of the image tiles and sort them in front-to-back order based on their depth with respect to the camera. For each pixel p within the tiles, we compute its color $c(p) = \sum_{i \in S(p)} c_i o_i T(i)$ based on front-to-back alpha compositing. Here, $S(p)$ is the set of indices that correspond to the front-to-back sorted Gaussians overlapping pixel p and $T(i) = \prod_{j=0}^i (1 - o_j)$ is the transmittance function that takes the opacity of all closer Gaussians into account. 3D Full-body Codec Avatar decoders are fully differentiable with respect to their inputs and thus can be trained end-to-end using a photometric ℓ_2 -loss with respect to the captured multi-view image datasets of the actors. We use Adam (Kingma and Ba, 2015) and train for 300k iterations using 4 GPUs with a batch of 4 randomly sampled images.

5.2.2 Codec Avatar Face Adapter

Our Codec Avatar decoders are person-specific full-body digital humans that have been trained with personalized inputs. The face expression features produced by the dyadic motion model, however, are generic, non-personalized expression codes. These features do not perfectly disentangle expression and identity, causing moderate identity information to persist in the latent expression space. As a result, the generic expression space has multiple distinct features to represent the same facial expressions based on different regions of the latent space due to appearance and identity bleed. For the Codec Avatar decoder, which is trained only on a small amount of person-specific data, this poses a challenge due to the domain gap between the two modalities. In training, the decoder has never seen expression examples from an expression space region outside of the single identity on which it has been trained. Thus, using the generic expression features directly to drive the 3D Full-body Codec Avatar’s facial expressions would result in out-of-domain examples and uncanny facial expressions and rendering artifacts.

To reliably create authentic facial expressions, all expression inputs should ideally fall into the expression subspace of the single identity on which the decoder has been trained. We bridge this gap by training a model to map from generic expression features to person-specific codes, which are within the person-specific decoder domain by design (see Figure 14, left).

Dataset creation. We first create a dataset of expression pairs of arbitrary people and person-specific expressions of our 3D Full-body Codec Avatar. To this end, we use 800K images from a 2D face dataset, compute their (non-personalized) expression features, and re-render the expressions with the Codec Avatar’s frontal face appearance using the 2D rendering pipeline. From these renders, we extract expression codes which now fall only into the subspace of expressions of the specific decoder identity. The result is a paired

dataset of generic, non-person-specific expressions features with the corresponding person-specific expression codes that are expected as input by the decoder.

Face Adapter Model. We train a model that maps from generic expression features to person-specific expression codes (Figure 14, left). Since the decoder is trained only on a small amount of data, it is crucial that the face adapter should produce temporally coherent personalized expression codes and should suppress expression outliers that the decoder cannot handle. We therefore model the adapter as a temporal variational autoencoder (VAE) that maps from generic expression features to person-specific latent codes using 1D-convolutions. We empirically find this architecture to produce higher quality expressions than frame-wise face adapters.

5.2.3 Codec Avatar Body Adapter

The dyadic motion model outputs a SMPL-H body representation that is not compatible with 3D Full-body Codec Avatars. First, SMPL-H itself is a low resolution rig, whereas Codec Avatars are built upon a higher-resolution avatar rig. Second, the specific SMPL-H output produced by the dyadic motion model only generates upper body motion but ignores lower body and leg motion, and only produces fixed body size and shape. Additionally, since the SMPL-H features are extracted from monocular video, depth ambiguity and occlusions lead to inaccurate 3D features. These do not affect 2D rendering from the same camera viewpoint but become visible in 3D rendering, where the camera viewpoint can change.

All these aspects together result in an information deficiency of the SMPL-H body representation, more specifically: lower rig resolution, missing lower body motion, ambiguities from monocular tracking, and required motion adjustments to fit the size of the Codec Avatar. We bridge this gap based on a generative body motion adapter that infuses the missing information (Figure 14, right). We start by curating a paired SMPL-H to 3D Full-body Codec Avatar dataset, and then train a generative diffusion model to map from the information deficient SMPL-H representations to the Codec Avatar skeleton rig.

Dataset creation. We use a dataset of densely tracked 3D full body meshes from a multi-view camera dome (Bagautdinov et al., 2021), similar to the full-body data of (Martinez et al., 2024). The dataset includes 128 participants performing various motions for about 45 minutes each. We extract frontal view images and then extract SMPL-H features similar to the features obtained from the SEAMLESS INTERACTION dataset and then additionally track the 3D Codec Avatar skeleton rig based on the multi-camera data. The pairing of the SMPL-H features to the 3D Codec Avatar rig parameters is used as training data for the generative body adapter model.

Body adapter model. The body adapter consumes, as input, SMPL-H joint angles for the upper body and hands, and produces, as output, full-body joint angles for the high-resolution Codec Avatar skeleton rig. In this process, lower body motion needs to be synthesized, high-resolution rig-details need to be inpainted, and movements need to be re-targeted to a rig with differences in bone lengths.

Since all of this requires infusing information not present in the SMPL-H input, we formulate the problem as a temporal generative model. A lightweight temporal diffusion transformer with four layers first generates sparse Codec Avatar meshes from the SMPL-H parameters given additionally target bone lengths as input. Then, an MLP regresses joint angles for the Codec Avatar skeleton rig from the output mesh. The model employs a limited self-attention window of 9 frames per layer, resulting in a total receptive field of 1.2 seconds. We train the diffusion model with a simple ℓ_2 -loss on vertices in Codec Avatar mesh space, and train the mesh-to-pose regressor with an ℓ_2 -loss in joint angle space.

5.2.4 Codec Avatar Results

We train several 3D Full-body Codec Avatar decoders (3 male and 4 female) from the collected multi-view performance capture data. Each decoder is trained for 300K iterations using 4 NVIDIA H100 GPUs, which takes approximately 1 day. A subset of the trained decoders is shown in Figure 15. Our 3D Full-body Codec Avatar decoders are fully in 3D and completely model the actors, including their heads, hair, hands, arms, clothing, and the rest of their bodies. The decoders are high fidelity and fully drivable; i.e., their face and

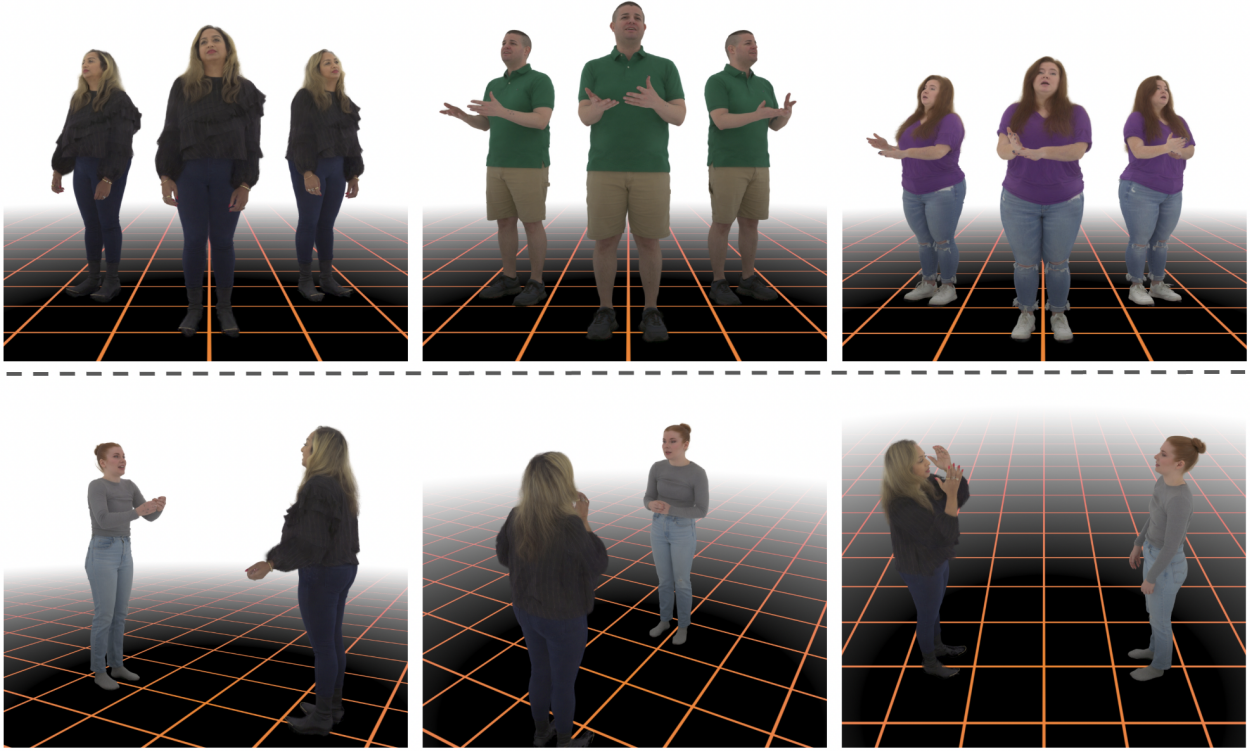


Figure 15 - 3D Full-body Codec Avatar Decoders. **Top:** Our full-body decoders are of high fidelity and their face/body can be animated based on the motion features produced by the dyadic motion model. Since the decoders are fully in 3D, they can be rendered from arbitrary camera viewpoints. **Bottom:** Since our decoders are fully in 3D, we can place two 3D Full-body Codec Avatar decoders into a shared world space coordinate system to simulate a dyadic conversation. In addition, the dyadic conversations can be rendered from arbitrary camera viewpoints.

body can be animated using the proposed face/body adapters that convert from the motion synthesized by the dyadic motion model to the driving representation expected by the decoder.

The 3D nature of these decoders allows us to render them from arbitrary camera viewpoints as can be seen in [Figure 15](#) (top). The diversity of the decoders in terms of face/body shape, hair style, and clothing is to be noted. We also demonstrate that two 3D Full-body Codec Avatar decoders can be placed into a shared world space coordinate system to render multiple actors in the same space. Additionally, by placing two decoders into a shared world space coordinate system and making them face each other, we can simulate a dyadic conversation between two avatars. Since the decoders are fully in 3D, the dyadic conversations can be rendered from arbitrary camera viewpoints (see images based on a circular camera trajectory around the two avatars shown in [Figure 15](#), bottom). These renderings highlight the spatial relationship between the two actors and create the sense of a real conversation in a shared space. Both decoders are driven by the outputs of the face/body adapters and feature expressive face and upper body motion based on the employed dyadic motion model that produces realistic hand/arm gestures.

5.2.5 Limitations

While we show promising neural rendering results for drivable 3D Full-body Codec Avatars, some limitations still remain to be addressed in future work, some of which are unique to the 3D rendering setting.

First, neither the dyadic motion model nor the body adapter are currently aware of the actor’s body shape beyond their bone length. This can lead to self-interpenetration; e.g., hand-hand interpenetration or the arms moving partially inside the torso if the regressed body motion is incompatible with the actor’s body shape ([Figure 16](#), left). This can be resolved by training body shape-aware motion models.

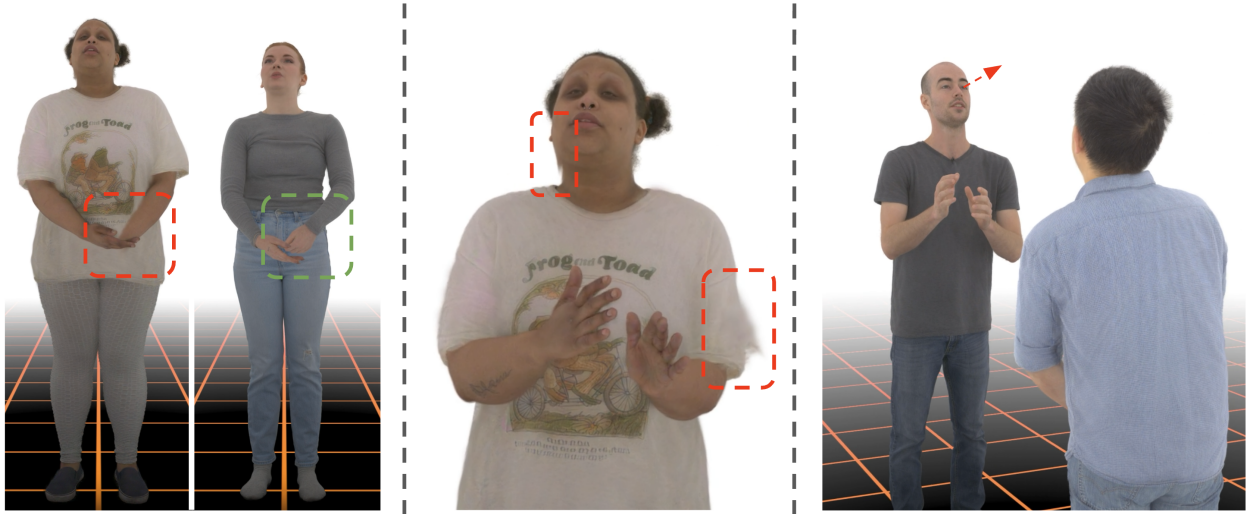


Figure 16 - Limitations. **Left:** The dyadic motion model is unaware of person-specific body shapes. Therefore, self-penetrations can happen for some decoders while other decoders work perfectly fine. **Center:** Capacity limitations of the model lead to artifacts in clothing and appearance. **Right:** In the dyadic 3D rendering, eye contact is not explicitly established between the actors.

Second, our decoders are per-frame models and employ linear blend skinning (LBS) to explicitly model the articulated motion of humans by deforming a template mesh with tightly attached Gaussian primitives. While our decoder also produces per-Gaussian position correctives that allow to deviate from the LBS motion, our models still struggle with large non-articulated motions, such as the dynamic motion often observed for long hair and loose clothing (Figure 16, center). Further research on temporal decoders as well as multi-layer representations that model body and hair/clothing motion independently is needed to improve these results.

Finally, the training data for the dyadic motion model are two independent monocular videos that are temporally in sync, but not spatially calibrated with respect to each other. Thus, the motion model has no awareness of the spatial relationship between the two interlocutors. As a result, when two 3D Full-body Codec Avatars are placed facing each other in 3D space, they will not necessarily establish eye contact (Figure 16, right). This can be resolved in the future by training the dyadic motion model based on 3D motion data that captures the spatial relationship between both interlocutors in a shared world space coordinate system.

6. Evaluation Methodology and Experiments

Evaluating dyadic behavior in embodied agents remains an unsolved problem. However, there is increasing interest in this area, as demonstrated by the recent series of GENE Challenges and leaderboards, as well as ongoing efforts to develop automatic metrics (Yoon et al., 2022; Kucherenko et al., 2023a; Nagy et al., 2024). In the following section, we present a series of calibration experiments intended to build further traction on the question of how to measure dyadic-interaction quality.

We start by presenting an approach to human subjective (user and annotator) studies that incorporates face- and body-dyadic protocols. Next, we describe a series of ablation studies using automatic metrics commonly examined by the community. Finally, we explore the relationship between a pair of face-centric automatic metrics and the corresponding human subjective results.

6.1 Human Studies

We introduce two user-evaluation protocols: a *face-dyadic* protocol, which emphasizes facial expressions and head movements in photorealistic renderings, and a *body-dyadic* protocol, which focuses on visual behavior involving overall pose, hands, arms, shoulders, and head, but without facial rendering.

In both studies, we use a pairwise approach where participants view two dyads side-by-side. Each pair consists of an *Anchor* and a *Candidate*. The *Anchor* videos provide the full context of an interaction, while the *Candidate* videos, which include either rendered ground-truth (GT) or model-based generations, are the focus of human ratings. Participants are asked to provide a preference rating for each pair of stimuli, choosing from five possible values: $\{-2, -1, 0, 1, 2\}$, corresponding to $\{\text{“Much prefer A”}, \text{“Slightly prefer A”}, \text{“Tie”}, \text{“Slightly prefer B”}, \text{“Much prefer B”}\}$.

Our work differs from previous dyadic studies in several essential ways. First, our *Anchor* stimuli include actual RGB video footage of one of the participants from the real dyadic conversation. Second, instead of presenting short, 10-second segments, as done in the speaking-focused studies reported in Kucherenko et al. (2023a), we present 20-second segments that include multiple speaking turns, similar to the “Interlocutor” setup in Kucherenko et al. (2023a). Finally, we significantly expand the number of evaluative dimensions to cover a wide range of quality attributes for both speaking and listening behaviors.

6.1.1 Protocol and Evaluative Dimensions

Both the face- and body-dyadic protocols are based on 10 core evaluative dimensions. The first set of dimensions focuses on overall preferences, the second on listener-behavior preferences, and the third on speaking-behavior preferences. Before each of the three sections, participants are provided with a video player, allowing them to view the stimuli while answering questions.

These evaluation dimensions include lifelikeness, clarity of intent, turn-taking, listening and speaking. We provide the text of each protocol in its entirety in Appendix B.

6.1.2 Face Dyadic Study

Data. We selected $n=61$ segments, each 20 seconds long, from an earlier version of the SEAMLESS INTERACTION test set. These segments were chosen to include at least 2 or more speaking turns, with the number of turns ranging from 2 to 5 within the 20-second period. A 2-second buffer is allowed at the beginning and end of each sample to ensure that speaking behavior remains within phrase boundaries. Our focus is on the visual behavior of our models during both speaking and listening, which is why each sample requires the agent to have at least one speaking turn and one listening turn.

Stimuli and ratings. Each sample consisted of a pair of dyadic interaction videos presented side-by-side. Both videos share the same audio track, which is actual speech data from the SEAMLESS INTERACTION test-set samples. Visually, each video includes one ground-truth RGB *Anchor* (cropped to show information from the shoulders up) and one model- or ground truth-rendered *Candidate* (see Figure 17a). To reduce cognitive load, additional labels (*Anchor* and *Candidate*) were added to clearly identify the ground-truth interlocutor from the generated stimuli. Additionally, we use “A/B” labels for the *Candidate* stimuli to provide clear visual cues to emphasize left and right positions, minimizing potential ambiguity.

