

EditYourself: Audio-Driven Generation and Manipulation of Talking Head Videos with Diffusion Transformers

John Flynn^{1,*} Wolfgang Paier^{1,*} Dimitar Dinev¹ Sam Nhut Nguyen¹
Hayk Poghosyan¹ Manuel Toribio¹ Sandipan Banerjee^{2,†} Guy Gafni^{1,‡}

¹ Pipio AI, ² Amazon

¹firstname.lastname@pipio.ai, ²sandgban@amazon.com

Project page: edit-yourself.github.io

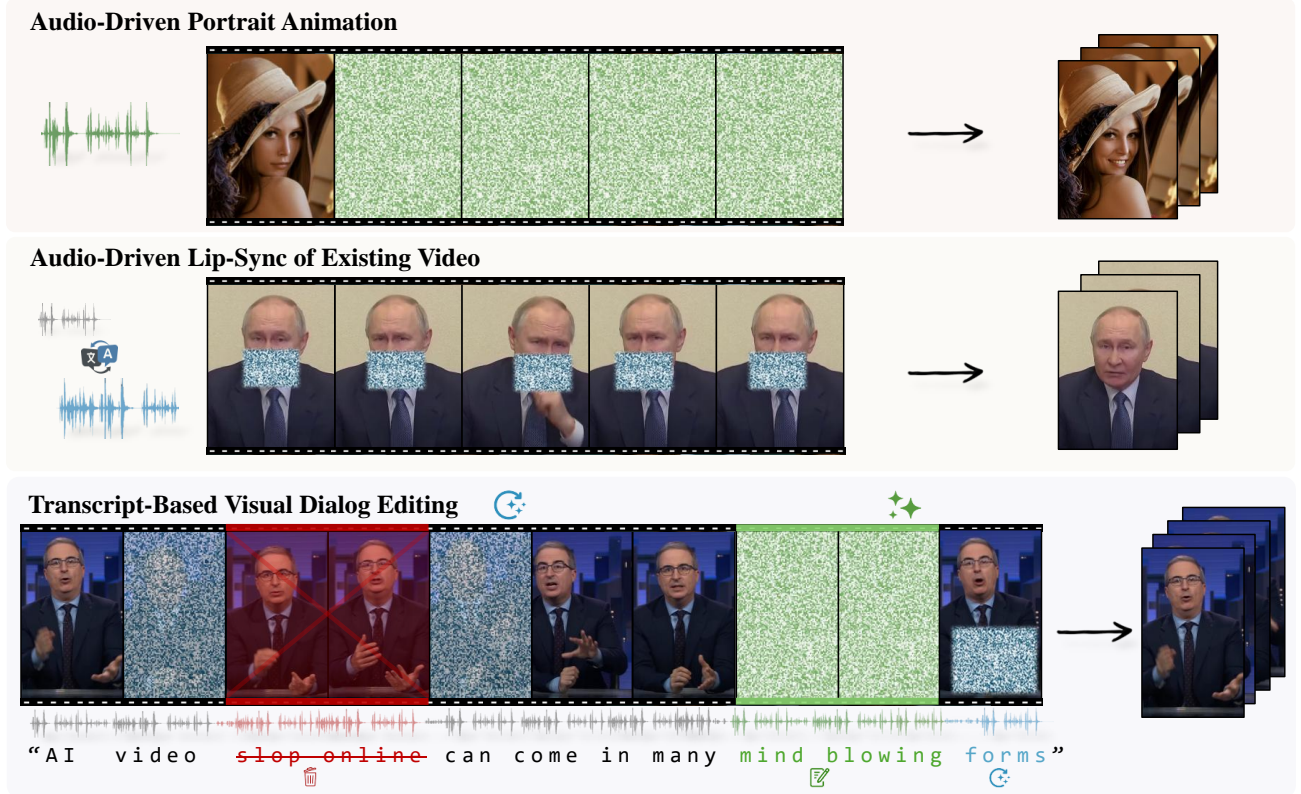


Figure 1. EditYourself is a multipurpose lip-syncing video diffusion model designed for transcription-based dialog editing, capable of lip-syncing from a single frame or an existing video, and seamlessly editing the video to match the new script.

Abstract

Current generative video models excel at producing novel content from text and image prompts, but leave a critical gap in editing existing pre-recorded videos, where minor alterations to the spoken script require preserving motion, temporal coherence, speaker identity, and accurate lip synchronization. We introduce EditYourself, a DiT-based framework for audio-driven video-to-video (V2V) editing that enables transcript-based modification of talking head videos, including the seamless addition, removal, and re-timing of visually spoken content. Building on a general-purpose video diffusion model, EditYourself augments its

V2V capabilities with audio conditioning and region-aware, edit-focused training extensions. This enables precise lip synchronization and temporally coherent restructuring of existing performances via spatiotemporal inpainting, including the synthesis of realistic human motion in newly added segments, while maintaining visual fidelity and identity consistency over long durations. This work represents a foundational step toward generative video models as practical tools for professional video post-production.

*Authors contributed equally.

†Work done while at Pipio AI.

‡Project Lead.

1. Introduction

A growing share of modern video content is human-centric, including movies, online courses, corporate communications, interviews, and short-form social media uploads. In these videos, creators often need to revise the spoken content post-recording to fix fumbled lines, update facts, remove filler words, tighten interviews or localize content across languages. Traditional non-linear video editing tools, however, provide limited support for such edits, as operations like inserting, removing or retiming speech typically introduce visible jump cuts or unnatural motion. Moreover, selectively re-rendering only parts of a human performance requires extensive manual intervention within existing post-production workflows [25, 46, 60]. Recent advances in video diffusion models suggest a promising alternative. These models can synthesize high-quality, temporally coherent human videos from text, images or audio, demonstrating an ability to model complex appearance, motion, and facial dynamics [35, 38]. This capability makes generative models a good candidate for not only content creation, but as editing engines that can repair, extend or reshape existing videos in a content-aware manner [80].

However, research on editing existing human-centric content remains far less mature than work on end-to-end generation. Most current approaches focus on Image-to-Video (I2V) generation from a single portrait image [11, 16, 47, 122]. While impressive in their realism, these methods frequently suffer from identity drift over time and incorrectly reproduce a subject’s likeness. A single image cannot capture the full range of facial details and speaking style present in a real performance. As a result, generated videos often hallucinate details such as teeth, wrinkles, facial hair or gestures, producing outputs that feel incorrect, especially when users are generating videos of themselves. On the other hand, V2V lip-sync models [62, 85, 112, 119] adhere closely to an input video to preserve visual fidelity and identity, but offer limited flexibility for editing. By operating under a fixed temporal structure that preserves the original frame count and timing, these methods make it difficult to insert or remove speech segments while maintaining temporal continuity. Consequently, existing V2V and I2V methods do not adequately support precise edits required for real-world post-production workflows.

Our work tackles this fundamental problem of temporal manipulation of existing talking-head videos, which we refer to as *visual dialog editing*: V2V editing driven by changes to the spoken dialog [4]. This setting goes beyond simple lip synchronization to completely new audio, and supports core post-production operations such as inserting, removing and retiming video segments while preserving visual continuity. Editing videos directly through their textual transcript provides an intuitive and expressive interface for creators, enabling precise word-level modifications such as

filler-word removal and post-shoot script revisions. More broadly, this transcript-centric workflow shifts video production from a “script-perfect-before-shooting” paradigm toward a “shoot once, refine later” model, enabling rapid updates, personalized variants and integration with higher-level control systems such as LLM-based AI agents for automated video editing [46, 60, 102].

In this work, we address this gap by re-framing talking-head video synthesis as a problem of visual dialog editing. We introduce *EditYourself*, a diffusion-based framework designed specifically for transcript-driven editing of talking head videos. By adapting a pre-trained general-purpose video diffusion model into a flexible, audio-driven V2V editor, our approach enables precise modification of existing videos, including addition, removal, and retiming of spoken segments, while maintaining accurate lip synchronization, visual identity, and temporal coherence over long videos.

In summary, our work makes the following contributions:

- **Lip-sync on a pretrained video diffusion model:** We introduce a two-stage training scheme that enables inference on speech audio across varying text, image, and video inputs, while maintaining accurate lip synchronization, together with a windowed audio conditioning strategy for precise speech-video alignment that does not require audio feature downsampling and remains robust across varying video frame rates.
- **Latent-space visual dialog editing:** We formulate transcript-driven video editing directly in latent space, supporting seamless addition, removal, and retiming of spoken segments.
- **Identity-preserving long video generation:** We introduce a reference-based identity conditioning mechanism, *Forward-Backward RoPE Conditioning*, together with TeaCache-aware inference, to stabilize appearance and temporal coherence over long videos.

Evaluations against recent I2V and V2V lip-sync benchmarks demonstrate that our method achieves SOTA visual quality and synchronization accuracy. In addition to offering competitive performance, our approach represents a foundational step toward utilizing video diffusion models as capable tools for editing human-centric video content.

2. Related Works

With the advent of diffusion models [59, 71], the field of video generation [78] has proliferated in recent years. Coupled with powerful 3D VAEs [56], these models have the capability of reconstructing the details and dynamics of an entire frame (instead of a small crop), opening the door to generating novel frames that are coherent with the rest of the video. There are several possible input modalities, which can be combined together, that define the task of the model. The common modalities are: (i)

Text-to-Video (T2V) synthesize the video from a textual input [38, 40, 107], **(ii)** *Image-to-Video* (I2V) animate a single image into a video [3], **(iii)** *First-Last-frame-to-Video* (FL2V) guide video generation between the given first and last frames [61, 106], and **(iv)** *Video-to-Video* (V2V) edit or transform video content while maintaining temporal consistency and structure. [67, 112]. The latest video generation works [30, 35, 73] focus on synthesizing clips that adhere to provided prompts by leveraging diffusion transformer blocks (DiTs) [83] as the main computational units in their models. A newer version of this, multi-modal diffusion transformer blocks (MM-DiT) [22] allow multiple input modalities to be represented in a common token space, facilitating joint attention across them.

