# StreamAvatar: Streaming Diffusion Models for Real-Time Interactive Human Avatars

Zhiyao Sun[1*†]   Ziqiao Peng[2*]   Yifeng Ma[3*]   Yi Chen[3]   Zhengguang Zhou[3]   Zixiang Zhou[3]
Guozhen Zhang[4]   Youliang Zhang[1]   Yuan Zhou[3‡]   Qinglin Lu[3§]   Yong-Jin Liu[1§]

[1]Tsinghua University   [2]Renmin University of China   [3]Tencent Hunyuan   [4]Nanjing University

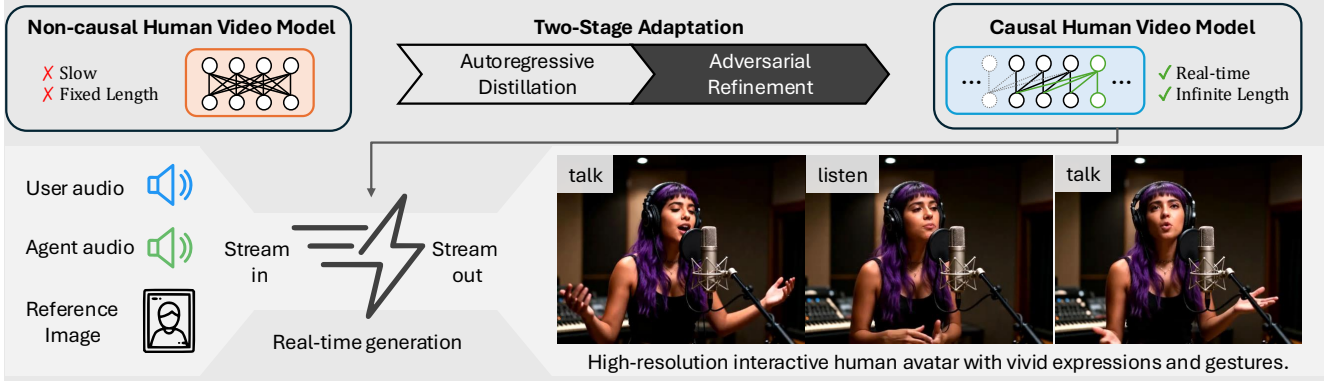Project Page: https://streamavatar.github.io

Figure 1. We propose StreamAvatar, which adapts human video diffusion models for real-time, streaming, and interactive video generation through a two-stage framework. Given a reference image and user/agent audio streams as input, StreamAvatar generates high-resolution streaming human video in real-time, producing vivid talking/listening expressions and gestures.

## Abstract

*Real-time, streaming interactive avatars represent a critical yet challenging goal in digital human research. Although diffusion-based human avatar generation methods achieve remarkable success, their non-causal architecture and high computational costs make them unsuitable for streaming. Moreover, existing interactive approaches are typically limited to head-and-shoulder region, limiting their ability to produce gestures and body motions. To address these challenges, we propose a two-stage autoregressive adaptation and acceleration framework that applies autoregressive distillation and adversarial refinement to adapt a high-fidelity human video diffusion model for real-time, interactive streaming. To ensure long-term stability and consistency, we introduce three key components: a Reference Sink, a Reference-Anchored Positional Re-encoding (RAPR) strategy, and a Consistency-Aware Discriminator. Building on this framework, we develop a one-shot, inter-active, human avatar model capable of generating both natural talking and listening behaviors with coherent gestures. Extensive experiments demonstrate that our method achieves state-of-the-art performance, surpassing existing approaches in generation quality, real-time efficiency, and interaction naturalness.*

## 1. Introduction

Human avatar generation holds substantial value, enabling a wide range of applications such as education, entertainment and virtual assistants. Real-time interactive human avatars offer even greater benefits by supporting fluid, dynamic communication with users. Diffusion models have achieved remarkable success in human avatar generation [4, 7, 9, 15, 19, 24, 35, 38]. However, in real-time interactive settings, diffusion-based human avatar generation methods still face three key challenges: *(1) Real-time streaming:* Most methods are inherently unsuitable for real-time streaming. Firstly, the iterative denoising process and long-context attention incur a prohibitive amount of computation. Secondly, the non-casual bidirectional attention

---

arXiv:2512.22065v1 [cs.CV] 26 Dec 2025

mechanism requires the entire video sequence to be processed at once, which is incompatible with streaming generation. *(2) Long-duration generation stability:* Streaming interaction requires generating long-duration videos; however, existing methods suffer from error accumulation and drifting when extended over time, resulting in degraded performance. *(3) Talking–listening interaction:* Current methods focus on *one-way talking* generation while neglecting the avatar's *listening* state. In conversational scenarios, not modeling the listening state makes the interaction feel unnatural. Although some studies explore listening video generation [22, 27, 39, 52] or interactive video generation [11, 55], they model only the facial or head region and fail to capture full-body expressiveness and reactivity.

To address these challenges, we propose StreamAvatar, which achieves real-time streaming interactive human video generation. StreamAvatar not only produces realistic results during speaking, but also generates rich and responsive listening motions conditioned on the interlocutor's audio. Specifically, we first train a powerful yet slow non-casual diffusion model for interactive human avatar generation as the teacher model, capable of producing both *talking and listening* behaviors. To achieve such interaction capability, we utilize a talking-listening audio mask to extract distinct talking and listening audio features, which are then injected through dedicated talking/listening modules of the model. This design allows the model to generate expressive co-speech or listening expressions and gestures.

Then, we propose a two-stage autoregressive adaptation and acceleration framework for transforming this teacher model into a real-time, streaming student model. In the *first* stage, we adapt the attention mechanism of the model from bidirectional attention into block-wise causal attention, and then perform distillation from the teacher model to obtain a few-step autoregressive generator, which significantly accelerates inference, reducing the DiT denoising process to 1/40 of its original runtime. To solve the critical challenges of long-video stability and identity preservation, we also introduce novel mechanisms for attention: a *Reference Sink*, which enforces persistent attention to the reference frame, and a *Reference-Anchored Positional Re-encoding (RAPR)* strategy, which resolves the training-inference mismatch and attention decay in long-sequence inference. In the *second* stage, we propose an adversarial refinement process to resolve the quality degradation such as distortion and blur caused by distillation, and improve long video consistency. Specifically, we introduce a *multitask consistency-aware discriminator*, which assesses both the realism of local frames and the consistency over all generated frames. The adversarial refinement stage significantly improves the generation quality and stability of our causal generator.

Our main contributions are summarized as follows:

- We introduce a two-stage autoregressive adaptation and acceleration framework for human video diffusion models. The first stage performs autoregressive distillation to convert a bidirectional diffusion models into real-time, streaming generator. The second stage performs adversarial refinement, which significantly improves generation quality and stability.
- To improve consistency and stability for long video generation, we propose three techniques, namely Reference Sink, Reference-Anchored Positional Re-encoding, and a consistency-aware discriminator.
- Building on this framework, we develop StreamAvatar, a one-shot, real-time, streaming interactive human video model that can generate both talking and listening behaviors with coherent expressions, gestures, and transitions.
- Extensive experiments demonstrates that our method outperforms or performs competitively with existing state-of-the-art approaches, while operating significantly faster.

## 2. Related Work

### 2.1. Audio-Driven Avatar Video Generation

**Traditional and 3D-Based Methods.** Early efforts relies on image-to-image translation [30, 51] or image warping [54] to generate lip-synced videos, but are often limited to the mouth region and fail to produce natural head motion. Later works, such as SadTalker [48] and VASA-1 [42], map audio to intermediate representations (e.g., 3DMMs or latent codes) and render videos with lightweight decoders, improving fidelity but still limited by the expressiveness of these intermediates. Another line of method directly model the subject with audio-conditioned 3D representations such as Neural Radiance Fields [10, 44] or 3D Gaussian Splatting [29]. However, they typically require per-subject training, making them unable to generalize in a one-shot manner.

