

TV2TV: A Unified Framework for Interleaved Language and Video Generation

Xiaochuang Han^{*,†}, Youssef Emad^{*,◇}, Melissa Hall^{*,◇}, John Nguyen^{*,◇}, Karthik Padthe^{*,◇}, Liam Robbins^{*,◇}, Amir Bar, Delong Chen, Michal Drozdal, Maha Elbayad, Yushi Hu, Shang-Wen Li, Sreya Dutta Roy, Jakob Verbeek, XuDong Wang, Marjan Ghazvininejad, Luke Zettlemoyer, Emily Dinan^{*}

Meta FAIR

^{*}Core contributors, [◇]Ordered alphabetically

Video generation models are rapidly advancing, but can still struggle with complex video outputs that require significant semantic branching or repeated high-level reasoning about what should happen next. In this paper, we introduce a new class of omni video-text models that integrate ideas from recent LM reasoning advances to address this challenge. More specifically, we present TV2TV, a unified generative modeling framework which decomposes video generation into an interleaved text and video generation process. TV2TV jointly learns language modeling (next-token prediction) and video flow matching (next-frame prediction) using a Mixture-of-Transformers (MoT) architecture. At inference time, TV2TV decides when to alternate between generating text and video frames, allowing the model to “think in words” about subsequent content before “acting in pixels” to produce frames. This design offloads much of the responsibility for deciding what should happen next to the language modeling tower, enabling improved visual quality and prompt alignment of generated videos. It also enables fine-grained controllability, allowing users to modify the video generation trajectory through text interventions at any point in the process. In controlled experiments on video game data, TV2TV demonstrates substantial improvements in both visual quality (preferred 91% of the time in human evaluations vs. a comparable text-to-video model) and controllability (19 point improvement in fine-grained instruction following accuracy vs. a “think-then-act” approach). TV2TV also scales to natural videos, as we show by augmenting sports videos with interleaved natural language action descriptions using vision-language models (VLMs). Training TV2TV on this corpus yields strong visual quality and prompt alignment, showcasing the model’s ability to reason about and generate complex real-world action sequences. Together, these results highlight TV2TV as a promising step toward video generation with open-ended textual reasoning and control.

[†]Correspondence: Xiaochuang Han xhan77@meta.com



1 Introduction

Despite incredible progress in visual quality, video generation models can still struggle with complex outputs that require significant semantic branching or repeated high-level reasoning about what should happen next. In this paper, we introduce a new class of omni video-text models that integrate ideas from recent LM reasoning advances to address this challenge. Our approach generalizes previous omni models that have focused on text and image modalities (Chameleon Team, 2024; Zhou et al., 2024; Wu et al., 2024; Deng et al., 2025; Li et al., 2025) as well as interactive video generations models like Genie (Bruce et al., 2024) that require explicit user input at each step. We instead show that it is possible to train an omni model which automatically decomposes video generation into an interleaved text and video generation process, thereby significantly improving quality and controllability.

We present TV2TV, which is a Transfusion-style modeling approach (Zhou et al., 2024) that jointly learns language modeling (next-token prediction) and video flow matching (Liu et al., 2022; Esser et al., 2024) (next-frame prediction). At inference time, TV2TV dynamically alternates between generating text and generating chunks of video frames, allowing the model to *think in words* about the content of the subsequent frames before *acting in pixels* to produce those frames. This approach offloads much of the semantic decision-making to the

