**I·I ByteDance**

# OmniHuman-1.5: Instilling an Active Mind in Avatars via Cognitive Simulation

**Jianwen Jiang** [*†]   **Weihong Zeng** [*]   **Zerong Zheng** [*]   **Jiaqi Yang** [*]
**Chao Liang** [*]   **Wang Liao** [*]   **Han Liang** [*]   **Yuan Zhang**   **Mingyuan Gao**

Intelligent Creation Lab, ByteDance

## Abstract

Existing video avatar models can produce fluid human animations, yet they struggle to move beyond mere physical likeness to capture a character's authentic essence. Their motions typically synchronize with low-level cues like audio rhythm, lacking a deeper semantic understanding of emotion, intent, or context. To bridge this gap, **we propose a framework designed to generate character animations that are not only physically plausible but also semantically coherent and expressive.** Our model, **OmniHuman-1.5**, is built upon two key technical contributions. First, we leverage Multimodal Large Language Models to synthesize a structured textual representation of conditions that provides high-level semantic guidance. This guidance steers our motion generator beyond simplistic rhythmic synchronization, enabling the production of actions that are contextually and emotionally resonant. Second, to ensure the effective fusion of these multimodal inputs and mitigate inter-modality conflicts, we introduce a specialized Multimodal DiT architecture with a novel Pseudo Last Frame design. The synergy of these components allows our model to accurately interpret the joint semantics of audio, images, and text, thereby generating motions that are deeply coherent with the character, scene, and linguistic content. Extensive experiments demonstrate that our model achieves leading performance across a comprehensive set of metrics, including lip-sync accuracy, video quality, motion naturalness and semantic consistency with textual prompts. Furthermore, our approach shows remarkable extensibility to complex scenarios, such as those involving multi-person and non-human subjects.

**Date:** August 27, 2025
**Project Page:** https://omnihuman-lab.github.io/v1_5

## 1 Introduction

> "System 1 operates automatically and quickly, with little or no effort and no sense of voluntary control. System 2 allocates attention to the effortful mental activities that demand it, including complex computations."
>
> — Daniel Kahneman, *Thinking, Fast and Slow*

---

[*]Equal contributions
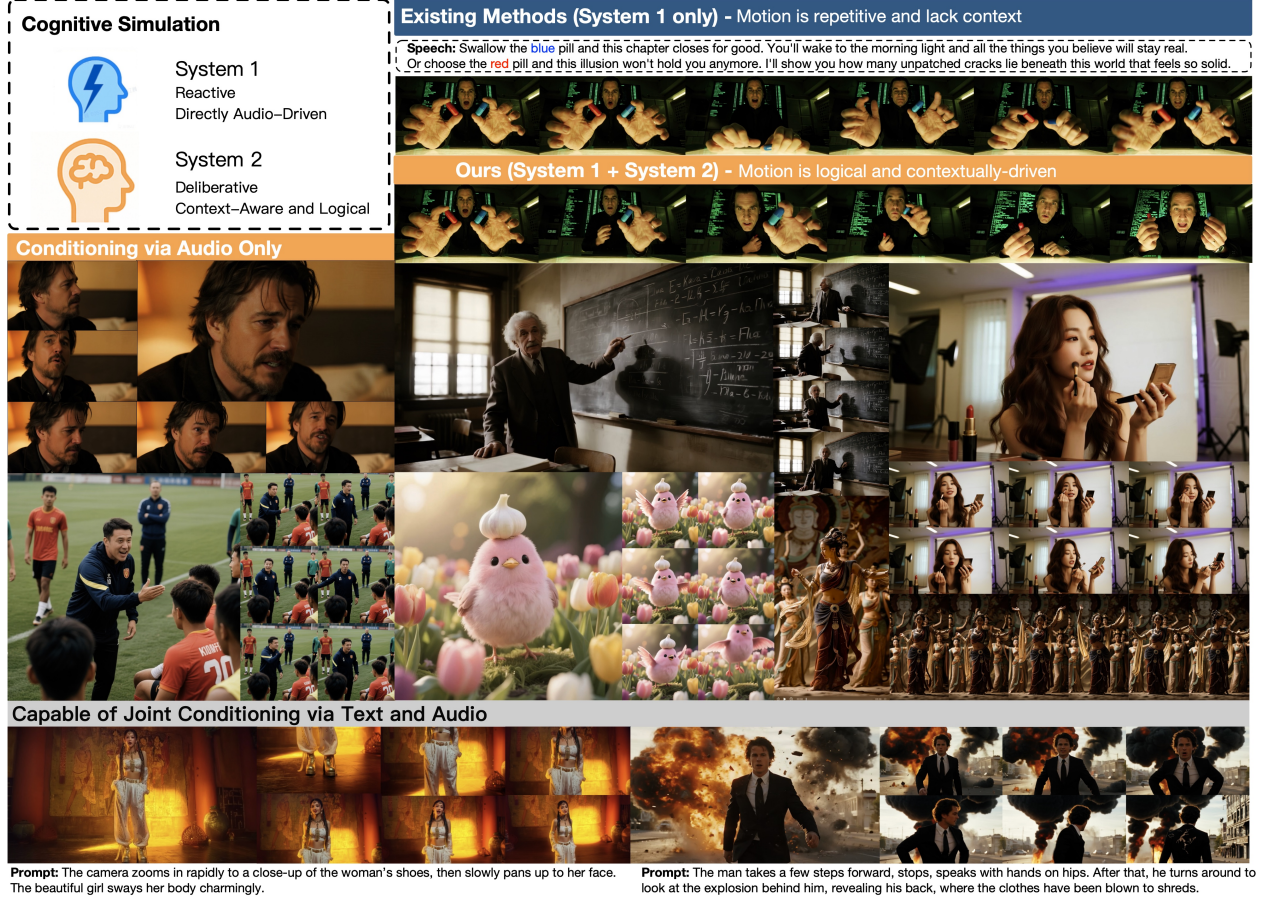[†]Project lead and corresponding author: jianwen.alan@gmail.com

**Figure 1  Simulating a Mind for Avatars.** We model avatar behavior by drawing on dual-system theory, which distinguishes between reactive System 1 and deliberative System 2 cognition. Top Left: Our framework combines System 1 actions (e.g., lip-sync, idle motions) with System 2 reasoning (e.g., logical gestures). Top Right: Conventional methods, analogous to System 1, excel at lip-sync but often produce repetitive, non-contextual motions. Bottom: In contrast, our method simulates both systems, generating diverse and naturally coherent behaviors that are semantically aligned with the provided audio and text.

The field of video avatars [8, 15, 19, 25, 30, 36, 39, 40, 53, 62, 67, 71, 74, 85, 86] aims to construct models capable of synthesizing realistic videos of characters from human-centric driving signals. The ultimate goal can be seen as creating lifelike avatars that are virtually indistinguishable from humans, capable of demonstrating both reasoned action and authentic emotion. This domain has witnessed rapid advancements over the past few years, with model capabilities evolving from early lip movement synthesis [29, 58, 59, 75, 93, 95] and portrait animation [10, 30, 67, 85, 86, 97] to more recent half-body [39, 65] and full-body generation [40]. As the scope of controllable generation expands, these models are increasingly expected to simulate a broader spectrum of human behaviors, moving beyond mere physical likeness to capture a character's authentic essence, as illustrated in Figure 1.

Recently, a series of audio-driven methods [15, 36, 40, 53, 74] based on Diffusion Transformers (DiT) [13, 35, 51, 55, 70] have emerged, enabling the generation of human motion videos that synchronize with audio in specific contexts. These approaches typically follow a similar paradigm: they first pre-train a foundational video generation model [55, 70] and then introduce conditional audio inputs to achieve animation. While this yields reasonable results, a closer inspection reveals that these models often only capture the direct and simplistic correlations between the audio signal and the resulting motion. Consequently, they tend to generate only synchronized lip movements and simple, repetitive accompanying gestures. When judged on their potential to function as convincing avatars, these outputs still exhibit a significant gap in naturalness

and plausibility when compared to authentic human behavior.

To understand the source of this gap, it is instructive to consider a prominent theory in human cognition, which posits that behavior is governed by two distinct systems [31, 32], often termed System 1 and System 2. System 1 is described as fast, unconscious, and reactive. This operational mode is analogous to that of current video avatar models, which excel at mapping input signals like audio directly to corresponding lip movements and simple gestures. In contrast, System 2 is characterized as deliberative, analytical, and effortful, involving reasoning about goals and context to deduce a course of action. This latter capability, the hallmark of intelligent behavior, remains challenging for existing methods to emulate robustly. Inspired by this dual-process model, we identify a clear path forward: leveraging the powerful reasoning capabilities [50, 54, 77, 80, 84, 89, 92] of modern Multimodal Large Language Models (MLLMs) [46, 64] to explicitly simulate the deliberative processes of System 2. We therefore propose a novel framework that integrates the principles of both systems, using MLLMs to simulate the goal-oriented aspects of System 2 while preserving the reactive capabilities analogous to System 1.

However, integrating the textual guidance from an MLLM into an avatar generation framework is non-trivial. Leveraging Chain-of-Thought (CoT) [80] prompting, MLLMs articulate their reasoning process as explicit text, a natural medium for expressing logic. Consequently, introducing an MLLM to enhance the high-level semantic coherence of motion inevitably strengthens the role of text as a conditional input. This creates a potential conflict of modalities within existing avatar models, where the audio signal already governs lip-sync and rhythmic motion, and the reference image constrains the character's appearance and potential range of motion. These distinct inputs are all intricately linked to the final generated motion. Therefore, we argue that a novel framework must be designed to mitigate modal conflicts while effectively leveraging the interdependencies among these diverse conditions. Developing such a framework is the key to simultaneously simulating both "System 1" and "System 2".