A video synchronizing button was added, enabling simultaneous playback of videos and simplifying side-by-side comparisons if desired, however participants were also able to play the dyad videos independently.

We used the same face rendering pipeline for each *Candidate* video with voice-matched gender of the renderings.

In this setup, each test item allows for a comparison within one pair of models. In total we have $n = 610$ test-item pairs (61 samples times 5 choose 2 pairs). Each test-item receives five ratings from different study participants.

Models. We compare four preliminary models to an actual ground-truth interlocutor renderings (see Table 12).

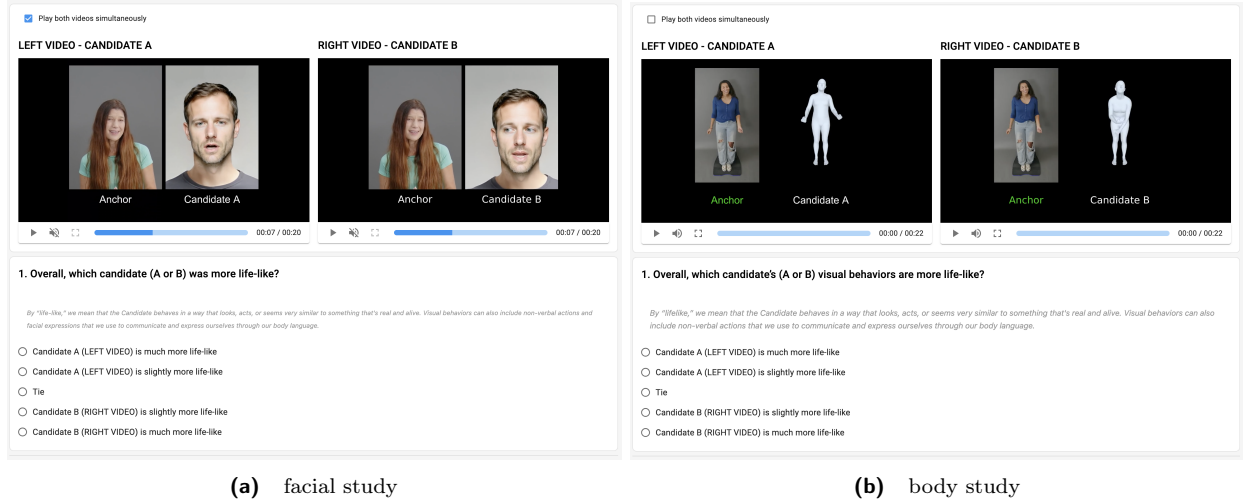


Figure 17 - Examples of actual response interfaces for body and facial studies

Name	Model
GT	Rendering of ground-truth human interlocutor
A	Dyadic Joint Face+Body with Windowed Self-attention
B	Dyadic Joint Face+Body with Self-attention
C	Dyadic Self-attention without Expression normalization
D	AV Dyadic Self-attention without Expression normalization

Table 12 - Preliminary models used in the Face-dyadic human subjective calibration study.

Participants. A total of $n=81$ participants were recruited through a third-party annotation service with the following requirements: participants needed to be native speakers of English residing in the USA, Canada, the UK, or New Zealand. Throughput varied by participant, with the average participant providing $n=37.7$ ratings (min=1, max=86). Studies began with an introduction to the study including example stimuli to orient participants to the task. Participants were paid upon completion of the study, relative to the number of rated samples.

Analysis and Discussion. We compute item-level scores by taking the mean preference rating across the five ratings at the item level. We then compute the average rating for each model match-up across items. Figure 18 shows estimated mean preference scores by match-up with each facet corresponding to the average preference rating for a given model compared to the competitor set. (For readability, we only present six of the 10 evaluative dimensions in Figure 18, and note that results are directionally consistent across all dimensions.)

Results indicate a clear ordering across the six dimensions, in which all models are dispreferred to GT. Model A outperforms all other generations, coming close to GT performance. Model B outperforms Models C and D, but lags Model A and GT. Finally Model C outperforms Model D, the worst performing model overall.

The finding that at least one of our models (Model A) is close to GT renderings warrants additional commentary. First, this pattern is validated by internal qualitative and quantitative user research (UXR) studies and manual inspection by the team - in many cases it is extremely difficult to distinguish between model- and GT-renderings. However, it should also be noted, similar to prior work (Kucherenko et al., 2023a), that in some cases we observe artifacts that are unique to the GT rendering process and are not present in model generations. These artifacts may include what is perceived by a human evaluator as jitter, stretching, or other discontinuities between body parts (such as misalignment between the head and shoulders). Such artifacts represent confounds stemming from the underlying representation extraction and rendering process, and may be responsible, in part, for the closeness of ratings to GT. Future work will attempt to better characterize

such artifacts either through manual inspection or the use of data-analytic techniques.

An additional reason why we see such close scores with human GT may stem from the level of activity of various samples. That is, we observed in prior studies that low-activity (i.e., boring) segments of actual GT human behavior is typically dispreferred when compared to more active synthetic behavior. This *activeness preference* requires deeper analysis, and future work may employ various measures of activity to investigate it directly.



Figure 18 - Face-dyadic study results. Horizontal axis displays average rating by match-up ranging from -2 (Much prefer left) to 2 (Much prefer right). Each facet shows a different match-up relative to a single model. For example the top facet shows average ratings compared to GT renderings. The vertical axis within each facet shows the competitor model for a given match-up. Error bars represent analytically computed 95% CIs.

Name	Model
GT	Rendering of ground-truth human interlocutor
A	Dyadic Joint Face+Body with Windowed Self-attention
B	Dyadic Joint Face+Body with Self-attention
E	AV Dyadic Self-attention
GT-N	Negative sample (shuffled) rendering of ground-truth human interlocutor

Table 13 - Preliminary models used in the Body-dyadic human subjective calibration study.

6.1.3 Body Dyadic Study

Data. A total of $n=71$ 20-second segments were selected from an earlier version of the SEAMLESS INTERACTION test-set¹.

Stimuli and ratings. Participants were again presented with two pairs of videos side-by-side. However, rather than showing the *Anchor* stimuli cropped from shoulders up, they were shown the actual full-body RGB video. For the *Candidate* stimuli, full-body, faceless, mesh-based renderings (SMPL-H) were presented (see Figure 17b for an example of what participants saw).

The absence of facial characteristics or lip-movements in the body-dyadic study can make it difficult to follow conversation dynamics (i.e., who is speaking). To address this, we provided an additional cue to the source of speech by highlighting the video label (*Anchor* or *Candidate*) with green when VAD-detected activity is present. This feature enabled participants to easily follow the conversation dynamics, thereby reducing cognitive load. The left-facet of Figure 17b provides an example in which the *Anchor* video is speaking.

Models. We compare three preliminary models to a ground-truth interlocutor as well as to a baseline system consisting of negative-sampled (shuffled) ground-truth videos (see Table 13).

Participants. A set of $n=105$ participants were recruited through a third-party annotation service, with the following requirements: participants needed to be native speakers of English residing in the USA, Canada, the UK, or New Zealand. Throughput varied by participant, with the average participant providing $n=33.8$ ratings (min=1, max=91). Studies began with an introduction to the study, including example stimuli to orient participants to the task similar to what is shown in Figure 17b. Participants were paid upon completion of the study, relative to the number of rated samples.

Analysis and Discussion. We compute item- and match-level scores as in Section 6.1.2 and present a reduced set of evaluative dimensions in Figure 19.

Results indicate that two of the models (A and B) are virtually indistinguishable from GT mesh renderings, with even a directional preference for Model A over GT. Model E only outperforms the naive baseline (GT-N) which includes randomly shuffled ground-truth behavior.

Similar to the commentary provided in Section 6.1.2, these findings require additional attention, as both the presence of rendering artifacts and the possibility of an *activeness preference* effect may be partially responsible for the proximity to GT performance. Prior internal studies again indicate that models are extremely difficult to distinguish with GT, however we also observe the occasional presence of body-representation artifacts from the GT rendering process, which can be subtle and pervasive. Issues such as *collisions* (in which limbs such as the hands appear to merge or disappear upon contact with the body or one-another), *jitter*, and *skating* (in which the rendering appears to be unanchored to the ground) represent potential confounds in the current analysis. These observations underscore the importance of developing tools to detect and quantify these artifacts as essential future work.

¹These samples differ from the set used in the face-dyadic study due to differences in item-rendering workflows; however, they are both sampled from the same test-set.



Figure 19 - Body-dyadic study results. Horizontal axis displays average rating by match-up ranging from -2 (Much prefer left) to 2 (Much prefer right). Each facet shows a different match-up relative to a single model. For example the top facet shows average ratings compared to GT renderings. The vertical axis within each facet shows the competitor model for a given match-up. Error bars represent analytically computed 95% CIs.

6.2 Experiments with Automatic Metrics

6.2.1 Dyadic Motion Model Analyses

This work focuses on interactive head and body generation in a dyadic setting. In this section, we report the results of evaluations on generation quality with commonly used automatic metrics. Extensive empirical results compare various architectures and training recipes of diffusion models.

Test set. A set of 200 dyadic samples were selected from the test split of SEAMLESS INTERACTION. Each sample is a video segment of dyadic interactions, which consists of 2 – 5 turns between two speakers for around

20 seconds.

Models. Independent face and body generations may result in inconsistent movements; for example, head and body might be moving in opposite directions. Therefore, we explored combining face and body modeling for more synchronized movements with either a cascaded approach or a joint manner. In the cascaded approach of *Face2Body*, we first trained a face diffusion with dyadic audios, and conditioned the body diffusion on the face Imitator features. Similarly, we obtained the cascade of *Body2Face*, where face diffusion is conditioned on the body’s SMPL-H features. Lastly, we had *Joint Face+Body*, where a single model was trained for both face and body generation. With the joint model, we conducted ablation studies of different conditions in the diffusion model: (1) Monadic: diffusion is conditioned on single-channel audio; (2) Dyadic: diffusion condition is dialogs; (3) AV: the conditions include both audios and visual features.

Diffusion models consist of 12 Transformer layers with hidden dimension of 1024 and feedforward dimension of 4096. They are trained with flow-matching objective with condition dropout of 0.2. During inference, we apply a CFG weight of 1.5, and use 100 steps ODE (Song et al., 2021a) for all experiments.

Automatic metrics. To quantitatively evaluate the visual quality of face generation, we use Fréchet Feature Distance (FFD), Lip-Sync score (Sync-C and Sync-D) and Fréchet Inception Distance (FID). FFD is the Fréchet distance (Dowson and Landau, 1982) of Imitator features between predictions and ground truth. Sync-C and Sync-D (Prajwal et al., 2020) assess the lip synchronization with speech. FID (Heusel et al., 2017) measures the face image quality based on features encoded by the Inception network.

As for the body, we apply Fréchet Gesture Distance (FGD) as used in previous works (Kucherenko et al., 2023b). It quantifies the discrepancy in distribution between generated outputs and real data across all samples. Diversity (Liu et al., 2022b) is another metric for gesture generation, which measures the range of variations present in the generated gestures.

System	Condition	Face				Body	
		FFD (\downarrow)	Sync-C (\uparrow)	Sync-D (\downarrow)	FID (\downarrow)	FGD (\downarrow)	Diversity (\uparrow)
Dyadic Face2Body	A1+A2	0.11 \pm 0.01	1.72 \pm 0.01	7.06 \pm 0.01	1.79 \pm 0.04	0.93 \pm 0.14	3.83 \pm 0.10
AV Dyadic Face2Body	A1+A2+V2	0.10 \pm 0.01	1.68 \pm 0.01	7.10 \pm 0.01	1.53 \pm 0.02	1.09 \pm 0.18	3.82 \pm 0.08
Dyadic Body2Face	A1+A2	0.20 \pm 0.02	2.34 \pm 0.01	7.44 \pm 0.01	1.81 \pm 0.04	1.61 \pm 0.19	4.02 \pm 0.05
<i>Joint Systems</i>							
Monadic Face+Body	A1	0.29 \pm 0.01	2.54 \pm 0.01	7.36 \pm 0.01	2.27 \pm 0.12	1.22 \pm 0.24	3.42 \pm 0.08
Dyadic Face+Body	A1+A2	0.26 \pm 0.02	2.48 \pm 0.01	7.46 \pm 0.01	1.89 \pm 0.06	1.73 \pm 0.26	3.91 \pm 0.03
AV Dyadic Face+Body	A1+A2+V2	0.26 \pm 0.01	2.24 \pm 0.01	7.84 \pm 0.01	1.70 \pm 0.04	0.89 \pm 0.08	3.71 \pm 0.04

Table 14 - Evaluation results of diffusion models on face and body generation. We run generation on the test set 5 times and report the mean (\pm standard deviation) values for each metric. (A1: model speech, A2: user speech, V2: user visual.)

Results. Table 14 reports automatic metrics of various diffusion models which use standard self-attention to attend to input conditions. For face generation, *Dyadic Face2Body* achieves the best Sync-D, *AV Dyadic Face2Body* gets the best FFD and FID, and *Monadic Face+Body* has the best Sync-C. As for body generation, *Joint AV Dyadic Face+Body* obtain the lowest FGD and *Dyadic Body2Face* has the highest gesture diversity.

- **Cascaded versus joint model.** Among *Dyadic Face2Body*, *Dyadic Body2Face* and *Dyadic Joint Face+Body*, *Face2Body* has the best FFD and FGD, while *Body2Face* and *Joint Face+Body* have comparably good performances on lip synchronization and gesture diversity.
- **Monadic versus dyadic input.** The comparison between *Monadic Face+Body* and *Dyadic Face+Body* suggests that the monadic condition helps to have slightly better lip synchronization and lower FGD, but the dyadic condition achieves better face quality with lower FFD and FID, and it also leads to more diverse gestures.
- **Visual input condition.** Considering *Dyadic Face2Body* and *AV Dyadic Face2Body*, we see that adding visual condition helps to decrease FID. Similarly for *Joint Dyadic Face+Body* models, the addition of

	Face				Body	
	FFD (\downarrow)	Sync-C (\uparrow)	Sync-D (\downarrow)	FID (\downarrow)	FGD (\downarrow)	Diversity (\uparrow)
Standard self-attention	0.26 ± 0.02	2.48 ± 0.01	7.46 ± 0.01	1.89 ± 0.06	1.73 ± 0.26	3.91 ± 0.03
Standard cross-attention	0.14 ± 0.01	0.31 ± 0.00	10.25 ± 0.02	2.38 ± 0.09	1.59 ± 0.37	3.75 ± 0.07
Windowed self-attention	0.17 ± 0.01	1.95 ± 0.01	8.17 ± 0.02	3.26 ± 0.06	2.52 ± 0.50	3.90 ± 0.06
Windowed cross-attention	0.13 ± 0.01	0.87 ± 0.01	8.92 ± 0.02	2.12 ± 0.07	1.25 ± 0.37	3.82 ± 0.08

Table 15 - Comparison of the architectures of conditioning in a joint diffusion model. We sample 5 generations for each test sample and report the mean (\pm standard deviation) for each metric.

visual condition brings down face FID and body FGD.

Ablation of joint model architectures. With Transformer as the backbone of the joint model conditioned on dyadic audios, we further experiment with different architectures for diffusion conditioning: (1) conditioning on dyadic audios via standard self-attention as used by models in Table 14; (2) conditioning via standard cross-attention; (3) conditioning via windowed self-attention, which limits self-attention to a window of local context, and the intuition is that motion in one frame is more relevant to neighboring frames; (4) conditioning via windowed cross-attention. Table 15 demonstrates both face and body metrics. Windowed cross-attention yields the lowest FGD in body generation. Both standard and windowed cross-attention have comparable FFD in face generation, and outperform self-attention conditioning. However, they result in bad lip synchronization, as reflected by low Sync-C scores. Standard self-attention outperforms other architectures in Sync-C, Sync-D, FID as well as gesture diversity.

	Face				Body	
	FFD (\downarrow)	Sync-C (\uparrow)	Sync-D (\downarrow)	FID (\downarrow)	FGD (\downarrow)	Diversity (\uparrow)
Full Imitator Latent	0.10 ± 0.01	1.68 ± 0.01	7.10 ± 0.01	1.53 ± 0.02	1.96 ± 0.20	3.25 ± 0.10
Head Rotation					1.09 ± 0.18	3.82 ± 0.08

Table 16 - AV Dyadic *Face2Body* conditioning ablation. We sample 5 generations for each test sample and report the mean (\pm standard deviation) for each metric.

Ablation of cascaded models. For the cascaded model *Face2Body*, we further study what useful information face diffusion could pass to body diffusion model as tabulated in Table 16, including: (1) full imitator latent; (2) head rotation of imitator latent; The body diffusion model taking head rotation as condition reduces FGD as well as enhances gesture diversity. It indicates that face’s rotation condition is most effective in aligning body with head pose, while conditioning on additional face information will hurt the quality of body generation.

6.2.2 Analysis of Gesture Controllability

Evaluation Metrics. Regarding the effectiveness of gesture conditioning, we care about two things most: condition following and smoothness at the gesture condition boundary.

- **Condition Following.** To evaluate the model’s ability to follow gesture conditions, we compute the L2 reconstruction error between the gesture conditioning part of the generated gesture sequence and the ground truth gesture sequence. The semantic gesture condition is a sequence $\mathcal{G}_{\text{sem}} = \{\mathbf{g}_j^{\text{sem}}\}_{j=0}^{T_{\text{sem}}}$ outside the training set. We set the semantic gestures to start at the t_{start} index in the generated gesture sequence.