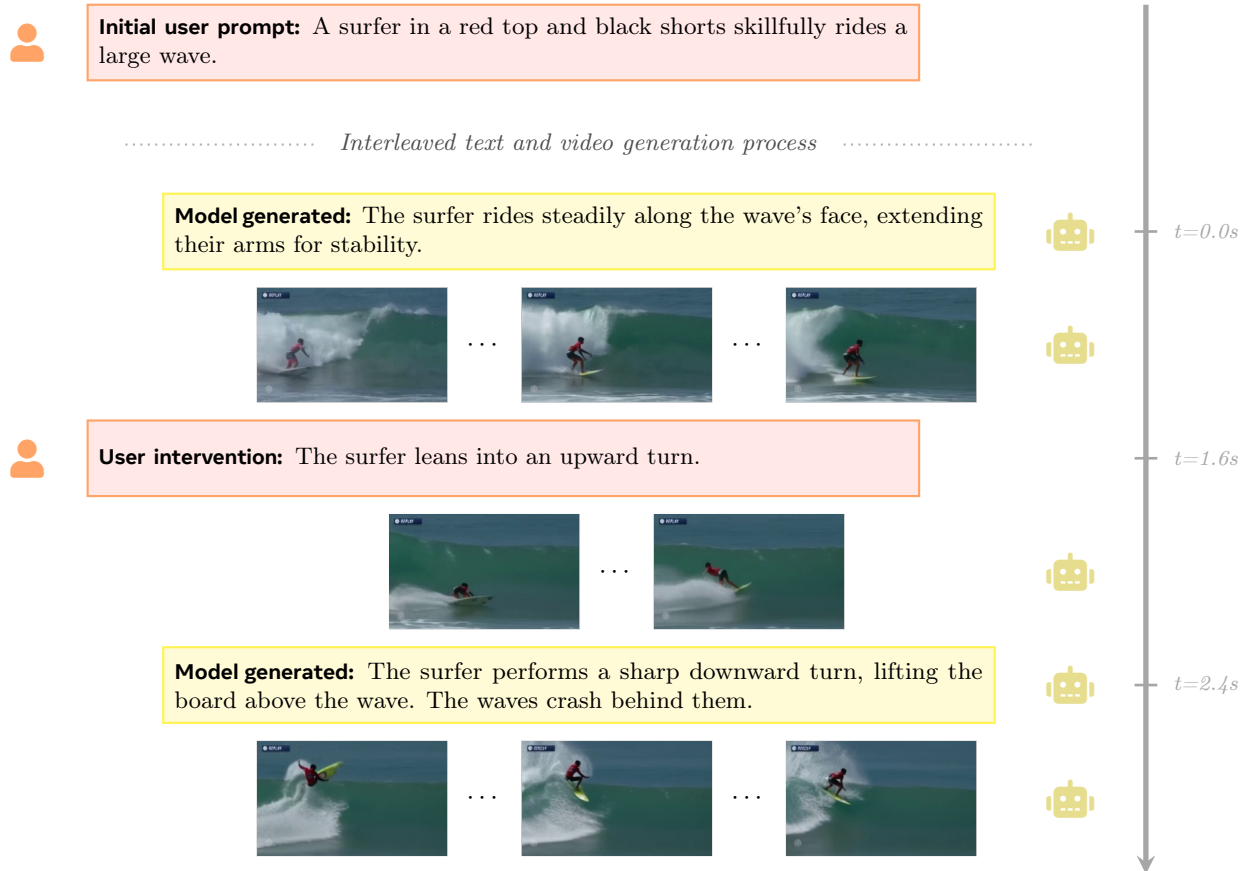


Figure 1 Overview of TV2TV interleaved text and video generation. TV2TV is a unified generative modeling framework which decomposes video generation into an interleaved text and video generation process. During inference, TV2TV dynamically alternates between autoregressively generating plans in text and semi-autoregressively generating chunks of video frames, allowing the model to *think in words* about the content of the subsequent frames before *acting in pixels* to produce those frames. This framework enables fine-grained and flexible control during video generation, allowing users to potentially intervene and modify the video generation trajectory at any point through textual prompting.

language modeling component of the model, capitalizing on recent advances in LLM reasoning capabilities and reducing the entropy of the video generation process. The design also affords strong controllability: users can inspect, edit, or steer the textual plan to modify the video generation trajectory at any timestep.

TV2TV adopts a Mixture-of-Transformers (MoT) (Liang et al., 2024) architecture with dedicated towers for video and text modalities, enabling modality-specific processing while maintaining a global self-attention over the entire multimodal input sequence. Following LMFusion (Shi et al., 2024), we initialize the text tower from a pre-trained language model. TV2TV is trained on interleaved sequences of text and temporally segmented chunks of video frames. We employ bi-directional attention within video frame chunks and causal attention otherwise.

We evaluate performance in two domains. First, we use video game data from *Counter Strike; Global Offensive* (CS:GO) curated and open-sourced by Pearce and Zhu (2022), which provides strongly correlated interleaved data via controller actions (represented as text) and resulting video gameplay. Controlled experiments on this data show that TV2TV substantially outperforms competing approaches in both visual quality – preferred 92% of the time over a comparable T2V baseline – and controllability, with a 19 point improvement in fine-grained instruction-following accuracy compared to a “think-then-act” (Think2V) method which generates a detailed text action plan prior to generating video. Relatedly, action-conditioned world models (e.g., Bruce et al., 2024) can also generate game or synthetic video-world data. However, having them automatically generate videos without dense human control requires either a separate controller model (e.g., Raad et al., 2024) or

a costly planning algorithm (e.g., Bar et al., 2025). Our TV2TV approach performs action generation and video-world generation flexibly and end-to-end.