Based on the preceding analysis, we propose a video avatar framework built upon a Multimodal Diffusion Transformer that features two key designs. First, to ensure a comprehensive simulation of both System 1 and System 2, we introduce an agentic system powered by MLLMs. These agents reason over multimodal inputs (text, reference image, audio) to generate high-level semantic context, providing a long-term, logically coherent signal (for "System 2") that complements the short-term, reactive audio signals (for "System 1"). Second, to effectively fuse these input conditions while mitigating modality interference, we employ dedicated Multimodal Branches for audio and text feature extraction, along with a Multimodal Attention mechanism for joint modeling of audio, text and video. This design iteratively updates and aligns audio and text features into a common semantic space. Furthermore, we introduce a distinct identity preservation strategy that avoids conditioning on the reference image during training. Instead, we guide the model by probabilistically conditioning on start/end frames during training and treating the reference image as a pseudo last frame at inference, which prevents the static image from interfering with dynamic, content-driven motion. As a result of these designs, our model generates vivid and lifelike human motion from audio and a single image. The resulting videos exhibit remarkable contextual and semantic coherence, effectively simulating both the reactive and deliberative aspects of human behavior.

We summarize our main contributions as follows:

- **A New Perspective on Avatar Modeling:** We introduce a new perspective for analyzing video avatars, framing the problem through the cognitive science lens of System 1 and System 2 thinking. Observing that current models primarily simulate System 1, we are the first to propose a holistic approach that models both.

- **A Framework for Dual-System Simulation:** To implement this vision, we propose a novel framework featuring two core components. First, MLLM-based agents generate deliberative guidance ("System 2"). Second, a specialized MMDiT architecture, equipped with a symmetric audio branch and our pseudo-last-frame strategy, synergistically fuses this guidance with reactive signals ("System 1"), thereby resolving critical modal conflicts.

- **Strong Empirical Performance and Generalization:** Our method not only achieves highly competitive results on standard benchmarks but is also significantly preferred in user studies for its contextual

naturalness and plausibility. Its versatility is further demonstrated by its successful extension to complex multi-person and non-human scenarios.

## 2 Related Work

**Video Generation.** The field of video generation has seen rapid advancements, largely building upon the success of diffusion models in visual synthesis [21, 61]. Current approaches can be broadly categorized based on their underlying architecture. The first category consists of methods based on pre-trained text-to-image U-Net models [6, 12]. These approaches typically insert temporal modules, such as attention or convolution layers, into the frozen U-Net backbone and then fine-tune the model on video data [18, 73]. This allows them to leverage powerful image priors, but their capabilities can be constrained by the original image-centric architecture. The second category is represented by models that adopt a Diffusion Transformer (DiT) architecture [3, 7, 35, 42, 43, 52, 88, 96]. These methods treat video as a sequence of spatiotemporal patches, processing them in a unified manner with a Transformer. This approach has demonstrated superior scalability and flexibility, enabling the generation of high-resolution videos with variable durations and aspect ratios, especially when trained on massive datasets. The third category explores the integration of Large Language Models (LLMs) with diffusion models. In the image generation domain, works such as [34, 49, 57, 81, 83, 91]have shown that LLMs can enhance compositional understanding and planning. For video generation, this area is still in a nascent stage but holds significant promise for improving the logical coherence and narrative structure of generated content.

**Video Avatar Model.** Video avatars aims to create realistic human videos from various driving signals, with methods typically categorized by their driving source and synthesis pipeline. One category is pose-driven animation, which generates video based on an explicit, externally provided motion sequence, such as skeletal poses [5, 33, 56, 68, 76, 87, 94]. As the motion generation step is largely bypassed in these methods, their focus shifts to achieving high-fidelity rendering. More central to this paper is the second key category: audio-driven animation. This task requires a model to first generate plausible human motion from an audio signal and subsequently synthesize the final video. Consequently, the model must handle the dual challenges of both motion generation and rendering. Within this category, a common approach is a two-stage pipeline, where a motion generation model first translates audio into an intermediate representation, such as 2D/3D keypoints or 3D mesh sequences [11, 44, 79, 100] or discrete motion tokens [22, 66]. A subsequent rendering model then synthesizes the video conditioned on this motion. More recently, end-to-end approaches have emerged [38–40, 78], which directly generate video from audio in a single step, aiming for improved audio-motion synchronization. Despite their architectural differences, these audio-driven methods predominantly treat motion generation as a direct mapping process, a fast, reactive function from audio features to motion output. This process does not explicitly model the high-level cognitive phase where humans' planning and reasoning ultimately determine the resulting physical actions. We believe that incorporating such a reasoning phase is essential for generating more plausible and intelligent human behaviors, and our work aims to provide a foundational exploration in this direction.

**Large Language Models on Cognitive Simulation.** Large Language Models (LLMs) have reshaped machine reasoning with a clear upward trajectory in cognitive capability. Foundational models like GPT-3 [4] and ChatGPT [1, 48] laid the groundwork, while recent iterations such as GPT-4o [27], OpenAI o-series [28] and DeepSeek-R1 [17] enhance contextual and multi-modal reasoning, with prompting being key to unlocking this potential—Chain-of-Thought (CoT) [80] revolutionized problem-solving by decomposing tasks into logical steps. Derivatives like Self-Consistency [77] using aggregating reasoning paths and Tree-of-Thought (ToT) [89] with branching exploration further improved accuracy in arithmetic, logic, and commonsense tasks, confirming LLMs can simulate human-like deductive/inductive thinking via structured prompting. This reasoning capability has further empowered autonomous agents as their core engine. Classic practical frameworks like Toolformer [54], MetaGPT [23], and AutoGPT [60] with autonomous goal refinement enable agents to transcend pre-defined rules, exhibiting intent understanding and error correction. Voyager [72] leverages LLMs for open-world strategies and adaptive responses, outperforming rule-based agents. Furthermore, the renowned Generative Agents [50] enable LLMs to simulate daily human behaviors (e.g., social interaction, scheduling), marking a pivotal shift in human-like agent simulation. Recent advancements like OpenAI's Deep

Research [47] and Monica's Manus [45] further showcase LLM-driven proficiency in complex task planning and cross-domain execution. Beyond agent systems, LLM reasoning also steers generative models across domains to solve the "controllable generation" challenge. For image editing, InstructPix2Pix [2] relies on LLM-generated instructions to refine image modifications. MetaQueries [49] uses learnable queries to prompt vision-language models to guide image generation with enhanced semantic alignment. Besides, recent LLM-driven video generation agents [26, 37, 84, 92] enable more controllable long video synthesis through collaborative agentic workflows. These works collectively position LLMs as universal planners that bridge high-level user intent and the execution logic of low-level generative models. Yet notably, the integration of LLM-driven reasoning and planning into more fine-grained intelligent human/avatar behavior generation remains an underexplored area.

# 3 Approach

## 3.1 Overview

Our goal is to generate character animations that are both visually realistic and logically coherent with multimodal inputs. To achieve this, we introduce a framework designed to simulate both System 1 (reactive) and System 2 (deliberative) cognitive processes. Our model is built upon a Diffusion Transformer (DiT) backbone [13, 16, 51, 55], which is first pre-trained [16] on general video generation tasks to acquire foundational video generation capabilities. We then transform this base model into a logical and expressive avatar through two critical designs, which are detailed in the following sections and illustrated in Figure 2.

Agentic Reasoning for Deliberative Control (Sec. 3.2): We first employ MLLM-based agents to reason about the input context and generate high-level semantic guidance. This step provides the deliberative control necessary to simulate "System 2".

Multimodal Fusion for Reactive Rendering (Sec. 3.3): Next, our specialized MMDiT [13] architecture fuses this semantic guidance with reactive signals like audio to simulate "System 1". To resolve modal conflicts, this architecture incorporates a pseudo-last-frame identity strategy, which prevents the static reference image from interfering with dynamic motion during training.

Beyond these core designs, our framework aligns with common practices [15, 36, 40] for simplicity. To support long-form video synthesis, our framework can operate autoregressively, continuing generation by using the final frames of a previously generated clip as the initial frames for a new segment [62]. Our framework operates in the compact latent space of a pre-trained 3D VAE [90] and is trained with a flow matching [41] objective. We omit further discussion of these standard components to focus on our main contributions.

## 3.2 Agentic Reasoning for Deliberative Control

**Module Inputs and Outputs.** To model the deliberative nature of "System 2", our Agentic Reasoning module reasons over the input conditions to generate high-level, logically coherent guidance. The inputs are the character's reference image and the corresponding audio clip, supplemented by an optional text prompt describing the desired character behavior. The module is designed to process these inputs to produce semantic conditions in two forms: Reasoning Text, the explicit chain-of-thought text from the agents used directly as a condition, and Reasoning Latents, intermediate MLLM features extracted to serve as an additional conditioning signal and integrated via a dedicated attention mechanism. Our primary approach utilizes the former, while we have also investigated the integration of the latter.