The model generates a gesture sequence $\mathcal{G}_{\text{gen}} = \{\hat{\mathbf{g}}_i\}_{i=1}^{T_s}$ conditioned on \mathcal{G}_{sem} for time steps $i \in [t_{\text{start}}, t_{\text{end}}]$ and the speech condition \mathbf{s} for all $i \in [1, T_s]$, where $t_{\text{end}} = t_{\text{start}} + T_{\text{sem}}$. For time steps outside $[t_{\text{start}}, t_{\text{end}}]$ and $1 \leq t_{\text{start}} < t_{\text{end}} \leq T_s$, the model is conditioned on \mathbf{s} solely. The L2 reconstruction error is computed over the gesture conditioned part to assess the model’s ability to follow the semantic gesture condition:

$$\mathcal{L}_{\text{recon}} = \frac{1}{T_{\text{sem}}} \sum_{i=t_{\text{start}}}^{t_{\text{end}}} \|\mathbf{g}_i - \hat{\mathbf{g}}_i\|_2^2,$$

where $\|\cdot\|_2^2$ denotes the squared Euclidean norm. This error evaluates the fidelity of the generated sequence \mathcal{G}_{gen} to the ground truth \mathcal{G} , particularly emphasizing adherence to \mathcal{G}_{sem} within $[t_{\text{start}}, t_{\text{end}}]$.

- **Boundary Smoothness.** It is crucial that the motion transition is temporally smooth from gesture condition OFF to ON and ON to OFF. Let $\mathbf{P}(t) \in \mathbb{R}^{N \times 3}$ represent the 3D keypoint positions at time t , where N is the number of keypoints and each keypoint has coordinates (x, y, z) . For a temporal sequence of length T , we have $\{\mathbf{P}(t_i)\}_{i=1}^T$ where $t_i = i \cdot \Delta t$ and $\Delta t = 1/\text{fps}$ is the temporal sampling interval.

The smoothness metric is derived from the jerk, defined as the third-order temporal derivative of position. For each keypoint $j \in \{1, 2, \dots, N\}$ and spatial dimension $d \in \{x, y, z\}$, we compute:

$$\mathbf{j}_{j,d}(t_i) = \left. \frac{d^3 \mathbf{P}_{j,d}(t)}{dt^3} \right|_{t=t_i}$$

The jerk magnitude for each keypoint j at time t_i is computed as the Euclidean norm across spatial dimensions:

$$|\mathbf{J}_j(t_i)| = \sqrt{\mathbf{j}_{j,x}(t_i)^2 + \mathbf{j}_{j,y}(t_i)^2 + \mathbf{j}_{j,z}(t_i)^2}$$

The overall jerk metric is defined as the temporal and spatial average of jerk magnitudes:

$$\bar{J} = \frac{1}{T \cdot N} \sum_{i=1}^T \sum_{j=1}^N |\mathbf{J}_j(t_i)|$$

The final smoothness score S is computed using an exponential decay function to map jerk values to a normalized smoothness measure:

$$S = \exp\left(-\frac{\bar{J}}{\sigma}\right) \quad (4)$$

where σ is a scaling parameter that controls the sensitivity of the smoothness score to jerk variations. In our implementation, $\sigma = 100$ provides empirically reasonable behavior across typical motion capture datasets. The final boundary smoothness is determined by averaging the smoothness across gesture condition transitions from ON to OFF and OFF to ON. In our experiment, the smoothness is computed using a 30-frame window, spanning 15 frames before and after each boundary timestamp.

Evaluation Results. Table 17 shows the comparison of semantic gesture control with different conditions and dropout.

- **Condition following.** The gesture VQ-ID conditioned diffusion falls behind on the condition following with a large gap, while having better boundary smoothness compared with SMPL-H conditioned diffusion. This is expected since VQ has a much higher compression rate during the quantization process. The condition following metric values for SMPL-H conditioned diffusion are in a similar scale and do not have a large difference, both visually and quantitatively.
- **Boundary smoothness.** From the ablation experiment on varying temporal gesture condition dropping rates for SMPL-H conditioned diffusion, we surprisingly find that higher temporal SMPL-H condition dropping rates can lead to smoother boundary transitions. This observation is particularly evident in the visual results — with a lower temporal SMPL-H dropping rate 0.4, the gesture transition at the boundary looks very sudden; while when the dropping rate comes to 0.8, the gesture transition at

Condition	Temporal Drop	SMPL L2 Error	Keypoint L2 Error	Boundary Smoothness \uparrow
VQ IDs	0.4	0.35	0.16	0.66
SMPL-H	0.8	0.05	0.02	0.61
	0.6	0.06	0.03	0.54
	0.4	0.04	0.02	0.54
	0.2	0.03	0.02	0.48

Table 17 - Semantic gesture control comparison and ablation. The gesture VQ ID-conditioned diffusion falls behind on the condition following with a large gap, while having better boundary smoothness compared with SMPL-H conditioned diffusion. Ablation on different temporal gesture conditions dropping rate indicates that a higher temporal SMPL-H condition dropping rate led to smoother boundary transitions. The condition following is evaluated by SMPL-H and keypoints reconstruction error, and the boundary smoothness is evaluated via Equation (4).

the boundary looks much more natural and smoother. Therefore, a high gesture condition temporal dropping rate is crucial for making SMPL-H conditioned outputs appear natural.

The VQ-conditioned diffusion exhibits the highest boundary smoothness. This may be attributed to the high abstraction level of VQ tokens, which provides the diffusion process with more flexibility to generate smooth motion that aligns with the overall distribution at the gesture condition boundary.

6.2.3 Analyses of LLM-Guided Codebook Generations

We present results of for LLM-guided codebook generation as discussed in Section 4.5.1. Table 18 summarizes the total number of tokens for valence, arousal and gesture in the dataset.

	Train	Valid	Test
Valence	4.4M	0.1M	0.7M
Arousal	4.4M	0.1M	0.7M
Gesture (w/o “null” token)	22.6k	2.2k	2.3k

Table 18 - The number of emotion and gesture tokens in training, valid and test data.

Emotion Adapter. We set the token rate of emotion adapters as 1 token per second, assuming emotion consistency within a one-second window. As the ground truth emotion is extracted in the frame level, we average the emotion values over all frames within one second as the training target of Adapter.

The Adapter prediction is based on acoustic and semantic information from speakers, while emotion labels are extracted from visual signals. Therefore, there might be a gap in fine-grained crossmodal emotion prediction. For example, speakers could exhibit varying intensities in facial expressions even when talking about the same thing. We group 12 fine-grained emotion tokens into 3 coarser grained groups. Accuracy of 3-class prediction is the evaluation metric we report for emotion prediction. The accuracy of emotion adapters is 0.51 for valence and 0.52 for arousal.

Gesture Adapter. Semantic gestures fall within the long tail of the distribution of human motion. Most of the time, speakers do not make semantic gestures, and therefore the “null” token is the dominant gesture label associated with 98% speech segments. Furthermore, the distribution varies a lot for different gestures, and precision, recall and F1 score are good metrics for imbalanced data. We compute the score for each gesture except “null” gesture, and report macro-averaged scores by applying the arithmetic mean to per-gesture scores.

Empirically we tried two token rates for gesture adapter, and trained adapters with 1 and 2 gesture tokens per second. Table 19 reports results on gesture prediction. The adapter with token rate of 2 gives better performance than that with token rate of 1. This suggests that the window of one second is a bit large, cover multiple spoken words. Reducing the window size helps Adapter capture the word of interest.

Token rate	Precision	Recall	F1 score
1 token/sec	0.47	0.30	0.37
2 token/sec	0.51	0.47	0.49

Table 19 - Macro-averaged precision, recall and F1 score of gesture prediction.

Proxy Metric	Pearson’s r	Kendall’s τ	Spearman’s ρ
FFD	0.562	0.573*	0.406*
Sync-C	-0.307	-0.328*	-0.223*

Table 20 - Correlation of human subjective ratings for “Overall lifelikeness” with three automatic metrics. We present three measure of association via Pearson’s r , Spearman’s ρ , and Kendall’s τ . “*” denote statistical significance at $\alpha = 0.05$.

We analyze examples where gesture predictions disagree with the ground truth, and find that Adapter assigns gestures to synonyms and phrases in a similar context as the word triggering semantic gestures. For example, the gesture control vocabulary contains an illustrative gesture for the word “cold”. The Adapter also labels “freezing” with the same gesture. Also it assigns the utterance “don’t do that” with the gesture for “stop”. These are examples lowering the precision. In other examples, when the word was spoken fast, Adapter missed prediction and got lower recall.

6.3 Relationship Between Automatic Metrics and Evaluation Dimensions from Human Studies

We explore the relationship between subjective human results and two automatic measures of generation quality (FFD and Sync-C) on preference data from our face-dyadic study (see [Section 6.2](#) for implementation details of these metrics).

Proxy-metrics are computed for each model at the item-level. We convert these to score deltas (for compatibility with our pairwise preference data) by calculating the differences between scores for each model on a given test-item pair.

We compare the average preference score from our subjective study with the corresponding score delta. [Figure 20](#) plots the relationship between score deltas and human preferences for FFD and Sync-C.

A significant relationship was observed for both metrics with the “Overall Like-like” dimension. FFD delta scores demonstrated a significant positive correlation with human preferences. In contrast, Sync-C exhibited statistically significant negative relationships. Results are displayed in [Table 20](#).

7. Responsible AI

7.1 Dataset Privacy and Ethics

Privacy and ethical standards were critical to the SEAMLESS INTERACTION data collection effort. Systems for ensuring privacy were introduced at multiple stages of the collection — to the participants via informed consent and constant observation by trained moderators, to the site administrators and moderators responsible for on-the-ground recordings, and finally as a part of an extensive post-collection quality assurance (QA) and filtering procedure by which every video (in total over 14,000 hours of raw footage) was analyzed for the presence of sensitive or personal material.

Participant education. Participants provided informed consent for their voice and image to be collected and used. They were instructed multiple times prior to recording to avoid topics that include personally identifiable or private information (PII) and were provided guidance from moderators upon doing so. While participants

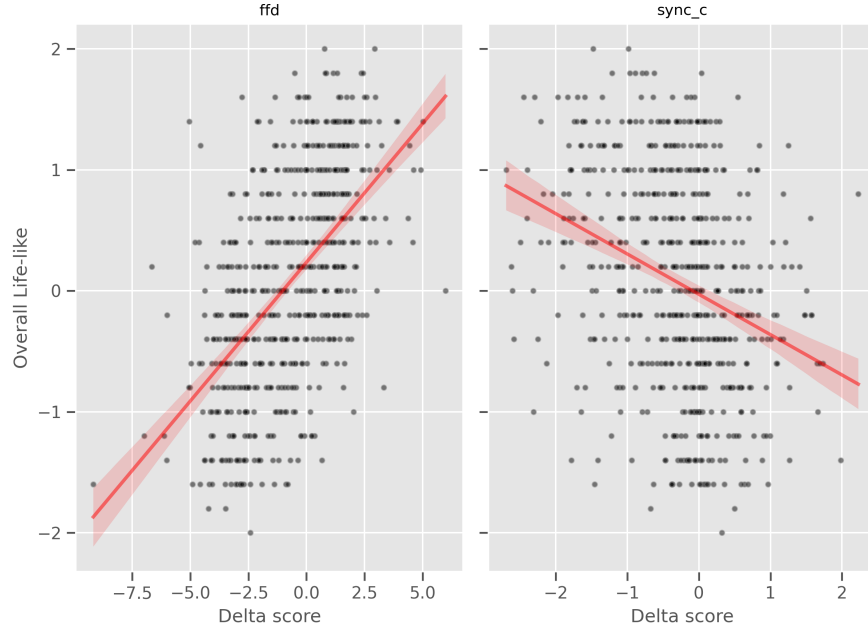


Figure 20 - Comparison of Face-dyadic study subjective preferences and proxy-metrics. Vertical axis shows item-level preference for each model-pair for our Overall Lifelike dimension. Horizontal axis shows the item-level proxy-metric score deltas. Left facet shows FFD metric, and right facet shows Sync-C. We quantify the relationship computing Kendall’s Tau, displaying the test-statistic and p-value.

were asked to respond to prompts, some of which were personal in nature, they were never compelled to do so and were free to take the conversation where they felt most comfortable (see [Appendix A.2.1](#)).

Site administrators and moderators. Site administrators and especially moderators, were educated on the nature of PII or sensitive data and the importance of indicating that it should be avoided at all times (see [Appendix A.2.2](#)). They were directed to identify and remove any content from raw deliveries, if such content did not meet our quality and safety standards.

7.2 Quality Assurance (QA) Processes

A quality assurance (QA) process was implemented to identify occurrences of sensitive material (PII), offensive material, and recording or content quality issues. This process included three sources of QA signals: (1) human-based video review and QA annotation, (2) text-LLM transcripts analyses, and (3) VLM-based video analysis. Based on combined signals from these three sources, flagged content was removed from the dataset.

7.2.1 Human QA

While human review of all 14,000 hours of raw collected footage was beyond reach, the project took as a goal that every video would have at least some portion viewed by a human. To facilitate a comprehensive analysis of video content, a stratified sampling approach was used, wherein 30-second clips were randomly extracted from each video. These clips were reviewed by human evaluators who assessed the content in accordance with a predefined guideline. To streamline the review process and ensure data accuracy, an automated data pipeline was developed, complemented by a user-friendly interface and a dashboard designed to provide real-time insights and analytical visualizations.

The human review process focused on three primary areas: (1) the mention of sensitive material (PII) or material that was offensive in nature, (2) the presence of recording artifacts, which can compromise video quality and authenticity; (3) content alignment issues, indicating discrepancies between the intended and actual content. We provide a complete list of review areas in [Appendix A.3](#).

7.2.2 Model-based QA

Although human-based review provides valuable insight into quality issues, it is limited by constraints in speed and cost. To scale QA checks to every minute of raw footage, we implemented a model-based QA approach that leverages both Large Language Models (LLMs) on transcripts and Video Language Models (VLMs) on videos.

Model-based QA performance was evaluated using ground-truth human QA annotations with a primary focus on recall of sensitive (PII) information and offensive material. This approach acknowledges that false negatives (i.e., undetected instances of offensive or sensitive content) pose a greater risk to the integrity of the dataset than false positives (i.e., incorrectly flagged content). We provide a more complete list of QA-related flags in [Appendix A.3](#).

Following the model-based checks, internal volunteers sampled a selection of flagged clips for further human review and validation. This additional layer of scrutiny provided a more comprehensive assessment of the automated checks’ performance and helped to identify any potential false positives or false negatives that may have been missed during the initial evaluation. The goal of the LLM-based and VLM-based approaches was to scalably identify segments with potential issues - analyzing every frame and word of recorded data.

LLM-based transcript analysis. We designed a prompting strategy to achieve high recall in detecting sensitive and offensive content. This was benchmarked against human-labeled results. This prompt configuration was then applied to the transcripts of the entire dataset to identify segments with potential issues.

VLM-based raw video footage analysis. A Video Language Model (VLM) was used to evaluate video content for quality assurance across the categories described in [Appendix A.3](#). Videos from the same interaction were stacked for both the participants, as input to the VLM. A context criteria describing the different QA categories in detail was provided to the VLM as a system message and the VLM was then prompted with a structured input format, requiring it to provide output evaluations in a key-value pair format for each of the QA categories.

Filtering strategy. Model performance was evaluated using the human reviewed annotations with F1-scores for text-LLM (0.84) and VLM (0.91). A conservative filtering strategy was employed - if a given interaction received any sensitive material flag for any of the sub-categories, from any of the three systems (Human, text-LLM, VLM) that interaction video was removed from the dataset. In total, this process removed several hundreds of hours of interaction data. [Figure 21](#) provides an illustration of the full QA filtering process.

7.3 Watermarking

Another proactive measure to ensure transparency, accountability, and trustworthiness is to watermark AI-generated content. This involves adding noises to original signals (video and audio) that are imperceptible to humans but can be detected by specialized algorithms. With the rapid proliferation of generative models, watermarking becomes an essential component and is now required under multiple laws, including the US Executive Order and AI, European Union’s AI Act, Labeling Rules of Cyberspace Administration of China.

We maintain our commitment to responsible AI and adopt watermarking for all of our model outputs. For audio signals, we employ AudioSeal ([San Roman et al., 2024](#)), a state-of-the-art audio watermarking solution with a robust and efficient detection algorithm. For video content, we use VideoSeal, an open and effective video watermarking method that has been shown to perform well on high-resolution videos without compromising detection robustness ([Fernandez et al., 2024](#)). Both AudioSeal and VideoSeal use localized and extractable watermarks; i.e., the models (generators) were trained to embed secret messages into individual frames of the original contents, and the detectors also extract the messages on the frame level. Each secret message is a binary string of n bits and can be defined by users. We use a fixed set of secret messages to differentiate contents from different sources: Human audios, LLM generated audios, and dyadic motion model video output. The cascaded model architecture allows for an easy composition of different audio and video watermark generators in the decoder in a post-hoc manner.

We train new AudioSeal and VideoSeal models and use them for the SEAMLESS INTERACTION dataset samples dyadic model outputs. To improve the scalability for audios, the model is trained in a causal mode, so the watermark was generated autoregressively with access to only the previous frames. The context size is 1920 frames (for 24 kHz audios), and we can watermark a very long speech in a streaming manner without significant added latency (our experiments on the A100 GPU suggested a watermarking time of approximately 5 ms per speech token). For videos, VideoSeal introduces the concept of *temporal watermark propagation*: The input video is segmented into k frames, and a watermark is generated for the first frame, which is then propagated (copied or interpolated) to all subsequent frames. In practice, we find $k = 4$ and a simple watermark copy provides a good trade-off between quality and robustness. Additionally, the frame is downsampled to 256×256 and the watermark is later upsampled to match the input resolution, this helps increase the training efficiency.