Next, we demonstrate how to scale this modeling paradigm to real-world videos which typically lack such temporally-aligned text data, by augmenting video data with interleaved natural language action descriptions using vision-language models (VLMs). Using such a pipeline, we curate 8K hours of interleaved text and video content in the sports domain, chosen for its dynamic motion and rich action content, and demonstrate that training TV2TV on this corpus yields strong visual quality and prompt alignment relative to both established external video generation models and controlled T2V and Think2V baselines. In holistic preference evaluations, TV2TV is favored 54.0% of the time compared to T2V (34.7% unfavorable) and 53.3% of the time compared to Think2V (41.3% unfavorable). We also present qualitative examples showcasing how users can dynamically steer video generation through intermediate textual prompting.

Together, these experiments highlight TV2TV as a promising step towards unifying advances in language model reasoning with highly controllable video generation systems, leveraging natural language not merely as input conditioning, but as an active reasoning mechanism for decomposing complex visual and temporal tasks.

Our contributions are summarized as follows:

- **TV2TV (§2)**: We introduce TV2TV, a unified generative modeling framework capable of decomposing video generation into an interleaved text and video generation process.
- **Controlled Experiments with Video Game Data (§3)** We validate this approach through controlled experiments with video game (CS:GO) data, demonstrating that TV2TV outperforms T2V and Think2V baselines in both visual quality and controllability. Specifically, in pairwise human evaluations, videos generated by TV2TV are preferred to those generated by a T2V baseline 92% of the time, and TV2TV shows a 19 point improvement in fine-grained instruction following accuracy compared to a Think2V baseline.
- **Scaling TV2TV to Real World Data (§4)** Finally, we scale this paradigm to real world video data by synthetically augmenting sport video data with interleaved captions. Training TV2TV on this corpus yields strong prompt alignment and visual quality relative to both external and controlled baselines (wins 54.0% vs. 34.7% and 53.3% vs. 41.3% against comparable T2V and Think2V baselines, respectively), showcasing TV2TV’s ability to seamlessly reason about and generate complex visual action sequences.

2 TV2TV

We present TV2TV, a unified modeling framework which decomposes video generation into an interleaved text and video generation process.

TV2TV builds on the Transfusion-style (Zhou et al., 2024) approach, jointly learning language modeling (next-token prediction) and video flow matching (next-frame prediction). At a global level, TV2TV autoregressively generates interleaved chunks of video frames and text tokens, maintaining strict temporal causality: each token or chunk of frames conditions only on preceding tokens or chunks. At the local level, within each chunk comprising one or several video frames, the model is non-autoregressive and uses a flow matching objective. This autoregressive video generation component, when considered independently of the interleaved text, follows a design similar to MAGI-1 (Teng et al., 2025).

TV2TV adopts a Mixture-of-Transformers (MoT) (Liang et al., 2024) architecture with dedicated towers for the video and text modalities, enabling modality-specific processing while maintaining a global self-attention over the entire multimodal input sequence. Following Shi et al. (2024), we initialize the text tower with a pretrained language model. At inference, TV2TV flexibly alternates between text and video generation, factorizing the video generation into interleaved textual planning and video segment generation. See Figure 2 for an overview.

In the following section, we detail TV2TV’s interleaved data representation (§2.1), architecture and optimization (§2.2), inference procedure (§2.3), and task formulation (§2.4).