**Multi-Step Reasoning Pipeline.** As shown in the top-right portion of Figure 2, this deliberative guidance is generated using two MLLMs that play distinct roles in a collaborative process. The first MLLM, the Analyzer, receives the reference image, a corresponding text caption (produced by an auxiliary model to ensure accurate image interpretation), the audio clip, and an optional user-provided text prompt. Guided by a chain-of-thought prompt, the Analyzer MLLM performs an iterative reasoning process on the inputs. It systematically infers the character's persona, language style, speech content, emotion, intent, and environmental context, consolidating these insights into a single structured representation, typically a JSON object. This context-rich information is then passed to the second MLLM, the Planner. The Planner receives the Analyzer's output, along with the original character image for visual context. Guided by a new instructional prompt, its task is to devise an
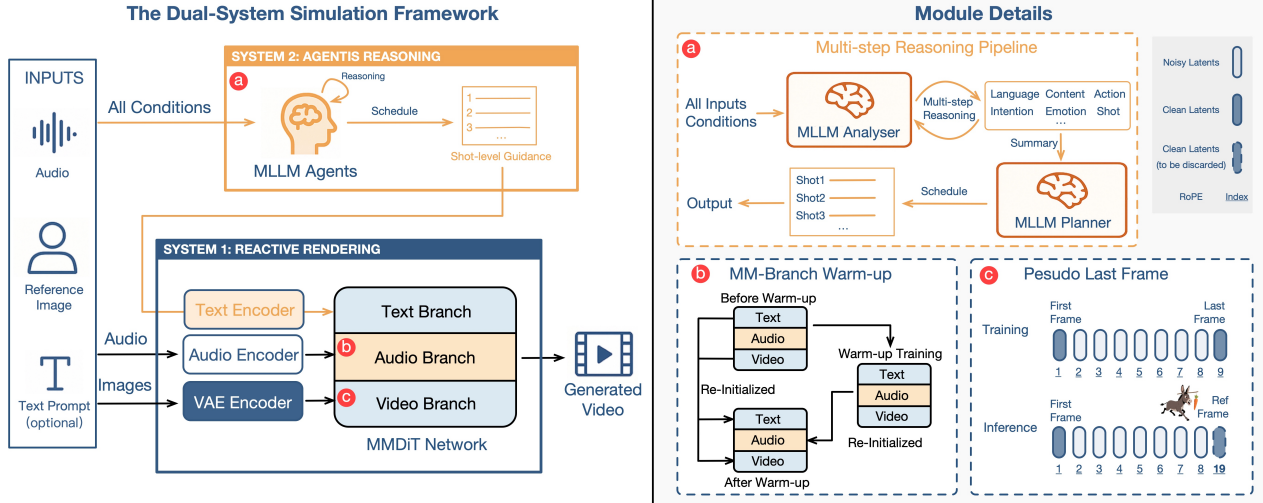
**Figure 2  The Dual-System Simulation Framework.** Our framework models avatar behavior by integrating a deliberative **System 2** for planning with a reactive **System 1** for rendering. **Left:** The overall pipeline. System 2 uses an MLLM Agent to reason over all inputs (audio, image, text) and generate a high-level "schedule". This schedule guides System 1's MMDiT network, which synthesizes the final video by fusing information within its dedicated text, audio, and video branches. **Right:** Key module details. (a) The reasoning pipeline consists of an MLLM Analyser and Planner that work together to create the schedule. (b,c) Our proposed **MM-Branch Warm-up** and **Pseudo Last Frame** methods mitigate multimodal conflicts during training.

action plan structured as a sequence of shots, where each shot defines the character's expressions and actions for a duration corresponding to a single generation pass of our diffusion model. This collaborative reasoning yields a comprehensive motion schedule that preserves a coherent character persona through consistent actions across the entire video.

**Reflective Re-planning.** To maintain logical coherence in long-form video generation, our agentic framework incorporates an optional "reflection" process. During autoregressive synthesis, the generation schedule is dynamically updated by re-evaluating the most recently generated output. This process mitigates a common challenge in diffusion-based synthesis, where subtle execution deviations can accumulate and degrade logical coherence over time, particularly in longer videos. In practice, the Planner takes the last generated frames and the original reference image as new inputs to re-assess its plan. This reflective loop corrects for semantic drift and helps maintain the video's logical consistency.

**Investigation of Latent Feature Conditioning.** We also investigated the aforementioned approach of utilizing latent features from a MLLM as a semantic conditioning signal. The approach directly utilized the audio tokens within the transformer of the Analyzer agent. The rationale was that cross-modal attention in the transformer layers would enrich these audio token representations with high-level semantics (e.g., inferred emotion and intent) while preserving their original temporal structure. Based on this consideration, we selected these "reasoning-infused" audio latents from the final transformer layer and concatenated them with the raw audio features. This combined signal then replaced the original audio input for the DiT network.

The above designs enable our agent to formulate a global, coherent plan for the entire scene. Unlike purely reactive methods that merely emulate "System 1", our approach further integrates deliberative reasoning from "System 2" to provide thoughtful, top-down guidance.

## 3.3  Reactive Rendering via Multimodal Diffusion

In this subsection, we describe how our diffusion model synthesizes the final video. It synergistically combines the high-level reasoning from the agents (primarily represented as text) with the low-level, reactive signals of audio inputs (primarily represented as audio features).

**Rethinking Reference Conditioning.** Before detailing our driving condition modeling, we must first analyze a critical input in video avatar models: the conditioning image, which serves two distinct purposes. The first is providing initial frames for autoregressive continuity, a standard practice we adopt by concatenating ground-truth (GT) frames. The second, more problematic purpose, is using a reference image for identity preservation. While early works used dedicated networks [24, 67, 99] and recent methods reuse model parameters [36, 40], both approaches inject a reference image sampled from the training video. As illustrated in Figure 3, we argue that this creates a critical artifact: the model learns a spurious correlation that the reference image should appear literally within the generated sequence, severely restricting motion dynamics and conflicting with the audio and text driving signals. While this artifact can be partially mitigated by probabilistically sampling reference images from outside the training video clip, this approach may introduce new problems, causing the model to learn that the generated output should exhibit significant variation from the reference image.

The root cause of this issue is that the reference image is an artificial construct, not a condition native to the video data itself. It is for this reason that our solution is to discard it entirely during training and, in its place, introduce a novel guidance mechanism. As shown in the bottom right of Figure 2, during training, we probabilistically condition the model on both the GT first and last frames of the video clip, as these are both native signals. During inference, we repurpose this mechanism by placing the user's reference image in the last frame's position, creating a pseudo last frame. Crucially, we shift its positional encoding (e.g., RoPE [63]) to maintain a fixed temporal distance from the generated content. This pseudo-frame, which is dropped after rendering, functions as a "carrot on a stick": it guides the model toward the reference identity without ever forcing it to replicate the image. As our experiments show, this approach eliminates training artifacts and mitigates autoregressive error, achieving a superior trade-off between motion dynamics and stability.



**Figure 3 Rationale for the Pseudo Last Frame. Left:** Reference-conditioning has trended toward simplification. **Right:** The dilemma of reference image sampling. Sampling from **within** the target video segment ensures high relevance but restricts motion diversity. Conversely, sampling from **outside** the segment, a scenario that becomes more frequent as datasets grow larger and more dynamic, leads to a drop in content relevance, causing inconsistencies that undermine the reference's purpose.

**Symmetric Fusion and Warm-Up.** With all training conditions now native and compatible, we address the challenge of joint modeling. We adopt an MMDiT backbone but depart from prior work in our approach to audio conditioning. Instead of injecting audio features via additional cross-attention layers, we introduce a dedicated audio branch, architecturally symmetric to the video and text branches. All three modalities are then fused at each layer through a **shared** multi-head self-attention mechanism. This symmetric design offers two key advantages. First, it allows audio features to be iteratively refined alongside video and text, ensuring deep semantic alignment. Second, it enables true joint modeling, as tokens from all three modalities mutually attend to one another, facilitating a more effective mapping to a shared semantic space. While this new branch adds parameters, the computational overhead is negligible due to the low ratio of audio to video tokens.

This symmetric architecture, while beneficial, introduces a training challenge. Naively training the entire model jointly leads to modality conflict, where the model learns to over-rely on the temporally dense audio signal for all predictions, thereby ignoring or washing out the more abstract guidance from the text branch. Freezing the pre-trained branches is also suboptimal, as it causes the audio branch to overfit and erroneously learn non-audio-related attributes like lighting and camera motion. To resolve this, we propose a two-stage warm-up strategy. In Stage 1, we train the full three-branch model jointly, forcing the model to learn an optimal division of labor: the text and video branches handle high-level semantics, compelling the audio branch to specialize in its core competencies (e.g., lip sync, speech mannerisms). For Stage 2, we construct
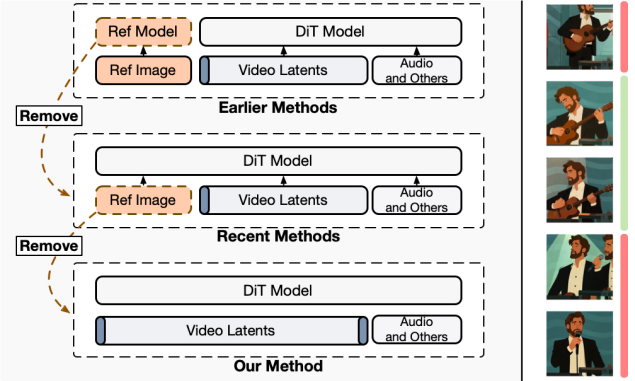
the final model. The text and video branches are initialized with their original pre-trained weights, while the audio branch is initialized using the warmed-up weights obtained from Stage 1. This model is then fine-tuned. This strategy ensures each branch begins with its own strong, specialized prior, mitigating modality conflict and allowing each input to retain its distinct conditioning power.

Ultimately, the proposed architecture executes the deliberative plan. By redesigning the reference conditioning and equipping the model with an audio conditioning branch, our rendering process faithfully translates the high-level guidance from System 2 while maintaining the reactive fidelity of System 1.

## 4 Experiments

### 4.1 Experimental Setup

**Implementation Details.** Our model is based on the MMDiT architecture and was pre-trained on a large-scale dataset of text-video/image pairs. For most experiments, the model generates 120-frame clips at 24 fps with a 480p resolution (short side). A separate super-resolution model of the same architecture is used to upscale outputs to 720p or 1080p. Longer videos are generated autoregressively. We used the AdamW optimizer with a learning rate of 5e-5, a global batch size of 256, and gradient clipping at a norm of 1.0. The training was conducted on 256 compute nodes and consisted of three stages: a 3-day audio branch warm-up, a 7-day main training phase, and a 1-day fine-tuning phase on high-quality data.