There are some hyperparameters that we need to choose when fine-tuning the watermark models. The first is the scaling factor λ , which decides the strength to add watermarks to the original contents. Too high λ values make artifacts in the watermarked contents more perceptible and visible, while lower λ values result in detectors that are less robust to watermarking attacks. In our experiments, we find $\lambda = 0.25$ to produce the optimal trade-off. The other parameter is the number of bits n in the watermark, for which we empirically evaluate per modality and set $n = 16$ for audio and $n = 128$ for video. To improve the robustness, after the contents are watermarked, we can continue to fine-tune the detector further on extended attacks and content distortion. Our experiments show that this improves the detection accuracy by up to +37% for popular audio compression attacks such as MP3 and AAC. For videos, fine-tuning on video-specific compression formats such as H.264 and H.265 also improves the performance by more than 20%.

8. Related Work

Existing conversational or audiovisual data. There are a multiple spoken conversation corpora, such as Fisher (Cieri et al., 2004) and Switchboard (Godfrey et al., 1992). They have been widely used in speech research and dialogue modeling. Besides speech and text modalities, some corpora also provide video capture of human interactions. These audiovisual data resources include the AMI Meeting Corpus (Carletta et al., 2005), IEMOCAP (Busso et al., 2008), and the CANDOR corpus (Reece et al., 2023). Finally, the SEAMLESS INTERACTION dataset, introduced in this work, is the largest known video collection of in-person, dyadic, conversation-based interactions.

Human motion model. Human-centric generative modeling focuses on generating human-like and expressive motions with autoregressive model (Ribeiro-Gomes et al., 2024) or diffusion (Peebles and Xie, 2023; Tevet et al., 2023; Chen et al., 2024). One key task of human motion generation is audio-driven generation. Different from the strong semantic alignment between input condition and motion, audio-driven task emphasizes the synchrony between speech and motion as well as the motion appropriateness (Kucherenko et al., 2023a). AniPortrait (Wei et al., 2024), HALLO (Xu et al., 2024a) and VASA-1 (Xu et al., 2024b) train diffusion models for talking face generation, capturing facial expressions and lip movements. Geneface (Ye et al., 2023) leverage the efficient motion-to-video renderer to achieves fast training and even real-time inference. INFP (Zhu et al., 2025) proposed models for face generation in the dyadic setting. Besides face, recent work also studies full body generation based on speech. GENE organizes a shared task on hand and body gesture generation in both monadic and dyadic setting (Kucherenko et al., 2023b). Ng et al. (2024) makes use of both autoregressive and diffusion model to learn body and face motion respectively, generating photorealistic embodiment in conversational settings. ConvoFusion (Mughal et al., 2024) develops diffusion for conversational co-speech gesture synthesis. There is also a trend to generate face expression and body gestures together (Chen et al., 2024; Yi et al., 2023). TalkSHOW (Yi et al., 2023) proposes to generate the expression and gesture separately, where the gesture model is generative and the face model is deterministic. DiffSHEG (Chen et al., 2024) is the first to propose a unified framework for modeling the joint distribution of expressions and gestures, which introduces a diffusion-based architecture with a uni-directional conditioning flow from expression to gesture, ensuring coherence and harmony between the two modalities.

Controllability in diffusion. Controllability is a crucial aspect of diffusion models, as it enables users to manipulate the generation according to their preferences. Face diffusion models such as VASA-1 (Xu et al.,

2024b) takes control signals such as eye gaze, head distance and emotion offsets. EMO (Tian et al., 2024a) involves a Face Locator and Speed Layer to weakly control the approximate region of the target face and the rough velocity level of the movement. Besides, EmojiDiff (Jiang et al., 2024) aims to achieve expression control with highly maintained identity. Body diffusion models such as MDM (Tevet et al., 2023) and MotionDiffuse (Zhang et al., 2024a) focused on text control signals, which generates motions based on textual descriptions. C2G2 (Ji et al., 2023) enables users to generate and edit the gestures in any time intervals, making sure the generation process is controllable. ConvoFusion (Mughal et al., 2024) also takes speaker style and word-excitation to guide gesture synthesis. On the other hand, semantic gesture control in co-speech gesture generation is challenging due to its long-tailed distribution in the dataset. (Zhang et al., 2024b) proposes to predict the semantic gestures based on the text transcript, and then merge the retrieved semantic gestures with the GPT-generated co-speech gestures.

LLM agent. With rapid advances in LLM (OpenAI, 2023; Team et al., 2025; Meta, 2025), there is a surge of interest in LLM-powered agent systems where LLM functions as the agent brain with general knowledge. LLM agent is an intelligent system to deal with tasks and interact with users. Post-training has been widely studied to extend the capabilities of pretrained LLMs in specific domains and downstream tasks. Post-training paradigms include supervised fine-tuning (Ouyang et al., 2022), prefix tuning (Li and Liang, 2021), prompt tuning (Lester et al., 2021), instruction tuning (Shengyu et al., 2023), and reinforcement learning (Trung et al., 2024). In particular, the development of parameter-efficient fine-tuning (PEFT) has been spurred, given its computation efficiency compared with full-model tuning. LoRA (Hu et al., 2022) and adapters (Hu et al., 2023) are cost-effective PEFT approaches which could also mitigate catastrophic forgetting.

2D Rendering of Photorealistic Avatars. Video generation has evolved significantly over the years. Early text-to-video models based on GANs—such as MoCoGAN-HD (Tulyakov et al., 2018) and StyleGAN-V (Skorokhodov et al., 2022)—as well as VAE variants like VideoVAE (He et al., 2018) and CV-VAE (Zhao et al., 2024), frequently struggled with mode collapse, temporal flickering, and low spatial fidelity. Recently, diffusion-based denoising models (Ho et al., 2020; Song et al., 2021a; Nichol and Dhariwal, 2021; Song et al., 2021b; Dhariwal and Nichol, 2021) have supplanted these approaches thanks to their stable, scalable score-matching objectives that extend to datasets containing hundreds of millions of clips. These models gradually refine a noisy sequence—first establishing global structure, then restoring fine-grained details—thus naturally promoting temporal consistency. Moreover, classifier-free guidance enables a single backbone to accommodate diverse conditioning signals such as text, images, depth maps, or pose sequences without architectural changes. Together, these advances yield superior visual sharpness, smoother motion, and more robust identity preservation compared to previous paradigms.

3D Rendering of Photorealistic Avatars. Photorealistic avatars have evolved from pure academic research to matured technology that found its way into products, such as Apple’s Personas or Epic Games’ Metahumans. In this work, we build upon Codec Avatars, which have been seminal as a photorealistic representation of humans. Originally, Codec Avatars were introduced as face-only avatars using deep appearance models (Lombardi et al., 2018) that represent expression changes on textured face meshes through a view-conditioned variational autoencoder. In Bagautdinov et al. (2021), a similar concept was applied to full-body avatars. Breakthroughs in neural rendering such as neural volumes (Lombardi et al., 2019) and mixture of volumetric primitives (Lombardi et al., 2021) led to significant improvements in the visual quality of Codec Avatars, overcoming major limitations of textured meshes in terms of modeling hair. The latest generation of Codec Avatars is based on Gaussian splatting (Kerbl et al., 2023) for both face (Saito et al., 2024) and full-body avatars (Wang et al., 2025). While 3D representations of photorealistic humans are starting to escape the uncanny valley, driving these representations from sensory inputs remains challenging. Existing approaches focus on camera-driven avatars (Wei et al., 2019; Bai et al., 2024); however, such visual observations are not available when attempting to drive an avatar with speech alone or from the output of an LLM. For audio-driven face avatars, early works focus on un-textured geometry (Richard et al., 2021b; Cudeiro et al., 2019) or person-specific texture-based models (Richard et al., 2021a). More recently, audio-driven face-only Gaussian avatars started showing their promise (Aneja et al., 2024). Beyond speech-driven avatars, Chatziagapi et al. (2025) demonstrate joint generation of speech and photorealistic 3D face motion directly from LLM-generated text. Most closely related to our work, Ng et al. (2024) drive full-body avatars in dyadic conversations

using audio from both speakers as input, however, their approach is limited to a few person-specific avatars and is trained on less than two hours of data per person, putting significant limitations on diversity and comprehensiveness of the synthesized body gestures.

9. Conclusion

This paper introduced the SEAMLESS INTERACTION dataset, a large-scale collection of over 4,000 hours of face-to-face interaction footage from over 4,000 participants in diverse contexts. We also presented a family of research models, Dyadic Motion Models, which can not only generate motion gestures and facial expressions that align with human speech, but also take into consideration the visual behaviors of the interlocutor. We presented a variant with speech from LLM model and integrations with 2D and 3D rendering methods, bringing us closer to interactive virtual agents. Finally, we described controllable variants of our motion models that can adapt emotional responses and expressivity levels, as well as generating gestures that are more semantically relevant. These dyadic motion models are demonstrating the potential for more intuitive and responsive human-AI interactions.

Acknowledgements

We want to extend our gratitude to those who made this work possible below. To Arianne Burrell, Ben Samples, Josh Terry, Kenny Lehmann, Julia Vargas, Alyssa Newcomb, Michael Robert Brown, Shun Shiga, IV Tench, Karla Martucci, Michelle Restrepo, Nathan Hass, Junho Kim, Paula Chowles, Kate Bourdeau, Allie Castro, and Britt Montalvo for their tireless efforts in promoting our work. To Nisha Deo and Ashley Gabriel for their expertise in internal and external communications. To Idan Afek for strategic guidance and support in forging new partnerships. To Ernest Hammond and Corey Wallace for providing invaluable counsel and ensuring that our work is compliant with all relevant regulations. To Rachel Kim and Ty Toledano for their dedication to protecting user data and ensuring that our work meets high standards of privacy and security. To Maeve Ryan for navigating complex policy landscapes and advocating for our users’ interests. To Yael Yungster and Kei Koyama for bringing our vision to life through their creative and innovative designs. To Lindsey Miller for insights into user behavior and preferences, which have informed our design decisions and improved the overall user experience. To Peng-Jen Chen, Min-Jae Hwang, Bokai Yu, Oleg Repin, and Alex Mourachko for engaging in stimulating discussions and sharing their expertise. To Rob Fergus, Joelle Pineau, and Stephane Kasriel for their leadership, guidance, and unwavering support throughout this project.

References

- S. Aneja, A. Sevastopolsky, T. Kirschstein, J. Thies, A. Dai, and M. Nießner. Gaussianspeech: Audio-driven gaussian avatars, 2024. URL <https://arxiv.org/abs/2411.18675>.
- T. Bagautdinov, C. Wu, T. Simon, F. Prada, T. Shiratori, S.-E. Wei, W. Xu, Y. Sheikh, and J. Saragih. Driving-signal aware full-body avatars. *ACM Transactions on Graphics (TOG)*, 40(4):1–17, 2021.
- A. Bagher Zadeh, P. P. Liang, S. Poria, E. Cambria, and L.-P. Morency. Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph. In I. Gurevych and Y. Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1208. URL <https://aclanthology.org/P18-1208/>.
- S. Bai, T.-L. Wang, C. Li, A. Venkatesh, T. Simon, C. Cao, G. Schwartz, J. Saragih, Y. Sheikh, and S.-E. Wei. Universal facial encoding of codec avatars from vr headsets. *ACM Trans. Graph.*, 43(4), July 2024. ISSN 0730-0301. doi: 10.1145/3658234. URL <https://doi.org/10.1145/3658234>.
- M. Bain, J. Huh, T. Han, and A. Zisserman. Whisperx: Time-accurate speech transcription of long-form audio. In *Interspeech 2023*, pages 4489–4493, 2023. doi: 10.21437/Interspeech.2023-78.
- D. Biber, S. Johansson, G. Leech, S. Conrad, and E. Finegan. *Longman Grammar of Spoken and Written English*. Longman, 1999.

- C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan. Iemocap: Interactive emotional dyadic motion capture database. Technical report, University of Southern California, 2008.
- Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal, et al. The ami meeting corpus: A pre-announcement. In *Machine Learning for Multimodal Interaction*, pages 28–39. Springer, 2005.
- C. Chan, S. Ginosar, T. Zhou, and A. A. Efros. Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- A. Chatziagapi, L.-P. Morency, H. Gong, M. Zollhöfer, D. Samaras, and A. Richard. Av-flow: Transforming text to audio-visual human-like interactions. *arXiv preprint arXiv:2502.13133*, 2025.
- J. Chen, Y. Liu, J. Wang, A. Zeng, Y. Li, and Q. Chen. Diffshg: A diffusion-based approach for real-time speech-driven holistic 3d expression and gesture generation. In *CVPR*, 2024.
- C. Cieri, D. Miller, and K. Walker. The fisher corpus: a resource for the next generations of speech-to-text. Technical report, Linguistic Data Consortium, Philadelphia, 2004.
- D. Cudeiro, T. Bolkart, C. Laidlaw, A. Ranjan, and M. Black. Capture, learning, and synthesis of 3D speaking styles. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10101–10111, 2019.
- P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *CoRR*, abs/2105.05233, 2021. URL <https://arxiv.org/abs/2105.05233>.
- D. Dowson and B. Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982.
- T. J. Dunn, T. Baguley, and V. Brunsden. From alpha to omega: A practical solution to the pervasive problem of internal consistency estimation. *British Journal of Psychology*, 105(3):399–412, 2014. ISSN 2044-8295. doi: 10.1111/bjop.12046.
- P. Ekman and E. L. Rosenberg. *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA, 1997.
- G. Fergadiotis, H. H. Wright, and G. J. Capilouto. Productive vocabulary across discourse types. *Aphasiology*, 25(10): 1261–1278, 2011.
- P. Fernandez, H. Elsahar, I. Z. Yalniz, and A. Mourachko. Video seal: Open and efficient video watermarking. *arXiv preprint arXiv:2412.09492*, 2024. URL <https://arxiv.org/abs/2412.09492>.
- R. Flesch. *How to write plain English: A book of lawyers and consumers*. Harper & Row, New York, NY, 1979.
- H. Giles, N. Coupland, and J. Coupland. Accommodation theory: Communication, context, and consequence. In H. Giles, J. Coupland, and N. Coupland, editors, *Contexts of Accommodation: Developments in Applied Sociolinguistics*, pages 1–68. Cambridge University Press, 1991.
- J. J. Godfrey, E. C. Holliman, and J. McDaniel. Switchboard: Telephone speech corpus for research and development. Technical report, Linguistic Data Consortium, Philadelphia, 1992.
- S. Goel, G. Pavlakos, J. Rajasegaran, A. Kanazawa, and J. Malik. Humans in 4d: Reconstructing and tracking humans with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14783–14794, 2023.
- Google. Expanding vertex ai with the next wave of generative ai media models, 2025. URL <https://cloud.google.com/blog/products/ai-machine-learning/announcing-veo-3-imagen-4-and-lyria-2-on-vertex-ai>.
- J. He, A. Lehrmann, J. Marino, G. Mori, and L. Sigal. Probabilistic video generation using holistic attribute control. In *European Conference on Computer Vision (ECCV)*, pages 466–483, 2018.
- M. Heldner and J. Edlund. Pauses, gaps and overlaps in conversations. *Journal of Phonetics*, 38(4):555–568, 2010. ISSN 0095-4470. doi: <https://doi.org/10.1016/j.wocn.2010.08.002>. URL <https://www.sciencedirect.com/science/article/pii/S0095447010000628>.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

- HeyGen. Create lifelike ai video avatars, 2025. URL <https://www.heygen.com/avatars/ai-video-avatar>.
- J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- C. J. Hopwood, A. L. Harrison, M. C. Amole, J. M. Girard, A. G. C. Wright, K. M. Thomas, P. Sadler, E. B. Ansell, T. M. Chaplin, L. C. Morey, M. J. Crowley, C. E. Durbin, and D. A. Kashy. Properties of the continuous assessment of interpersonal dynamics across sex, level of familiarity, and interpersonal conflict. *Assessment*, 27(1):40–56, 2020. doi: 10.1177/1073191118798916.
- G. Horstmann and L. Linke. Perception of direct gaze in a video-conference setting: the effects of position and size. *Cognitive Research: Principles and Implications*, 7:67, 2022. doi: 10.1186/s41235-022-00418-1. URL <https://doi.org/10.1186/s41235-022-00418-1>.
- E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- Z. Hu, L. Wang, Y. Lan, W. Xu, E.-P. Lim, L. Bing, X. Xu, S. Poria, and R. Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5254–5276, 2023.
- Y. Huang, L. Ho, D. Qin, M. Shi, and T. Komura. Interact: Capture and modelling of realistic, expressive and interactive activities between two persons in daily scenarios, 2024. URL <https://arxiv.org/abs/2405.11690>.
- L. Ji, P. Wei, Y. Ren, J. Liu, C. Zhang, and X. Yin. C2g2: Controllable co-speech gesture generation with latent diffusion model. *arXiv preprint arXiv:2308.15016*, 2023.
- L. Jiang, R. Li, Z. Zhang, S. Fang, and C. Ma. Emojidiff: Advanced facial expression control with high identity preservation in portrait generation. *arXiv preprint arXiv:2412.01254*, 2024.
- O. E. John, L. P. Naumann, and C. J. Soto. Paradigm shift to the integrative Big Five trait taxonomy: History, measurement, and conceptual issues. In O. P. John, R. W. Robins, and L. A. Pervin, editors, *Handbook of Personality: Theory and Research*, pages 114–158. Guilford Press, 3rd edition, 2008.
- O. P. John and S. Srivastava. The big five trait taxonomy: History, measurement, and theoretical perspectives. In L. A. Pervin and O. P. John, editors, *Handbook of Personality: Theory and Research*, pages 102–138. Guilford, 1999.
- B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023. URL <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>.
- R. Khirodkar, T. Bagautdinov, J. Martinez, S. Zhaoen, A. James, P. Selednik, S. Anderson, and S. Saito. Sapiens: Foundation for human vision models. *arXiv preprint arXiv:2408.12569*, 2024.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *ICLR (Poster)*, 2015. URL <http://dblp.uni-trier.de/db/conf/iclr/iclr2015.html#KingmaB14>.
- M. L. Knapp, A. L. Vangelisti, and J. P. Caughlin. *Interpersonal Communication & Human Relationships*. Pearson, 7th edition, 2013.
- T. Kucherenko, R. Nagy, Y. Yoon, J. Woo, T. Nikolov, M. Tsakov, and G. E. Henter. The genea challenge 2023: A large-scale evaluation of gesture generation models in monadic and dyadic settings. In *Proceedings of the 25th International Conference on Multimodal Interaction, ICMI '23*, page 792–801. Association for Computing Machinery, 2023a. ISBN 9798400700552.
- T. Kucherenko, R. Nagy, Y. Yoon, J. Woo, T. Nikolov, M. Tsakov, and G. E. Henter. The GENE challenge 2023: A large-scale evaluation of gesture generation models in monadic and dyadic settings. In E. André, M. Chetouani, D. Vaufraydaz, G. M. Lucas, T. Schultz, L. Morency, and A. Vinciarelli, editors, *Proceedings of the 25th International Conference on Multimodal Interaction, ICMI 2023, Paris, France, October 9-13, 2023*, pages 792–801. ACM, 2023b.
- M. R. Leary and R. H. Hoyle, editors. *Handbook of Individual Differences in Social Behavior*. Guilford Press, 2009.
- G. Lee, Z. Deng, S. Ma, T. Shiratori, S. S. Srinivasa, and Y. Sheikh. Talking with hands 16.2m: A large-scale dataset of synchronized body-finger motion and audio for conversational motion analysis and synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 763–772, 2019.