2.1. Audio-Driven Talking Head Generation

Early Methods. Audio-driven facial animation, in particular lip-syncing, has been an active research topic with a variety of methods explored. GAN-based methods [8, 32, 53, 54, 88] achieved early success with appropriate audio representations, such as Wav2Lip and Wav2Vec [2, 85] for conditioning. These methods can indeed lip-sync a video, however are unable to make larger changes like head motion. Adding 3D Morphable Models [7] from traditional graphics as an intermediate representation allows enhanced control over the subject in the video. Lip-sync and head control can be added using the parameterizations (*e.g.* blendshapes) of the models [15, 89, 109]. Volumetric rendering techniques from graphics [74, 79] have found success in representing human avatars [26, 34, 113, 115]. More recently, 3D Gaussian Splatting [55] techniques have also been used to successfully lip-sync videos [12, 63].

Diffusion based Methods. In the last couple of years, latent diffusion models [90] have become the backbone of choice for generating talking head videos from a single source image or video. The earlier set of these models typically use a pre-trained 2D/3D VAE [56] to encode the source and a trainable UNet-style module [91] for denoising. A trainable copy of the UNet acts as a reference net to inject control signals into the denoiser’s feature space. These signals can be audio representations [2, 85], emotion embeddings [122], face and body keypoints [5, 11, 42] or identity information [6, 110]. However, recent models replace the denoising UNet with a diffusion transformer (DiT) [83] for improved scalability and global context handling [16, 30, 57, 64, 84, 101, 105], and focus on multi-stage training [49, 70] where the model is incrementally trained on a higher data dimensionality of the source (*e.g.* audio, then image, then video). The final model can then be controlled by only audio [29, 47] or combining it with blendshapes [104], pose information [27, 77, 84], external embeddings [10, 33, 101] that is injected into the la-

tent model via cross attention or a multi-modal block [22]. The majority of these models [20, 30, 97, 105] use a flow-matching objective [71] rather than denoising diffusion due to its faster sampling and straightforward noise to data path.

2.2. Video Manipulation

Video-to-Video Editing. As a consequence of the above line of research, V2V editing and manipulation applications have gained considerable popularity of late. Some of these models, like Runway’s Gen-3 Alpha/Gen-4.5¹, perform style transfer [51] on existing videos while others (*e.g.* [99], [103], [42], [66]) focus on direct motion transfer from conditioning signals. For explicit content insertion/removal from video frames, inpainting models leveraging optical flow have been explored [117, 123], with recent works using diffusion [3, 121] and flow matching [50, 108].

The video editing problem can be formulated as a collection of image editing steps (or “slices” [14]) directly using a pretrained T2I model. Finetuning the model [24, 80, 93] enables text-based editing, while enhancements such as feature banks and optical flow [58, 68, 100, 121] can improve restyling quality and object removal. Latent diffusion models [90] are suitable for surgically editing specific regions of the video, as there is a clear mapping between the space and time coordinates of any given video pixel to the latent tokens it generates [96]. Keyframes, a concept from traditional video editing, can also be used to loosen the one-to-one correspondence between the input and output video [5, 6, 112] and produce videos that match the subject, but could have different head or hand movements.

Transcript-Based Editing. A problem closely related to our work is how to handle changes in the spoken script without requiring re-shooting of the whole segment. Early methods [25] presented a dynamic programming-based synthesis strategy to assemble new speech videos combining visemes, 3DMM-based blending, and a recurrent video generation network, while [114] used a fast phoneme search and neural re-targeting to transfer mouth movements from the source to a target. The talking-head editing process can also be broken down into audio-to-dense-landmark motion and motion-to-video stages [36, 111]. Although these methods enable transcript-based editing, they are limiting in that they either require subject-specific data or they struggle to generalize to diverse videos.

2.3. Identity-Preserving Long Video Generation

Generating longer videos with diffusion models remains a technical challenge. The spatio-temporal dimensions of the output video are determined by those of the noise tensor (*i.e.* the sequence length of the noise tokens), practically limited

¹<https://runwayml.com/research/introducing-runway-gen-4.5>

by GPU memory. Naive auto-regressive techniques experience drastic reductions in video quality and identity preservation [43]. Recent techniques mitigate video quality degradation [43, 64, 81, 118], but these are not sufficient to prevent identity drift in human faces (*i.e.* loss of facial details, over-smoothing, over-saturation of the skin, changes in facial hair). A simple approach is to leverage a reference subject image encoding, as done in [10, 16, 48, 70, 92, 99, 125]. As the reference image can contain background information not related to the subject’s identity, embeddings from CLIP [86] or the face-specific ArcFace model [19] can also be used [28, 39, 110, 116, 120]. These features can be integrated into the DiT via cross-modal adapters.

3. Method

We base our model on LTX-Video [35], a general-purpose video diffusion model that supports text, image, and video-conditioned generation (T2V, I2V, and V2V), which we introduce in Section 3.1. Building on this backbone, we introduce a set of extensions that specialize the model for audio-driven and transcript-based video editing. Specifically, we include (i) cross-modal audio conditioning and a V2V lip-sync training strategy (3.2), (ii) a latent-space formulation of visual dialog editing that supports transcript-driven addition, removal, and retiming of speech (3.3), (iii) a caching-aware long-inference strategy for temporally consistent generation over long durations (3.4), and (iv) reference-based identity conditioning with a novel Forward-Backward Rotary Positional Embedding (RoPE) [95] mechanism to stabilize appearance across both edited and fully synthesized segments (3.4).

3.1. Preliminaries

Baseline Network. We use the LTX-0.9.7 [35] DiT and the associated Video-VAE as our baseline, which follows the common 3D causal VAE and flow-matching DiT pattern for video generation, including 3D RoPE [95] for spatio-temporal positions and adaptive normalization for timestep conditioning. With 14B parameters, the DiT model operates in a highly compressed latent space using a separately pre-trained Video-VAE. The VAE encoder’s rather aggressive compression rate ($32 \times 32 \times 8$) results in a significantly lower token count, aimed at increasing performance towards interactive applications. Videos are generated in a two-pass fashion: (i) denoising is first performed on a coarser, lower-resolution representation of the video, (ii) followed by a learned upsampling of the latents and a second, higher-resolution denoising pass. Crucially, LTX-Video was pre-trained with a flexible multi-task objective spanning T2V, I2V, keyframe generation, and various forms of spatial and temporal inpainting. This is achieved by masking tokens and assigning them distinct conditioning timesteps, regardless of the global diffusion step.

Flow Matching Training Objective. We adopt the Flow Matching [71] paradigm, following LTX-Video [35]. Formally, video samples are encoded into a latent representation $\mathbf{x}_0 \sim p_{\text{data}}$ using the LTX-Video Video-VAE. We define a linear probability path to interpolate between latent video representation \mathbf{x}_0 and a noise distribution $\mathbf{x}_1 \sim \mathcal{N}(0, \mathbf{I})$ via the displacement flow $\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\mathbf{x}_1$ with a continuous time step $t \in [0, 1]$. The DiT, v_θ , is trained to predict the velocity field that transforms noise back into data by minimizing our base training objective

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t, \mathbf{x}_1, \mathbf{x}_0, \mathbf{c}} [\|v_\theta(\mathbf{x}_t, t, \mathbf{c}) - (\mathbf{x}_1 - \mathbf{x}_0)\|_2^2], \quad (1)$$

where \mathbf{c} denotes the available input conditions (text prompt, in the LTX-Video base model). In the subsequent subsections, we modify this objective to include audio and identity conditioning. Please see Equation (8) for the expanded training loss formulation.

At inference, new videos can be generated by solving the probability flow ODE, $\frac{d\mathbf{x}_t}{dt} = v_\theta(\mathbf{x}_t, t, \mathbf{c})$, which requires integrating the velocity field from $t = 1$ to $t = 0$:

$$\mathbf{x}_0 = \mathbf{x}_1 - \int_0^1 v_\theta(\mathbf{x}_t, t, \mathbf{c}) dt \quad (2)$$

In practice, this integration is discretized using a first-order Euler solver over 40 steps following the update rule $\mathbf{x}_{t_{i-1}} = \mathbf{x}_{t_i} - \Delta t \cdot v_\theta(\mathbf{x}_{t_i}, t_i, \mathbf{c})$. For further details, please refer to the original LTX paper [35] and repository [69].

3.2. Cross-Modal Audio & Video Conditioning

Audio Conditioning Strategy. Inspired by DiT-based portrait animation methods [20, 47, 70, 84, 101, 105], we extend a pre-trained video diffusion model with an audio conditioning modality by introducing additional cross-attention layers into the transformer blocks. Specifically, we insert one such layer into each DiT block, positioned between the text cross-attention and the FFN. As keys and values, we use pre-extracted Whisper-small [87] features $\mathbf{c}_{\text{audio}} \in \mathbb{R}^{L \times B \times C}$ with L the sequence length, B the number of encoder block outputs and C the channel dimension. The proposed conditioning mechanism however is agnostic to the choice of audio representation. The audio features are processed by a learned projection and pooling module (Audio Projection) to produce lip-sync embeddings at the latent video frame rate. These embeddings are then shared across all DiT blocks, allowing audio information to modulate video features at every layer while preserving the pretrained DiT’s token structure.

To minimize disruption to the pretrained DiT at the start of training, we initialize the Audio Projection module’s convolution layers as average pooling operators and set the audio cross-attention output projections to zero. During training, we randomly drop audio conditioning with probability

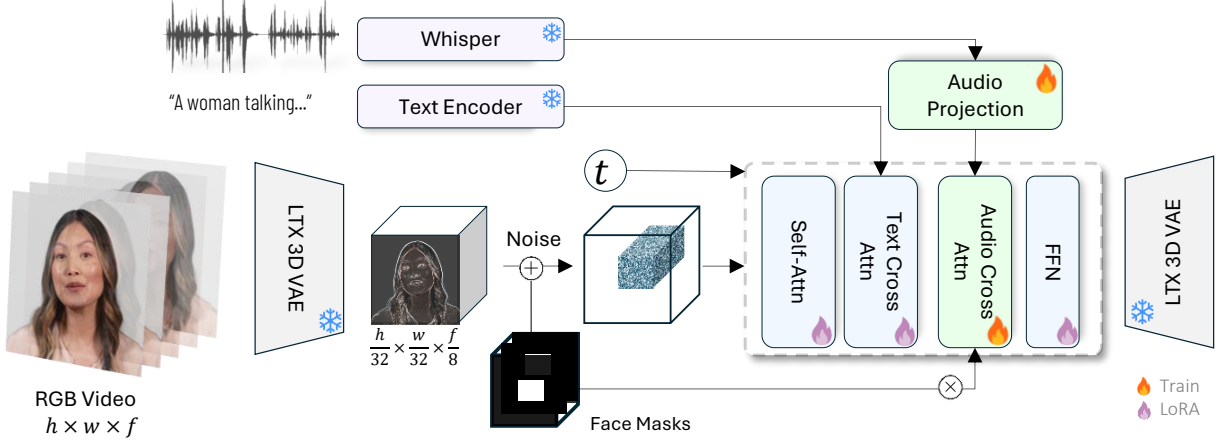


Figure 2. Our proposed pipeline. A global audio projection layer and audio cross-attention layers are added to the network’s architecture. For V2V lip syncing, we apply noise to tokens corresponding to the mouth area and task the model with spatio-temporally inpainting them.