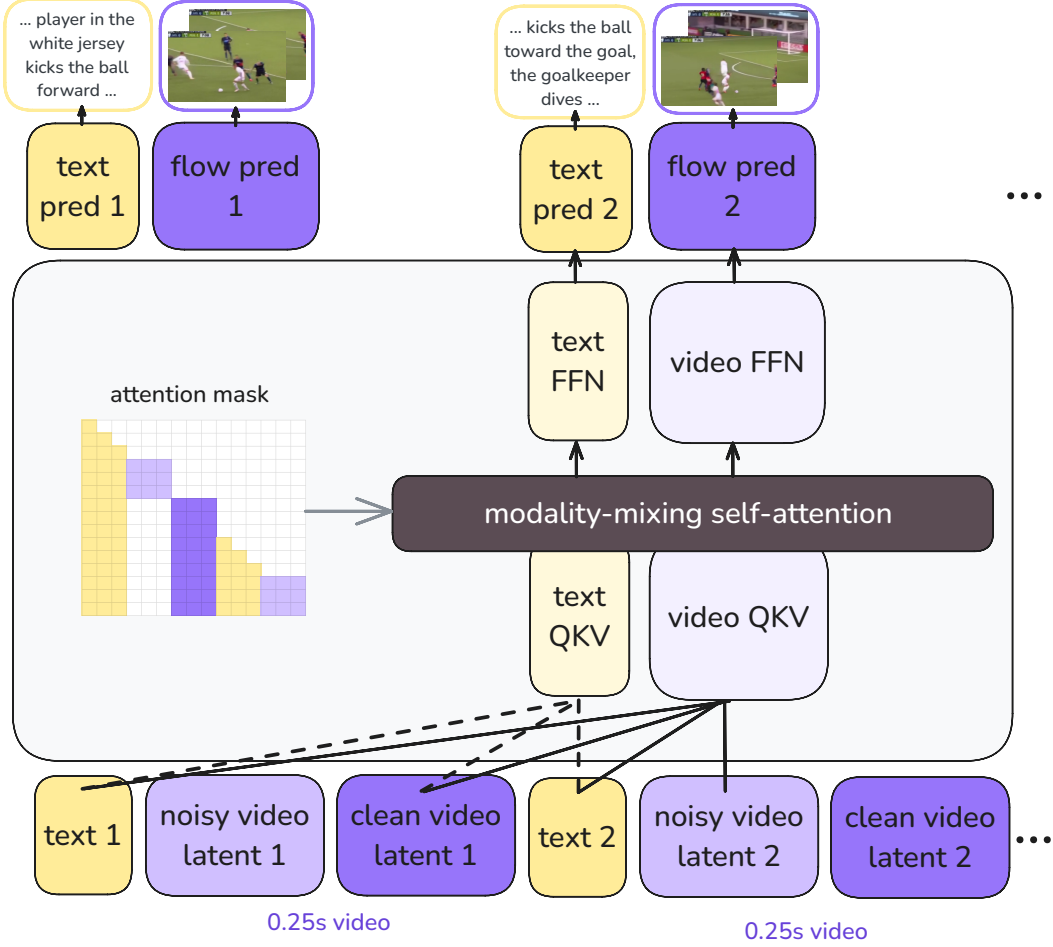


Figure 2 Overview of TV2TV architecture. TV2TV builds on the Transfusion (Zhou et al., 2024) modeling approach, jointly learning language modeling and video flow matching. TV2TV autoregressively generates interleaved chunks of video frames and text tokens, maintaining strict temporal causality: each token or chunk of frames can only attend to preceding tokens or chunks. TV2TV adopts a Mixture-of-Transformers (MoT) (Liang et al., 2024) architecture with dedicated towers for the video and text modalities.

2.1 Data representation

TV2TV models both discrete text and continuous videos in an interleaved fashion. Below, we describe how we prepare these interleaved sequences.

Discrete text tokens. To represent language, we use a regular BPE tokenizer to obtain discrete text tokens.¹

Continuous video tokens. To represent video, we use a VAE tokenizer with a causal 3D CNN backbone to obtain continuous video tokens.² Our video tokenizer has a temporal compression factor of 4, so every 4 frames are grouped together as an atomic *chunk* of frames, except for the first frame. The video tokenizer was trained with sequences of 49 frames, so we chunk and pad longer videos to multiple 49 frames, pass them through the tokenizer, and obtain a series of *latent frames* in the latent space, one for each atomic chunk of frames. We work with 16 FPS videos throughout this work, so each latent frame corresponds to 0.25 seconds of video. For simplicity, we use *frame chunk* to refer to such 0.25-second video in the latent space from now.

¹Specifically, we use the `tiktoken`-based tokenizer from Dubey et al. (2024).

²Specifically, we use the `Cosmos-Tokenize1-CV4x8x8-360p` tokenizer from Agarwal et al. (2025).