**Diffusion-Based Avatars.** Recently, diffusion models have set a new standard for generation quality and one-shot generalization in the avatar space. A large body of work has focused on integrating audio features into large-scale video diffusion models to generate high-fidelity, expressive portraits [5, 7, 15, 33, 38] and semi/full-body videos [4, 9, 18, 19, 24, 25, 35]. While these methods produce state-of-the-art visual quality, their reliance on an iterative, bidirectional denoising process makes them computationally prohibitive and fundamentally unsuitable for real-time, streaming applications. Our work bridges this gap, targeting the generation of high-fidelity human avatars in real-time.

### 2.2. Streaming Video Diffusion Models

A large number of video diffusion models are limited to fixed-length generation due to their use of bidirectional attention. Common workarounds, such as motion or overlapping frames, often create unnatural transitions, identity drift, and cumulative errors. Moreover, these meth-

ods cannot reduce generation latency. To enable efficient, streaming generation, one prominent line of work focuses on distilling bidirectional models into few-step causal autoregressive systems. CausVid [46] pioneered this by re-architecting a bidirectional DiT with block-causal attention, using Diffusion Forcing (DF) [2] and distribution matching distillation (DMD) [45]. Self-Forcing [14] later improved upon DF by identifying a critical train-test mismatch and introducing a "student-forcing" scheme, where the student model conditions on its own prior outputs. Seaweed-APT2 [21] integrates adversarial post-training into autoregressive distillation to further boost generation quality.

However, a key "train-short-test-long" challenge persists for these autoregressive methods: quality degrades when extrapolating to sequences far beyond the training horizon. LongLive [43] and Self-Forcing++ [8] both address this by applying DMD to segments sampled from long videos generated by the student model. In contrast, we identify that this extrapolation failure also stems from an out-of-distribution (OOD) issue induced by Rotary Positional Embeddings (RoPE). We propose a re-encoding mechanism to resolve this instability without the need of generating long videos. Similar to LongLive, we also employ an attention sink to mitigate identity drift.

## 2.3. Interactive and Responsive Video Avatars

Natural human interaction is a dyadic process of speaking and listening. Recent work on avatar generation [3, 23, 47] achieves high-speed performance and satisfies the real-time requirements of interactive scenarios. However, these models process only talking audio and therefore cannot produce listener responses that are temporally aligned with the interlocutor's speech. Research dedicated to listener behaviors [22, 27, 34, 39, 52] focuses on the listening state and can generate listener responses aligned with the interlocutor's audio, yet these methods do not model how an avatar should transition smoothly between speaking and listening. Methods targeting unified speaking–listening generation [11, 41, 55] are capable of responding to the interlocutor's audio and exhibiting natural transitions between speaking and listening, but their motion space is largely restricted to the head-and-shoulder region, limiting their ability to produce rich hand and body motions that are essential for expressive interaction.

Our approach is designed to overcome these limitations: it runs in real time, generates motion beyond the head-and-shoulder region, produces listener responses synchronized to the interlocutor's audio, and supports smooth, unified transitions between speaking and listening.

## 3. Method

In this section, we first introduce our two-stage autoregressive adaptation and acceleration framework (Sec. 3.1), and

then develop a real-time, streaming interactive human video model (Sec. 3.2) based on this framework.

## 3.1. Streaming Human Video Diffusion Model

Recent human video generation approaches are generally built on video models constructed with diffusion transformers [28]. They typically adopt bidirectional attention, and take a reference image $x_0^0$, a text prompt $y$, and an audio clip $\{a_t\}$ as input conditions. In the following, we describe our two-stage autoregressive adaptation and acceleration framework using this class of models as a representative example.

### 3.1.1. Stage 1: Autoregressive Distillation.

This stage has two goals: (1) re-architect the model for autoregressive generation, and (2) distill its bidirectional, iterative denoising process into a causal, few-step model.

**Re-architecture.** The distillation process operates on the Diffusion Transformer (DiT) within the VAE latent space. The original DiT is bidirectional over a fixed generation window $\{x_t^n\}_{t=0}^T$, where $n = 0, 1, \ldots, N$ denotes the diffusion timestep and $T + 1$ is the window size. Inspired by prior works [14, 46], we split the generation window into smaller chunks: a clean 1-frame Reference chunk ($\{x_0^0\}$) and subsequent Generation chunks $\{\{x_t^n\}_{t=s_i}^{e_i}\}_{i=1}^{T/C}$ of size $C$, where $s_i = (i - 1) \cdot C + 1$ and $e_i = s_i \cdot C$. See Fig. 2 for an example with $C = 3$. We enforce causal attention between chunks and bidirectional attention within chunks, allowing the model to better capture local dynamics while enabling autoregression. To allow for efficient inference for long videos, a rolling KV cache is adopted to store a fixed window of context information. This modification preserves the original network weights, providing a strong starting point for distillation.

**Distillation Pipeline.** We first follow CausVid [46] and introduce an initialization stage to stabilize subsequent distillation. Specifically, we use the teacher model to generate avatar videos from a small set of avatar images and audio clips and record the corresponding denoising trajectories to construct a dataset of Ordinary Differential Equation (ODE) solution pairs $(\{x_t^n\}, \{x_t^0\})$. Since our goal is to distill a *few-step* autoregressive student generator, we retain only the timesteps used by the student model and train it to predict $\{x_t^0\}$ from $\{x_t^n\}$ for each chunk using a regression loss. Next, we apply Score Identity Distillation (SiD) [53] to train the student model by distilling from the teacher. Importantly, we adopt the student-forcing scheme from Self Forcing [14], where the student model predicts the next chunk $\{\hat{x}_t\}_{t=s_j}^{e_j}$ conditioned on its own output from the previous chunks $\{\{\hat{x}_t\}_{t=s_i}^{e_i}\}_{i<j}$, thereby mitigating the training–test mismatch. Note that in the original Self Forcing, the next chunk's denoising process is conditioned on clean (fully denoised) previous chunks $\{\hat{x}_t^0\}$, which requires an additional forward pass to update the KV-cache after denoising. How-
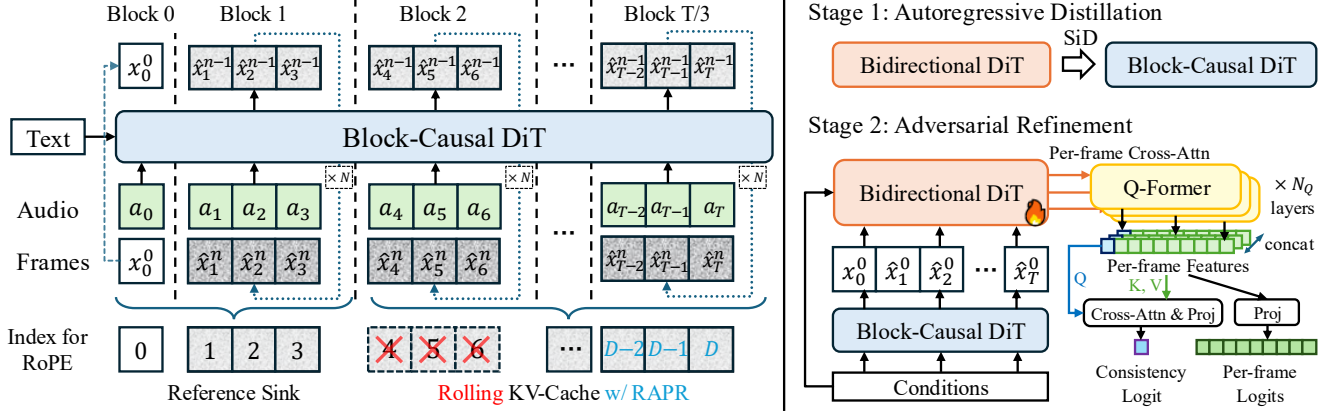
Figure 2. Overview of the two-stage autoregressive adaptation and acceleration framework. The original bidirectional DiT is first transformed into a block-causal DiT with block size $C = 3$. Then, in stage 1, we apply Score Identity Distillation to distill from the bidirectional teacher into a block-causal student. A Reference Sink and Reference-Anchored Positional Re-encoding is introduced to improve long-term stability and consistency. In stage 2, we apply an adversarial refinement process guided by a Consistency-Aware Discriminator, to further improve generation quality, consistency, and stability.

ever, we empirically find that omitting this update step does not noticeably degrade generation quality. This means the next chunk can instead be conditioned on noisy previous chunks $\{\hat{x}_t^1\}$, which reduces one forward pass per chunk and improves efficiency.