\bar{p}_{audio} by detaching these layers. Overall, the Audio Projection module and associated cross-attention layers introduce approximately 2B additional learnable parameters. The resulting architecture is illustrated in Figure 2.

To restrict attention to temporally local audio context, we associate each video frame index i with a window of W audio features $\tilde{\mathbf{c}}_{\text{audio}}^i \in \mathbb{R}^{W \times B \times C}$. Because the sampling rate of audio features (e.g. Whisper embeddings) f_a typically differs from the video frame rate f_v , naively selecting the nearest audio features can introduce sub-frame audio-video misalignment. Prior approaches often address this mismatch by interpolating audio features to f_v . However, this strategy is fragile for two reasons: (1) downsampling can discard high-frequency information present in modern speech embeddings, and (2) fixed-size windows correspond to different temporal durations across videos with varying frame rates.

To address these issues, we sample audio features on a phase-shifted grid that preserves the original audio feature rate f_a while aligning audio windows to video frames. Specifically, for each video frame index i , we extract a center-aligned window of W audio features from $\mathbf{c}_{\text{audio}}$ at fractional audio indices u_n using linear interpolation, yielding $\tilde{\mathbf{c}}_{\text{audio}}^i[n]$, where n denotes the index within the window.

$$\begin{aligned} u_n &= i \frac{f_a}{f_v} + \left(n - \frac{W-1}{2}\right), n = 0, \dots, W-1, \\ k_n &= \lfloor u_n \rfloor, \\ \alpha_n &= u_n - k_n, \end{aligned}$$

$$\tilde{\mathbf{c}}_{\text{audio}}^i[n] = (1 - \alpha_n) \mathbf{c}_{\text{audio}}[k_n] + \alpha_n \mathbf{c}_{\text{audio}}[k_n + 1] \quad (3)$$

This design decouples audio temporal resolution from video frame rate, ensuring consistent window semantics across videos with arbitrary frame rates.

To encode relative position within the audio window, we introduce a learned, fixed-size positional embedding tensor

$\mathbf{P} \in \mathbb{R}^{W \times B \times C}$, where each slice $\mathbf{P}[n]$ corresponds to a window index.

$$\tilde{\mathbf{c}}_{\text{audio+pos}}^i[n] = \tilde{\mathbf{c}}_{\text{audio}}^i[n] + \mathbf{P}[n] \quad (4)$$

V2V Lip-Sync. LTX-Video [35] is a general-purpose video diffusion model supporting text, image, and video-conditioned generation. In its standard I2V usage, the model conditions on clean latents for an initial frame and noisy latents for subsequent frames, with self-attention propagating information to guide temporal generation. We build on this mechanism to specialize the model for audio-driven V2V lip synchronization by selectively regenerating the mouth region in talking-head videos. For each source video, we detect lower-face bounding boxes using Mediapipe [76] and compute an enclosing box over groups of eight consecutive frames, yielding a binary mask \mathbf{M} per latent frame. During training, noise ϵ is applied at timestep t only within the masked region \mathbf{M} , and the model is trained to inpaint the corresponding tokens over space and time while preserving unmasked content. See Fig. 3.

$$\mathbf{x}_t = \mathbf{M} \odot [(1-t)\mathbf{x}_0 + t\epsilon] + (1-\mathbf{M}) \odot \mathbf{x}_0 \quad (5)$$

We also restrict the audio cross-attention layers to only update tokens belonging to the face region by multiplying the cross-attention output with the face mask:

$$\mathbf{z}_{\text{out}} = \mathbf{z}_{\text{in}} + \mathbf{M} \odot \text{AudioAttn}(\mathbf{z}_{\text{in}}, \mathbf{c}_a) \quad (6)$$

where \mathbf{z} denote the hidden latents in the DiT.

We further apply random conditioning dropout to ensure that the model remains robust to missing spatial and temporal inputs, and can operate under different combinations of conditioning signals. Following the original training procedure for LTX-Video, we randomly drop first-frame condi-

tioning with probability \bar{p}_{ff} , reducing the objective to text-to-video generation. Similarly, we randomly drop video-to-video conditioning with probability \bar{p}_{v2v} to preserve the model’s image-to-video generation capability. When video-to-video conditioning is absent, we set the spatial mask to $\mathbf{M} = 1$, allowing audio cross-attention to update all tokens, enabling unconstrained latent generation. Tab. 1 reports the values we chose for p_{ff} , p_{v2v} and p_{audio} .

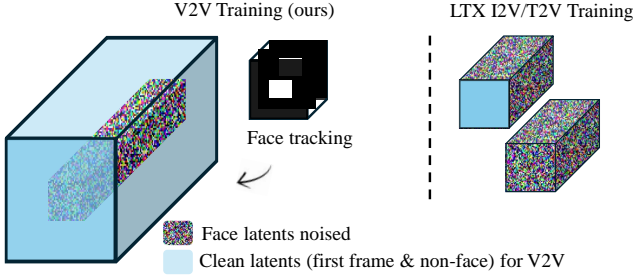


Figure 3. In order to train the audio attention layers, we fully noise the tokens corresponding to mouth region throughout the training sample. We retain clean latents of the first frame, similar to image-to-video training in LTX. The model learns to in-paint the mouth through time and space using the audio, and the initial mouth shape as conditions.

Our masked training strategy affords substantial flexibility at V2V inference time. By selectively scaling and positioning the mask \mathbf{M} (see Fig. 4), the model can be configured to synchronize only the lips, the face, or the entire head. In the *Head* mask mode, the model leverages its generative prior over head motion learned during T2V and I2V lip-sync training, selectively re-synthesizing head dynamics to match the timing and prosody of the new speech. In contrast, the *Face* and *Mouth* modes progressively constrain generation to smaller spatial regions, producing new content within the masked area while increasingly adhering to the original video outside the mask.

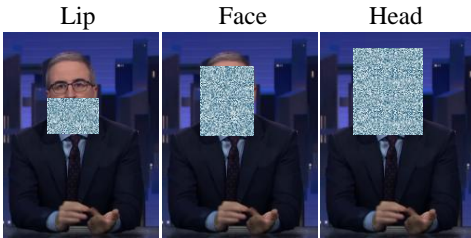


Figure 4. V2V Inference Modes. Adjusting the mask \mathbf{M} in inference enables different synchronization levels: *Lip* for mouth-only sync, *Face* for expressions, and *Head* to synthesize new head dynamics matching the audio prosody.

While noise levels are sampled individually for each token, we must bias the sampling of noise timesteps towards

the high-noise region. If noise levels in the mouth region are too low during training, the pre-trained DiT’s inherent ability to denoise video tokens will allow it to infer the mouth correctly without relying on the audio signal, leading to a trivialization of the cross-attention layers in training, causing a collapse of the lip-sync task when inferencing on novel audio. This noise bias aligns with the observation of previous works that crucial lip-sync details are determined primarily in earlier stages of denoising [66, 84, 99]. We use a shifted log-normal distribution with shift $\mu = 2.05$ which places 90% of timesteps in the range [0.60, 0.98]. The tendency of the lip shapes and mouth movements to be determined during the early noise stages is further validated during inference.

3.3. Visual Dialog Editing

Modification of the transcript in text domain admits changes to both audio and video domains. We address the visual synthesis challenge and use commercially available solutions for zero-shot voice cloning [17, 21]. We target the generation of new frames, re-generation of existing frames, eliminating discontinuity artifacts across edit boundaries (jump-cuts), all while ensuring accurate lip synchronization to the target audio. Specifically, we define the following operations for a video segment:

1. **Addition:** Insertion of new content at arbitrary timestamps, seamlessly adhering to surrounding boundary frames (when present).
2. **Removal:** Deletion of existing content while smoothing the resulting temporal discontinuity to avoid visible jump cuts.
3. **Re-render:** Selective inpainting of video content over specified spatial and temporal regions (e.g. correcting an awkward facial expression or replacing a hand gesture).
4. **Retime:** Altering the total duration of a video segment to match changes in script duration (e.g. for language localization), implemented via evenly-distributed additions/removals.

Unlike Addition and Removal, which are localized operations, Retime applies distributed temporal adjustments across the entire segment. This distinction is particularly important for dubbing, since translated speech often involves changes in word order, density, and duration; strict correspondence with the video timeline is lost, requiring adjustments to the overall duration of the segment rather than at specific words. We motivate the operations with the following examples in Figs. 5 and 6.

We obtain word-level timestamps using an automated transcription tool [1, 18]. The user-edited transcript is then diff-checked against the original to identify changed spans and map them to their corresponding audio timestamps. This results in a finalized set of Addition, Removal, and Retime operations.

Original: “This feature rocks and we will most likely launch it.”

Revised: “This awesome new_(+0.9s) feature rocks and we will ~~most likely~~_(-0.5s) launch it next week_(+0.6s).”

Figure 5. Example of a script-driven temporal edit, illustrating the complexity of V2V operations. New content is highlighted in green, and a redaction is shown with a red strike-through, accompanied by the required duration change for each operation. Two addition operations and a removal operation are needed to account for these edits.

EN: “It’s stated in our terms and conditions.”

DE: “Das ist in unseren Allgemeinen Geschäftsbedingungen festgelegt. (+1.1s)”

Figure 6. Example of a *Retime* operation needed for language localization/dubbing. The English phrase expands significantly when translated into German, requiring the model to expand the duration of the entire segment by approximately +1.1s to maintain natural speech.