Interleaved text and video sequences. In the interleaved sequence, we organize text segments and frame chunks chronologically according to their timestamps. Video frame chunks are timestamped by their start time, while text segments are assigned timestamps as defined in §3 and §4. Our generation process respects temporal causality: each text segment or video frame chunk is conditioned only on content from earlier timestamps. When a text segment and video frame chunk share the same timestamp, we place the text before the frame chunk in the sequence such that the video generation can condition on the associated text (i.e., allowing the model to *think* in text before *acting* to generate those frames). Additionally, we introduce two special discrete tokens – *beginning-of-frame* (BOF) and *end-of-frame* (EOF) – that delimit the continuous tokens of each video frame chunk. This design enables the model to automatically transition between text generation mode and video generation mode during inference (see §2.3 for details).

Clean and noisy latents. Though across the interleaved sequence we predict text and video frame chunks autoregressively, within each frame chunk we perform flow matching on the continuous video tokens. Flow matching and diffusion methods require noisy, interpolated input representations for training. However, autoregressive generation under teacher-forcing requires access to clean representations from previous sequence elements to maintain proper conditioning context. To resolve this conflict, we maintain two copies of each video frame’s representation in the input sequence: a noisy frame chunk followed immediately by a clean frame chunk. This design allows the model to condition on clean historical context while learning to denoise current frames. This contrasts with MAGI-1 (Teng et al., 2025), which addresses this challenge by enforcing monotonicity in noise, ensuring earlier video chunks are cleaner than later ones. While this maintains a single representation per frame, it does not work well in our interleaved text-video setting as it limits the model’s ability to effectively interact with and condition on the textual components of the sequence.

Notation. In summary, the text and video data are tokenized and prepared as an interleaved sequence of discrete and continuous tokens. We refer to the text, noisy video frame chunk, and clean video frame chunk representations as \mathbf{x}^{txt} , $\mathbf{x}^{\text{noisy-vid}}$, and $\mathbf{x}^{\text{clean-vid}}$, respectively. Thus, each sequence for our TV2TV training is in the form:

$$[\mathbf{x}_1^{\text{txt}}, \mathbf{x}_1^{\text{noisy-vid}}, \mathbf{x}_1^{\text{clean-vid}}, \mathbf{x}_2^{\text{txt}}, \mathbf{x}_2^{\text{noisy-vid}}, \mathbf{x}_2^{\text{clean-vid}}, \dots, \mathbf{x}_N^{\text{txt}}, \mathbf{x}_N^{\text{noisy-vid}}, \mathbf{x}_N^{\text{clean-vid}}] \quad (1)$$

where N is the total number of latent video frames. As a shorthand for indexing specific modalities, we use $\mathbf{x}^{\text{txt}} = \oplus_{i=1}^N \mathbf{x}_i^{\text{txt}}$ for the text tokens in the sequence, $\mathbf{x}^{\text{vid}} = \oplus_{i=1}^N (\mathbf{x}_i^{\text{noisy-vid}}, \mathbf{x}_i^{\text{clean-vid}})$ for the video tokens in the sequence, and $\mathbf{x}^{\text{all}} = \oplus_{i=1}^N (\mathbf{x}_i^{\text{txt}}, \mathbf{x}_i^{\text{noisy-vid}}, \mathbf{x}_i^{\text{clean-vid}})$ for the full interleaving sequence. \oplus indicates concatenation along the sequence dimension.

2.2 Architecture and optimization

TV2TV adopts an MoT (Liang et al., 2024) architecture with dedicated transformer towers for each video and text. Following Shi et al. (2024), the text tower is initialized from a pretrained LM. During training, TV2TV jointly learns language modeling (next-token prediction) and video flow matching (next-frame prediction). We provide further architectural and optimization details below.

Noise interpolation. Flow matching on our continuous video frame representations requires interpolating them with pure noise samples to produce the noisy latents. Specifically, we use rectified flow (Liu et al., 2022) in a configuration similar to Esser et al. (2024). Let $\mathbf{x}^{\text{clean-vid}}$ be the video latents encoded by the continuous video tokenizer.

$$\mathbf{x}^{\text{noisy-vid}} = t\mathbf{x}^{\text{clean-vid}} + (1 - t)\epsilon \quad (2)$$

where t is sampled from a logit-normal distribution, $t \sim \text{logistic}(\mathcal{N}(\mu, \sigma^2))$,³ and ϵ is a Gaussian noise, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

³We use $\mu = 0$ and $\sigma = 1.4$ in this work.