**Reference Sink.** To avoid the context window growing infinitely, prior works [14, 46] introduce a rolling KV cache. However, as new chunks are generated, the KV pairs for the reference frame are eventually evicted, leading to severe identity drift in long videos. Therefore, we introduce an attention sink to the KV cache, where the KV pairs of the reference frame $x_0$ are permanently held and never evicted. This ensures the model can always attend to the original identity. We also find that extending the sink to include the first generation chunk $\{x_t\}_{t=s_1}^{e_1}$ further improves consistency.

**Reference-Anchored Positional Re-encoding.** However, adding Reference Sink alone is insufficient for maintaining consistency in long video generation. We identify that this limitation primarily arises from two critical issues induced by the standard Rotary Positional Embedding (RoPE) [31]: 1) *Training-test Mismatch*: Vanilla RoPE assigns global frame indices as positional indices. As the model is trained only on short clips (e.g., $T$ frames), it never encounters the large positional indices required for long-duration streaming, leading to catastrophic out-of-distribution (OOD) issues. 2) *Attention Decay*: The long-term decay property inherent in RoPE causes the attention scores toward the Reference Sink to diminish as the generation window moves further away, exacerbating identity drift. To address this, we propose *Reference-Anchored Positional Re-encoding (RAPR)* (See Fig. 3). RAPR changes how positional indices are managed in the KV cache. The mecha-

nism is as follows: 1) We store *non-encoded* keys in the KV cache. 2) When generating the current frame $x_t$, we calculate its capped distance to the reference $x_0$ with a maximum limit $D$ (i.e., $\min(t, D)$), which serves as its RoPE index. 3) We synchronously shift the indices of all other cached keys to maintain their correct relative positions based on this capped distance. RoPE is then applied to these re-calculated positions. RAPR provides two crucial benefits. First, by capping the maximum distance $D$, it prevents attention decay for distant frames, ensuring the model consistently attends to the reference. Second, it mitigates the training-test mismatch problem. By enabling RAPR during both training and inference, the model learns to operate within a finite positional-encoding space (defined by $D < T$). This allows the model to simulate long-video positional shifts during training using only short clips, guaranteeing stability during extended inference.

### 3.1.2. Stage 2: Adversarial Refinement

After the Stage 1 distillation, we obtain a few-step causal generator capable of real-time, streaming generation. However, due to the aggressive reduction in diffusion steps and architectural modifications, the distilled model may exhibit visual artifacts (e.g., blurring in hands or teeth) and temporal inconsistencies. To address these issues, we introduce an adversarial refinement stage featuring a novel *consistency-aware* discriminator (see Fig. 3). We first follow prior works [20] to initialize the discriminator from the pretrained teacher model's backbone, and insert $N_Q$ Querying Transformers (Q-Formers) [17] with learnable *per-frame* queries into the intermediate layers to extract deep features for each frame. The discriminator's logit output adopts a dual-branch design: 1) *Local Realism Branch:* applies a linear
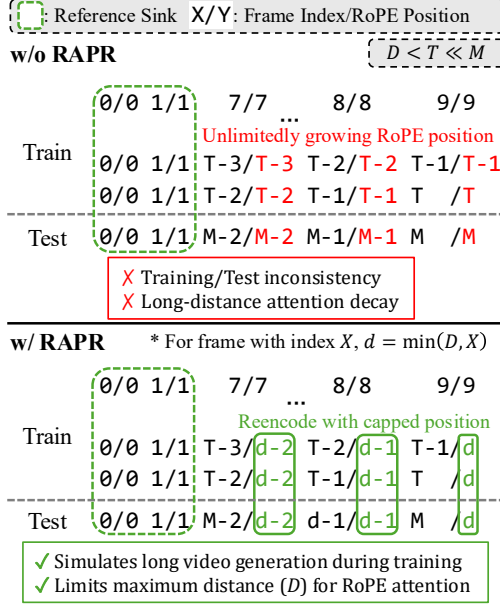
Figure 3. Vanilla RoPE vs RoPE with Reference-Anchored Positional Re-encoding (RAPR). RAPR improves long video generation without the need for long video training.
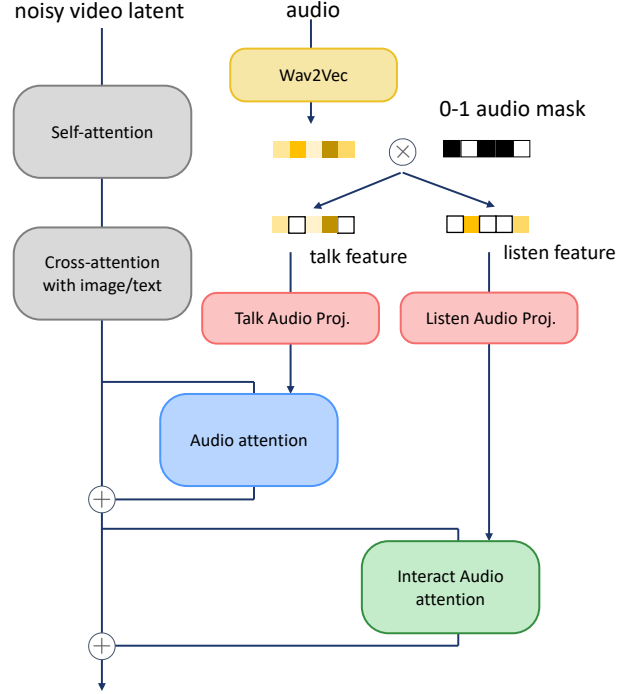


Figure 4. The architecture of the transformer block from our interactive human generation model. We extend the original block with audio-related attention modules to support talking and listening condition.

projection directly to the per-frame features, producing per-frame logits that assess the realism of individual generated frames; 2) *Global Consistency Branch:* enforces identity and temporal consistency by applying cross-attention between the reference frame's feature and those of all subsequent frames, followed by a linear projection to yield a single logit. This global branch explicitly penalizes undesired deviation from the reference. We train the student generator against this discriminator using a relativistic adversarial loss [16] and R1/R2 gradient penalty [26], as proposed in [13, 20]. Note that different from the distillation stage, the adversarial stage is trained with real video data, which directly optimize the generation distribution towards real distribution.

### 3.2. Interactive Human Generation

Our aim is to develop a real-time, streaming human video model that support *natural interactive human behaviors* in realistic scenes. We first train a bidirectional teacher model, and then distill it into a real-time causal model using our proposed framework. To start with, we adopt the Wan2.2-TI2V-5B [37] as the backbone, and extend its transformer block with two audio-related attention modules. Fig. 4 shows the extended transformer block architecture.