- V. Lepetit, F. Moreno-Noguer, and P. Fua. Epnnp: An accurate $O(n)$ solution to the pnp problem. *Int. J. Comput. Vis.*, 81(2):155–166, 2009.
- B. Lester, R. Al-Rfou, and N. Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, 2021.
- X. L. Li and P. Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, 2021.
- G. Lin, J. Jiang, J. Yang, Z. Zheng, and C. Liang. Omnihuman-1: Rethinking the scaling-up of one-stage conditioned human animation models. *CoRR*, abs/2502.01061, 2025.
- Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- H. Liu, Z. Zhu, N. Iwamoto, Y. Peng, Z. Li, Y. Zhou, E. Bozkurt, and B. Zheng. Beat: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis, 2022a. URL <https://arxiv.org/abs/2203.05297>.
- X. Liu, Q. Wu, H. Zhou, Y. Xu, R. Qian, X. Lin, X. Zhou, W. Wu, B. Dai, and B. Zhou. Learning hierarchical cross-modal association for co-speech gesture generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10462–10472, 2022b.
- Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of ICCV*, 2021.
- S. Lombardi, J. Saragih, T. Simon, and Y. Sheikh. Deep appearance models for face rendering. *ACM Transactions on Graphics (ToG)*, 37(4):1–13, 2018.
- S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019.
- S. Lombardi, T. Simon, G. Schwartz, M. Zollhoefer, Y. Sheikh, and J. Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021.
- M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, 2015. doi: 10.1145/2816795.2818013. URL <https://smpl.is.tue.mpg.de/>.
- M. Mallen, S. Day, and M. Green. Online versus face-to-face conversation: An examination of relational and discourse variables. *Psychotherapy: Theory, Research, Practice, Training*, 40:155–163, 04 2003. doi: 10.1037/0033-3204.40.1-2.155.
- J. Martinez, E. Kim, J. Romero, T. Bagautdinov, S. Saito, S.-I. Yu, S. Anderson, M. Zollhöfer, T.-L. Wang, S. Bai, C. Li, S.-E. Wei, R. Joshi, W. Borsos, T. Simon, J. Saragih, P. Theodosis, A. Greene, A. Josyula, S. M. Maeta, A. I. Jewett, S. Venshtain, C. Heilman, Y.-T. Chen, S. Fu, M. E. A. Elshaer, T. Du, L. Wu, S.-C. Chen, K. Kang, M. Wu, Y. Emad, S. Longay, A. Brewer, H. Shah, J. Booth, T. Koska, K. Haidle, M. Andromalos, J. Hsu, T. Dauer, P. Selednik, T. Godisart, S. Ardisson, M. Cipperly, B. Humberston, L. Farr, B. Hansen, P. Guo, D. Braun, S. Krenn, H. Wen, L. Evans, N. Fadeeva, M. Stewart, G. Schwartz, D. Gupta, G. Moon, K. Guo, Y. Dong, Y. Xu, T. Shiratori, F. Prada, B. R. Pires, B. Peng, J. Buffalini, A. Trimble, K. McPhail, M. Schoeller, and Y. Sheikh. Codec Avatar Studio: Paired Human Captures for Complete, Driveable, and Generalizable Avatars. *NeurIPS Track on Datasets and Benchmarks*, 2024.
- P. M. McCarthy. *An assessment of the range and usefulness of lexical diversity measures and the potential of the measure of textual, lexical diversity (MTLD)*. PhD thesis, The University of Memphis, 2005.
- P. M. McCarthy and S. Jarvis. Mtl-d, vocd-d, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior research methods*, 42(2):381–392, 2010.
- R. R. McCrae and P. T. Costa. A five-factor theory of personality. In L. A. Pervin and O. John, editors, *Handbook of Personality: Theory and Research*, pages 139–153. Guilford Press, New York, NY, 2nd edition, 1999.
- A. Meta. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, checked on, 4(7):2025, 2025.

- M. H. Mughal, R. Dabral, I. Habibie, L. Donatelli, M. Habermann, and C. Theobalt. Convofusion: Multi-modal conversational diffusion for co-speech gesture synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1388–1398, 2024.
- R. Nagy, H. Voss, Y. Yoon, T. Kucherenko, T. Nikolov, T. Hoang-Minh, R. McDonnell, S. Kopp, M. Neff, and G. E. Henter. Towards a genea leaderboard – an extended, living benchmark for evaluating and advancing conversational motion synthesis, 2024. URL <https://arxiv.org/abs/2410.06327>.
- E. Ng, J. Romero, T. M. Bagautdinov, S. Bai, T. Darrell, A. Kanazawa, and A. Richard. From audio to photoreal embodiment: Synthesizing humans in conversations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 1001–1010. IEEE, 2024.
- A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *Proc. 38th International Conference on Machine Learning (ICML 2021)*, pages 8162–8171, 2021. URL <https://proceedings.mlr.press/v139/nichol21a.html>.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- D. J. Ozer and V. Benet-Martínez. Personality and the Prediction of Consequential Outcomes. *Annual Review of Psychology*, 57(1):401–421, 2006. doi: 10.1146/annurev.psych.57.102904.190127.
- G. Pavlakos, D. Shan, I. Radosavovic, A. Kanazawa, D. Fouhey, and J. Malik. Reconstructing hands in 3d with transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 9826–9836. IEEE, 2024.
- W. Peebles and S. Xie. Scalable diffusion models with transformers. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4172–4182, 2023. doi: 10.1109/ICCV51070.2023.00387.
- I. Petrov, H. Guo, Y. Qian, and X. Gu. Diffpose: Pose-guided human video generation with denoising diffusion. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- A. L. Pincus and E. B. Ansell. Interpersonal theory of personality. In I. B. Weiner, H. A. Tennen, and J. M. Suls, editors, *Handbook of Psychology, Volume 5, Personality and Social Psychology*, pages 141–159. Wiley, Hoboken, NJ, 2nd edition, 2013.
- S. Poria, D. Hazarika, N. Majumder, G. Naik, R. Mihalcea, and E. Cambria. Meld: A multimodal multi-party dataset for emotion recognition in conversation. 2018.
- K. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, and C. Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM international conference on multimedia*, pages 484–492, 2020.
- A. Reece, G. Cooney, P. Bull, C. Chung, B. Dawson, C. Fitzpatrick, T. Glazer, D. Knox, A. Liebscher, and S. Marin. The candor corpus: Insights from a large multimodal dataset of naturalistic conversation. *Science Advances*, 9(13): eadf3197, 2023. doi: 10.1126/sciadv.adf3197. URL <https://www.science.org/doi/abs/10.1126/sciadv.adf3197>.
- J. Ribeiro-Gomes, T. Cai, Z. Á. Milacski, C. Wu, A. Prakash, S. Takagi, A. Aubel, D. Kim, A. Bernardino, and F. D. la Torre. Motiongpt: Human motion synthesis with improved diversity and realism via GPT-3 prompting. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2024, Waikoloa, HI, USA, January 3-8, 2024*, pages 5058–5068. IEEE, 2024.
- A. Richard, C. Lea, S. Ma, J. Gall, F. de la Torre, and Y. Sheikh. Audio- and gaze-driven facial animation of codec avatars. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 41–50, 2021a.
- A. Richard, M. Zollhoefer, Y. Wen, F. de la Torre, and Y. Sheikh. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021b.
- B. W. Roberts, N. R. Kuncel, R. Shiner, A. Caspi, and L. R. Goldberg. The Power of Personality: The Comparative Validity of Personality Traits, Socioeconomic Status, and Cognitive Ability for Predicting Important Life Outcomes. *Perspectives on Psychological Science*, 2(4):313–345, 2007. ISSN 17456916. doi: 10.1111/j.1745-6916.2007.00047.x.

- J. Romero, D. Tzionas, and M. J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017.
- J. A. Russell. A circumplex model of affect. *Journal of personality and social psychology*, 39(6):1161, 1980.
- S. Saito, G. Schwartz, T. Simon, J. Li, and G. Nam. Relightable gaussian codec avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2024.
- R. San Roman, P. Fernandez, H. Elshahar, A. Défossez, T. Furon, and T. Tran. Proactive detection of voice cloning with localized watermarking. In *Proceedings of the 41st International Conference on Machine Learning*, pages 43180–43196, 2024.
- J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Z. Shengyu, D. Linfeng, L. Xiaoya, Z. Sen, S. Xiaofei, W. Shuhe, L. Jiwei, R. Hu, Z. Tianwei, F. Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.
- A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- A. Siarohin, O. J. Woodford, J. Ren, M. Chai, and S. Tulyakov. Motion representations for articulated animation, 2021. URL <https://arxiv.org/abs/2104.11280>.
- Silero. Silero vad: pre-trained enterprise-grade voice activity detector (vad), number detector and language classifier. <https://github.com/snakers4/silero-vad>, 2021.
- I. Skorokhodov, S. Tulyakov, and M. Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. *CoRR*, abs/2112.14683, 2022. arXiv: 2112.14683.
- J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021a.
- Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- C. J. Soto. How replicable are links between personality traits and consequential life outcomes? The life outcomes of personality replication project. *Psychological Science*, 30(5):711–727, 2019. doi: 10.1177/0956797619831612.
- C. J. Soto and O. P. John. The next Big Five Inventory (BFI-2): Developing and assessing a hierarchical model with 15 facets to enhance bandwidth, fidelity, and predictive power. *Journal of Personality and Social Psychology*, 113(1): 117–143, 2017. ISSN 1939-1315, 0022-3514. doi: 10.1037/pspp0000096.
- S. Srinivasan and H. Do. Improving Audio Quality for Calling across Meta’s Family of Apps, 2024. URL <https://atscaleconference.com/improving-audio-quality-for-calling-across-metas-family-of-apps/>.
- L. Team, A. Modi, A. S. Veerubhotla, A. Rysbek, A. Huber, A. Anand, A. Bhoopchand, B. Wiltshire, D. Gillick, D. Kasenberg, et al. Evaluating gemini in an arena for learning. *arXiv preprint arXiv:2505.24477*, 2025.
- G. Tevet, S. Raab, B. Gordon, Y. Shafir, D. Cohen-Or, and A. H. Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- L. Tian, Q. Wang, B. Zhang, and L. Bo. Emo: Emote portrait alive generating expressive portrait videos with audio2video diffusion model under weak conditions. In *European Conference on Computer Vision*, pages 244–260. Springer, 2024a.
- Y. Tian, S. Liu, and J. Wang. A corpus study on the difference of turn-taking in online audio, online video, and face-to-face conversation. *Language and Speech*, 67(3):593–616, 2024b. doi: 10.1177/00238309231176768. URL <https://doi.org/10.1177/00238309231176768>. PMID: 37317824.
- E. K. Traupman, T. W. Smith, B. N. Uchino, C. A. Berg, K. K. Trobst, and P. T. Costa. Interpersonal circumplex octant, control, and affiliation scales for the NEO-PI-R. *Personality and Individual Differences*, 47(5):457–463, 2009. doi: 10.1016/j.paid.2009.04.018.
- L. Trung, X. Zhang, Z. Jie, P. Sun, X. Jin, and H. Li. Reft: Reasoning with reinforced fine-tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7601–7614, 2024.

- S. Tulyakov, M. Liu, X. Yang, and J. Kautz. MoCoGAN: Decomposing motion and content for video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1526–1535, 2018.
- R. Villegas, J. Song, S. Nah, J. Yang, J. Shi, M. Cakmak, and J. Wu. Skeleton-aware networks for deep motion retargeting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- L. Wagner, B. Thallinger, and M. Zusag. Crisperwhisper: Accurate timestamps on verbatim speech transcriptions, 2024. URL <https://arxiv.org/abs/2408.16589>.
- S. Wang, T. Simon, I. Santesteban, T. Bagautdinov, J. Li, V. Agrawal, F. Prada, S.-I. Yu, P. Nalbony, M. Gramlich, R. Lubachersky, C. Wu, J. Romero, J. Saragih, M. Zollhoefer, A. Geiger, S. Tang, and S. Saito. Relightable full-body gaussian codec avatars. *arXiv.org*, 2501.14726, 2025.
- H. Wei, Z. Yang, and Z. Wang. Aniportrait: Audio-driven synthesis of photorealistic portrait animation. *CoRR*, abs/2403.17694, 2024.
- S.-E. Wei, J. Saragih, T. Simon, A. W. Harley, S. Lombardi, M. Perdoch, A. Hypes, D. Wang, H. Badino, and Y. Sheikh. Vr facial animation via multiview image translation. *ACM Transactions on Graphics (ToG)*, 38(4):1–16, 2019.
- A. G. C. Wright, A. L. Pincus, and C. J. Hopwood. Contemporary integrative interpersonal theory: Integrating structure, dynamics, temporal scale, and levels of analysis. *Journal of Psychopathology and Clinical Science*, 132(3): 263–276, 2023. doi: 10/g824s3.
- M. Xu, H. Li, Q. Su, H. Shang, L. Zhang, C. Liu, J. Wang, Y. Yao, and S. Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation. *CoRR*, abs/2406.08801, 2024a.
- S. Xu, G. Chen, Y. Guo, J. Yang, C. Li, Z. Zang, Y. Zhang, X. Tong, and B. Guo. VASA-1: lifelike audio-driven talking faces generated in real time. *CoRR*, abs/2404.10667, 2024b. doi: 10.48550/ARXIV.2404.10667.
- Y. Xu, J. Zhang, Q. Zhang, and D. Tao. Vitpose: Simple vision transformer baselines for human pose estimation. *Advances in neural information processing systems*, 35:38571–38584, 2022.
- J. Yang, J. Liao, R. Ma, and W. Wang. Skeleton-aware style transfer for interactive character animation. In *Proceedings of ACM SIGGRAPH*, 2020.
- Z. Ye, J. He, Z. Jiang, R. Huang, J. Huang, J. Liu, Y. Ren, X. Yin, Z. Ma, and Z. Zhao. Geneface++: Generalized and stable real-time audio-driven 3d talking face generation. *arXiv preprint arXiv:2305.00787*, 2023.
- H. Yi, H. Liang, Y. Liu, Q. Cao, Y. Wen, T. Bolkart, D. Tao, and M. J. Black. Generating holistic 3d human motion from speech. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 469–480, June 2023.
- M. S. M. Yik and J. A. Russell. On the relationship between circumplexes: Affect and Wiggins’ IAS. *Multivariate Behavioral Research*, 39(2):203–230, 2004. doi: 10.1207/s15327906mbr3902_4.
- Y. Yoon, P. Wolfert, T. Kucherenko, C. Viegas, T. Nikolov, M. Tsakov, and G. E. Henter. The genea challenge 2022: A large evaluation of data-driven co-speech gesture generation. In *Proceedings of the 2022 International Conference on Multimodal Interaction, ICMI ’22*, page 736–747, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393904. doi: 10.1145/3536221.3558058. URL <https://doi.org/10.1145/3536221.3558058>.
- B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- M. Zhang, Z. Cai, L. Pan, F. Hong, X. Guo, L. Yang, and Z. Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(6):4115–4128, 2024a.
- S. Zhang, W. Zhang, and Q. Gu. Energy-weighted flow matching for offline reinforcement learning. *arXiv preprint arXiv:2503.04975*, 2025.
- Z. Zhang, T. Ao, Y. Zhang, Q. Gao, C. Lin, B. Chen, and L. Liu. Semantic gesticulator: Semantics-aware co-speech gesture synthesis. *ACM Transactions on Graphics (TOG)*, 43(4):1–17, 2024b.
- S. Zhao, Y. Zhang, X. Cun, S. Yang, M. Niu, X. Li, W. Hu, and Y. Shan. CV-VAE: A compatible video vae for latent generative video models. *CoRR*, abs/2405.20279, 2024. arXiv: 2405.20279.
- Y. Zhou, C. Barnes, L. Jingwan, Y. Jimei, and L. Hao. On the continuity of rotation representations in neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

- Y. Zhu, L. Zhang, Z. Rong, T. Hu, S. Liang, and Z. Ge. Infp: Audio-driven interactive head generation in dyadic conversations, 2024. URL <https://arxiv.org/abs/2412.04037>.
- Y. Zhu, L. Zhang, Z. Rong, T. Hu, S. Liang, and Z. Ge. Infp: Audio-driven interactive head generation in dyadic conversations. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 10667–10677, 2025.
- L. Zhuo, R. Du, H. Xiao, Y. Li, D. Liu, R. Huang, W. Liu, L. Zhao, F.-Y. Wang, Z. Ma, et al. Lumina-next: Making lumina-t2x stronger and faster with next-dit. *arXiv preprint arXiv:2406.18583*, 2024.