**Training Data.** Our training set consists of 15,000 hours of filtered video data. Following prior work, we used a lip-sync model to identify and discard the audio from videos with poor lip-audio correlation. These samples, comprising 70% of the data, were used with audio-dropout during training. For the final fine-tuning stage, we ranked our training data by quality metrics and selected the top 100 hours.

**Evaluation Datasets.** To rigorously evaluate our model, we noted that current DiT-based methods already perform well in standard human speaking scenarios. To test the true generalization limits of our approach, we therefore constructed two novel and highly challenging test sets. Our first custom benchmark is a diverse single-subject set of 150 cases, including real-world human portraits, AIGC figures, anime characters and animals. Each image was manually paired by experts with a corresponding audio track, such as speech, singing or theatrical performances, to create a demanding generalization test. To assess performance in more complex scenes, we also built a multi-subject set of 57 cases, featuring the same visual diversity and expert-paired audio for multi-character interactions. Furthermore, to evaluate our model's text-conditioning, experts wrote descriptive prompts for all 150 single-subject cases, allowing us to measure adherence to textual guidance. Finally, for fair comparison with prior work, we adopted their experimental settings, using 100 videos from CelebV-HQ [98] for the talking-head task and the CyberHost [39] test set (269 videos, 119 identities) for evaluating performance in full-body scenarios.

**Evaluation Metrics.** To comprehensively assess our method, we employ a multi-faceted evaluation protocol that includes both objective and subjective metrics. For objective evaluation, we measure generation quality using the Fréchet Inception Distance (FID) [20] and Fréchet Video Distance (FVD) [69], alongside no-reference Image Quality (IQA) and Aesthetics (ASE) scores [82]. We also evaluate audio-visual synchronization with Sync-C [9], hand quality using Hand Keypoint Confidence (HKC) and Hand Keypoint Variance (HKV) [39].

However, as these objective metrics often fail to capture higher-level semantic qualities and overall perceptual realism, we also conducted a comprehensive subjective user study with 40 participants. This study involved two main protocols. The first was a **pairwise comparison**, where participants viewed two videos from different methods in a randomized order. This comparison was twofold: from a positive perspective, users selected the video with the best overall quality, from which we calculated the Good/Same/Bad (GSB) score, defined as $(\text{Wins} - \text{Loses})/(\text{Wins} + \text{Loses} + \text{Ties})$; from a critical perspective, they identified specific flaws: Lip-sync Inconsistency (LSI), Motion Unnaturalness (MU), and Image Distortion (ID), allowing us to compute a defect rate for each. The second protocol was a **best-choice selection task**, where participants selected the single best video from all competing methods. This yielded a Top-1 selection rate, providing a direct measure of overall appeal.

| Method | IQA ↑ | ASE ↑ | Sync-C ↑ | HKC ↑ | HKV ↑ |
|---|---|---|---|---|---|
| *Ablation on Agentic Reasoning* | | | | | |
| Ours w/o Multi-Step Reasoning | 4.795 | 3.901 | 3.853 | 0.576 | 157.638 |
| Ours w/o Analyzer | 4.793 | 3.910 | 4.278 | 0.572 | 148.381 |
| Ours w/o Reasoning (System 1 Only) | 4.784 | 3.885 | 3.507 | 0.544 | 122.376 |
| *Ablation on Conditioning Modules* | | | | | |
| Ours w/ Cross-Attention | 4.745 | 3.856 | 3.263 | 0.558 | 116.317 |
| Ours w/o MM-Warmup | 4.752 | 3.866 | 3.993 | 0.549 | 164.080 |
| Ours w/ Ref. Image | 4.772 | 3.896 | 3.982 | 0.559 | 160.889 |
| Ours w/o Ref. & Pseudo Frame | 4.682 | 3.878 | 4.141 | 0.564 | 160.986 |
| **Ours (Full Model)** | 4.790 | 3.901 | 4.087 | 0.571 | 168.912 |

**Table 1  Ablation studies on our proposed framework.**

**(a)** Ablation on agentic reasoning.

| Method | LSI↓ | MU↓ | ID↓ | GSB↑ |
|---|---|---|---|---|
| Ours (w/o Reasoning) | 0.12 | 0.58 | 0.11 | −0.29 |
| Ours (Full Model) | 0.12 | 0.37 | 0.04 | +0.29 |

**(b)** Comparison of conditioning.

| Conditioning Method | LSI↓ | MU↓ | ID↓ | GSB↑ |
|---|---|---|---|---|
| Previous Work [40] | 0.21 | 0.39 | 0.17 | −0.23 |
| Ours (Proposed) | 0.03 | 0.25 | 0.07 | +0.23 |

**(c)** GSB Score Comparisons

| GSB Comparisons | TA ↑ | Mot ↑ | VQ ↑ |
|---|---|---|---|
| Ours vs. Base Model | -0.02 | +0.18 | +0.14 |

**Table 2  Pairwise subjective ablation study and component comparison.** We report Lip-sync Inconsistency (LSI), Motion Unnaturalness (MU), Image Distortion (ID), and an overall Good/Same/Bad preference score (GSB). Lower is better for LSI, MU, and ID.

## 4.2  Ablation Studies

In this section, we conduct a series of ablation studies to rigorously validate the contributions of our proposed components. The experiments are performed on a custom, single-subject test set of 150 video clips. Our analysis systematically isolates the impact of two key elements: (1) the agentic reasoning module and (2) the proposed conditioning architecture within our diffusion model. For a comprehensive assessment, we employ both quantitative and subjective evaluations. The quantitative metrics provide an objective measure of performance, while the subsequent user studies assess perceptual quality in terms of lip-sync consistency, motion naturalness, image quality, and overall user preference.

**The Effectiveness of Agentic Reasoning.** Here, we analyze the contribution of the Agentic Reasoning module by dissecting its intermediate steps. We conduct experiments by first removing the multi-step reasoning process, then ablating the entire Analyzer, and finally, removing the reasoning module altogether, which results in a "System 1 only" model. As shown in Table 1, standard metrics for image quality (IQA) and lip-sync (Sync-C) show only minor variations across these ablations. This is expected, as these metrics primarily evaluate low-level fidelity, which remains high in all diffusion-based variants. However, they are not designed to measure higher-level semantic qualities like logical coherence. A more telling objective trend emerges with the HKV metric, which progressively decreases as reasoning capabilities are weakened, indicating that the generated animations become more static and less expressive.

To truly evaluate the core contributions of our reasoning module to these semantic qualities, we therefore turn to subjective evaluation. The results, presented in (a) of Table 2, offer a direct comparison between the models with and without the agentic reasoning module. The overall user preference (GSB score) immediately reveals a substantial advantage for our full model. More specifically, the introduction of reasoning leads to a significant reduction in perceived motion unnaturalness (MU), with over a 20% improvement in pairwise comparisons. Furthermore, it maintains or slightly improves lip-sync consistency and image quality, as reflected by the LSI and ID metrics. These findings support the effectiveness of our Agentic Reasoning module, particularly in its ability to enhance the plausibility and semantic motion naturalness of the generated animations, which are qualities not fully captured by objective metrics.

| Method | Subjective Evaluation | | | | Quantitative Metrics | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | DA↑ | LSI↓ | MU↓ | GSB↑ | IQA↑ | ASE↑ | Sync-D↓ | HKC↑ | HKV↑ |
| InterActHuman | - | - | - | - | 4.574 | 3.643 | 8.163 | 0.553 | 103.91 |
| Ours w/o Reasoning | 0.88 | 0.13 | 0.63 | -0.26 | 4.576 | 3.631 | 7.541 | 0.611 | 138.43 |
| **Ours (Full Model)** | 0.94 | 0.04 | 0.12 | +0.26 | 4.529 | 3.653 | 6.904 | 0.614 | 158.36 |

**Table 3  Comparison with existing methods on multi-person animation.** We report quantitative metrics and pairwise subjective evaluation results, including Driving Accuracy (DA), Lip-sync Inconsistency (LSI), Motion Unnaturalness (MU) and an overall user preference score derived from a Good/Same/Bad (GSB) evaluation.

**The Effectiveness of Proposed Conditioning Modules.** We now ablate our core architectural designs, with results presented alongside the previous study in Table 1 and 2. In these experiments, the Agentic Reasoning module remains fixed, providing identical inputs to all model variants. We test several key variations: using standard cross-attention for audio integration instead of our MM-Attention, removing the MM-Warmup strategy, conditioning on a reference image and omitting the pseudo-last-frame at inference. As shown in Table 1, our full model again leads in most objective metrics, with its superior HKC and HKV scores highlighting enhanced motion dynamics. To further validate our approach, (b) of the Table 2 presents a direct subjective comparison against OmniHuman-1 [40], a state-of-the-art method that utilizes a reference attention mechanism and standard cross-attention for audio injection instead of our proposed conditioning implementation. The results show that our method achieves a significant advantage not only in the overall GSB score but also across multiple fine-grained dimensions, including lip-sync accuracy, motion naturalness, and visual quality. This clearly demonstrates the effectiveness of our proposed conditioning technique, which in turn provides a robust foundation for executing the plans generated by the agentic reasoning module.

In addition to the ablation studies, Table 2 presents a direct comparison with the base model under text-only conditioning. For this evaluation, the audio component was disregarded to isolate the assessment of visual fidelity and character motion. We conducted a GSB pairwise comparison to analyze performance across three key areas: text alignment (TA), motion naturalness (Mot), and visual quality (VQ). The results reveal that our model successfully integrates multi-modal inputs while preserving a text-prompt-following capability on par with the pre-trained general model. Critically, our approach also demonstrates a significant lead in both motion naturalness and overall visual quality, as reflected by the Mot. and VQ metrics.