**Audio Mask for Interaction Phase Identification.** To distinguish between speaking and listening phases, we adopt an *audio mask* rather than *audio track separation* used in previous methods [55]. The audio mask is obtained through TalkNet [32], a joint audio–video detection method. We ob-

served that TalkNet provides more accurate temporal segmentation than audio separation methods. Audio track separation modifies waveforms, often producing signals that deviate from the data distribution used in Wav2Vec [1] pre-training. Such deviations degrade the quality of extracted Wav2Vec features. The audio mask avoids this issue by preserving the original waveform while marking each frame as either speaking (1) or listening (0). It is applied *after* Wav2Vec feature extraction, where it multiplicatively modulates the extracted representations rather than altering the raw waveform. This design maintains high-fidelity audio features while providing precise temporal control over speaking and listening intervals.

**Integration within the Generation Framework.** The masked Wav2Vec features are injected into the video latent through two audio-related attention modules:

- *Audio Attention*, which introduces talking cues to drive expressive human motion during speaking segments.
- *Interact Audio Attention*, which introduces listening cues to generate natural reactive behaviors during listening intervals.

All other layers in the framework, such as self-attention and cross-attention, remain audio-independent and focus on modeling visual and contextual dependencies. Note that our method disables text-based control and instead allows the

model to generate natural motions purely from the audio. Specifically, we fix the text prompt to *a person is speaking and listening.*

**Listen-and-Talk Generation Ability.** By combining the base video model's expressive generative capacity with phase-specific audio conditioning enabled by the audio mask, the system achieves *natural interactive human generation*. It produces smooth transitions between speaking and listening states and generates appropriate gestures and reactions, demonstrating realism and coherence far beyond traditional talk-only approaches.

## 4. Experiment

### 4.1. Experimental Setups

**Dataset.** Our data comes from a combination of SpeakerVid-5M [49] and self-collected videos. We sample and filter $\sim 200$ hours of 720P videos for training. To construct a balanced training set for both speaking and listening behaviors, we leverage the audio masks produced by TalkNet [32], which label each video frame as either speaking or listening. For every clip, we compute the ratio of listening frames in the detected audio mask. Clips with a high listening ratio are categorized as listening data, whereas clips with a low listening ratio are categorized as speaking data. By controlling the proportion of clips selected from each category, we ensure a balanced distribution of speaking and listening samples in the final dataset. In addition, we filter out samples with extreme head orientations by detecting the face pose and discarding videos in which the subject's head deviates excessively from a forward-facing direction.

**Implementation Details.** For bidirectional teacher model training, we finetune the base model for 20000 iterations with a batch size of 32 and learning rate of 5e-6. For the causal student model, we adopt the following set of parameters: denoising steps $N = 3$, chunk size $C = 3$, maximum distance for RAPR $D = 10$, and the number of Q-Formers $N_Q = 3$. The attention window is comprised of a 4-frame reference sink and a 6-frame rolling KV-cache. In Stage 1, we train the ODE initialization for 5000 iterations with a batch size of 8 and learning rate of 2e-6, followed by SiD distillation for 6000 iterations with a batch size of 16 and learning rate of 3e-6. In Stage 2, the adversarial refinement is trained for 1400 iterations with a batch size of 32 and learning rate of 5e-6. When streamlining the DiT denosing and VAE decoding on two H800 GPUs, our model achieves *real-time* generation with a 1.2-second latency.

**Evaluation Metrics.** We comprehensively evaluate our model's performance using multiple metrics. The Fréchet Inception Distance (FID) [12] and Fréchet Video Distance (FVD) [36] assess the distributional similarity between generated and real videos. Q-Align [40] evaluates image quality (IQA) and aesthetic score (ASE), Sync-C and Sync-D [6] measure audio–lip synchronization, while Hand Keypoint Variances (HKV) [25] quantifies gesture dynamics. We also adopt the Human Anomaly (HA) score from VBench-2.0 [50] to assess distortion in body, hands, and faces. To evaluate the motion richness during listening, we compute the variances of body, hand, and face keypoints observed during the speaker's listening phase, denoted as LBKV, LHKV, and LFKV, respectively.

**Compared Baselines.** We compare our model with state-of-the-art avatar video generation methods on both talking and interactive tasks. For talking video generation, we evaluate against Hallo3 [7], HunyuanVideo-Avatar (HY-Avatar) [4], OmniAvatar [9], EchoMimicV3 [24], and StableAvatar [35]. To the best of our knowledge, existing full-body avatar generation methods primarily focus on the talking phase and have not yet achieved unified 2D full-body listening–talking video generation. Therefore, to evaluate the effectiveness of our method in the listening phase, we compare it with a representative talking-only baseline, OmniAvatar. When generating test samples, feeding multi-speaker audio into a talking-only method would cause it to produce mouth movements even during the listening segments. Therefore, when generating samples with OmniAvatar, we mute the test audio during the listening periods.

### 4.2. Qualitative Results

For talking avatar video generation, qualitative comparisons are shown in Fig. 5. EchoMimicV3 and Hallo3 exhibit noticeable distortions, particularly in the hand regions, and suffer from error accumulation that leads to identity drift over long sequences. OmniAvatar produces stable results with fewer distortions but shows the least motion dynamics, largely preserving the posture of the reference frame. HunyuanVideo-Avatar and StableAvatar generate more dynamic motions but occasionally introduce distortions; StableAvatar further produces artifacts such as subtitles and jitters at sliding-window boundaries. In contrast, our method achieves realistic and consistent avatars with rich motion dynamics while exhibiting fewer distortions.

For interactive avatar generation, qualitative results are shown in Fig. 6. During the talking phase, our method produces accurate lip movements accompanied by vivid gestures, while during the listening phase, it responds naturally to auditory cues. It also generates smooth and realistic transitions between talking and listening. Compared to the baseline, which remains largely static during listening, our method exhibits greater expressiveness and realism.

Please refer to our appendix and demo video for more comparisons and results.
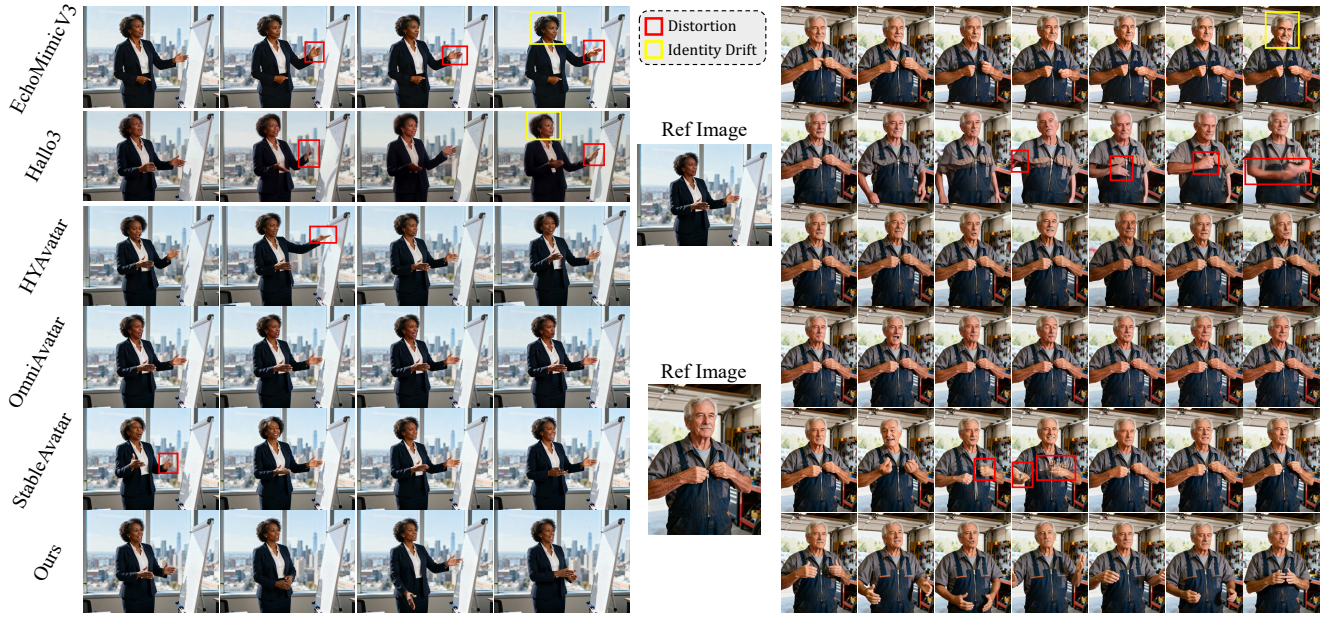
Figure 5. **Qualitative comparison with SoTA talking avatar video generation methods.** Please zoom in for details.