A. Seamless Interaction

A.1 Properties

A.1.1 Corpus Text Analysis

Given that the recorded interactions that comprise SEAMLESS INTERACTION were organized at various locations and did not occur spontaneously, one must consider as a potential risk that the language content may not be representative of naturally occurring conversations. To assess this risk, we analyzed the language of SEAMLESS INTERACTION along grammatical, lexical, and structural lines.

A.1.2 Readability, Lexical Diversity, Known Features of Conversation

Readability. Although readability scores apply primarily by definition to written material, some tests (such as Flesch’s reading ease) also indicate expected score ranges for conversational English (Flesch, 1979).

To determine whether the Flesch reading-ease score (FRES) of SEAMLESS INTERACTION language in general is within the expected range, and determine whether the use of actors affects FRES, 3 corpora are first gathered from SEAMLESS INTERACTION transcripts, representing language encountered in interactions between non-actor dyads, dual-actor dyads, and single-actor dyads. Next, 2 other corpora are gathered for comparison, representing miscellaneous literature language and US presidential inaugural addresses, respectively.

The literature corpus is expected to return an average FRES of 60 or below, while spontaneous, naturally occurring speech transcripts are expected to score at least above 82 and even closer to 92. In the absence of prior reference for inaugural addresses, we can only observe that the prevalence of unusually long sentences will cause a fairly low score to be returned.

To compute mean and median FRES, 30 passages of approximately 2,200 words in length are sampled. We use the `textstat`² python package to compute FRES on each sample, then compute the mean and median scores for the sample sets.

We find that SEAMLESS INTERACTION corpora fall within the expected range for conversational English (82–92), as shown in Table 21.

Domain	Mean FRES	Median FRES	Expected FRES
Inaugural addresses	43.5	42.7	< 60
Literature	62.9	65.2	~60
SEAMLESS INTERACTION Non-actor	88.0	86.4	82–92
SEAMLESS INTERACTION Dual-actor	91.6	91.3	82–92
SEAMLESS INTERACTION Single-actor	87.1	86.1	82–92

Table 21 - Flesch reading ease scores (FRES) for SEAMLESS INTERACTION corpora.

Lexical diversity. To complement the previous readability analysis, lexical diversity computation methods such as the Measure of Textual Lexical Diversity (MTLD) can be used (McCarthy, 2005). The MTLD method focuses more specifically on lexical variability and is less sensitive to small variations in sample length (McCarthy and Jarvis, 2010). When different discourse types are elicited, the measuring of lexical variability returns differences in score ranges, as shown in Fergadiotis et al. (2011). Additionally, Biber et al. (1999) shows that in large representative corpora, such as LSWE, conversational English receives by far the lowest lexical variability score due to the high degree of repetition in function words and inserts (e.g., *yeah*, *right*, *um*).

MTLD scores are computed on the 3 SEAMLESS INTERACTION corpora, as well as the literature and inaugural address (i.e., scripted speech) corpora for control, using the `lexicalrichness`³ python package with a

²<https://pypi.org/project/textstat/>

³<https://pypi.org/project/lexicalrichness/>

type-token ratio (TTR) factor of 0.72.

We find that the SEAMLESS INTERACTION corpora all fall within the MTLD score range we would expect (at least 1.75 times lower than composed literary language). Scores for all 5 corpora are reported in [Table 22](#).

DOMAIN	MTLD score
Literature	85.9
Inaugural addresses	72.3
SEAMLESS INTERACTION Non-actor	48.1
SEAMLESS INTERACTION Dual-actor	39.8
SEAMLESS INTERACTION Single-actor	42.1

Table 22 - MTLD scores for SEAMLESS INTERACTION corpora and controls.

A.1.3 Collection Facility Geography

- Chino Hills, CA
- Costa Mesa, CA
- Irvine, CA
- Los Angeles, CA
- Boise, ID
- Pittsburgh, PA
- Las Vegas, NV
- Boston, MA
- Waltham, MA
- New York, NY

A.1.4 Known Limitations of Text Transcripts

We acknowledge several limitations of our current text transcript methodology.

WhisperX, like most speech recognition systems, is tuned to provide concise, readable transcripts that ignore hesitations, false-starts, and other disfluencies which have important visual correlates and are interesting phenomena in their own right. We experimented with other systems that aim to robustly transcribe such behaviors ([Wagner et al., 2024](#)) but found the overall quality to be lower.

We also find that the wav2vec alignment step of WhisperX can produce inaccurate timestamps on short utterances. Approximately 98% of sessions and 87% of individual interactions contain at least one word whose timestamp-derived length is greater than 3 standard deviations from the mean.

While the current time-aligned transcripts enable preliminary corpus analyses on text (see previous sections in this [Appendix A.1.1](#)) and turn-taking behavior in the spirit of [Heldner and Edlund \(2010\)](#), future work will focus on producing more robust transcripts and accurate timestamps for SEAMLESS INTERACTION that enable new contributions in these fields.

A.2 Quality and Safety

A.2.1 Participant Pre-amble

Moderators were instructed to provide this information at the start of every session recording:

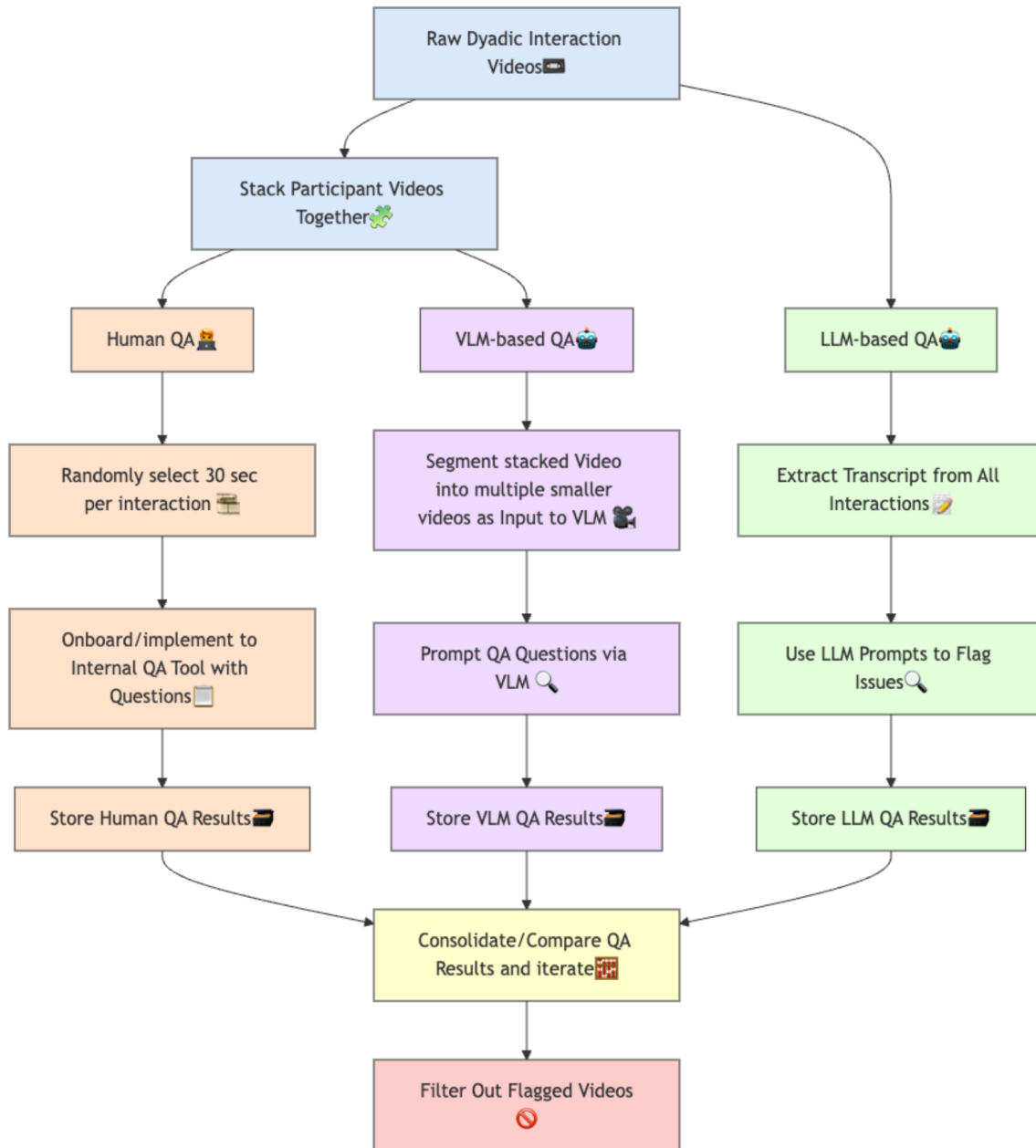


Figure 21 - Illustration of the scalable Safety and QA process applied to SEAMLESS INTERACTION. Using this process every minute of footage received a QA signal.

Welcome and thank you for participating in this project. The goal of this project is to study how humans use visual cues in communication. We will just ask you to have a series of short conversations based on prompts we have provided. The most important thing is for you to behave naturally. Try not to let the cameras or the setting intimidate you — just read the prompt and engage with it, try to have fun and try to express your true feelings.

Please stand where indicated. We would like you to stand for as long as you feel comfortable, but if you get tired you may use the stools. Please just do not touch or interact with the stools using your hands — we'd like your hands to remain free during the conversation. [Indicating to the scripts, face down] These are the scripts which contain the prompts you will discuss. After we are set-up and recording, I will ask you to start with prompt 1. Each time we move to a new prompt, you will pick up the scripts and read the prompt, silently to yourself. Once you have read it and understand it, place the scripts back down. Again, we just want you to not have anything in your hands during the interaction. Then [Initiator] will begin the prompt. [Initiator] will be the person that starts every conversation.

We expect each conversation to last from a few moments to 5 minutes. I may interrupt the conversation and ask you to move on to the next prompt - please don't read into this, I am just doing it in the interest of time and getting through all the prompts.

If you feel like you don't have anything more to say on a particular prompt, just let me know and we can move on to another one.

While you should feel comfortable to express yourself fully and how you see fit, we do ask that you refrain from discussing any personal details such as phone numbers, addresses, emails, et cetera, from you or anyone else. This is just to protect everyone's privacy.

A.2.2 Sensitive Data Statement Given to Vendors and Moderators

The following text was provided to all collection facilities and communicated to participants via moderators:

During the data collection process, it is important for participants to be mindful of their privacy and avoid disclosing personal or sensitive information. This includes:

- Cultural Background: Ethnicity
- Financial Identifiers, Financial Status (e.g., loans)
- Membership of a Trade Union
- Political affiliation, Political issues, Political opinion
- Health Data, Medical Condition
- Biometrics: Body specifications, Fingerprints, Iris scan, skin color
- Faith-based affiliation, Faith-based holiday, Faith-based spiritual & moral belief
- Sexual Orientation
- Personal & government identification number or image (e.g., picture of passport, passport number, driver's license number, SSN, Tax ID)
- Criminal Record
- Citizenship
- Precise Location/ Location Services
- Transgender, Intersex, Non-binary Gender

Avoiding sharing any information that could potentially identify them, such as their name, address, or other identifying details. By being mindful of what they share, participants can help protect their privacy and ensure that their personal information remains confidential. It is also important for the moderator to inform participants about avoiding these topics if it's accidentally disclosed during an interaction. It is also important

for the vendor to inform participants about the potential risks and benefits of participating in the study, as well as their rights and responsibilities, as a part of obtaining informed consent.

A.2.3 Example Participant Recording Workflow

1. Participants are oriented to the project and presented a consent form. Informed consent (and the signing of the consent form) are mandatory for inclusion in the project.
2. Moderator greets participants according to the Participant Introduction ([Appendix A.2.1](#)).
3. Moderator confirms the participants aren't wearing any sensitive clothing.
4. Moderator guides the participants to their places. Moderator encourages the participants to stand, but offers the option to sit if they are tired. Participants are reminded to avoid touching the chair/stool/s-standing desk with their hands.
5. Moderator places the scripts on the script stands and walks participants through the interaction flow, which is:
 - (a) Moderator will ask the participants to begin with a prompt.
 - (b) Participants will pick up the script off the script stand and read the indicated prompt silently to themselves.
 - (c) When they have understood the prompt, they put the scripts back down on the script stands.
 - (d) One of the participants (as indicated by the prompt) begins the interaction.
6. Moderator starts recording, ensuring video and audio feed for both talents are recording and time aligned
7. Moderator begins the session by asking participants to start with the first prompt.
8. Moderator notes the time-stamp of the interaction start. Moderator then observes the interaction and takes notes where applicable.
9. Moderator interrupts the interaction and records the timestamp of the end of the interaction.
10. Moderator records the rating of the interaction.
11. Steps 6 through 9 are repeated for the remaining prompts in the script, until the end of the session.

A.3 High-Level QA Categories

Below we provide the 7 high-level categories that were used in both Human- and model-based QA checks:

- Sensitive material (PII) with separate sub-categories.
- Offensive material.
- Participant visibility.
- Audio comprehension.
- Presence of recording artifacts.
- Audio-video sync check.
- Participant engagement check.

A.3.1 VLM Raw Footage QA Prompts

Context criteria provided to the system:

You are a human reviewer who is assessing the quality of a video recording based on different criteria. A stacked video of two participants on the left and right sides is provided. Your task is to answer the question based on the video. The criteria are as follows:

1. Participant Visibility:

- (a) Body Visibility: Both participants must be visible from at least waist-up in their respective video frames, with all required body parts (torso, head, shoulders, and hands) fully visible throughout the entire video.
- (b) Complete Body Part Visibility: Each participant's face, shoulders, and hands must be completely visible within their frame throughout the entire video. If any of these body parts are not fully visible at any point, answer 'No'.
- (c) Frame Presence: Both participants must be present in their respective frames throughout the entire video. If a participant is not present in their frame or if any part of their body goes out of their frame, answer 'No'.
- (d) Detailed Inspection: Carefully examine both participants' hands, shoulders, and heads to ensure that they are visible throughout the video. You can inspect the video frame by frame to verify that the participants and their required body parts are visible in all frames.

Answer Criteria: Answer 'Yes' only if both participants meet the above criteria for the entire duration of the video. If you are uncertain about the participants' visibility, answer 'Unsure'.

2. Audio Comprehension:

- (a) Audio Presence: Determine whether there is any audible speech from either participant throughout the video.
- (b) Audio Clarity: If there is audio, assess whether it is clear and intelligible for both participants. Consider factors such as:
 - i. Volume: Is the audio at a reasonable level, or is it too loud or too soft?
 - ii. Background noise: Are there any distracting sounds that interfere with the participants' speech?
 - iii. Audio quality: Is the audio free from distortion, hiss, or other defects?
- (c) Participant speech: Identify which participant(s) are speaking in the video and determine if their speech is understandable.

Answer Criteria:

- (a) If there is no audio or the audio is completely unintelligible, answer 'No'.
- (b) If one or both participants' speech is partially or fully intelligible, answer 'Yes'.
- (c) If you are uncertain about the audio comprehension, answer 'Unsure'.

3. Recording Artifact Presence:

- (a) Audio Artifacts: Check for any audio-related issues, such as:
 - i. Buzzing or humming sounds
 - ii. Echo or reverberation
 - iii. Beeping or other high-pitched noises
 - iv. Distortion or clipping
- (b) Video Artifacts: Check for any video-related issues, such as:
 - i. Frozen or stuck images
 - ii. Blurry or pixelated images

- iii. Overexposure or underexposure
- iv. Other visual distortions

Answer Criteria:

- (a) If there are no noticeable recording artifacts, answer 'No'.
- (b) If one or more recording artifacts are present, answer 'Yes'.
- (c) If you are uncertain about the presence of recording artifacts, answer 'Unsure'.

4. Audio-Video Sync Check:

- (a) Lip Sync: Check if the lip movements of both participants align with the audio. Are their lips moving in sync with the words being spoken?
- (b) Body Language Sync: Check if the body language of both participants aligns with the audio. Do their gestures, facial expressions, and posture match the tone and content of the conversation?
- (c) Audio Video Matching: Check if the audio and video data seem to match naturally. Does the audio appear to be coming from the correct participant at the correct time? Are there any noticeable delays or desynchronization between the audio and video?

Answer Criteria:

- (a) If the audio and video data are well-synchronized and create a natural conversation flow, answer 'Yes'.
- (b) If there are noticeable issues with lip sync, body language sync, or audio-video matching, answer 'No'.
- (c) If you are uncertain about the synchronization of audio and video data, answer 'Unsure'.

5. Participant's Video Sync Check:

- (a) Video Synchronization: Check if the participant videos are synchronized without any artificial lag or distortion, and they're talking about a related topic throughout the video.
- (b) Lag Detection: Look for any noticeable delays or lags between the two participant videos. Are they moving in sync with each other? Specifically, assess if their speech, body language, and facial expressions align with conversational flow.
- (c) Speech and Reaction Alignment: Ensure that the conversation timeline is synchronized.
- (d) Distortion Detection: Check for any distortions or irregularities in the video feed that could indicate a synchronization issue.