## 4.3  Further Exploration on Applications

**Applications on Diverse Inputs.** As depicted in Figure 4, we also explore the generalization capabilities of our model on non-human subjects, including anthropomorphic and animal characters. The results demonstrate remarkable robustness, which we attribute to our dual-system framework that effectively integrates high-level understanding with low-level synthesis. Furthermore, the fourth row of the Figure 4 highlights the model's capacity for conversational understanding, enabled by the agentic reasoning module. For the top and bottom images, we provide the same conversational dialogue but input the audio track corresponding to a different speaker for each. As shown, the characters seamlessly transition between speaking and idle states, correctly reflecting the turn-taking in the dialogue. This capability, when combined with efficient model acceleration techniques, showcases the potential of our framework for real-time, interactive conversational agent applications.

**Applications on Multi-person Scenarios.** To enable multi-person animation, we extend our model with two modifications. First, we condition the synthesis on a speaker-specific mask that directs audio feature injection exclusively to the masked regions during the multimodal attention process. Following InterActHuman [78], we employ a lightweight, plug-and-play predictor to dynamically generate these masks, ensuring robust speaker tracking through movement and occlusion without affecting the baseline single-person model. Second, we leverage our framework's inherent agent-based design by augmenting the Planner to also accept this mask to identify the active speaker. With the rest of the reasoning pipeline remaining identical, this simple extension enables the model to generate logically consistent and coordinated actions for all individuals
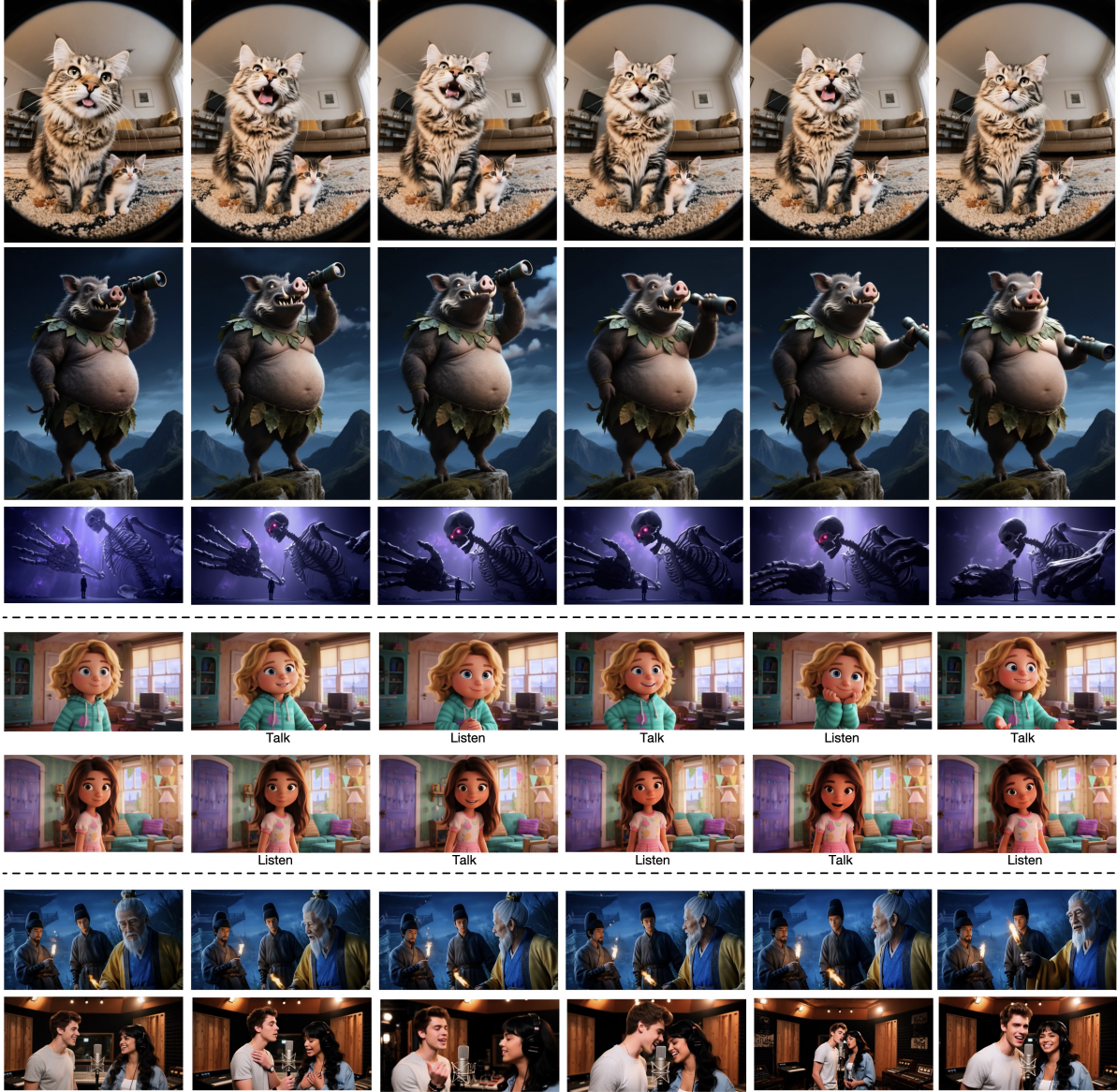
**Figure 4  Generalization and Multi-Person Results.** The top rows show the model's generalization across various non-human subjects. The fourth row presents a dialogue scenario, where characters correctly respond to conversational audio by switching between speaking and idle states. The bottom rows showcase performance in multi-person scenes, with coordinated behavior for both speakers and listeners.

in the scene.

As shown in Table 3, we present a quantitative and subjective comparison on our multi-person test set. Our full model, equipped with the agentic reasoning module, demonstrates significant improvements over two baselines that lack this capability: our model without agentic reasoning (ablation) and InterActHuman. Specifically, our method shows a clear advantage on metrics measuring gesture motion dynamics (HKC and HKV) and achieves better lip-sync accuracy. It is worth noting that we use Sync-D [9] for evaluation, as the original Sync-C is effective for single-person lip-sync, and is less reliable for non-speaking individuals in multi-person scenarios. Furthermore, in pairwise subjective evaluations against the ablation model, our full model achieves higher driving accuracy (DA), defined as the ratio of correctly animated individuals (either speaking or silent) to the total number of people present. It also produces fewer instances of lip-sync

| Method | IQA ↑ | ASE ↑ | Sync-C ↑ | FID ↓ | FVD ↓ | | Method | IQA ↑ | ASE ↑ | Sync-C ↑ | FID ↓ | FVD ↓ | HKC ↑ | HKV ↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SadTalker | 2.953 | 1.812 | 3.843 | 36.648 | 171.848 | | Skyreel-A1 | 3.889 | 2.525 | 2.983 | 69.619 | 70.678 | 0.786 | 28.840 |
| Hallo | 3.505 | 2.262 | 4.130 | 35.961 | 53.992 | | FantasyTalking | 3.892 | 2.738 | 3.548 | 52.332 | 47.052 | 0.838 | 18.845 |
| EchoMimic | 3.307 | 2.128 | 3.136 | 35.373 | 54.715 | | OmniAvatar | 3.871 | 2.728 | 6.589 | 42.163 | 43.998 | 0.795 | 56.574 |
| Loopy | 3.780 | 2.492 | 4.849 | 33.204 | 49.153 | | MultiTalk | 3.822 | 2.681 | 6.868 | 37.308 | 32.783 | 0.817 | 62.753 |
| Hallo-3 | 3.451 | 2.257 | 3.933 | 38.481 | 42.125 | | OmniHuman-1 | 4.142 | 3.024 | 7.443 | 31.641 | 27.031 | 0.898 | 47.561 |
| OmniHuman-1 | 3.875 | 2.656 | 5.199 | 31.435 | 46.393 | | | | | | | | | |
| Ours | 3.817 | **2.663** | 5.053 | **31.320** | 45.771 | | Ours | **4.144** | **3.030** | 7.243 | **31.160** | 27.642 | 0.875 | **72.113** |

**Table 4  Quantitative comparison with audio-conditioned animation baselines. (Left)** Portrait animation on the CelebV-HQ test set. **(Right)** Full-body animation on the CyberHost test set.

| Method | Top-1 (%) |
|---|---|
| Skyreel-A1 | 5% |
| FantasyTalking | 8% |
| OmniAvatar | 14% |
| MultiTalk | 18% |
| OmniHuman-1 | 22% |
| **Ours** | 33% |

**(a)** Best-Choice Selection.



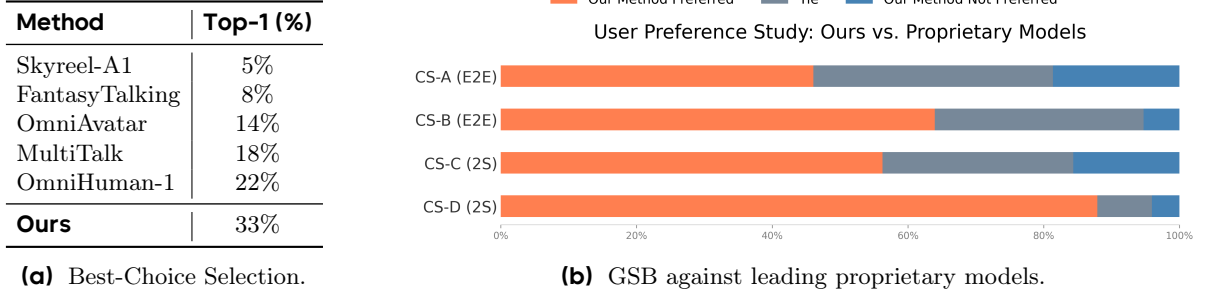**(b)** GSB against leading proprietary models.