We leverage LTX-Video’s architectural flexibility and formulate visual dialog editing as a specialized inpainting task. To realize these edits, we modify the latent video frame sequence directly along the spatial and temporal axes. Frame addition is implemented by inserting fully noised latent frames at the corresponding locations, while frame removal deletes existing latent frames from the sequence. Exploiting the causality of the VAE encoder, we define a mapping between each latent frame at index n and its corresponding range of input video frames indexed from $8(n-1)+1$ to $8n+1$ exclusive, with latent frame 0 mapped to video frame 0. This mapping provides a clean proxy for temporal editing in latent space at an 8-frame resolution.

To mitigate visual artifacts introduced by frame removal, we apply additional noise to latent frames adjacent to removed segments, allowing the diffusion process to smoothly regenerate motion. Selective re-rendering is implemented by setting the spatial mask $\mathbf{M} = 1$ at arbitrary regions—such as the face, head or hands—so that only those regions can be noised and regenerated while the remainder of the video remains unchanged. Since editing alters the sequence length, temporal rotary positional embeddings (RoPE) are computed on the edited latent sequence. Finally, newly inserted frames are fully unmasked ($\mathbf{M} = 1$), while existing frames retain face-region masking during lip-sync generation.

The overall latent-space editing process is illustrated in Figure 7.

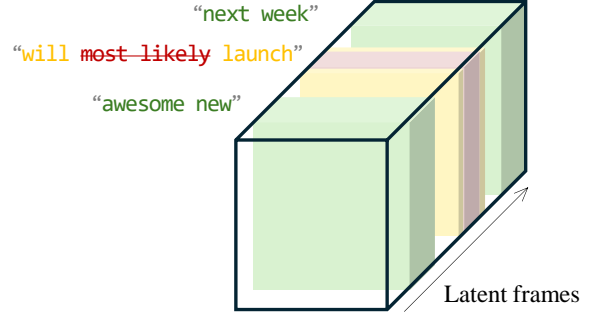


Figure 7. Timeline editing via manipulating latent frames. **Addition:** Insert new latents with full noise. **Removal:** Discard corresponding latents and noise adjacent latents to **re-render** a smooth transition.

3.4. Identity-Preserving Long Inference

Long Inference. For long video generation, the sequence of latent frames is processed in blocks. Rather than fully denoising each block before proceeding to the next one (autoregressive long inference), we adapt the *Time-aware position shift fusion* (TAPSF) long inference strategy proposed in Sonic [10, 47]. We first encode the entire video into latent space and logically partition the full latent video into non-overlapping inference blocks. We choose a block width of 17 latent frames (136 video frames). For each timestep, we iteratively perform a **single** denoising step on each block of frames. For the next denoising timestep, the partition of frame latents into blocks is offset, such that the next denoising step will integrate context from adjacent frames, sharing longer-form context over many such denoising and offsetting steps (Fig. 8). The model then naturally bridges context between adjacent blocks throughout the entire denoising process, increasing inter-block stability. This strategy switches the order of looping between denoising steps and frame blocks compared to typical autoregressive-style inference. Further details can be found in [47].

At the video boundaries, frames that fall outside a block are handled by evaluating the overlapping regions twice: once from each neighboring block and averaging the resulting predicted velocities. RoPE are computed in the global frame coordinate space, ensuring consistent temporal positioning across shifts.

One disadvantage of TAPSF is its incompatibility with popular cache-based acceleration techniques such as [72, 124]. These methods accelerate inference by re-using the DiT blocks’ outputs across previous timesteps if the inputs are “similar” enough. Similarity of the hidden states can be determined by, *e.g.* a rescaled difference norm of the timestep-embedding-modulated inputs. With TAPSF, the previously computed block outputs correspond to a temporally *shifted* set of features rather than a stationary representation of the same frames at a different noise level.

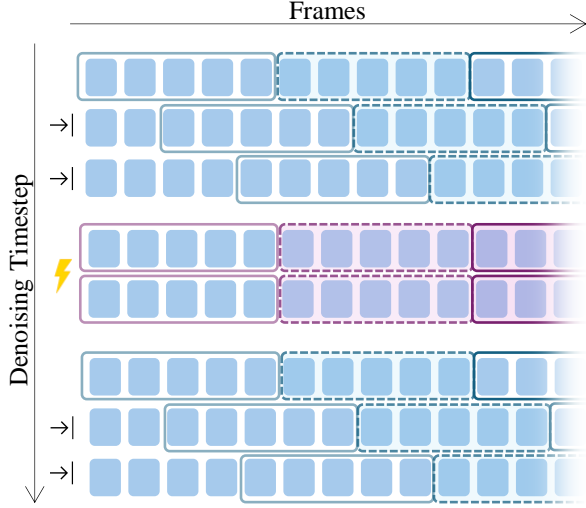


Figure 8. Long inference strategy: video latent frames are grouped into inference blocks and denoised iteratively. We apply a position shift to the blocks after each denoising step (blue inference blocks) to propagate context over longer windows throughout denoising. During medial timesteps (purple inference blocks) we disable the shift to benefit from TeaCache.

Observing that lip synchronization, identity cues, and large-scale motion are primarily determined during early denoising timesteps, we adjust TAPSF to not shift blocks during the middle steps of denoising, and maintain a shift of 5 latent frames during the early denoising steps (which are heavy on lip-sync, identity, and large motion) and late denoising steps (fine details). This modification allows us to apply adaptive caching during the middle 75% of denoising steps, resulting in a speedup of approximately $1.6\times$ while preserving the benefits of TAPSF for long-range temporal coherence.

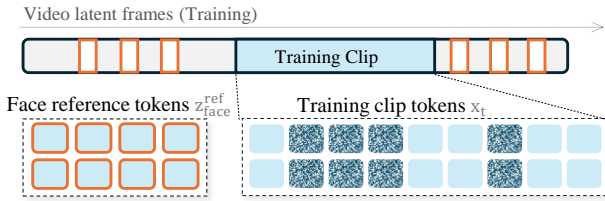


Figure 9. Identity conditioning: we train the DiT to use unnoised face tokens from outside the training clip to better preserve subject identity. These reference tokens are taken from a temporal neighborhood of the training clip and randomly added to the DiT’s input sequence.

Identity Conditioning. Popular portrait animation methods [57, 97, 101, 110] address drift in subject identity by injecting facial features that capture the speaker’s visual char-

acteristics, such as CLIP [86], DiNoV2 [82], or face embeddings [19] into dedicated cross-attention layers. For I2V models, the image prompt corresponding to the first frame of the video also implicitly serves as the identity reference. In this setup, the first frame conditions the model along two paths: through self-attention, as its clean tokens are present in the sequence of tokens entering the DiT and all the noisy tokens attend to it. In addition, features from the reference image can be injected through cross-attention.

In the V2V setting, a full reference *video* of the subject is available during both training and inference. InfiniteTalk [112] dynamically swaps the single reference frame for each inference block, with a frame from the video, to preserve appearance and coarse temporal progression in “sparse video-to-video dubbing.” However, we aim to leverage the subject’s identity and speaking style present throughout the entire video, rather than reducing conditioning to a single frame at a time.

To this end, we fine-tune the self-attention mechanism to condition on reference frames. During training, we randomly sample 6 latent frames (corresponding to 64 video frames) from a temporal window of ± 5 s around the target clip, encode them, and retain only tokens corresponding to the lower face region. These *face reference tokens* $z_{\text{face}}^{\text{ref}}$ are kept un-noised and concatenated to the video tokens along the sequence dimension.

OmniHuman-1 [70] also conditions on a reference frame as concatenation to the token sequence, specifically by zeroing the temporal component of its 3D Rotary Positional Embedding (RoPE), effectively removing temporal ordering and motion information while still providing appearance cues. While this design suffices for a single reference frame, it fails to extend to our reference-video setting, where multiple reference frames correspond to the same spatial region (e.g., the mouth) but capture different temporal states. Zeroing their temporal embeddings would align all reference tokens on the same RoPE phase, leading to aggregation bias—the model tends to average or equally attend to all reference tokens, despite each representing visuals for distinct phonemes or lip positions. To mitigate this, we assign unique sentinel temporal indices ($t = -1, -2, \dots$) to reference tokens from different frames. This preserves their distinct temporal identities while keeping them separable from the generated video’s temporal sequence. In inference, we sample face reference tokens from the block processed if available, and optionally from adjacent blocks.

For inference on fully synthetic blocks lacking clean video latents, as in long I2V generation or extended additions, we must also prevent drift in global frame appearance. To this end, We propose *forward-backward RoPE conditioning* during inference, assigning full-frame reference tokens $z_{\text{frame}}^{\text{ref}}$ from the last-seen and first-available future clean latent frames (see Fig. 10). These tokens are

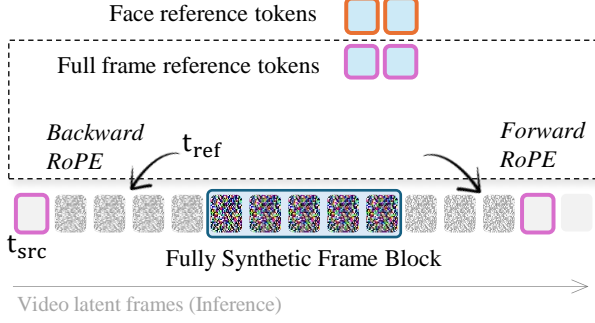


Figure 10. During inference of fully-synthetic blocks (i.e., blocks without V2V or first-frame condition), we prevent global appearance drift by adding full-frame reference tokens (the closest past and future latent frames) to the input sequence. These reference frames are not noised and their temporal indices are adjusted such that the temporal distance between the reference frames and the block is not greater than 3.

assigned “fake” temporal indices t_{ref} to provide appearance cues that are aligned enough in the temporal phase of RoPE without forcing exact replication of those frames. i.e. $\text{RoPE}(\mathbf{z}_{\text{frame}}^{\text{ref}}, t_{\text{ref}})$ with

$$t_{\text{ref}} = \begin{cases} t_{\text{source}} & \text{if } \Delta t \leq 3 \\ t_{\text{block (end)}} + 3 & \text{if } \Delta t > 3 \text{ (forward)} \\ t_{\text{block (start)}} - 3 & \text{if } \Delta t > 3 \text{ (backward)} \end{cases} \quad (7)$$

where $\Delta t = |t_{\text{source}} - t_{\text{block}}|$ is the temporal distance between the source reference frame and the frame block boundary. Similar to the face reference tokens, the frame reference tokens $\mathbf{z}_{\text{frame}}^{\text{ref}}$ are also concatenated with the video tokens along the sequence dimension for inference.