Dropout. To enable text classifier-free guidance (CFG) during inference, we drop out $\mathbf{x}_i^{\text{txt}}$ randomly with a small probability $p_{\text{txt-drop}}$. To alleviate exposure bias in the teacher-forcing sequential training, a popular strategy is to inject a small amount of noise into clean data (Teng et al., 2025). In our case, we apply a soft dropout for $\mathbf{x}^{\text{clean-vid}}$, flipping $\mathbf{x}_i^{\text{clean-vid}}$ to $\mathbf{x}_i^{\text{noisy-vid}}$ with a small probability $p_{\text{clean-vid-flip}}$.

Input projection. The input text tokens \mathbf{x}^{txt} are projected by a linear embedding layer to a sequence of text hidden states $\mathbf{h}_{\text{in}}^{\text{txt}}$. The video frame latents \mathbf{x}^{vid} are projected to a sequence of video hidden states $\mathbf{h}_{\text{in}}^{\text{vid}}$ via a U-Net downsampler (Ronneberger et al., 2015). The timestep t is integrated in the downsampler via adding a time embedding. For each token in the video latents \mathbf{x}^{vid} corresponding to different spatial patches, we also add an absolute 2D position embedding through the downsampler.

$$\mathbf{h}_{\text{in}}^{\text{txt}} = \text{Proj}_{\text{txt}}(\mathbf{x}^{\text{txt}}) \quad (3)$$

$$\mathbf{h}_{\text{in}}^{\text{vid}} = \text{UNet-Down}_{\text{vid}}(\mathbf{x}^{\text{vid}}, t) \quad (4)$$

Modality-specific self-attention. Following the original MoT design (Liang et al., 2024), the text hidden states $\mathbf{h}_{\text{in}}^{\text{txt}}$ and video hidden states $\mathbf{h}_{\text{in}}^{\text{vid}}$ are transformed to their respective queries, keys, and values via separate Q, K, V matrices. The pre-attention layer normalization is also modality-specific and is folded into the QKV functions in the equations below for simplicity.

$$\mathbf{h}_Q^{\text{txt}} = Q_{\text{txt}}(\mathbf{h}_{\text{in}}^{\text{txt}}), \quad \mathbf{h}_K^{\text{txt}} = K_{\text{txt}}(\mathbf{h}_{\text{in}}^{\text{txt}}), \quad \mathbf{h}_V^{\text{txt}} = V_{\text{txt}}(\mathbf{h}_{\text{in}}^{\text{txt}}) \quad (5)$$

$$\mathbf{h}_Q^{\text{vid}} = Q_{\text{vid}}(\mathbf{h}_{\text{in}}^{\text{vid}}), \quad \mathbf{h}_K^{\text{vid}} = K_{\text{vid}}(\mathbf{h}_{\text{in}}^{\text{vid}}), \quad \mathbf{h}_V^{\text{vid}} = V_{\text{vid}}(\mathbf{h}_{\text{in}}^{\text{vid}}) \quad (6)$$

Attention is then computed across all tokens in the interleaving sequence. The attention-weighted values are projected back to the hidden state dimension using modality-specific O matrices.

$$\mathbf{h}_O^{\text{txt}} = O_{\text{txt}} \left(\text{softmax} \left(\frac{\text{mask}(\mathbf{h}_Q^{\text{txt}} \mathbf{h}_K^{\text{all}T})}{\sqrt{d}} \right) \mathbf{h}_V^{\text{all}} \right) \quad (7)$$

$$\mathbf{h}_O^{\text{vid}} = O_{\text{vid}} \left(\text{softmax} \left(\frac{\text{mask}(\mathbf{h}_Q^{\text{vid}} \mathbf{h}_K^{\text{all}T})}{\sqrt{d}} \right) \mathbf{h}_V^{\text{all}} \right) \quad (8)$$

where mask denotes a hybrid attention mask—applying a causal mask to the positions of text tokens and a block-causal mask to the positions of noisy and clean video tokens. An additional principle for masking is that noisy video tokens cannot be attended by any future tokens in the sequence. A global 1D RoPE is also applied here to all positions for all modalities.