**Figure 5  Subjective User Preference Study.** We present results from two evaluation settings: (Left) a best-choice selection task comparing our method against academic baselines, and (Right) a GSB pairwise comparison against leading proprietary models.

inconsistency (LSI) and motion unnaturalness (MU). These subjective advantages are reflected in the overall win rate (GSB score), collectively validating the effectiveness of our proposed method.

## 4.4  Comparison among Recent Methods

**Comparisons with State-of-the-Art Methods.** We conduct a comprehensive evaluation of our method against leading academic baselines across two distinct scenarios: portrait and full-body generation. For the portrait scenario, we compare our model against both specialized talking-head methods and state-of-the-art DiT-based approaches, including SadTalker [93], EchoMimic [8], Hallo [85], Hallo3 [10], Loopy [30] and OmniHuman-1 [40], on the CelebV-HQ test set. For the more challenging task of full-body synthesis, our evaluation on the CyberHost test set includes a strong suite of recent DiT-based models: Skyreels-A1 [14], FantasyTalking [74], OmniAvatar [15], MultiTalk [36] and OmniHuman-1 [40].

The quantitative results in Table 4 show our method consistently ranking in the top two across most metrics. In the portrait scenario, our model performs on par with the strong OmniHuman-1 baseline. We attribute this to the limited motion range in portrait videos, which challenges objective metrics in capturing subtle facial expressiveness. Our advantages become more pronounced in the full-body scenario. While leading in image quality and lip-sync, our model excels in generating dynamic, large-scale movements, evidenced by a high HKV score. Crucially, it achieves this without sacrificing local detail, maintaining a competitive HKC score. Taken together, these evaluation results demonstrate the clear advantages of our method over existing approaches.

To provide a more holistic assessment of perceptual quality, we conducted user studies to supplement our quantitative metrics. First, we performed a user preference evaluation against the top academic baselines from our full-body comparison, with results shown in Figure 5a. Additionally, we benchmarked our method against four leading proprietary models, which we categorized as either end-to-end (E2E) or two-stage (2S) systems (combining I2V and video dubbing). To comply with EULAs and avoid conflicts of interest, these models were anonymized as CS-A (E2E), CS-B (E2E), CS-C (2S), and CS-D (2S). The results of this comparison, shown in Figure 5b. These user studies demonstrate the superiority of our method, which is particularly evident in
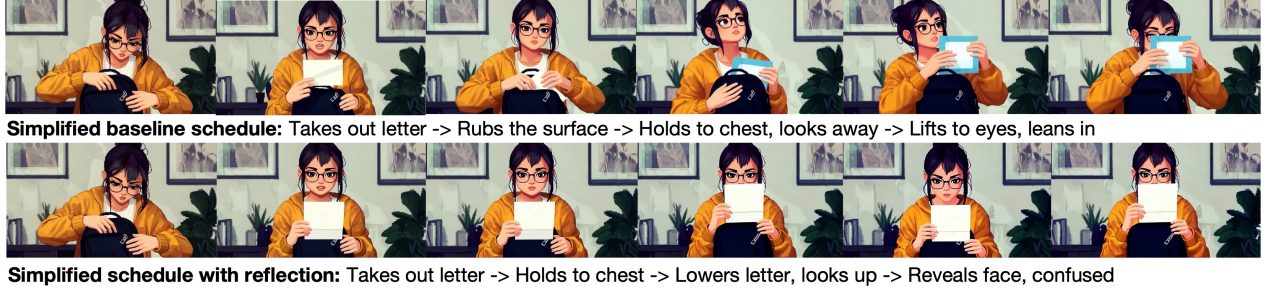
**Simplified baseline schedule:** Takes out letter -> Rubs the surface -> Holds to chest, looks away -> Lifts to eyes, leans in



**Simplified schedule with reflection:** Takes out letter -> Holds to chest -> Lowers letter, looks up -> Reveals face, confused

**Figure 6  Qualitative results of the reflection process.** Without reflection (first row), an ill-planned action ("Rubs the surface") causes object inconsistency. With reflection (second row), the model revises its plan to a more logical action, ensuring consistency.
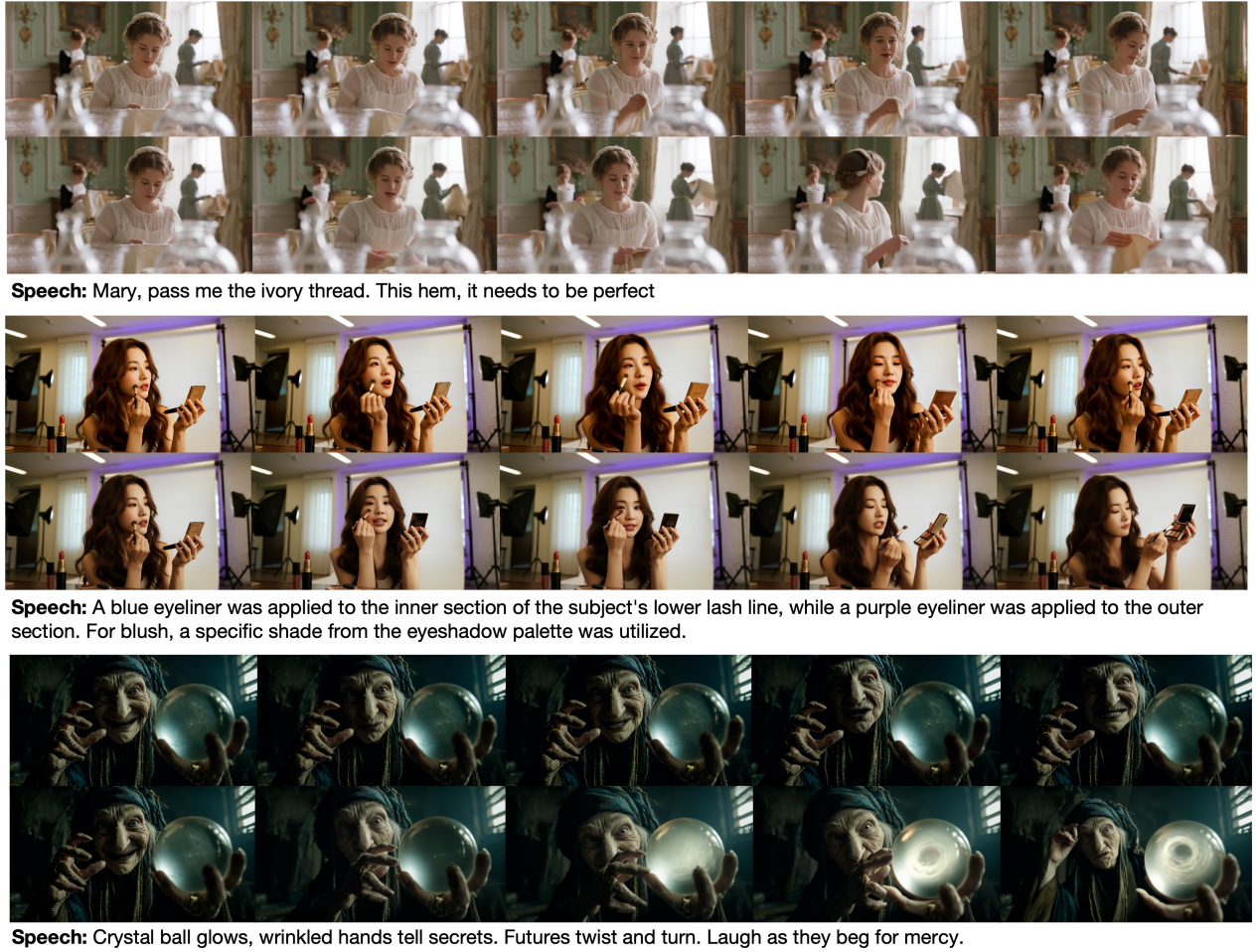


**Speech:** Mary, pass me the ivory thread. This hem, it needs to be perfect



**Speech:** A blue eyeliner was applied to the inner section of the subject's lower lash line, while a purple eyeliner was applied to the outer section. For blush, a specific shade from the eyeshadow palette was utilized.



**Speech:** Crystal ball glows, wrinkled hands tell secrets. Futures twist and turn. Laugh as they beg for mercy.

**Figure 7  Qualitative comparison of our model against OmniHuman-1 [40]** For each pair of examples, our model (bottom row) generates actions with higher semantic consistency to the speech prompt than the baseline (top row). For example, our model correctly depicts a character applying makeup and a glowing crystal ball as described in the speech, actions which are absent in the baseline's results.

its handling of contextual coherence, a factor to which human users are highly sensitive yet one that objective metrics often fail to capture.

## 4.5 Extended Experimental Results

**Ablation Study on the Reasoning Process.** We visualize the impact of our reflection process in Figure 6. This optional module is designed to correct errors from the initial action plan. Without reflection (top row), the model generates an action schedule in a single pass. This can lead to logical inconsistencies; for instance, after "Takes out letter," it generates "Rubs the surface," causing the letter to vanish and breaking semantic continuity. In contrast, our model with reflection (bottom row) revises the plan after generating the first segment. It observes the outcome of "Takes out letter" and corrects subsequent actions to be relevant to the theme of letter-reading, ensuring logical progression and mitigating error accumulation. However, as this reflection process introduces additional inference overhead, it was disabled for the quantitative comparisons in this paper. We also investigated injecting reasoning latents into the synthesis model. While this encouraged more nuanced facial expressions and subtle movements, it also suppressed large, dynamic actions. As this appeared to be more of an aesthetic trade-off than a clear improvement and was not a significant factor in user evaluations, we excluded it from our final model configuration.