We note that adapting the temporal embedding of frames as a method to address consistency is also proposed in concurrent work [44, 92].

Combined with our long-inference approach, our model can output minutes-long videos without noticeable identity drift (see [project page](#)).

3.5. Training Loss

Our final training loss becomes

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t, \epsilon, \mathbf{x}_0, \mathbf{c}} [\|\mathbf{M} \cdot v_{\theta}(\mathbf{z}_{\text{in}}, t, \mathbf{c}_{\text{audio}}, \mathbf{c}_{\text{text}}, \mathbf{M}) - \mathbf{u}_t\|_2^2] \quad (8)$$

where the target velocity \mathbf{u}_t is masked to focus the learning signal on the mouth region:

$$\mathbf{u}_t = \mathbf{M} \odot (\epsilon - \mathbf{x}_0)$$

with the input tokens $\mathbf{z}_{\text{in}} = [\mathbf{z}_{\text{face}}^{\text{ref}}, \mathbf{z}_{\text{frame}}^{\text{ref}}, \mathbf{x}_t]$, containing the face reference tokens, frame reference tokens, and noisy video tokens, all concatenated along the sequence dimension.

4. Experiments

4.1. Training

We base our model on the Lightricks LTX-0.9.7 architecture and open-sourced weights [35, 69].

Dataset. We collect a total of 1,070 hours of talking-head footage. This includes 70 hours of proprietary, high-quality frontal recordings, along with 1,000 hours of shorter user-generated content gathered from YouTube, exhibiting substantial variability in appearance, identity, pose, background, gestural dynamics, and composition. The videos span a broad range of resolutions (0.25 to 2.0 MP), aspect ratios ($2 : 1, 1.78 : 1, \dots, 1 : 1, \dots, 0.56 : 1, 0.5 : 1$), and frame rates (24–60 fps). This enables support for diverse input formats during inference. Additionally, we filter the dataset for scene cuts, lip-sync confidence score ($\text{SyncC} \geq 3$ [85]), and temporal offset ($\leq 40\text{ms}$) [13], bitrate ($\geq 2000\text{Kbps}$), frame-rate ($\in [24 - 60]$), number of frames (≥ 121) and corrupted videos, finally yielding 475 hours of video cut into 121-frame video clips. The data filtration process is depicted in Fig 11. We generate video captions with CogVLM2 [41].

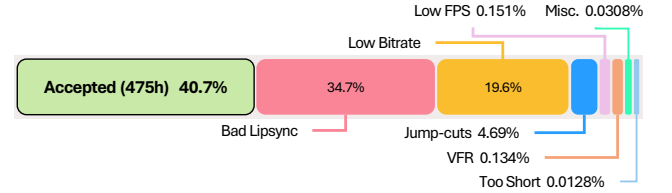


Figure 11. Data filtering results: a significant share of in-the-wild short-form videos exhibited low visual quality or bad lip-sync. We removed these from the dataset to achieve optimal training performance.

Model Training. Training proceeds in two stages. In stage 1, the new audio projection module and cross-attention layers are optimized (with all other weights frozen) for $20k$ steps, encouraging the model to learn lip synchronization without compromising the generative prior of the pre-trained DiT. In Stage 2, training transitions to a 128-rank low-rank adaptation (LoRA) of the full DiT for another $10k$ iterations, which enables identity conditioning, further improves lip synchronization, and is instrumental to the visual fidelity of the resulting videos.

We train our model to generate videos with any combination of input conditions and at a range of video frame rates and resolutions. To support varied input combinations, we randomly enable the video-to-video, first-frame, and audio conditions with probabilities p_{v2v} , p_{ff} and p_{audio} , respectively (see Table 1). Note that audio conditioning is

Table 1. Hyperparameters across training stages.

Parameter	Stage 1: Audio-Only	Stage 2: LoRA
Steps	20k	10k
Base LR: lr_{adam}	1e-5	1e-4
Audio Cond: p_{audio}	1.0	0.9
FF Cond: p_{ff}	0.9	0.9
V2V Cond: p_{v2v}	0.9	0.9
ID Cond: p_{id}	0.0	0.5

always enabled in the first training stage since the unconditional case reduces to the base DiT model. We enable identity reference conditioning with probability p_{id} only during the second stage. To support multiple video resolutions, we randomly resize each video clip to a resolution between 0.25 and 2.0 megapixels, then trim the number of frames to maintain an equal number of video tokens per batch. The original frame rate of each video clip is encoded in the 3D RoPE, which is computed separately for each batch element.

We use the Muon Optimizer [52] for matrix-shaped parameters and AdamW [75] for the remaining parameters which improved lip-sync and lowered our flow matching loss substantially in the first stage. We use learning rate lr_{adam} for AdamW and $lr_{\text{muon}} = 100 \times lr_{\text{adam}}$ for Muon and increase lr_{adam} and lr_{muon} by $10\times$ in the second stage. We train for 42 hours on $8\times$ H100 GPUs with a batch size of 4 for each GPU.

To improve training efficiency and stability, we incorporate two complementary techniques: (1) immiscible diffusion with KNN-based noise selection ($k = 4$) [65], which reduces diffusion trajectory mixing and accelerates convergence, and (2) contrastive flow matching [94], which enforces uniqueness across conditional flows to enhance audio-visual correspondence and identity preservation. These techniques work synergistically - immiscible diffusion reduces trajectory miscibility while contrastive flow matching explicitly maximizes dissimilarities between flow from different audio conditions, helping the model better distinguish between different audio features and their corresponding visual representations across the entire denoised region.

4.2. Evaluation

We evaluate our model’s video generation performance in terms of lip synchronization and visual fidelity across both I2V and V2V settings. Extensive qualitative comparisons are provided on our [project page](#), which we recommend viewing to assess results in their native video format. We therefore focus on quantitative evaluation below.

Video-to-Video. To evaluate a re-render of an existing video with a new audio track, we first consider a controlled reconstruction (self-reenactment) setting where ground truth video is available. This setup measures the model’s ability to preserve visual fidelity and speaker identity, including mouth shape, facial details, and temporal dynamics. We evaluate on a subset of 100 videos² from the TalkVid dataset [9], and compare against several state-of-the-art V2V lip-sync systems. These baselines include both open-source research models and widely deployed commercial solutions.

In addition to self-reenactment, we evaluate a re-render with novel audio by pairing each source video with audio from a different TalkVid video, measuring the model’s ability to preserve visual fidelity while accurately synchronizing unseen speech.

We report standard image and video quality metrics FID [37] and FVD [98] as well as identity preservation (CSIM [31]) and lip-sync accuracy (Sync-C and Sync-D [13]). Results are summarized in Table 2.

Image-to-Video. To evaluate performance in the I2V setting, we condition on the first frame of each TalkVid video and the first four seconds of the corresponding audio track. We compare against a range of recent I2V talking-head generation methods, including both open-source research models and commercial systems.

In addition to lip-sync and fidelity metrics, we report VBench [45] evaluation scores for Subject Consistency, Background Consistency, Aesthetic Quality, and Motion Smoothness on a subset of 30 videos using a one-minute audio track.

4.3. Performance Optimizations

With all of the following optimizations enabled, our model renders a 10-second 1080p video in 225 seconds on a single H100 GPU. In comparison, InfiniteTalk, while comparable in image fidelity, requires approximately 10,000 seconds to generate the same video. These results demonstrate that our system achieves practical inference speeds suitable for real-world, long-form video editing workflows.

VAE Tiling and Latent Frame Blocking. For memory optimization, we implement temporal tiling with an overlap of 16 video frames for the VAE encoding and decoding. For denoising, we choose a block size of 17 latent frames (corresponding to 136 video frames), which balances memory efficiency with sufficient temporal context for stable generation. This block structure is consistent with the shifting strategy described in Section 3.4 and enables scalable inference on long sequences.

²For InfiniteTalk we reduce to 20 videos due to long runtimes.

Table 2. Quantitative results on **Video-to-Video** lip-syncing evaluated on the **TalkVid** [9] dataset. We compare methods on **Novel Audio** (audio from a different video) and **Self-Reenactment** (audio from the source video). Metrics include **FID** and **FVD** (image/video fidelity ↓), **CSIM** (identity preservation ↑), and **Sync-C/D** (lip-sync confidence ↑ and distance ↓). **Pose Preservation** indicates if the method retains the original head pose. We highlight **best** and **second best** performance.

Method	FID ↓	FVD ↓	CSIM ↑	Sync-C ↑	Sync-D ↓	Pose Preservation
<i>Video-to-Video (Self-Reenactment)</i>						
<i>Open Source</i>						
LatentSync [62]	46.22	148.19	0.88	7.11	1.559	✓
InfiniteTalk [112]	53.12	304.04	0.86	7.31	1.538	✗
MuseTalk [119]	47.78	134.46	0.82	5.43	1.954	✓
<i>Commercial</i>						
Pixverse V5 Lipsync	45.94	116.54	0.89	5.67	1.880	✓
Sync.so React 1	48.07	161.04	0.86	6.59	1.604	✓
Sync.so V2 Pro	35.51	109.88	0.91	7.43	1.547	✓
Veed Lipsync	53.02	171.94	0.86	6.74	1.637	✓
Creatify Lipsync	64.88	341.47	0.69	7.04	1.585	✓
Ours	37.10	109.04	0.92	7.50	1.480	✓
<i>Video-to-Video (Novel Audio)</i>						
LatentSync [62]	56.84	157.73	0.88	6.90	1.857	✓
InfiniteTalk [112]	42.93	286.03	0.88	6.89	1.568	✗
MuseTalk [119]	49.59	140.49	0.82	5.04	1.922	✓
Ours	41.25	104.18	0.89	7.36	1.502	✓

Table 3. Quantitative results on **Image-to-Video** lip-syncing evaluated on **TalkVid** [9] and **VBench** [45]. We report standard lip-sync metrics (definitions follow Table 2) and general video quality metrics: **Subj./Back.** (subject/background consistency ↑), **Aesth.** (aesthetic quality ↑), and **Motion** (smoothness ↑). We highlight **best** and **second best** performance.