Modality-specific feed-forward network. Again, following the original MoT (Liang et al., 2024) design, after self-attention, we use modality-specific FFNs to further transform text and video representations separately. The pre-FFN layer normalization is also modality-specific and is folded in the FFN function for simplicity.⁴

$$\mathbf{h}_{\text{FFN}}^{\text{txt}} = \text{FFN}_{\text{txt}}(\mathbf{h}_O^{\text{txt}}) \quad (9)$$

$$\mathbf{h}_{\text{FFN}}^{\text{vid}} = \text{FFN}_{\text{vid}}(\mathbf{h}_O^{\text{vid}}) \quad (10)$$

$$(11)$$

Output projection. After L layers of self-attention and FFNs, the resulting hidden states are projected either to logits in text via an output embedding or to a predicted flow via a U-Net upsampler.

$$\mathbf{s}_{\text{logits}}^{\text{txt}} = \text{LM-Head}_{\text{txt}}(\mathbf{h}_{\text{FFN}}^{\text{txt}}) \quad (12)$$

$$\mathbf{v}_{\text{pred}}^{\text{noisy-vid}} = \text{UNet-Up}_{\text{vid}}(\mathbf{h}_{\text{FFN}}^{\text{noisy-vid}}, \mathbf{h}_{\text{in}}^{\text{noisy-vid}}, t) \quad (13)$$

⁴We also do not show residual connections for simplicity of notation.

Training objective. The training objective is a combination of text cross entropy loss and video flow loss. All model parameters are optimized jointly using both losses.

$$\mathcal{L}_{\text{txt}} = \text{CE}(\mathbf{s}_{\text{logits}}^{\text{txt}}, \mathbf{x}^{\text{txt}}) \quad (14)$$

$$\mathcal{L}_{\text{vid}} = \text{MSE}(\mathbf{v}_{\text{pred}}^{\text{noisy-vid}}, (\mathbf{x}^{\text{clean-vid}} - \epsilon)) \quad (15)$$

$$\mathcal{L} = \lambda_{\text{txt}} \mathcal{L}_{\text{txt}} + \lambda_{\text{vid}} \mathcal{L}_{\text{vid}} \quad (16)$$

Parameter initialization. Due to the modality-specific design above, we have a separate set of parameters for each modality. The text tower’s parameters θ_{txt} include the parameters from Proj_{txt} , \mathbf{Q}_{txt} , \mathbf{K}_{txt} , \mathbf{V}_{txt} , \mathbf{O}_{txt} , FFN_{txt} , and $\text{LM-Head}_{\text{txt}}$. The video tower’s parameters θ_{vid} include the parameters from UNet-Down, \mathbf{Q}_{vid} , \mathbf{K}_{vid} , \mathbf{V}_{vid} , \mathbf{O}_{vid} , FFN_{vid} , and UNet-Up_{vid}. Following Shi et al. (2024), we initialize θ_{txt} from a pre-trained Llama model.

2.3 Inference

TV2TV’s key innovation lies in its ability to dynamically switch between text and video generation during inference. This is achieved by using a special *beginning-of-frame* (BOF) token which controls the transition from text to video generation. By default, TV2TV operates in text mode, generating tokens autoregressively like standard LLMs. At each autoregressive step $i - 1$, the model samples a next token x_i^{txt} and proceeds as follows:

- **Text token ($x_i^{\text{txt}} \in V$):** If the current token x_i^{txt} is within the regular language vocabulary V (not special token BOF or EOS), the model continues autoregressive text generation.
- **BOF token ($x_i^{\text{txt}} = \text{BOF}$):** If the current token is the BOF token, the model
 1. Extends the sequence with tokens representing one video frame chunk ($\mathbf{x}_i^{\text{noisy-vid}}$).
 2. Initializes these tokens from a normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$.
 3. Runs m steps of an ODE solver (e.g., Euler) using $\mathbf{x}_i^{\text{noisy-vid}}$ and the KV cache of the previous tokens (optionally applying CFG).⁵
 4. Run forward pass with the final output of the solver $\mathbf{x}_i^{\text{clean-vid}}$ and updates the KV cache, and then resumes autoregressive text generation.
- **End of sequence ($x_i^{\text{txt}} = \text{EOS}$ or context length exceeded):** The generation process terminates.

Extending sequences beyond the trained context length Because TV2TV is globally autoregressive, it can natively generate videos longer than its trained context length using sliding windows. To extend an interleaved text-video sequence $\oplus_{i=1}^N(\mathbf{x}_i^{\text{txt}}, \mathbf{x}_i^{\text{vid}})$, we retain the second half $\oplus_{i=N/2}^N(\mathbf{x}_i^{\text{txt}}, \mathbf{x}_i^{\text{vid}})$ and use it as a condition $\oplus_{i=1}^{N/2}(\mathbf{x}_i^{\text{txt}}, \mathbf{x}_i^{\text{vid}})$ for the next generation window.