Figure 6. **Qualitative comparison with the baseline on interactive avatar video generation.** While the baseline remains almost static during the listening phase, our avatar reacts naturally to the listening audio, showing fluid and realistic transitions between talking and listening. Notice how it responds to laughter in the first example and how its expressions and gestures shift seamlessly between phases.

Table 1. Quantitative comparison with SoTA talking avatar video generation methods. Metrics are reported on both short and long datasets, separated by "/". Best in **bold** and second best underlined. FID/FVD for the long dataset is not applicable because the dataset is synthesized. We also report the number of denoising steps, resolution, and time taken to produce a 5-second video on a single H20 GPU for reference.

| Method | FID | FVD | ASE | IQA | Sync-C | Sync-D | HKV | HA | Steps | Res. | Time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| StableAvatar | 75.20/ - | 603.54/ - | 4.66/4.87 | 3.02/3.83 | 4.24/2.90 | 10.92/11.57 | 42.92/38.06 | 0.909/0.963 | 40 | 480p | 12min |
| OmniAvatar | 87.24/ - | 851.93/ - | 4.45/4.66 | 2.85/3.74 | **7.60**/6.47 | **7.99**/8.98 | 8.64/15.17 | **0.974**/0.974 | 25 | 480p | 36min |
| HY-Avatar | 76.49/ - | **557.46**/ - | 4.67/4.80 | 3.00/3.78 | 6.71/5.98 | 8.79/8.58 | 54.31/80.46 | 0.947/0.975 | 50 | 720p | 74min |
| Hallo3 | 117.41/ - | 1009.27/ - | 4.36/4.68 | 2.79/3.58 | 5.62/4.34 | 9.70/9.61 | 35.69/59.15 | 0.958/0.959 | 51 | 480p | 32min |
| EchoMimicV3 | 78.65/ - | 724.29/ - | 4.66/4.87 | 3.02/3.90 | 3.10/2.20 | 10.63/11.64 | 25.53/41.31 | 0.969/0.937 | 25 | 480p | 7min |
| Ours (baseline) | 96.58/ - | 885.97/ - | 4.29/4.01 | 2.75/3.13 | 7.04/5.99 | 8.34/8.36 | 67.65/111.69 | 0.948/0.971 | | | |
| +ref sink | 88.75/ - | 772.10/ - | 4.55/4.64 | 2.93/3.64 | 7.03/5.82 | 8.39/8.27 | 75.94/102.16 | 0.950/0.973 | | | |
| +RAPR | 81.63/ - | 753.46/ - | 4.64/4.84 | 2.98/3.81 | 7.06/5.99 | 8.34/8.18 | 48.29/61.07 | 0.956/0.961 | | | |
| +GAN w/o $D_{CA}$ | 79.68/ - | 741.30/ - | 4.65/4.91 | 3.02/3.80 | 7.05/6.38 | 8.31/8.15 | 35.50/29.85 | 0.947/0.993 | | | |
| **Ours** | **74.21**/ - | 707.34/ - | **4.68/4.94** | **3.03/3.91** | 7.06/**6.64** | 8.21/**8.14** | 48.35/65.77 | **0.974/0.993** | 3 | 720p | 20s |

Table 2. Quantitative results for interactive avatar generation.

| Method | LBKV | LHKV | LFKV |
|---|---|---|---|
| Baseline | 6.05 | 4.53 | 2.39 |
| Ours | **15.88** | **16.24** | **7.11** |

Table 3. Ablation on the audio mask position.

| Method | LBKV | LHKV | LFKV |
|---|---|---|---|
| Pre-Mask (undistilled) | 16.98 | 15.49 | 5.72 |
| Ours (undistilled) | **17.74** | **21.44** | **5.81** |

## 4.3. Quantitative Results

For talking video generation, we construct two evaluation datasets: 1) a short dataset consisting of 50 real avatar images paired with 5-second audio clips and corresponding ground-truth videos, and 2) a long dataset consisting of 25 synthesized avatar images paired with 20-second audio clips. Quantitative results for talking avatar video generation are reported in Table 1. Despite being only a 3-step causal model and requiring significant less time for generation, our method achieves superior or highly competitive performance across nearly all metrics, demonstrating strong visual quality and accurate lip synchronization. Notably, it achieves one of the largest motion magnitudes while maintaining a lower anomaly rate, which is challenging to accomplish simultaneously. We also report results on EMTD [25] in Appendix Sec. 6.1.

For interactive avatar generation, we sample 50 videos from SpeakerVid-5M for evaluation. Tab. 2 presents the comparison results. Our method significantly outperforms the talking-only baseline across all metrics, demonstrating its ability to produce rich motions during listening.

## 4.4. Ablation Study

We evaluate the key components of our autoregressive distillation framework and interactive model. Starting from Self Forcing [14] as the baseline, we incrementally add the Reference Sink (+ref sink), the RAPR strategy (+RAPR), the adversarial refinement stage with a standard discriminator (+GAN w/o $D_{CA}$), and finally replace it with our consistency-aware discriminator (Ours). As shown in

Tab. 1, adding each components consistently improves the performance. It should be noted that although some ablation variants exhibit larger motion dynamics, this often comes at the cost of distorted or blurry frames, as reflected by the ASE and IQA scores.

To evaluate the effectiveness of our design regarding where the audio mask is applied, we introduce a **Pre-Mask** variant, which applies the audio mask to the raw audio before feeding it into Wav2Vec for feature extraction. In contrast, **Ours** applies the audio mask to the features extracted by Wav2Vec. To reduce computational cost, we conduct experiments using the undistilled model. As shown in Tab. 3, **Ours** outperforms **Pre-Mask** across all metrics. This indicates that applying the audio mask directly to the audio degrades the resulting Wav2Vec features and ultimately harms model performance, thereby validating the effectiveness of our design. Please refer to the Appendix and demo video for more results.

## 5. Conclusion

In this paper, we address three major challenges in current human video avatars. For real-time streaming generation, we propose a two-stage autoregressive adaptation and acceleration framework with a distillation and an adversarial refinement process. For stable and consistent long video generation, we introduce the Reference Sink, Reference-Anchored Positional Re-encoding, and a consistency-aware discriminator. Finally, to enable natural interaction, we develop StreamAvatar, a real-time, streaming human video

model that generates both talking and listening behaviors with coherent expressions, gestures, and transitions. Extensive experiments demonstrate that StreamAvatar achieves state-of-the-art performance while operating orders of magnitude faster than competitive methods.

**Limitations and Future Work.** Despite impressive performance, our method exhibits certain limitations. Due to the limited temporal context, it may produce inconsistent content in regions that remain occluded for an extended period. Incorporating long-term memory mechanisms could help alleviate this issue. In addition, VAE decoding currently accounts for more than half of the total processing time in our pipeline, and future work could explore more efficient VAE decoding to further reduce the streaming latency.