Method	<i>Image-to-Video: TalkVid [9]</i>					<i>Image-to-Video: VBench [45]</i>			
	FID ↓	FVD ↓	CSIM ↑	Sync-C ↑	Sync-D ↓	Subj. ↑	Back. ↑	Aesth. ↑	Motion ↑
<i>Open Source</i>									
Hallo3 [16]	68.11	524.05	0.79	6.25	1.679	0.939	0.941	0.306	0.984
InfiniteTalk [112]	55.47	285.75	0.86	6.75	1.642	0.975	0.950	0.445	0.992
Sonic [47]	63.66	413.54	0.85	6.98	1.629	0.969	0.936	0.400	0.985
StableAvatar [97]	65.56	459.33	0.79	6.39	1.683	0.958	0.936	0.318	0.990
<i>Commercial</i>									
Creatify Aurora	82.83	332.77	0.74	6.53	1.606	0.972	0.937	0.419	0.991
Veed Fabric	68.59	439.78	0.82	6.03	1.898	0.973	0.948	0.495	0.991
Kling Pro V2	72.32	547.28	0.82	6.12	1.843	0.977	0.953	0.472	0.990
OmniHuman V1.5	76.68	363.13	0.78	5.27	2.307	0.960	0.939	0.451	0.989
Fal/HunyuanVideo-Avatar [10]	57.05	383.17	0.81	5.04	2.372	0.968	0.935	0.427	0.986
Fal/MultiTalk [57]	65.69	312.69	0.83	6.50	1.789	0.956	0.929	0.435	0.989
Fal/StableAvatar [97]	68.52	593.48	0.77	4.61	2.939	0.974	0.951	0.459	0.985
Ours	55.38	312.46	0.86	7.21	1.516	0.973	0.953	0.486	0.992

Quantization. To optimize inference speed without compromising quality improvements from the fine-tuning stage, we employ a hybrid FP8 quantization strategy. Specifically, we quantize the large pre-trained base weights in the FFN

and Attention layers, while preserving LoRA adapters and architectural bottlenecks *e.g.* final projection layers and embeddings in BF16 precision. This selective quantization yields a $\times 2.5$ speedup while maintaining output quality.

Hybrid Sequence Parallelism. To eliminate computational bottlenecks in the attention layers during high-resolution generation, we adopt Hybrid Sequence Parallelism by combining Ulysses and Ring Attention [23]. This approach distributes the attention computation across multiple GPUs, improving throughput without increasing memory pressure. When deployed on an $8\times H100$ GPU node, this strategy provides an aggregate $\times 8$ speedup within the denoising loop.

4.4. Ablation Study

What happens if Identity conditioning is dropped?

Without reference tokens, the subject’s appearance drifts away rather quickly from its original state since the DiT can access the first-frame condition only in the first and last block (due to circular padding [47]). All intermediate inference blocks have no direct access to clean lower face tokens as a reference. Our TAPSF long inference strategy softens the impact of drift as appearance information is shared between neighboring blocks over the course of the denoising process, which finally results in a smooth appearance drift that is most obvious in the middle of the video (Fig. 12, V2V-Frame72). In the I2V case (latent frames are fully noised), this drift becomes more severe, resulting in a complete change of appearance and scene (Fig. 12, I2V-Frame72). Fig. 12 visualizes the benefit of using face (FR) and full-frame (FF) reference tokens during inference. Pure V2V exhibits a clear identity drift (beard growth) after just 72 frames, whereas V2V+FR maintains the original identity well throughout the video (Fig. 12). While in the video-to-video case, only the person’s face and small background details are affected, fully synthetic sections (e.g., long additions or completely re-rendered sections) suffer from obvious discontinuities even at the scene level (Fig. 12, I2V-Frame72). I2V+FR maintains the appearance of the lower face, but shows a mild scene drift (curtains, hair, eyebrows) since the face reference tokens contain only the lower face and a small portion of the background. I2V+FR+FF maintains full scene and identity consistency while being able to fully re-render the captured performance with novel head poses and facial expressions (eyebrows, forehead), see Fig. 12 (bottom row).

Training with/without Identity Reference Condition.

While the DiT is capable of utilizing face reference tokens without having been exposed to them during training, incorporating reference tokens during training with a probability of 50% leads to improved rendering quality, particularly in complex scenes with highly structured and dynamic backgrounds (Fig. 13).

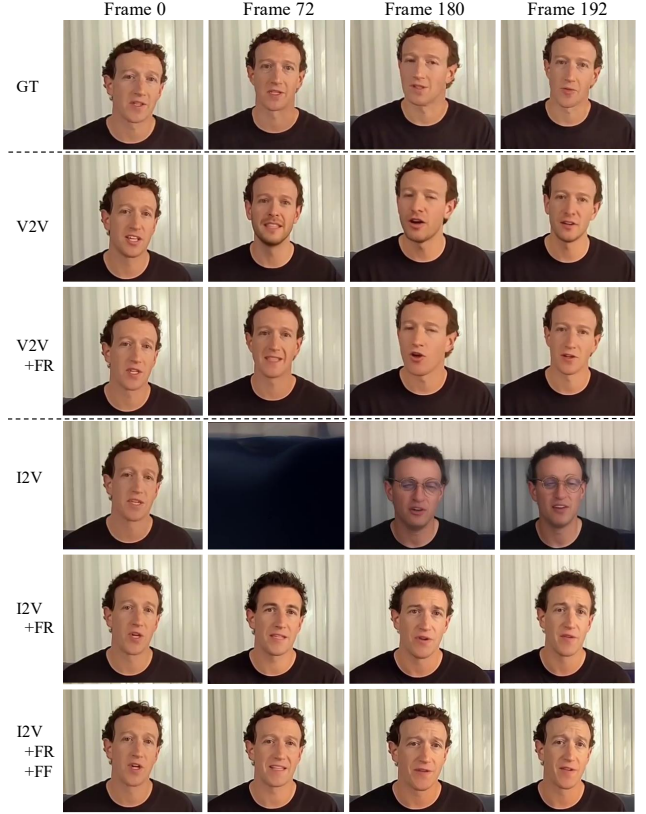


Figure 12. Ablation of the reference frame conditioning. From left to right, we show 4 representative frames from an 8-second clip rendered with different reference conditioning variants. From top to bottom: groundtruth (GT), pure video-to-video (V2V), video-to-video with face reference tokens (V2V+FR), pure image-to-video (I2V), image-to-video with face reference tokens (I2V+FR), image-to-video with face and full-frame reference tokens (I2V+FR+FF).

5. Conclusion

In summary, we present EditYourself, a diffusion-based framework for audio-driven talking head synthesis that extends a general-purpose video diffusion model with audio-driven V2V editing capabilities. Through a two-stage training scheme and a windowed audio conditioning strategy, our approach enables precise lip synchronization while preserving visual fidelity to the original video content. We further introduce Forward-Backward RoPE Conditioning to maintain stable identity and appearance over extended durations. By operating directly in latent space, EditYourself enables transcript-based modification of videos, including addition, removal and retiming, offering a practical step toward using generative video models as tools for professional post-production, alongside their role in end-to-end synthesis.

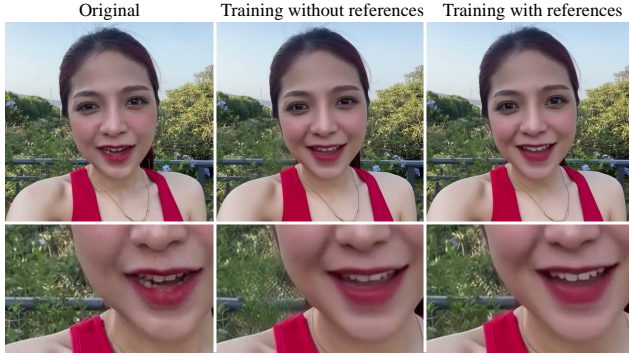


Figure 13. The difference in render quality with and without exposing the DiT to reference tokens during training. Left: ground truth, middle: no reference tokens during training, right: with reference tokens during training. Without training for reference tokens at sentinel timesteps, render artifacts appear, particularly with complex and dynamic backgrounds.

5.1. Ethical Considerations

EditYourself’s capacity for high-fidelity synthesis and the granular manipulation of existing footage necessitates careful consideration of potential misuse, particularly in the context of visual forgeries and the dissemination of misinformation. We emphasize that responsibility for the ethical use of these techniques is shared across the research community, including those who deploy or build upon methods introduced in this work. We therefore recommend a multi-layered approach to responsible deployment: (1) Establishing legal barriers such as explicit declarations of content ownership, and (2) implementing technical safeguards including celebrity detection, identity verification and robust digital watermarking. Furthermore, the authors strongly advocate for continued research into content provenance and synthetic media detection to mitigate risks associated with unauthorized generation and to ensure the ethical evolution of generative video tools.

5.2. Acknowledgements

The authors would like to thank the Lightricks LTX-Video team for open-sourcing their model weights, and specifically Ofir Bibi, Yoav HaCohen, and Nisan Chirput for their technical insights during the development of this work.

References

- [1] Amazon Web Services, Inc. Amazon transcribe. <https://aws.amazon.com/transcribe/>, 2025. Software. 6
- [2] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020. 3
- [3] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, Yuanzhen Li, Michael Rubinstein, Tomer Michaeli, Oliver Wang, Deqing Sun, Tali Dekel, and Inbar Mosseri. Lumiere: A space-time diffusion model for video generation, 2024. 3
- [4] Floraine Berthouzoz, Wilmot Li, and Maneesh Agrawala. Tools for placing cuts and transitions in interview video. *ACM Trans. Graph.*, 31(4), 2012. 2
- [5] Antoni Bigata, Rodrigo Mira, Stella Bounareli, Michał Stypułkowski, Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Keysync: A robust approach for leakage-free lip synchronization in high resolution, 2025. 3
- [6] Antoni Bigata, Michał Stypułkowski, Rodrigo Mira, Stella Bounareli, Konstantinos Vougioukas, Zoe Landgraf, Nikita Drobyshev, Maciej Zieba, Stavros Petridis, and Maja Pantic. Keyface: Expressive audio-driven facial animation for long sequences via keyframe interpolation, 2025. 3
- [7] Volker Blanz and Thomas Vetter. Face recognition based on fitting a 3d morphable model. *IEEE Transactions on pattern analysis and machine intelligence*, 25(9):1063–1074, 2003. 3
- [8] Lele Chen, Zhiheng Li, Ross K. Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *ECCV*, 2018. 3
- [9] Shunian Chen, Hejin Huang, Yexin Liu, Zihan Ye, Pengcheng Chen, Chenghao Zhu, Michael Guan, Rongsheng Wang, Junying Chen, Guanbin Li, Ser-Nam Lim, Harry Yang, and Benyou Wang. Talkvid: A large-scale diversified dataset for audio-driven talking head synthesis, 2025. 10, 11
- [10] 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, 2025. 3, 4, 7, 11
- [11] Zhiyuan Chen, Jiajiong Cao, Zhiqian Chen, Yuming Li, and Chenguang Ma. Echomimic: Lifelike audio-driven portrait animations through editable landmark conditioning, 2024. 2, 3
- [12] Kyusun Cho, Jounghbin Lee, Heeji Yoon, Yeobin Hong, Jaehoon Ko, Sangjun Ahn, and Seungryong Kim. Gaussian talker: Real-time talking head synthesis with 3d gaussian splatting. In *MM*, 2024. 3
- [13] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In *Workshop on Multi-view Lip-reading, ACCV*, 2016. 9, 10
- [14] Nathaniel Cohen, Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Slicedit: Zero-shot video editing with text-to-image diffusion models using spatio-temporal slices. In *ICML*, 2024. 3
- [15] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael Black. Capture, learning, and synthesis of 3D speaking styles. In *CVPR*, 2019. 3
- [16] 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*, 2025. 2, 3, 4, 11