**Visual Comparison with Baseline.** In Figure 7, we present a visual comparison with OmniHuman-1, a strong baseline identified in our preceding experiments. The corresponding speech content for each video is provided. As can be seen, our method demonstrates significantly stronger semantic relevance and logical correlation between the audio and the generated actions. For instance, in the first video, the character turns her head when calling out "Mary"; in the second, she performs the specific actions of applying eyeliner and gesturing towards the eyeshadow palette as described; and in the third, the crystal ball glows and changes in response to the wizard's incantation. These high-level, context-aware results are difficult to capture with objective metrics and remain a challenge for existing methods. This qualitative evidence further substantiates the superiority of our approach. We encourage readers to view the examples on our project page for a more intuitive understanding of the correlation between generated motion, speech content and character intent. This alignment is a key aspect often overlooked in prior work.

## 5 Conclusion

In this work, we introduced a new paradigm for human video generation inspired by the dual-system theory of human cognition. We argued that existing methods primarily simulate reactive "System 1" thinking, failing to align motion with high-level intent. We proposed OmniHuman-1.5, a framework that additionally models deliberative "System 2" processes through two key innovations: an MLLM-based agent for semantic planning and a specialized MMDiT architecture with a novel pseudo last frame strategy to fuse multimodal signals. Experiments show our approach generates more expressive and logically consistent results, which users significantly preferred for their naturalness and plausibility. By demonstrating this framework's effectiveness, even extending it to multi-person scenarios, we believe simulating cognitive agency offers a new perspective for creating the next generation of lifelike digital humans.

## 6 Broader Impact

Our core contribution in this work is to introduce a novel paradigm for video avatar generation. By simulating a dual-system cognitive framework, our model achieves a new level of expressive capability and logical coherence in motion, moving beyond the limitations of single-process generation. While this advancement opens up exciting possibilities for creative applications like AI-driven film production and music videos, we are acutely aware of the potential for misuse associated with highly realistic avatar technologies. To address these ethical concerns, we advocate for a robust framework of responsible deployment. Although current results may still bear subtle artifacts of AI generation, which can serve as a minor deterrent, proactive safeguards are essential. We strongly recommend the following measures: (1) applying prominent, visible watermarks to all generated content to clearly label it as AI-generated; (2) implementing filtering algorithms to reject inappropriate or malicious input prompts and to review output content; and (3) embedding traceable, invisible watermarks to ensure accountability and aid in source identification if misuse occurs. By integrating these safety protocols, we can help ensure that our technology fosters creativity while minimizing the risks of malicious applications such as fraud or disinformation.

# 7  Acknowledgment

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

[2] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 18392–18402, 2023.

[3] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. OpenAI Blog, 1(8):1, 2024.

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.

[5] Di Chang, Yichun Shi, Quankai Gao, Jessica Fu, Hongyi Xu, Guoxian Song, Qing Yan, Yizhe Zhu, Xiao Yang, and Mohammad Soleymani. Magicpose: Realistic human poses and facial expressions retargeting with identity-aware diffusion. arXiv preprint arXiv:2311.12052, 2023.

[6] Junsong Chen, Yue Wu, Simian Luo, Enze Xie, Sayak Paul, Ping Luo, Hang Zhao, and Zhenguo Li. Pixart-{\delta}: Fast and controllable image generation with latent consistency models. arXiv preprint arXiv:2401.05252, 2024.

[7] Shoufa Chen, Mengmeng Xu, Jiawei Ren, Yuren Cong, Sen He, Yanping Xie, Animesh Sinha, Ping Luo, Tao Xiang, and Juan-Manuel Perez-Rua. Gentron: Diffusion transformers for image and video generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6441–6451, 2024.

[8] Zhiyuan Chen, Jiajiong Cao, Zhiquan Chen, Yuming Li, and Chenguang Ma. Echomimic: Lifelike audio-driven portrait animations through editable landmark conditions. arXiv preprint arXiv:2407.08136, 2024.

[9] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In Computer Vision–ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part II 13, pages 251–263. Springer, 2017.

[10] Jiahao Cui, Hui Li, Yun Zhan, Hanlin Shang, Kaihui Cheng, Yuqi Ma, Shan Mu, Hang Zhou, Jingdong Wang, and Siyu Zhu. Hallo3: Highly dynamic and realistic portrait image animation with diffusion transformer networks. arXiv preprint arXiv:2412.00733, 2024.

[11] Xiang Deng, Youxin Pang, Xiaochen Zhao, Chao Xu, Lizhen Wang, Hongjiang Xiao, Shi Yan, Hongwen Zhang, and Yebin Liu. Stereo-talker: Audio-driven 3d human synthesis with prior-guided mixture-of-experts. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2025.

[12] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In Forty-first international conference on machine learning, 2024.

[13] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In Forty-first International Conference on Machine Learning, 2024.

[14] Zhengcong Fei, Debang Li, Di Qiu, Jiahua Wang, Yikun Dou, Rui Wang, Jingtao Xu, Mingyuan Fan, Guibin Chen, Yang Li, et al. Skyreels-a2: Compose anything in video diffusion transformers. arXiv preprint arXiv:2504.02436, 2025.

[15] Qijun Gan, Ruizi Yang, Jianke Zhu, Shaofei Xue, and Steven Hoi. Omniavatar: Efficient audio-driven avatar video generation with adaptive body animation. arXiv preprint arXiv:2506.18866, 2025.

[16] Yu Gao, Haoyuan Guo, Tuyen Hoang, Weilin Huang, Lu Jiang, Fangyuan Kong, Huixia Li, Jiashi Li, Liang Li, Xiaojie Li, et al. Seedance 1.0: Exploring the boundaries of video generation models. arXiv preprint arXiv:2506.09113, 2025.

[17] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025.

[18] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. arXiv preprint arXiv:2307.04725, 2023.

[19] Tianyu He, Junliang Guo, Runyi Yu, Yuchi Wang, Jialiang Zhu, Kaikai An, Leyi Li, Xu Tan, Chunyu Wang, Han Hu, et al. Gaia: Zero-shot talking avatar generation. arXiv preprint arXiv:2311.15230, 2023.

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017.

[21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.

[22] Steven Hogue, Chenxu Zhang, Hamza Daruger, Yapeng Tian, and Xiaohu Guo. Diffted: One-shot audio-driven ted talk video generation with diffusion-based co-speech gestures. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1922–1931, 2024.

[23] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. In The Twelfth International Conference on Learning Representations, 2023.

[24] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8153–8163, 2024.

[25] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8153–8163, 2024.

[26] Panwen Hu, Jin Jiang, Jianqi Chen, Mingfei Han, Shengcai Liao, Xiaojun Chang, and Xiaodan Liang. Storyagent: Customized storytelling video generation via multi-agent collaboration. arXiv preprint arXiv:2411.04925, 2024.

[27] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024.

[28] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. arXiv preprint arXiv:2412.16720, 2024.

[29] Jianwen Jiang, Gaojie Lin, Zhengkun Rong, Chao Liang, Yongming Zhu, Jiaqi Yang, and Tianyun Zhong. Mobileportrait: Real-time one-shot neural head avatars on mobile devices. arXiv preprint arXiv:2407.05712, 2024.

[30] Jianwen Jiang, Chao Liang, Jiaqi Yang, Gaojie Lin, Tianyun Zhong, and Yanbo Zheng. Loopy: Taming audio-driven portrait avatar with long-term motion dependency. In The Thirteenth International Conference on Learning Representations, 2025. URL https://openreview.net/forum?id=weM4YBicIP.

[31] Daniel Kahneman. Thinking, fast and slow. macmillan, 2011.

[32] Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In Handbook of the fundamentals of financial decision making: Part I, pages 99–127. World Scientific, 2013.

[33] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion image-to-video synthesis via stable diffusion. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pages 22623–22633. IEEE, 2023.

[34] Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal language models. Advances in Neural Information Processing Systems, 36:21487–21506, 2023.

[35] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. arXiv preprint arXiv:2412.03603, 2024.

[36] Zhe Kong, Feng Gao, Yong Zhang, Zhuoliang Kang, Xiaoming Wei, Xunliang Cai, Guanying Chen, and Wenhan Luo. Let them talk: Audio-driven multi-person conversational video generation. arXiv preprint arXiv:2505.22647, 2025.

[37] Yunxin Li, Haoyuan Shi, Baotian Hu, Longyue Wang, Jiashun Zhu, Jinyi Xu, Zhen Zhao, and Min Zhang. Anim-director: A large multimodal model powered agent for controllable animation video generation. In SIGGRAPH Asia 2024 Conference Papers, pages 1–11, 2024.

[38] Chao Liang, Jianwen Jiang, Wang Liao, Jiaqi Yang, Weihong Zeng, Han Liang, et al. Alignhuman: Improving motion and fidelity via timestep-segment preference optimization for audio-driven human animation. arXiv preprint arXiv:2506.11144, 2025.

[39] Gaojie Lin, Jianwen Jiang, Chao Liang, Tianyun Zhong, Jiaqi Yang, Zerong Zheng, and Yanbo Zheng. Cyberhost: A one-stage diffusion framework for audio-driven talking body generation. In The Thirteenth International Conference on Learning Representations, 2025. URL https://openreview.net/forum?id=vaEPihQsAA.

[40] Gaojie Lin, Jianwen Jiang, Jiaqi Yang, Zerong Zheng, and Chao Liang. Omnihuman-1: Rethinking the scaling-up of one-stage conditioned human animation models. arXiv preprint arXiv:2502.01061, 2025.