# References

[1] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *NeurIPS*, 2020. 5

[2] Boyuan Chen, Diego Marti Monso, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. In *NeurIPS*, 2024. 3

[3] Ming Chen, Liyuan Cui, Wenyuan Zhang, Haoxian Zhang, Yan Zhou, Xiaohan Li, Songlin Tang, Jiwen Liu, Borui Liao, Hejia Chen, Xiaoqiang Liu, and Pengfei Wan. MIDAS: multimodal interactive digital-human synthesis via real-time autoregressive video generation. *CoRR*, abs/2508.19320, 2025. 3, 1

[4] Yi Chen, Sen Liang, Zixiang Zhou, Ziyao Huang, Yifeng Ma, Junshu Tang, Qin Lin, Yuan Zhou, and Qinglin Lu. Hunyuanvideo-avatar: High-fidelity audio-driven human animation for multiple characters. *CoRR*, abs/2505.20156, 2025. 1, 2, 6

[5] Zhiyuan Chen, Jiajiong Cao, Zhiquan Chen, Yuming Li, and Chenguang Ma. Echomimic: Lifelike audio-driven portrait animations through editable landmark conditions. In *AAAI*, pages 2403–2410. AAAI Press, 2025. 2

[6] Joon Son Chung and Andrew Zisserman. Out of time: Automated lip sync in the wild. In *ACCV Workshops (2)*, pages 251–263. Springer, 2016. 6

[7] 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 video diffusion transformer. In *CVPR*, pages 21086–21095. Computer Vision Foundation / IEEE, 2025. 1, 2, 6

[8] Justin Cui, Jie Wu, Ming Li, Tao Yang, Xiaojie Li, Rui Wang, Andrew Bai, Yuanhao Ban, and Cho-Jui Hsieh. Selfforcing++: Towards minute-scale high-quality video generation, 2025. 3

[9] Qijun Gan, Ruizi Yang, Jianke Zhu, Shaofei Xue, and Steven Hoi. Omniavatar: Efficient audio-driven avatar video generation with adaptive body animation. *CoRR*, abs/2506.18866, 2025. 1, 2, 6

[10] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, pages 5764–5774. IEEE, 2021. 2

[11] Ying Guo, Xi Liu, Cheng Zhen, Pengfei Yan, and Xiaoming Wei. ARIG: autoregressive interactive head generation for real-time conversations. *CoRR*, abs/2507.00472, 2025. 2, 3, 1

[12] 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. In *NIPS*, pages 6626–6637, 2017. 6

[13] Nick Huang, Aaron Gokaslan, Volodymyr Kuleshov, and James Tompkin. The GAN is dead; long live the gan! A modern GAN baseline. In *NeurIPS*, 2024. 5

[14] Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the train-test gap in autoregressive video diffusion. *CoRR*, abs/2506.08009, 2025. 3, 4, 8

[15] Xiaozhong Ji, Xiaobin Hu, Zhihong Xu, Junwei Zhu, Chuming Lin, Qingdong He, Jiangning Zhang, Donghao Luo, Yi Chen, Qin Lin, Qinglin Lu, and Chengjie Wang. Sonic: Shifting focus to global audio perception in portrait animation. In *CVPR*, pages 193–203. Computer Vision Foundation / IEEE, 2025. 1, 2

[16] Alexia Jolicoeur-Martineau. The relativistic discriminator: a key element missing from standard GAN. In *ICLR (Poster)*. OpenReview.net, 2019. 5

[17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742. PMLR, 2023. 4

[18] Gaojie Lin, Jianwen Jiang, Chao Liang, Tianyun Zhong, Jiaqi Yang, and Yanbo Zheng. Cyberhost: Taming audio-driven avatar diffusion model with region codebook attention. *CoRR*, abs/2409.01876, 2024. 2

[19] Gaojie Lin, Jianwen Jiang, Jiaqi Yang, Zerong Zheng, and Chao Liang. Omnihuman-1: Rethinking the scaling-up of one-stage conditioned human animation models. *CoRR*, abs/2502.01061, 2025. 1, 2

[20] Shanchuan Lin, Xin Xia, Yuxi Ren, Ceyuan Yang, Xuefeng Xiao, and Lu Jiang. Diffusion adversarial post-training for one-step video generation. *CoRR*, abs/2501.08316, 2025. 4, 5

[21] Shanchuan Lin, Ceyuan Yang, Hao He, Jianwen Jiang, Yuxi Ren, Xin Xia, Yang Zhao, Xuefeng Xiao, and Lu Jiang. Autoregressive adversarial post-training for real-time interactive video generation. *CoRR*, abs/2506.09350, 2025. 3

[22] Xi Liu, Ying Guo, Cheng Zhen, Tong Li, Yingying Ao, and Pengfei Yan. Customlistener: Text-guided responsive interaction for user-friendly listening head generation. In *CVPR*, pages 2415–2424. IEEE, 2024. 2, 3

[23] Chetwin Low and Weimin Wang. Talkingmachines: Realtime audio-driven facetime-style video via autoregressive diffusion models. *CoRR*, abs/2506.03099, 2025. 3

[24] Rang Meng, Yan Wang, Weipeng Wu, Ruobing Zheng, Yuming Li, and Chenguang Ma. Echomimicv3: 1.3b parameters are all you need for unified multi-modal and multi-task human animation. *CoRR*, abs/2507.03905, 2025. 1, 2, 6

[25] Rang Meng, Xingyu Zhang, Yuming Li, and Chenguang Ma. Echomimicv2: Towards striking, simplified, and semi-body human animation. In *CVPR*, pages 5489–5498. Computer Vision Foundation / IEEE, 2025. 2, 6, 8, 1

[26] Lars M. Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *ICML*, pages 3478–3487. PMLR, 2018. 5

[27] Evonne Ng, Hanbyul Joo, Liwen Hu, Hao Li, Trevor Darrell, Angjoo Kanazawa, and Shiry Ginosar. Learning to listen: Modeling non-deterministic dyadic facial motion. In *CVPR*, pages 20363–20373. IEEE, 2022. 2, 3

[28] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4172–4182. IEEE, 2023. 3

[29] Ziqiao Peng, Wentao Hu, Junyuan Ma, Xiangyu Zhu, Xiaomei Zhang, Hao Zhao, Hui Tian, Jun He, Hongyan Liu, and Zhaoxin Fan. Synctalk++: High-fidelity and efficient synchronized talking heads synthesis using gaussian splatting. *CoRR*, abs/2506.14742, 2025. 2

[30] K. R. Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C. V. Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *ACM Multimedia*, pages 484–492. ACM, 2020. 2

[31] Jianlin Su, Murtadha H. M. Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 4

[32] Ruijie Tao, Zexu Pan, Rohan Kumar Das, Xinyuan Qian, Mike Zheng Shou, and Haizhou Li. Is someone speaking? exploring long-term temporal features for audio-visual active speaker detection. In *ACMMM*. 5, 6

[33] 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 *ECCV (83)*, pages 244–260. Springer, 2024. 2

[34] Minh Tran, Di Chang, Maksim Siniukov, and Mohammad Soleymani. DIM: dyadic interaction modeling for social behavior generation. In *ECCV (37)*, pages 484–503. Springer, 2024. 3

[35] Shuyuan Tu, Yueming Pan, Yinming Huang, Xintong Han, Zhen Xing, Qi Dai, Chong Luo, Zuxuan Wu, and Yu-Gang Jiang. Stableavatar: Infinite-length audio-driven avatar video generation. *CoRR*, abs/2508.08248, 2025. 1, 2, 6

[36] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. FVD: A new metric for video generation. In *DGS@ICLR*. OpenReview.net, 2019. 6

[37] Team Wan, 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. 5

[38] 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. *CoRR*, abs/2504.04842, 2025. 1, 2