- [17] Deepdub AI. Deepdub: The virtual ai studio. <https://deepdub.ai>, 2023. Software. 6
- [18] Deepgram. Deepgram: Ai speech-to-text. <https://deepgram.com>, 2025. Software. 6
- [19] Jiankang Deng, Jia Guo, Jing Yang, Niannan Xue, Irene Kotsia, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. *IEEE TPAMI*, 44(10):5962–5979, 2022. 4, 8
- [20] Fangyu Du, Taiqing Li, Ziwei Zhang, Qian Qiao, Tan Yu, Dingcheng Zhen, Xu Jia, Yang Yang, Shunshun Yin, and Siyuan Liu. Rap: Real-time audio-driven portrait animation with video diffusion transformer, 2025. 3, 4
- [21] ElevenLabs. Elevenlabs: Prime voice ai. <https://elevenlabs.io>, 2023. Software. 6
- [22] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 3
- [23] Jiarui Fang and Shangchun Zhao. Usp: A unified sequence parallelism approach for long context generative ai, 2024. 12
- [24] Ruoyu Feng, Wenming Weng, Yanhui Wang, Yuhui Yuan, Jianmin Bao, Chong Luo, Zhibo Chen, and Baining Guo. Ccredit: Creative and controllable video editing via diffusion models. In *CVPR*, 2024. 3
- [25] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. Text-based editing of talking-head video. *ACM Trans. Graph.*, 38(4):68:1–68:14, 2019. 2, 3
- [26] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Niessner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, 2021. 3
- [27] Qijun Gan, Yi Ren, Chen Zhang, Zhenhui Ye, Pan Xie, Xiang Yin, Zehuan Yuan, Bingyue Peng, and Jianke Zhu. Humandit: Pose-guided diffusion transformer for long-form human motion video generation, 2025. 3
- [28] 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. 4
- [29] Qijun Gan, Ruizi Yang, Jianke Zhu, Shaofei Xue, and Steven Hoi. Omniavatar: Efficient audio-driven avatar video generation with adaptive body animation, 2025. 3
- [30] Xin Gao, Li Hu, Siqi Hu, Mingyang Huang, Chaonan Ji, Dechao Meng, Jinwei Qi, Penchong Qiao, Zhen Shen, Yafei Song, Ke Sun, Linrui Tian, Guangyuan Wang, Qi Wang, Zhongjian Wang, Jiayu Xiao, Sheng Xu, Bang Zhang, Peng Zhang, Xindi Zhang, Zhe Zhang, Jingren Zhou, and Lian Zhuo. Wan-s2v: Audio-driven cinematic video generation, 2025. 3
- [31] Safouane El Ghazouali, Umberto Michelucci, Yassin El Hillali, and Hichem Noura. Csim: A copula-based similarity index sensitive to local changes for image quality assessment, 2024. 10
- [32] Jiazhi Guan, Zhanwang Zhang, Hang Zhou, Tianshu HU, Kaisiyuan Wang, Dongliang He, Haocheng Feng, Jingtuo Liu, Errui Ding, Ziwei Liu, and Jingdong Wang. Stylesync: High-fidelity generalized and personalized lip sync in style-based generator. In *CVPR*, 2023. 3
- [33] Jiazhi Guan, Zhiliang Xu, Hang Zhou, Kaisiyuan Wang, Shengyi He, Zhanwang Zhang, Borong Liang, Haocheng Feng, Errui Ding, Jingtuo Liu, Jingdong Wang, Youjian Zhao, and Ziwei Liu. Resyncer: Rewiring style-based generator for unified audio-visually synced facial performer. In *ECCV*, 2024. 3
- [34] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, 2021. 3
- [35] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, Poriya Panet, Sapir Weissbuch, Victor Kulikov, Yaki Bitterman, Zeev Melumian, and Ofir Bibi. Ltx-video: Realtime video latent diffusion, 2024. 2, 3, 4, 5, 9
- [36] Bo Han, Heqing Zou, Haoyang Li, Guangcong Wang, and Chng Eng Siong. Text-based talking video editing with cascaded conditional diffusion, 2024. 3
- [37] 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, 2018. 10
- [38] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models, 2022. 2, 3
- [39] Fa-Ting Hong, Zunnan Xu, Zixiang Zhou, Jun Zhou, Xiu Li, Qin Lin, Qinglin Lu, and Dan Xu. Audio-visual controlled video diffusion with masked selective state spaces modeling for natural talking head generation. *arXiv preprint arXiv:2504.02542*, 2025. 4
- [40] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. In *ICLR*, 2023. 3
- [41] Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024. 9
- [42] Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone: Consistent and controllable image-to-video synthesis for character animation, 2023. 3
- [43] Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the train-test gap in autoregressive video diffusion. In *NeurIPS*, 2025. 4
- [44] Yubo Huang, Hailong Guo, Fangtai Wu, Shifeng Zhang, Shijie Huang, Qijun Gan, Lin Liu, Sirui Zhao, Enhong Chen, Jiaming Liu, and Steven Hoi. Live avatar: Streaming real-time audio-driven avatar generation with infinite length, 2025. 9

- [45] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Vbench: Comprehensive benchmark suite for video generative models, 2023. 10, 11
- [46] Mina Huh, Ding Li, Kim Pimmel, Hijung Valentina Shin, Amy Pavel, and Mira Dontcheva. Videodiff: Human-ai video co-creation with alternatives. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, 2025. Association for Computing Machinery. 2
- [47] Xiaozhong Ji, Xiaobin Hu, Zhihong Xu, Junwei Zhu, Chuming Lin, Qingdong He, Jiangning Zhang, Donghao Luo, Yi Chen, Qin Lin, et al. Sonic: Shifting focus to global audio perception in portrait animation. In *CVPR*, 2025. 2, 3, 4, 7, 11, 12
- [48] Jianwen Jiang, Chao Liang, Jiaqi Yang, Gaojie Lin, Tianyun Zhong, and Yanbo Zheng. Loopy: Taming audio-driven portrait avatar with long-term motion dependency. *arXiv preprint arXiv:2409.02634*, 2024. 4
- [49] Jianwen Jiang, Weihong Zeng, Zerong Zheng, Jiaqi Yang, Chao Liang, Wang Liao, Han Liang, Yuan Zhang, and Mingyuan Gao. Omnihuman-1.5: Instilling an active mind in avatars via cognitive simulation, 2025. 3
- [50] Zeyinzi Jiang, Zhen Han, Chaojie Mao, Jingfeng Zhang, Yulin Pan, and Yu Liu. Vace: All-in-one video creation and editing. In *ICCV*, 2025. 3
- [51] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 3
- [52] Keller Jordan, Yuchen Jin, Vlado Boza, You Jiacheng, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. 10
- [53] Prajwal K R, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and C V Jawahar. Towards automatic face-to-face translation. In *MM*, 2019. 3
- [54] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020. 3
- [55] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 3
- [56] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2, 3
- [57] 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, 2025. 3, 8, 11
- [58] Juil Koo, Paul Guerrero, Chun-Hao Paul Huang, Duygu Ceylan, and Minhyuk Sung. Videohandles: Editing 3d object compositions in videos using video generative priors. In *CVPR*, 2025. 3
- [59] Chieh-Hsin Lai, Yang Song, Dongjun Kim, Yuki Mitsufuji, and Stefano Ermon. The principles of diffusion models, 2025. 2
- [60] Mackenzie Leake and Wilmot Li. Chunkyedit: Text-first video interview editing via chunking. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, 2024. Association for Computing Machinery. 2
- [61] Chenyang Lei, Yazhou Xing, Hao Ouyang, and Qifeng Chen. Deep video prior for video consistency and propagation, 2022. 3
- [62] Chunyu Li, Chao Zhang, Weikai Xu, Jingyu Lin, Jinghui Xie, Weiguo Feng, Bingyue Peng, Cunjian Chen, and Weiwei Xing. Latentsync: Taming audio-conditioned latent diffusion models for lip sync with syncnet supervision. *arXiv preprint arXiv:2412.09262*, 2024. 2, 11
- [63] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Talkinggaussian: Structure-persistent 3d talking head synthesis via gaussian splatting. In *ECCV*, 2024. 3
- [64] Wuyang Li, Wentao Pan, Po-Chien Luan, Yang Gao, and Alexandre Alahi. Stable video infinity: Infinite-length video generation with error recycling, 2025. 3, 4
- [65] Yiheng Li, Feng Liang, Dan Kondratyuk, Masayoshi Tomizuka, Kurt Keutzer, and Chenfeng Xu. Improved immiscible diffusion: Accelerate diffusion training by reducing its miscibility. *arXiv preprint arXiv:2505.18521*, 2025. 10
- [66] Chao Liang, Jianwen Jiang, Wang Liao, Jiaqi Yang, Zerong zheng, Weihong Zeng, and Han Liang. Alignhuman: Improving motion and fidelity via timestep-segment preference optimization for audio-driven human animation, 2025. 3, 6
- [67] Feng Liang, Bichen Wu, Jialiang Wang, Licheng Yu, Kunpeng Li, Yinan Zhao, Ishan Misra, Jia-Bin Huang, Peizhao Zhang, Peter Vajda, et al. Flowvid: Taming imperfect optical flows for consistent video-to-video synthesis. In *CVPR*, 2014. 3
- [68] Feng Liang, Akio Kodaira, Chenfeng Xu, Masayoshi Tomizuka, Kurt Keutzer, and Diana Marculescu. Looking backward: Streaming video-to-video translation with feature banks. In *ICLR*, 2025. 3
- [69] Lightricks. Ltx-video. <https://github.com/Lightricks/LTX-Video>, 2024. GitHub repository. 4, 9
- [70] Gaojie Lin, Jianwen Jiang, Jiaqi Yang, Zerong Zheng, and Chao Liang. Omnihuman-1: Rethinking the scaling-up of one-stage conditioned human animation models, 2025. 3, 4, 8
- [71] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023. 2, 3, 4
- [72] Feng Liu, Shiwei Zhang, Xiaofeng Wang, Yujie Wei, Haonan Qiu, Yuzhong Zhao, Yingya Zhang, Qixiang Ye, and Fang Wan. Timestep embedding tells: It's time to cache for video diffusion model. *arXiv preprint arXiv:2411.19108*, 2024. 7
- [73] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, Lifang He, and Lichao Sun. Sora: A review

on background, technology, limitations, and opportunities of large vision models, 2024. 3