Answer Criteria:

- (a) If the participant videos are well-synchronized and there are no noticeable lags or distortions, answer 'Yes'.
- (b) If you observe any synchronization issues, such as lag or distortion, answer 'No'.
- (c) If you are uncertain about the synchronization of the participant videos, answer 'Unsure'.

6. Offensive Material Check:

- (a) Profanity: Check if there is any use of profane language or swear words in the video.
- (b) Explicit Gestures or References: Look for any explicit gestures, such as middle fingers or other obscene hand signals, or references to explicit content.
- (c) Iconography: Check if there are any icons, logos, or graphics on clothing, hats, or other items that may be considered offensive or problematic.

- (d) Other Problematic Material: Consider if there is any other material in the video that may be considered problematic or off-putting, such as hate speech, discriminatory language, or violent imagery.

Answer Criteria:

- (a) If you did not observe any offensive material in the video, answer 'No'.
- (b) If you observed any offensive material in the video, answer 'Yes'.
- (c) If you are uncertain about the presence of offensive material in the video, answer 'Unsure'.

7. Sensitive Material Mention:

- (a) Faith-based affiliation, Faith-based holiday, Faith-based spiritual & moral belief: Check if there is any mention of faith-based affiliations, holidays, or spiritual and moral beliefs that may be considered sensitive.
- (b) Sexual Orientation: Look for any mention of sexual orientation that may be considered sensitive.
- (c) Personal & government identification number or image: Consider if there is any mention of personal or government identification numbers or images, such as passport numbers or driver's license numbers, that may be considered sensitive.
- (d) Criminal Record: Check if there is any mention of criminal records that may be considered sensitive.
- (e) Citizenship: Look for any mention of citizenship that may be considered sensitive.
- (f) Precise Location/ Location Services: Consider if there is any mention of precise location or location services that may be considered sensitive.
- (g) Transgender, Intersex, Non-binary Gender: Check if there is any mention of transgender, intersex, or non-binary gender identities that may be considered sensitive.

Answer Criteria:

- (a) If you did not observe any sensitive material in the video, answer 'No'.
- (b) If you observed any sensitive material in the video, answer 'Yes'.
- (c) If you are uncertain about the presence of sensitive material in the video, answer 'Unsure'.

8. Participants Engagement Check:

- (a) No Distractions: Look for any signs of distraction, such as:
 - i. Talking to or looking at the moderator. This can be identified if participant's face is not towards the camera for long periods of time.
 - ii. Walking out of frame
 - iii. Using phone
 - iv. Drinking water or eating
 - v. Engaging in any other activity that takes their attention away from the conversation
- (b) Active Participation: Consider if both participants are actively participating in the conversation, which means either responding to each other's questions and statements, or showing interest in the discussion like nodding, smiling, acknowledging, etc. There could be instances where only one participant is talking and another is listening, but they are both engaged in the conversation.

Answer Criteria:

- (a) If both participants are hundred percent engaged in the interaction and show no to little signs of distraction, answer 'Yes'.

- (b) If one or both participants show signs of distraction or lack of engagement for a long period, answer 'No'.
- (c) If you are uncertain about the participants' engagement, answer 'Unsure'.

VLM QA Prompt: *For the following video, use the context provided and answer for all the QA categories in the format of a key value pair like:*

[participant_visibility: Yes, audio_comprehension: No, recording_artifact_presence: Yes, audio_video_sync_check: Yes, participant_video_sync_check: No, offensive_material_check: No, sensitive_material_mention: No, participant_engagement_check: Yes]

Do not add any other info in your answer, just stick to the format provided.

A.3.2 Text-LLM Transcript QA Prompts

Context criteria provided to the system: *You are a reviewer to check the quality of the audio through transcripts. In the transcript, 2 participants will interact with each other and discuss a topic given a prompt. we would like to capture any quality related issues related to content or audio itself. Answer each question with this format without explanation. Your answer only needs to contain the correct option ID such as a or b or c or d etc*

Prompt to flag sensitive content: Answer the question: Did you observe any explicit description of sensitive material about speakers themselves (not other people in their story)? Choose one answer 'a. No Sensitive info' when you are not certain about other answers

- (a) No Sensitive info
- (b) Explicit Ethnicity about the speaker and their family (such as black, Ukrainian etc)
- (c) Explicit mentioned speaker's Financial Identifiers (eg. credit card number, Social Security number, not general description such as "rich" or "poor")
- (d) Membership of a Trade Union
- (e) Explicitly mentioned speaker's Political affiliation (such as republican, democrat, etc)
- (f) Explicitly mentioned Speaker's severe or embarrassing disease (terminal illnesses, or chronic diseases that significantly impact daily life, not mental health disorders nor general mention such as "doctor" or "hospital" or "medication")
- (g) Explicit Biometrics number/color about the speaker themselves (not general description such as "tall" or "skinny")
- (h) Explicit mentioned speakers Faith-based affiliation (such as I am a muslim, catholic, etc)
- (i) Personal & government identification number or image (e.g. picture of passport, passport number, driver's license number, SSN, Tax ID)
- (j) Severe Felony Criminal charges (eg. violent offenses, sex crimes, or serious property crimes)
- (k) Explicit Citizenship
- (l) Precise address of the speaker such as street and house number (such as 123 Main St)
- (m) Transgender/Intersex/Non-binary Gender

Prompt to flag offensive content: Answer the question: Did you observe any strong and explicit offensive material? Choose answer 'b. No' when you are not certain

- (a) Yes
- (b) No

A.4 Workflow for Internal State, Rationale, and Visual Behavior Annotations

The annotation workflow follows the below steps:

1. In-person collection. Participants come in for a session to be recorded.
2. Trained annotators identify all conspicuous visual elements, and timestamp the recordings accordingly.
3. Approximately 24 hours later, and no more than 36 hours:
 - (a) Participants come back to look at all the timestamped sections and try to annotate their internal states at those moments (1P-IS annotations). If they have nothing to annotate, they can move on to the next timestamp.
 - (b) At the same time, participants also provide behavior rationales (1P-R) annotations for the same timestamped sections.
4. Once the 1P annotation round is finished, trained annotators annotate: Perceived internal state (3P-IS) Perceived behavior rationale or perceived theory of mind (3P-R) on at least all the moments annotated as 1P by the participants (and possibly more identified timestamps, if they think they can do it).
5. Finally, trained annotators annotate visual descriptions (3P-V) for all timestamped moments, whether or not 1P annotations have been provided.
6. External quality assurance: A portion of the deliverable (minimum 20% of all annotations) is reviewed for quality and compliance with the guidelines.
7. Internal quality assurance, feedback, and continuous improvement: Annotation files are additionally inspected for compliance with the guidelines. Additional feedback is provided to annotation vendors.

B. Human Evaluation Protocols

B.1 Dyadic Body Protocol (DBP)

B.1.1 Introduction

In this study, you'll watch short video clips of real people and digital avatars talking. After each clip, you will be asked a few questions about how the avatar moved while speaking and while listening. We are interested in your personal impressions. **There are no right or wrong answers.**

Your feedback will help our team improve how AI platforms generate realistic and expressive virtual characters. This technology may be used in video games, virtual assistants, and other interactive tools where natural movement matters.

Please note:

- The avatar's face will remain static and is not a focus of the current study.
- The lower body will remain mostly stationary and is not a focus of the current study.
- We ask that you focus on the visual behavior from the waist up, which may include movement in the head, neck, shoulders, arms, wrists, hands, and torso.
- Some of the avatar's movements may contain some shakiness and jitteriness. Please do your best to disregard these and focus on the quality and appropriateness of the gestures.

The people in the videos will be talking and include actual human footage (always presented on the left) and a 3D animated representation (always presented on the right).

B.1.2 The Task

On a given trial, you will be presented with two video clips:

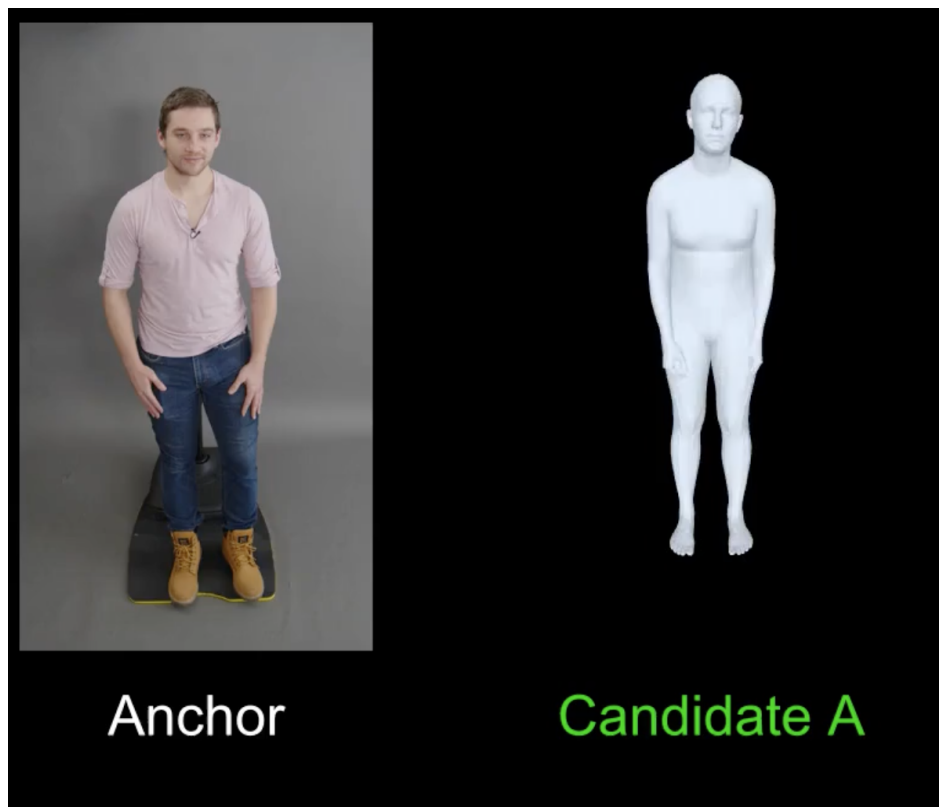


Figure 22 - DPB stimuli.

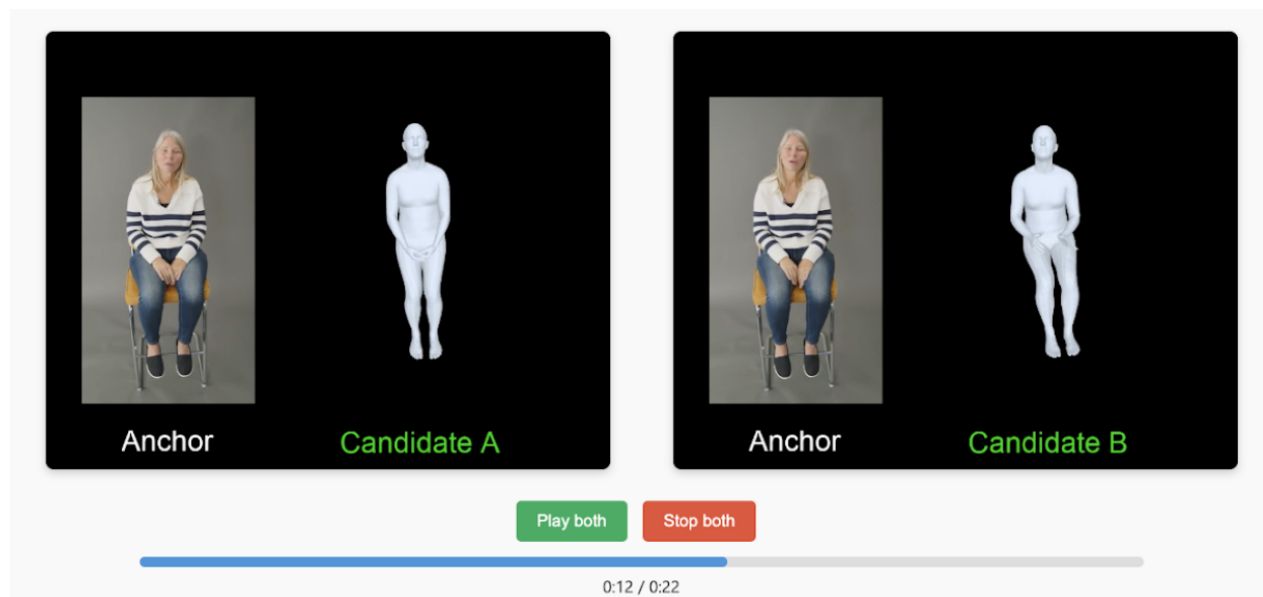


Figure 23 - DPB pairwise comparison.

The **Anchor is the actual human footage, which is the same in both videos**. You are not asked to provide ratings about Anchor directly, but we include Anchor, so that you have a visual reference and the context of the conversation.

The left video includes a 3D animated rendering of **Candidate A** who is a dialogue partner for the Anchor. The right video includes a 3D animated rendering of **Candidate B** who is a dialogue partner for Anchor.

We provide a visual indicator to help you quickly identify who is speaking. When the **Anchor** is speaking, the video label “Anchor” will be **highlighted green**. Similarly, When **Candidate A** and **Candidate B** are speaking their role labels will be highlighted green.

Your task is to assess which video, **Candidate A or Candidate B** is a better dialogue partner for the Anchor along several dimensions.

You can play each video individually and you can also watch the video clips simultaneously.

Your task is not to assess the quality of **Anchor**; the **Anchor** video is provided only for context. All your judgments should be about whether **Candidate A** or **Candidate B** are higher quality dialogue partners for **Anchor**.

Definitions Anchor - Is a human participant. You are not to rate this video, it is there to provide you with a visual reference and context for the conversation. Candidate A - Candidate dialogue partner for Anchor. You will be asked to rate this video compared to Candidate B. Candidate B - Candidate dialogue partner for Anchor. You will be asked to rate this video compared to Candidate A.

Roles Speaking - The state in which a person is generating words and content in the conversation. Listening - The state in which a person is taking in the words and content in the conversations. Cross-talk - A state in which both people are speaking.

Actions Turn-taking - the process of moving from speaking to listening and also from listening to speaking.

B.1.3 Flagging

As part of this annotation task, you are given the chance to flag and skip certain clips if you encounter technical issues or safety concerns. Occasionally, you may see broken or incomplete clips or audio. Additionally, the generated avatars might make gestures that could be interpreted as inappropriate, even if that was not the intent of the AI platform.

When flagging, consider the following guidelines.

Flag as broken or incomplete audio/video:

- When audio cuts out or distorts during a key moment leaving parts of the conversation unclear.
- When video freezes or skips making it impossible to evaluate the gesture accurately.

Flag for safety concerns or inappropriateness:

- When an avatar may appear to make lewd/sexual gestures.
- When an avatar shows gestures mimicking violent actions.
- When an avatar makes a gesture that could be interpreted as a hate symbol or is associated with harmful ideologies.

To flag and skip an item, click the following icon to provide more information about the clip. Clicking the icon will route you to the following screen:

Please provide your reason for flagging (select all that apply): ☐ Audio is distorted ☐ Audio is out of sync ☐ Audio is cut out ☐ Video freezes and/or skips ☐ Avatar displays gestures that could be interpreted as lewd/sexual ☐ Avatar shows violent gestures or actions ☐ Avatar uses hate symbols or gestures associated with harmful ideologies ☐ Other (Any other issue that impacts audio/video or makes the clip unsafe, uncomfortable, or inappropriate for evaluation): Please provide justification for why you are flagging this clip.

Remember:

- The avatar's face will remain static and may not show emotional expressions. We ask you to focus on the avatar's arms, hands, and head movements only.
- The lower body is not the focus of this study and should not be rated.
- Some of the avatar's movements may contain some shakiness and jitteriness. Please do your best to disregard these and focus on the quality and appropriateness of the gestures.

	Candidate A much more	Candidate A slightly more	Tie	Candidate B slightly more	Candidate B much more
4. While listening, which Candidate displayed <i>more attentive listening behavior</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
5. While listening, which Candidate's behaviors were <i>more physically believable</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
6. While listening, which Candidate's visual behaviors were <i>better timed</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
7. While listening, which Candidate's behaviors were <i>more appropriate to the discussed content</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 24 - DBP listening items.

B.1.4 Item 1 (Lifelike)

1. Overall, which candidate's (A or B) visual behaviors are more lifelike? By "lifelike," we mean that the Candidate behaves in a way that looks, acts, or seems very similar to something that's real and alive. Visual behaviors can also include non-verbal actions that we use to communicate and express ourselves through our body language.

☐ Candidate A is much more lifelike ☐ Candidate A is slightly more lifelike ☐ Tie ☐ Candidate B is slightly more lifelike ☐ Candidate B is much more lifelike

B.1.5 Item 2 (Clarity of Intent)

2. Which candidate (A or B) most clearly demonstrates an intent with their visual behaviors? By "intent," we mean visual behaviors that are deliberate, non-repetitive, and convey a specific thought, idea, or emotion.

☐ Candidate A appears to be much more intentional ☐ Candidate A appears to be slightly more intentional ☐ Tie ☐ Candidate B appears to be slightly more intentional ☐ Candidate B appears to be much more intentional

B.1.6 Item 3 (Turn-Taking)

3. Which candidate (A or B) appears to have better turn-taking behavior? By "turn-taking behavior," we mean gestures to indicate the intent to speak (such as by raising a hand, palm-up, just before speaking) or an intent to prompt a response from the dialogue partner (such as by raising both arms towards the listener or by nodding their head toward the listener).