[41] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. arXiv preprint arXiv:2209.03003, 2022.

[42] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi Wan, Ranchen Ming, Xiaoniu Song, Xing Chen, et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation model. arXiv preprint arXiv:2502.10248, 2025.

[43] Willi Menapace, Aliaksandr Siarohin, Ivan Skorokhodov, Ekaterina Deyneka, Tsai-Shien Chen, Anil Kag, Yuwei Fang, Aleksei Stoliar, Elisa Ricci, Jian Ren, et al. Snap video: Scaled spatiotemporal transformers for text-to-video synthesis. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7038–7048, 2024.

[44] Rang Meng, Xingyu Zhang, Yuming Li, and Chenguang Ma. Echomimicv2: Towards striking, simplified, and semi-body human animation. In Proceedings of the Computer Vision and Pattern Recognition Conference, pages 5489–5498, 2025.

[45] Monica. Manus: General ai agent that bridges mind and action. Website, 2025. URL https://manus.im/. Accessed: [2025-07-16].

[46] OpenAI. Chatgpt (gpt-4o version). https://chat.openai.com/, 2024.

[47] OpenAI. Deep research system card. Technical Report, 2025. URL https://cdn.openai.com/deep-research-system-card.pdf. Accessed: [2025-07-16].

[48] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744, 2022.

[49] Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang Wang, Zhiyang Xu, Jiuhai Chen, Kunpeng Li, Felix Juefei-Xu, et al. Transfer between modalities with metaqueries. arXiv preprint arXiv:2504.06256, 2025.

[50] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In Proceedings of the 36th annual acm symposium on user interface software and technology, pages 1–22, 2023.

[51] William Peebles and Saining Xie. Scalable diffusion models with transformers. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 4195–4205, 2023.

[52] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. arXiv preprint arXiv:2410.13720, 2024.

[53] Di Qiu, Zhengcong Fei, Rui Wang, Jialin Bai, Changqian Yu, Mingyuan Fan, Guibin Chen, and Xiang Wen. Skyreels-a1: Expressive portrait animation in video diffusion transformers. arXiv preprint arXiv:2502.10841, 2025.

[54] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. Advances in Neural Information Processing Systems, 36:68539–68551, 2023.

[55] Team Seawead, Ceyuan Yang, Zhijie Lin, Yang Zhao, Shanchuan Lin, Zhibei Ma, Haoyuan Guo, Hao Chen, Lu Qi, Sen Wang, et al. Seaweed-7b: Cost-effective training of video generation foundation model. arXiv preprint arXiv:2504.08685, 2025.

[56] Ruizhi Shao, Youxin Pang, Zerong Zheng, Jingxiang Sun, and Yebin Liu. Human4dit: Free-view human video generation with 4d diffusion transformer. arXiv preprint arXiv:2405.17405, 2024.

[57] Weijia Shi, Xiaochuang Han, Chunting Zhou, Weixin Liang, Xi Victoria Lin, Luke Zettlemoyer, and Lili Yu. Lmfusion: Adapting pretrained language models for multimodal generation. arXiv preprint arXiv:2412.15188, 2024.

[58] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. Advances in neural information processing systems, 32, 2019.

[59] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13653–13662, 2021.

[60] Significant-Gravitas. Autogpt: Build, deploy, and run ai agents. GitHub Repository, 2023. URL https://github.com/Significant-Gravitas/AutoGPT. Accessed: [2025-07-16].

[61] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502, 2020.

[62] Michal Stypulkowski, Konstantinos Vougioukas, Sen He, Maciej Zieba, Stavros Petridis, and Maja Pantic. Diffused heads: Diffusion models beat gans on talking-face generation. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 5091–5100, 2024.

[63] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. Neurocomputing, 568:127063, 2024.

[64] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.

[65] Linrui Tian, Siqi Hu, Qi Wang, Bang Zhang, and Liefeng Bo. Emo2: End-effector guided audio-driven avatar video generation. arXiv preprint arXiv:2501.10687, 2025.

[66] Linrui Tian, Siqi Hu, Qi Wang, Bang Zhang, and Liefeng Bo. Emo2: End-effector guided audio-driven avatar video generation. arXiv preprint arXiv:2501.10687, 2025.

[67] Linrui Tian, Qi Wang, Bang Zhang, and Liefeng 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, 2025.

[68] Shuyuan Tu, Qi Dai, Zihao Zhang, Sicheng Xie, Zhi-Qi Cheng, Chong Luo, Xintong Han, Zuxuan Wu, and Yu-Gang Jiang. Motionfollower: Editing video motion via lightweight score-guided diffusion. arXiv preprint arXiv:2405.20325, 2024.

[69] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. https://openreview.net/forum?id=rylgEULtdN, 2019.

[70] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. arXiv preprint arXiv:2503.20314, 2025.

[71] Cong Wang, Kuan Tian, Jun Zhang, Yonghang Guan, Feng Luo, Fei Shen, Zhiwei Jiang, Qing Gu, Xiao Han, and Wei Yang. V-express: Conditional dropout for progressive training of portrait video generation. arXiv preprint arXiv:2406.02511, 2024.

[72] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291, 2023.

[73] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. arXiv preprint arXiv:2308.06571, 2023.

[74] Mengchao Wang, Qiang Wang, Fan Jiang, Yaqi Fan, Yunpeng Zhang, Yonggang Qi, Kun Zhao, and Mu Xu. Fantasytalking: Realistic talking portrait generation via coherent motion synthesis. arXiv preprint arXiv:2504.04842, 2025.

[75] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 10039–10049, 2021.

[76] Xiang Wang, Shiwei Zhang, Changxin Gao, Jiayu Wang, Xiaoqiang Zhou, Yingya Zhang, Luxin Yan, and Nong Sang. Unianimate: Taming unified video diffusion models for consistent human image animation. arXiv preprint arXiv:2406.01188, 2024.

[77] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In The Eleventh International Conference on Learning Representations.

[78] Zhenzhi Wang, Jiaqi Yang, Jianwen Jiang, Chao Liang, Gaojie Lin, Zerong Zheng, Ceyuan Yang, and Dahua Lin. Interacthuman: Multi-concept human animation with layout-aligned audio conditions. arXiv preprint arXiv:2506.09984, 2025.

[79] Huawei Wei, Zejun Yang, and Zhisheng Wang. Aniportrait: Audio-driven synthesis of photorealistic portrait animation. arXiv preprint arXiv:2403.17694, 2024.

[80] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.

[81] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. arXiv preprint arXiv:2508.02324, 2025.

[82] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, et al. Q-align: Teaching lmms for visual scoring via discrete text-defined levels. arXiv preprint arXiv:2312.17090, 2023.

[83] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. In Forty-first International Conference on Machine Learning, 2024.

[84] Weijia Wu, Zeyu Zhu, and Mike Zheng Shou. Automated movie generation via multi-agent cot planning. arXiv preprint arXiv:2503.07314, 2025.

[85] Mingwang Xu, Hui Li, Qingkun Su, Hanlin Shang, Liwei Zhang, Ce Liu, Jingdong Wang, Luc Van Gool, Yao Yao, and Siyu Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation. arXiv preprint arXiv:2406.08801, 2024.

[86] Sicheng Xu, Guojun Chen, Yu-Xiao Guo, Jiaolong Yang, Chong Li, Zhenyu Zang, Yizhong Zhang, Xin Tong, and Baining Guo. Vasa-1: Lifelike audio-driven talking faces generated in real time. arXiv preprint arXiv:2404.10667, 2024.

[87] Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using diffusion model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1481–1490, 2024.

[88] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. arXiv preprint arXiv:2408.06072, 2024.

[89] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances in neural information processing systems, 36:11809–11822, 2023.

[90] Lijun Yu, Jos Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion–tokenizer is key to visual generation. arXiv preprint arXiv:2310.05737, 2023.

[91] Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. arXiv preprint arXiv:2309.02591, 2023.

[92] Zhengqing Yuan, Yixin Liu, Yihan Cao, Weixiang Sun, Haolong Jia, Ruoxi Chen, Zhaoxu Li, Bin Lin, Li Yuan, Lifang He, et al. Mora: Enabling generalist video generation via a multi-agent framework. arXiv preprint arXiv:2403.13248, 2024.

[93] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8652–8661, 2023.

[94] Yuang Zhang, Jiaxi Gu, Li-Wen Wang, Han Wang, Junqi Cheng, Yuefeng Zhu, and Fangyuan Zou. Mimicmotion: High-quality human motion video generation with confidence-aware pose guidance. arXiv preprint arXiv:2406.19680, 2024.

[95] Jian Zhao and Hui Zhang. Thin-plate spline motion model for image animation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 3657–3666, 2022.

[96] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. arXiv preprint arXiv:2412.20404, 2024.

[97] Tianyun Zhong, Chao Liang, Jianwen Jiang, Gaojie Lin, Jiaqi Yang, and Zhou Zhao. Fada: Fast diffusion avatar synthesis with mixed-supervised multi-cfg distillation. arXiv preprint arXiv:2412.16915, 2024.

[98] Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and Chen Change Loy. Celebv-hq: A large-scale video facial attributes dataset. In European conference on computer vision, pages 650–667. Springer, 2022.

[99] Luyang Zhu, Dawei Yang, Tyler Zhu, Fitsum Reda, William Chan, Chitwan Saharia, Mohammad Norouzi, and Ira Kemelmacher-Shlizerman. Tryondiffusion: A tale of two unets. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4606–4615, 2023.

[100] Shaobin Zhuang, Kunchang Li, Xinyuan Chen, Yaohui Wang, Ziwei Liu, Yu Qiao, and Yali Wang. Vlogger: Make your dream a vlog. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8806–8817, 2024.