[39] Yinuo Wang, Yanbo Fan, Xuan Wang, Yu Guo, and Fei Wang. Diffusion-based realistic listening head generation via hybrid motion modeling. In *CVPR*, pages 15885–15895. Computer Vision Foundation / IEEE, 2025. 2, 3

[40] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, Qiong Yan, Xiongkuo Min, Guangtao Zhai, and Weisi Lin. Q-align: Teaching lmms for visual scoring via discrete text-defined levels. In *ICML*. OpenReview.net, 2024. 6

[41] You Xie, Tianpei Gu, Zenan Li, Chenxu Zhang, Guoxian Song, Xiaochen Zhao, Chao Liang, Jianwen Jiang, Hongyi Xu, and Linjie Luo. X-streamer: Unified human world modeling with audiovisual interaction. *CoRR*, abs/2509.21574, 2025. 3, 1

[42] 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. In *NeurIPS*, 2024. 2

[43] Shuai Yang, Wei Huang, Ruihang Chu, Yicheng Xiao, Yuyang Zhao, Xianbang Wang, Muyang Li, Enze Xie, Yingcong Chen, Yao Lu, Song Han, and Yukang Chen. Longlive: Real-time interactive long video generation. *CoRR*, abs/2509.22622, 2025. 3

[44] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, Jinzheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. In *ICLR*. OpenReview.net, 2023. 2

[45] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Frédo Durand, William T. Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *CVPR*, pages 6613–6623. IEEE, 2024. 3

[46] Tianwei Yin, Qiang Zhang, Richard Zhang, William T. Freeman, Frédo Durand, Eli Shechtman, and Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. In *CVPR*, pages 22963–22974. Computer Vision Foundation / IEEE, 2025. 3, 4

[47] Haojie Yu, Zhaonian Wang, Yihan Pan, Meng Cheng, Hao Yang, Chao Wang, Tao Xie, Xiaoming Xu, Xiaoming Wei, and Xunliang Cai. LLIA - enabling low-latency interactive avatars: Real-time audio-driven portrait video generation with diffusion models. *CoRR*, abs/2506.05806, 2025. 3

[48] 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 *CVPR*, pages 8652–8661. IEEE, 2023. 2

[49] Youliang Zhang, Zhaoyang Li, Duomin Wang, Jiahe Zhang, Deyu Zhou, Zixin Yin, Xili Dai, Gang Yu, and Xiu Li. Speakervid-5m: A large-scale high-quality dataset for audio-visual dyadic interactive human generation. *CoRR*, abs/2507.09862, 2025. 6

[50] Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Yuanhan Zhang, Jingwen He, Wei-Shi Zheng, Yu Qiao, and Ziwei Liu. Vbench-2.0: Advancing video generation benchmark suite for intrinsic faithfulness. *CoRR*, abs/2503.21755, 2025. 6

[51] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *CVPR*, pages 4176–4186. Computer Vision Foundation / IEEE, 2021. 2

[52] Mohan Zhou, Yalong Bai, Wei Zhang, Ting Yao, Tiejun Zhao, and Tao Mei. Responsive listening head generation: A benchmark dataset and baseline. In *ECCV (38)*, pages 124–142. Springer, 2022. 2, 3

[53] Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation. In *ICML*. OpenReview.net, 2024. 3

[54] Yang Zhou, Dingzeyu Li, Xintong Han, Evangelos Kalogerakis, Eli Shechtman, and Jose Echevarria. Makeittalk: Speaker-aware talking head animation. *CoRR*, abs/2004.12992, 2020. 2

[55] Yongming Zhu, Longhao Zhang, Zhengkun Rong, Tianshu Hu, Shuang Liang, and Zhipeng Ge. INFP: audio-driven interactive head generation in dyadic conversations. In *CVPR*, pages 10667–10677. Computer Vision Foundation / IEEE, 2025. 2, 3, 5, 1

# Appendix

## 6. Additional Experiments

### 6.1. Quantitative Evaluation on EMTD

To further evaluate our approach, we compare it with baseline methods and ablation variants on the EchoMimicV2 Testing Dataset (EMTD) [25]. The EMTD dataset contains 110 front-facing, half-body speech videos. Quantitative results are presented in Tab. 4. Our method outperforms all comparison methods across almost all metrics, and the ablation results further demonstrate the effectiveness of our design.

### 6.2. Qualitative Ablation Results

We present qualitative ablation results in Fig. 8 and in the demo video. The baseline model (Baseline) fails to preserve identity and drifts away over time. Introducing the reference sink (+Ref Sink) improves identity preservation but remains insufficient for long-duration generation. Incorporating the RAPR strategy (+RAPR) further enhances temporal consistency and stability in long sequences, though occasional failures and blur/distortion persist. Applying adversarial distillation without the consistency-aware discriminator results in worse consistency. As we add these components, the performance improves consistently, which validates the effectiveness of our design.

### 6.3. User Study

We conduct a user study to comprehensively evaluate our method. Participants are shown paired video clips generated by our approach and a comparison method, and asked to assess the two along five dimensions: audio–lip synchronization (Sync), motion dynamics (Dynamics), temporal continuity and smoothness (Continuity), visual quality and naturalness (Quality), and identity preservation (Identity). For each pair, participants indicate whether they prefer our method, prefer the comparison method, or have no preference. In total, we collect 960 paired comparisons from 24 participants. As illustrated in Fig. 7, our method consistently outperforms the state-of-the-art baselines across almost all comparisons, which aligns closely with our quantitative evaluation.

### 6.4. Comparison with Interactive Head Generation

We compare our method with state-of-the-art interactive head generation methods, including INFP [55] and ARIG [11]. As these methods are not open-sourced, we conduct qualitative comparison with the results from their project pages, as shown in Fig. 9 and the demo video. Although our model is designed and trained for body video generation, it still performs on par with these dedicated head avatar methods, while delivering best visual quality.
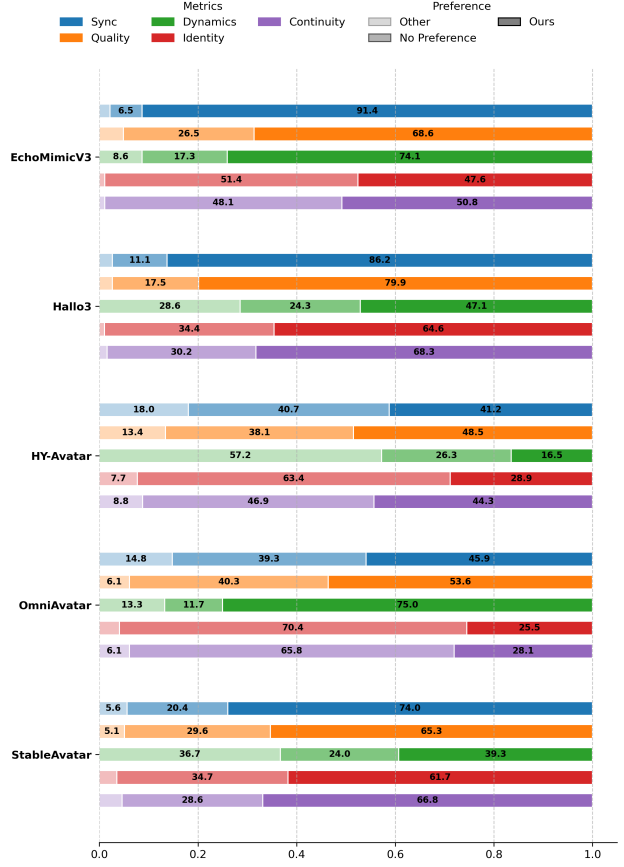


Figure 7. **User Study:** Ours vs all comparison methods.