- [74] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019. 3
- [75] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017. 10
- [76] Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays, Fan Zhang, Chuoling Chang, Ming Yong, Juhyun Lee, Wan-Teh Chang, Wei Hua, Manfred Georg, and Matthias Grundmann. Mediapipe: A framework for perceiving and processing reality. In *Third Workshop on Computer Vision for AR/VR at IEEE Computer Vision and Pattern Recognition (CVPR) 2019*, 2019. 5
- [77] Xingpei Ma, Jiaran Cai, Yuansheng Guan, Shenneng Huang, Qiang Zhang, and Shunsi Zhang. Playmate: Flexible control of portrait animation via 3d-implicit space guided diffusion. In *ICML*, 2025. 3
- [78] Yue Ma, Kunyu Feng, Zhongyuan Hu, Xinyu Wang, Yucheng Wang, Mingzhe Zheng, Xuanhua He, Chenyang Zhu, Hongyu Liu, Yingqing He, Zeyu Wang, Zhifeng Li, Xiu Li, Wei Liu, Dan Xu, Linfeng Zhang, and Qifeng Chen. Controllable video generation: A survey, 2025. 2
- [79] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 3
- [80] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha, Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid Hoshen. Dreamix: Video diffusion models are general video editors, 2023. 2, 3
- [81] Mang Ning, Mingxiao Li, Jianlin Su, Albert Ali Salah, and Itir Onal Ertugrul. Elucidating the exposure bias in diffusion models. In *ICLR*, 2024. 4
- [82] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023. 8
- [83] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023. 3
- [84] Ziqiao Peng, Jiwen Liu, Haoxian Zhang, Xiaoqiang Liu, Songlin Tang, Pengfei Wan, Di Zhang, Hongyan Liu, and Jun He. Omnisync: Towards universal lip synchronization via diffusion transformers. In *NeurIPS*, 2025. 3, 4, 6
- [85] 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 *MM*, 2020. 2, 3, 9
- [86] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. 4, 8
- [87] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. 4
- [88] Siddarth Ravichandran, Ondřej Texler, Dimitar Dinev, and Hyun Jae Kang. Synthesizing photorealistic virtual humans through cross-modal disentanglement. In *CVPR*, 2023. 3
- [89] Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando de la Torre, and Yaser Sheikh. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In *ICCV*, 2021. 3
- [90] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3
- [91] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 3
- [92] Junyoung Seo, Rodrigo Mira, Alexandros Haliassos, Stella Bounareli, Honglie Chen, Linh Tran, Seungryong Kim, Zoe Landgraf, and Jie Shen. Lookahead anchoring: Preserving character identity in audio-driven human animation. *arXiv preprint arXiv:2510.23581*, 2025. 4, 9
- [93] Uriel Singer, Amit Zohar, Yuval Kirstain, Shelly Sheynin, Adam Polyak, Devi Parikh, and Yaniv Taigman. Video editing via factorized diffusion distillation. In *ECCV*, 2024. 3
- [94] George Stoica, Vivek Ramanujan, Xiang Fan, Ali Farhadi, Ranjay Krishna, and Judy Hoffman. Contrastive flow matching. *arXiv preprint arXiv:2506.05350*, 2025. 10
- [95] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. 4
- [96] Xingye Tian, Tianxiong Zhong, Xuebo Wang, Boyuan Jiang, Pengfei Wan, Xin Tao, and Zhiwei Zhang. Vfrtok: Variable frame rates video tokenizer with duration-proportional information assumption. In *NeurIPS*, 2025. 3
- [97] 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, 2025. 3, 8, 11
- [98] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges, 2019. 10
- [99] Team Wan et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 3, 4, 6
- [100] Ge Wang, Songlin Fan, Hangxu Liu, Qianjian Song, Hewei Wang, and Jinfeng Xu. Consistent video editing as flow-driven image-to-video generation, 2025. 3
- [101] 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. In *ACM MM*, 2025. 3, 4, 8

- [102] Sitong Wang, Zheng Ning, Anh Truong, Mira Dontcheva, Dingzeyu Li, and Lydia B. Chilton. Podreels: Human-ai co-creation of video podcast teasers, 2024. 2
- [103] 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. 3
- [104] Zhongjian Wang, Peng Zhang, Jinwei Qi, Guangyuan Wang, Chaonan Ji, Sheng Xu, Bang Zhang, and Liefeng Bo. Omnitalter: One-shot real-time text-driven talking audio-video generation with multimodal style mimicking. In *NeurIPS*, 2025. 3
- [105] Cong Wei, Bo Sun, Haoyu Ma, Ji Hou, Felix Juefei-Xu, Zecheng He, Xiaoliang Dai, Luxin Zhang, Kunpeng Li, Tingbo Hou, et al. Mocha: Towards movie-grade talking character synthesis. *arXiv preprint arXiv:2503.23307*, 2025. 3, 4
- [106] Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners, 2025. 3
- [107] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation, 2022. 3
- [108] Xian Wu and Chang Liu. Ditpainter: Efficient video inpainting with diffusion transformers, 2025. 3
- [109] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. In *CVPR*, 2023. 3
- [110] Zunnan Xu, Zhentao Yu, Zixiang Zhou, Jun Zhou, Xiaoyu Jin, Fa-Ting Hong, Xiaozhong Ji, Junwei Zhu, Chengfei Cai, Shiyu Tang, et al. Hunyuanportrait: Implicit condition control for enhanced portrait animation. In *CVPR*, 2025. 3, 4, 8
- [111] Songlin Yang, Wei Wang, Jun Ling, Bo Peng, Xu Tan, and Jing Dong. Context-aware talking-head video editing. In *MM*, 2023. 3
- [112] Shaoshu Yang, Zhe Kong, Feng Gao, Meng Cheng, Xiangyu Liu, Yong Zhang, Zhuoliang Kang, Wenhan Luo, Xunliang Cai, Ran He, and Xiaoming Wei. Infinitetalk: Audio-driven video generation for sparse-frame video dubbing, 2025. 2, 3, 8, 11
- [113] Shunyu Yao, RuiZhe Zhong, Yichao Yan, Guangtao Zhai, and Xiaokang Yang. Dfa-nerf: Personalized talking head generation via disentangled face attributes neural rendering, 2022. 3
- [114] Xinwei Yao, Ohad Fried, Kayvon Fatahalian, and Maneesh Agrawala. Iterative text-based editing of talking-heads using neural retargeting. *ACM Trans. Graph.*, 40(3), 2021. 3
- [115] 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*, 2023. 3
- [116] Shenghai Yuan, Jinfa Huang, Xianyi He, Yunyang Ge, Yujun Shi, Liuhan Chen, Jiebo Luo, and Li Yuan. Identity-preserving text-to-video generation by frequency decomposition. In *CVPR*, 2025. 4
- [117] Kaidong Zhang, Jingjing Fu, and Dong Liu. Flow-guided transformer for video inpainting. In *ECCV*, 2022. 3
- [118] Lvmin Zhang, Shengqu Cai, Muiyang Li, Gordon Wetstein, and Maneesh Agrawala. Frame context packing and drift prevention in next-frame-prediction video diffusion models. In *NeurIPS*, 2025. 4
- [119] Yue Zhang, Zhizhou Zhong, Minhao Liu, Zhaokang Chen, Bin Wu, Yubin Zeng, Chao Zhan, Yingjie He, Junxin Huang, and Wenjiang Zhou. Musetalk: Real-time high-fidelity video dubbing via spatio-temporal sampling. *arXiv preprint arXiv:2410.10122*, 2024. 2, 11
- [120] Yuechen Zhang, Yaoyang Liu, Bin Xia, Bohao Peng, Zexin Yan, Eric Lo, and Jiaya Jia. Magic mirror: Id-preserved video generation in video diffusion transformers. In *ICCV*, 2025. 4
- [121] Zhixing Zhang, Bichen Wu, Xiaoyan Wang, Yaqiao Luo, Luxin Zhang, Yinan Zhao, Peter Vajda, Dimitris Metaxas, and Licheng Yu. Avid: Any-length video inpainting with diffusion model. In *CVPR*, 2024. 3
- [122] Longtao Zheng, Yifan Zhang, Hanzhong Guo, Jiachun Pan, Zhenxiong Tan, Jiahao Lu, Chuanxin Tang, Bo An, and Shuicheng Yan. Memo: Memory-guided diffusion for expressive talking video generation, 2024. 2, 3
- [123] Shangchen Zhou, Chongyi Li, Kelvin C.K. Chan, and Chen Change Loy. ProPainter: Improving propagation and transformer for video inpainting. In *ICCV*, 2023. 3
- [124] Xin Zhou, Dingkan Liang, Kaijin Chen, , Tianrui Feng, Xiwu Chen, Hongkai Lin, Yikang Ding, Feiyang Tan, Hengshuang Zhao, and Xiang Bai. Less is enough: Training-free video diffusion acceleration via runtime-adaptive caching. *arXiv preprint arXiv:2507.02860*, 2025. 7
- [125] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jia-ashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. *Advances in Neural Information Processing Systems*, 37: 110315–110340, 2024. 4