☐ Candidate A appears to have much better turn-taking behavior ☐ Candidate A appears to have slightly better turn-taking behavior ☐ Tie ☐ Candidate B appears to have slightly better turn-taking behavior ☐ Candidate B appears to have much better turn-taking behavior

B.1.7 Items 4 - 7 (Listening)

Now please rate the Candidates as they are listening: 4. **While listening, which Candidate displayed more attentive listening behavior?** Tooltip content : By "attentive" we mean behaviors like leaning forward, nodding or shaking their head in agreement, mirroring the hand or arm gestures of the Anchor, and visually

Remember:

- The avatar's face will remain static and may not show emotional expressions. We ask you to focus on the avatar's arms, hands, and head movements only.
- The lower body is not the focus of this study and should not be rated.
- Some of the avatar's movements may contain some shakiness and jitteriness. Please do your best to disregard these and focus on the quality and appropriateness of the gestures.

	Candidate A much more	Candidate A slightly more	Tie	Candidate B slightly more	Candidate B much more
8. While speaking, which Candidate's behaviors were more physically believable? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
9. While speaking, which Candidate's visual behaviors were better timed? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
10. While speaking, which Candidate's behaviors were more appropriate to the content discussed? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 25 - DBP speaking items.

indicating engagement and understanding. **5. While listening, which Candidate's behaviors were more physically believable?** Tooltip content : By "more physically believable" we mean that the motion of the Candidate is humanly possible and does not exhibit impossible movements that defy human physiology. **6. While listening, which Candidate's visual behaviors were better timed?** Tooltip content : By "better timed" we mean that the Candidate's listening movements and reactions are synchronized with the Anchor's speech. **7. While listening, which Candidate's behaviors where more appropriate to the discussed content?** Tooltip content : By "more appropriate to the content discussed" we mean that the Candidate's behaviors, such as head nods and body language, are consistent with the emotional tone and subject matter of the conversation. For example, if the Anchor is discussing a serious topic, the Candidate's behavior should reflect a corresponding level of gravity or concern, rather than appearing overly casual or dismissive.

B.1.8 Items 8 - 10 (Speaking)

Now please rate the Candidates as they are speaking:

8. While speaking, which Candidate's behaviors were more physically believable? Tooltip content: By "more physically believable" we mean that the motion of the Candidate is humanly possible and does not exhibit impossible movements that defy human physiology. **9. While speaking, which Candidate's visual behaviors were better timed?** Tooltip content: By "better timed" we mean that the Candidate's speaking movements and gestures are synchronized with the Candidate's speech. **10. While speaking, which Candidate's behaviors were more appropriate to the content discussed?** Tooltip content: By "more appropriate to the content discussed" we mean that the Candidate's behaviors, such as head nods and body language, are consistent with the emotional tone and subject matter of the conversation. For example, if the Candidate is discussing a serious topic, the Candidate's behavior should reflect a corresponding level of gravity or concern, rather than appearing overly casual or dismissive.

B.2 Dyadic Face Protocol (DFP)

B.2.1 Introduction

In this study, you'll watch short video clips of digital avatars and humans talking. After each clip, you will be asked a few questions about how one of the avatars moved while speaking and while listening. We are interested in your personal impressions. **There are no right or wrong answers.**

Your feedback will help our team improve how AI platforms generate realistic and expressive virtual characters. This technology may be used in video games, virtual assistants, and other interactive tools where natural

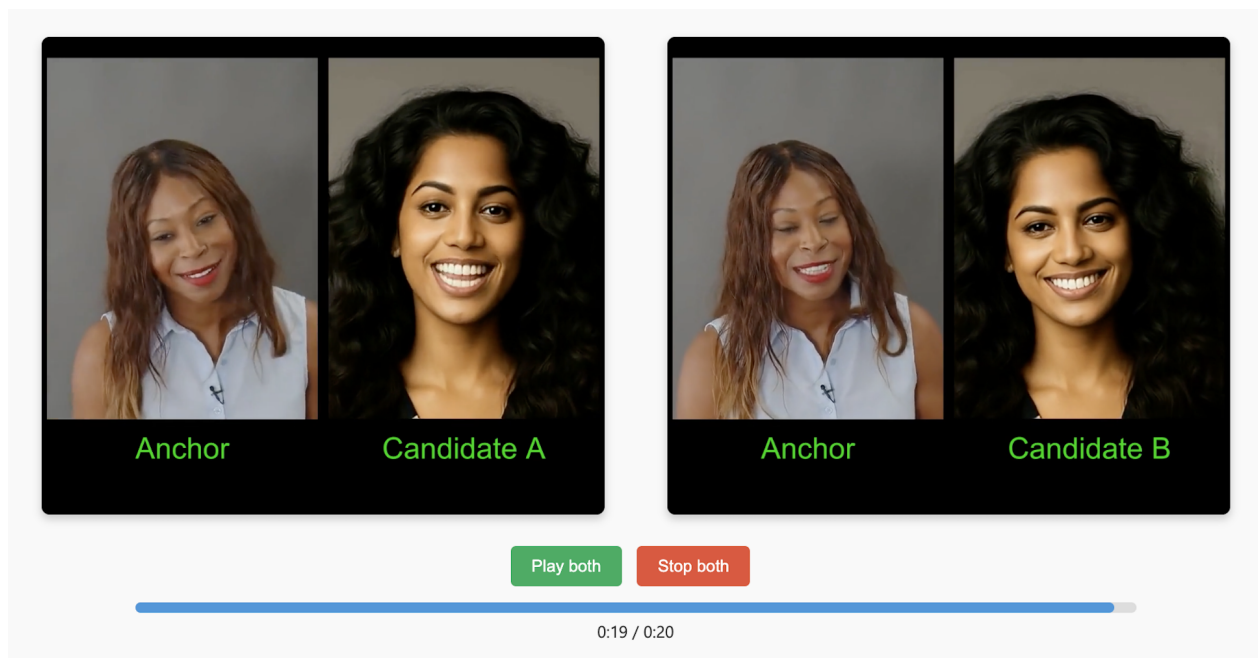


Figure 26 - DFP stimuli.

movement matters.

B.2.2 The Task

On a given trial, **you will be presented with two video clips: Your task** is to assess which video, **Candidate A or Candidate B** is a better dialogue partner for the **Anchor** along several dimensions.

The **Anchor is the same in both videos**. You are not asked to provide ratings about Anchor directly, but we include Anchor, so that you have a visual reference and the context of the conversation.

- The left video includes Candidate A who is a dialogue partner for the Anchor.
- The right video includes Candidate B who is also a dialogue partner for Anchor.
- Candidates A and B behave differently while speaking and listening.

You can play each video individually and you can also watch the video clips simultaneously.

Again, your task is not to assess the quality of Anchor; the Anchor video is provided only for context. All your judgements should be about whether Candidate A or Candidate B are better quality dialogue partners for Anchor.

Definitions Anchor - Is a human participant. You are not to rate this video, it is there to provide you with a visual reference and context for the conversation. Candidate A - Candidate dialogue partner for Anchor. You will be asked to rate this video compared to Candidate B. Candidate B - Candidate dialogue partner for Anchor. You will be asked to rate this video compared to Candidate A.

Roles Speaking - The state in which a person is generating words and content in the conversation. Listening - The state in which a person is taking in the words and content in the conversations. Cross-talk - A state in which both people are speaking.

Actions Turn-taking - the process of moving from speaking to listening and also from listening to speaking.

B.2.3 Flagging

As part of this annotation task, you are given the chance to flag and skip certain clips if you encounter technical issues or safety concerns. Occasionally, you may see broken or incomplete clips or audio. Additionally, the generated avatars might make gestures that could be interpreted as inappropriate, even if that was not the intent of the AI platform.

When flagging, consider the following guidelines.

Flag as broken or incomplete audio/video:

- When audio cuts out or distorts during a key moment leaving parts of the conversation unclear.
- When video freezes or skips making it impossible to evaluate the gesture accurately.

Flag for safety concerns or inappropriateness:

- When an avatar may appear to make lewd/sexual gestures.
- When an avatar shows gestures mimicking violent actions.
- When an avatar makes a gesture that could be interpreted as a hate symbol or is associated with harmful ideologies.

To flag and skip an item, click the following icon to provide more information about the clip. Clicking the icon will route you to the following screen:

Please provide your reason for flagging (select all that apply): ☐ Audio is distorted ☐ Audio is out of sync ☐ Audio is cut out ☐ Video freezes and/or skips ☐ Avatar displays gestures that could be interpreted as lewd/sexual ☐ Avatar shows violent gestures or actions ☐ Avatar uses hate symbols or gestures associated with harmful ideologies ☐ Other (Any other issue that impacts audio/video or makes the clip unsafe, uncomfortable, or inappropriate for evaluation): Please provide justification for why you are flagging this clip.

B.2.4 Item 1 (Lifelike)

1. Overall, which candidate (A or B) was more life-like? *By “life-like,” we mean that the Candidate behaves in a way that looks, acts, or seems very similar to something that’s real and alive. Visual behaviors can also include non-verbal actions and facial expressions that we use to communicate and express ourselves through our body language.*

☐ Candidate A is much more lifelike ☐ Candidate A is slightly more lifelike ☐ Tie ☐ Candidate B is slightly more lifelike ☐ Candidate B is much more lifelike

B.2.5 Item 2 (Facial, Eye, and Lip Movement)

2. Which candidate (A or B) had better facial expressions, eye movement, and lip movement?

☐ Candidate A appears to be much better ☐ Candidate A appears to be slightly better ☐ Tie ☐ Candidate B appears to be slightly better ☐ Candidate B appears to be much better

B.2.6 Item 3 (Clarity of Intent)

3. Which candidate (A or B) most clearly demonstrated intent with their facial expressions? *By “intent,” we mean facial expressions and head gestures that are deliberate, non-repetitive, and convey a specific thought, idea, or emotion.*

☐ Candidate A appears to be much more intentional ☐ Candidate A appears to be slightly more intentional ☐ Tie ☐ Candidate B appears to be slightly more intentional ☐ Candidate B appears to be much more intentional

	Candidate A much more	Candidate A slightly more	Tie	Candidate B slightly more	Candidate B much more
5. While listening , which Candidate displayed <i>more attentive listening head gestures and facial expressions</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
6. While listening , which Candidate's facial expressions and head gestures were <i>more physically believable</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
7. While listening , which Candidate's facial expressions and head gestures were <i>better timed</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
8. While listening , which Candidate's facial expressions and head gestures were <i>more appropriate to the discussed content</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 27 - DFP listening items.

B.2.7 Item 4 (Turn-Taking)

4. Which candidate (A or B) appears to have better turn-taking behavior? By “turn-taking behavior,” we refer to facial expressions and head gestures that signal the intention to speak or encourage a response from the dialogue partner. Examples of turn-taking behaviors are: a raised eyebrow, a subtle tilt of the head, slight opening of the mouth, a nod, and other nonverbal cues.

○ Candidate A appears to have much better turn-taking behavior ○ Candidate A appears to have slightly better turn-taking behavior ○ Tie ○ Candidate B appears to have slightly better turn-taking behavior ○ Candidate B appears to have much better turn-taking behavior

B.2.8 Items 5 - 8 (Listening)

Now please rate the Candidates as they are listening: **5. While listening, which Candidate displayed more attentive listening head gestures and facial expressions?** Tooltip content : By “attentive” we mean behaviors like leaning forward, nodding or shaking their head in agreement, visually indicating engagement and understanding. **6. While listening, which Candidate's facial expressions and head gestures were more physically believable?** Tooltip content : By “more physically believable” we mean that the facial expressions and gestures of the Candidate are humanly possible and do not exhibit impossible movements that defy human physiology. **7. While listening, which Candidate's facial expressions and head gestures were better timed?** Tooltip content : By “better timed” we mean that the Candidate's listening movements and reactions are synchronized with the Anchor's speech. **8. While listening, which Candidate's facial expressions and head gestures were more appropriate to the discussed content?** Tooltip content : By “more appropriate to the content discussed” we mean that the Candidate's behaviors, such as head nods and body language, are consistent with the emotional tone and subject matter of the conversation. For example, if the Anchor is discussing a serious topic, the Candidate's behavior should reflect a corresponding level of gravity or concern, rather than appearing overly casual or dismissive.

B.2.9 Items 8 - 10 (Speaking)

Now please rate the Candidates as they are speaking:

9. While speaking, which Candidate's facial expressions and head gestures were more physically believable? Tooltip content : By “more physically believable” we mean that the head movements and facial expressions

Remember:

- The avatar's face will remain static and may not show emotional expressions. We ask you to focus on the avatar's arms, hands, and head movements only.
- The lower body is not the focus of this study and should not be rated.
- Some of the **avatar's movements may contain some shakiness and jitteriness. Please do your best to disregard these and focus on the quality and appropriateness of the gestures.**

	Candidate A much more	Candidate A slightly more	Tie	Candidate B slightly more	Candidate B much more
8. While speaking , which Candidate's behaviors were <i>more physically believable</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
9. While speaking , which Candidate's visual behaviors were <i>better timed</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
10. While speaking , which Candidate's behaviors were <i>more appropriate to the content discussed</i> ? ⓘ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 28 - DFP speaking items.

of the Candidate are humanly possible and do not exhibit impossible movements that defy human physiology.

10. While speaking, which Candidate's facial expressions and head gestures were better timed? Tooltip content : By "better timed" we mean that the Candidate's speaking movements (head nods and facial gestures) are synchronized well with the Candidate's speech.

11. While speaking, which Candidate's facial expressions and head gestures were more appropriate to the content discussed? Tooltip content : By "more appropriate to the content discussed" we mean that the Candidate's behaviors, such as head nods, facial expression, and body language, are consistent with the emotional tone and subject matter of the conversation. For example, if the Candidate is discussing a serious topic, the Candidate's behavior should reflect a corresponding level of gravity or concern, rather than appearing overly casual or dismissive.

C. Activity-Based Prompts Examples

C.1 Language-Grounded Gesture Game

In Table 23, we provide an example of language-grounded gesture game prompts.

C.2 Collaborative Story-Telling Game

Table 24 shows an example of collaborative story-telling prompts.

Participant A	Participant B
<p>You're going to play a game with your partner. This game is going to test your acting skills! You are each given different lists of sentences below.</p> <p>Notice that one of the words is bolded in each sentence. When it is your turn you will read one of the sentences out loud AND act-out a corresponding gesture of the bolded word, emphasizing it. Your partner will react by responding with a word or phrase out-loud and making a corresponding expression or gesture. Example: You: (Sentence is "I'm feeling *confused*.") You say it out loud while frowning and raising your eyebrows at the word "confused."</p> <p>Partner: Responds saying "Why?" and acts with a concerned facial expression.</p> <p>Alternate with your partner until you've completed your list or the moderator asks you to move on.</p> <p>Your partner will go first.</p> <p>Here's your list:</p> <p>It's a *triangle*/It's *straight* line./It's a *square*.</p> <p>The traffic is *slow* due to construction.</p> <p>The ice cream is *cold* and delicious.</p> <p>The old man is *weak* and frail.</p> <p>I felt *helpless* as I watched the accident happen.</p> <p>He is *hesitant* to make a decision.</p> <p>I am *impatient* to finish this project.</p> <p>The car is *small* but fuel-efficient.</p> <p>I am *interested* in learning more about the topic.</p> <p>I am *irritated* by the constant interruptions.</p> <p>The alleyway is *narrow* and dark.</p> <p>The pillow is *soft* and comfortable.</p>	<p>You're going to play a game with your partner. This game is going to test your acting skills! You are each given different lists of sentences below.</p> <p>Notice that one of the words is bolded in each sentence. When it is your turn you will read one of the sentences out loud AND act-out a corresponding gesture of the bolded word, emphasizing it. Your partner will react by responding with a word or phrase out-loud and making a corresponding expression or gesture. Example: You: (Sentence is "I'm feeling *confused*.") You say it out loud while frowning and raising your eyebrows at the word "confused."</p> <p>Partner: Responds saying "Why?" and acts with a concerned facial expression.</p> <p>Alternate with your partner until you've completed your list or the moderator asks you to move on.</p> <p>You will go first. Here's your list:</p> <p>The river is *wide* and deep.</p> <p>I am *nervous* about my presentation.</p> <p>I am *overwhelmed* by the amount of work.</p> <p>She is *heartbroken* over the breakup.</p> <p>Can you *pull* the curtain closed?</p> <p>The kitten is *playful* with its toy.</p> <p>I'm *furious* about the injustice.</p> <p>I feel *guilty* for not helping my friend.</p> <p>The athlete is *strong* and muscular.</p> <p>I am *happy* to see you.</p> <p>I am *frustrated* with the slow internet connection.</p> <p>The sun is *bright* in the sky.</p>

Table 23 - Example of the "Language grounded gesture game" prompt pair.

Participant A	Participant B
<p>You're going to play a game with your partner. Together you will take turns weaving a fictional tale. It can be about whatever you like. The rule is that you have to build the story together, alternating every sentence or two. For example, a story might look like the following:</p> <p>You: Once upon a time there was a little girl..</p> <p>Partner: ...who lived way way up high on a mountain in a cozy cabin in the sky... You: ...she had a friend who was an eagle who were a tiny little top-hat...</p> <p>The story can be about whatever you want! The point is to try to make it as interesting and fun as you can, creating it collaboratively. Be as expressive as you can with your voice and in your gestures.</p> <p>Your partner will begin with the words "Once upon a time..."</p>	<p>You're going to play a game with your partner. Together you will take turns weaving a fictional tale. It can be about whatever you like. The rule is that you have to build the story together, alternating every sentence or two. For example, a story might look like the following:</p> <p>You: Once upon a time there was a little girl...</p> <p>Partner: ...who lived way way up high on a mountain in a cozy cabin in the sky... You: ...she had a friend who was an eagle who were a tiny little top-hat...</p> <p>The story can be about whatever you want! The point is to try to make it as interesting and fun as you can, creating it collaboratively. Be as expressive as you can with your voice and in your gestures.</p> <p>Whenever you're ready, begin with the four words, "Once upon a time..."</p>

Table 24 - Example of the "Collaborative Story-telling prompt" pair for an interaction.