### 6.5. Comparison with Streaming Interactive Avatar Generation

We further compare our method with current state-of-the-art streaming interactive avatar generation methods, including MIDAS [3] and X-Streamer [41]. As these methods are not open-sourced, we conduct qualitative comparison with the results from their project pages, as shown in Fig. 10 and the demo video. Our method produces more accurate lip synchronization and more vivid expressions than MIDAS. It is also worth noting that our method is *one-shot*, whereas MIDAS requires person-specific finetuning. Our method also exhibits more natural listening behaviors, more diverse motions, and higher visual quality than X-Streamer.

## 7. Runtime Details

Our model generates videos of 25 FPS. To enable real-time generation, we distribute the DiT denoising and VAE de-

Table 4. Quantitative comparison with SoTA talking avatar video generation methods on the EMTD dataset. Best in **bold** and second best underlined.

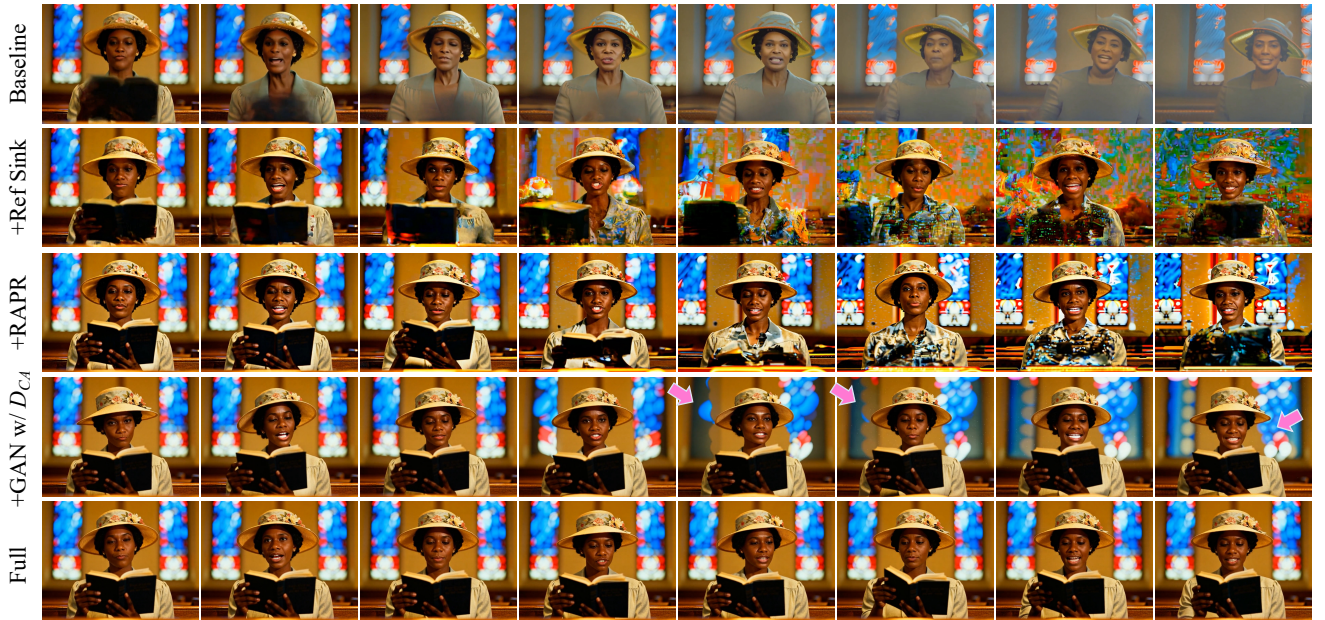| Method | FID | FVD | ASE | IQA | Sync-C | Sync-D | HKV | HA |
|---|---|---|---|---|---|---|---|---|
| StableAvatar | 91.63 | 840.86 | 3.67 | 2.37 | 3.04 | 12.16 | 57.99 | 0.794 |
| OmniAvatar | 75.20 | 982.09 | 3.72 | 2.45 | <u>7.68</u> | <u>7.97</u> | 29.04 | 0.889 |
| HY-Avatar | <u>63.09</u> | <u>765.05</u> | 3.91 | 2.57 | 7.35 | 8.36 | **66.07** | 0.880 |
| Hallo3 | 91.15 | 898.19 | 3.57 | 2.26 | 5.62 | 9.72 | 29.52 | 0.874 |
| EchoMimicV3 | 67.35 | 822.89 | <u>3.98</u> | <u>2.68</u> | 3.00 | 12.20 | 56.58 | <u>0.921</u> |
| Ours (baseline) | 107.50 | 1254.66 | 3.25 | 2.12 | 7.23 | 8.26 | 80.15 | 0.934 |
| +ref sink | 81.10 | 1060.86 | 3.69 | 2.36 | 7.60 | 8.05 | 79.67 | 0.908 |
| +RAPR | 63.71 | 801.68 | 4.03 | 2.66 | 7.45 | 8.04 | 62.09 | 0.925 |
| +GAN w/o global | 59.87 | 749.32 | 4.03 | 2.64 | 7.67 | 7.88 | 36.79 | 0.929 |
| **Ours** | **61.84** | **683.14** | **4.13** | **2.78** | **8.06** | **7.93** | <u>62.60</u> | **0.935** |



Figure 8. Qualitative Ablation Results.

coding processes across two H800 GPUs. We evaluate the performance using the Real Time Factor (RTF), which is the ratio between the inference time and the input segment length, and the First Frame Delay (FFD), when generating videos at a resolution of 928x704. The results are listed in Tab. 5. Since the RTF values of all modules are below 1, the system supports real-time generation. The overall latency is given by the sum of the FFD and the input chunk length (0.48s), resulting in a total of 1.20s.

## 8. Ethical Considerations

This work focuses on talking avatar generation for constructive, human-centered applications, and is not intended to

Table 5. Evaluation of the real-time performance of our model.

| Module | RTF | FFD |
|---|---|---|
| DiT | 0.69 | 0.33s |
| VAE | 0.82 | 0.39s |

support deceptive or harmful media. As with any generative technology, misuse is possible, such as creating fraudulent identities, fabricating false narratives, or generating avatars for harassment. To mitigate these risks, we commit to safeguards including embedding watermarks and clearly disclosing that all outputs are synthetic. We also aim to

Figure 9. Qualitative comparison with SoTA interactive head generation methods. Please ignore the arrows which come with the original video on ARIG's project page.
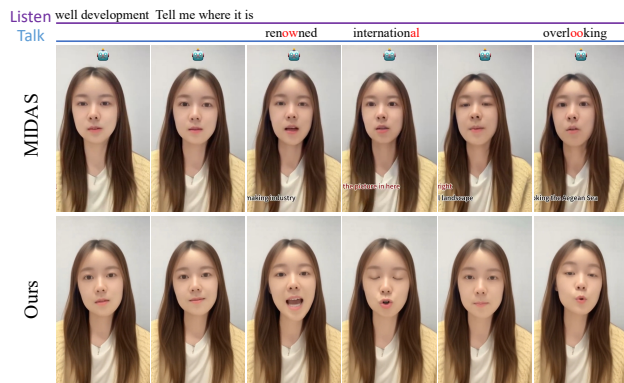


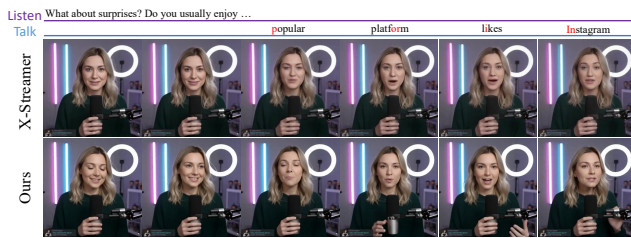Figure 10. Qualitative comparison with MIDAS [3].



Figure 11. Qualitative comparison with X-Streamer [41].

collaborate with the research community to develop improved deepfake detection tools and support efforts to establish standards for media provenance.