2.4 Task formulation

While thus far we have presented a general approach for training on and generating interleaved text and video sequences, in this paper, we focus on *video generation* as the primary objective while leveraging text generation as an auxiliary task that provides semantic guidance. Rather than generating complex video content directly, the TV2TV framework decomposes the video generation process into interleaved text and video sequences. This factorization offers two advantages. First, it offloads much of the semantic complexity to the text generation components of the model, reducing the burden on the video generation component. Second, it enables flexible user control during generation, allowing users to intervene and modify the video generation trajectory at any point through textual prompting.

⁵When CFG is enabled during inference, we maintain both a text-conditional sequence $[\mathbf{x}_1^{\text{txt}}, \mathbf{x}_1^{\text{clean-vid}}, \dots, \mathbf{x}_{i-1}^{\text{txt}}, \mathbf{x}_{i-1}^{\text{clean-vid}}, \mathbf{x}_i^{\text{txt}}, \mathbf{x}_i^{\text{noisy-vid}}]$ and a text-unconditional sequence $[\mathbf{x}_1^{\text{clean-vid}}, \dots, \mathbf{x}_{i-1}^{\text{clean-vid}}, \mathbf{x}_i^{\text{noisy-vid}}]$, and take ODE steps contrastively on the two $\mathbf{x}_i^{\text{noisy-vid}}$.

A central question is how to obtain such interleaved text and video data for training. Video game data with associated controller actions provides a natural source, as the controller action can act as the textual “plan” for the subsequent video frames. §3 presents controlled experiments with such data.

To extend this modeling paradigm to general domain video data, which largely lacks clean and temporally-aligned captions, in §4, we experiment with a methodology for synthetically augmenting real-world videos with interleaved natural language action descriptions using vision-language models (VLMs).

3 Controlled Experiments with Video Game Data

In this section, we evaluate TV2TV using video gameplay footage (video) paired with controller actions (text). Video games offer an ideal testbed for interleaved text and video approaches: controller inputs serve as textual “plans” that reflect the subsequent actions displayed in the gameplay footage. By using the same video data while varying the text representations, we can directly evaluate whether letting the model “think in text” *interleavingly* before “acting in pixels” improves video generation quality. Additionally, we investigate whether this approach enables fine-grained user control during inference.

Specifically, we design our experiments to address the following research questions:

- **Overall visual quality:** Does generating interleaved planning text improve overall video quality compared to non-interleaved baselines?
- **Fine-grained controllability:** Does training with interleaved text improve the user controllability of video generation, i.e., can a mid-sequence text intervention reliably steer the video?

We train TV2TV and two baselines on gameplay video, actions, and metadata from the *Counter-Strike: Global Offensive* (CS:GO) dataset curated and open-sourced by [Pearce and Zhu \(2022\)](#).

3.1 Modeling details

We compare TV2TV with two controlled baselines, which we refer to as T2V and Think2V, to isolate different aspects of our approach:⁶

- **T2V:** We adopt the same autoregressive video modeling framework as TV2TV, but the model is trained without any interleaved text, conditioning solely on the meta-prompt (and an initial starting frame).
- **Think2V:** Rather than generating text and video in an interleaved fashion, after receiving the meta-prompt, Think2V first generates a detailed roll-out of all subsequent text actions *before* producing any video frames. In other words, text generation (‘*thinking*’) is followed by video generation in a sequential, non-interleaved manner. Again, we adopt the same autoregressive video modeling framework as TV2TV.

See [Figure 3](#) for an illustration of the different sequence representations used for TV2TV, T2V, and Think2V.

For these experiments, all models adopt a 3B-MoT backbone with modality-specific 3B-parameter text and video towers. The text tower is initialized with Llama-3.2-3B ([Dubey et al., 2024](#)). We train for 50K steps with a batch size of 128, utilizing a cosine learning rate scheduler with a maximum learning rate of $3e-4$. Detailed model and training configurations are provided in [Table 6](#) in [§A.1](#).

3.2 Data details

All models are trained on 95 hours of video from the video game *Counter-Strike: Global Offensive* (CS:GO) curated and open-sourced by [Pearce and Zhu \(2022\)](#). Each video frame is associated with a controller action, but as we group each 4 frames into a single video latent, we concatenate 4 controller actions, stringify them as text, and pass them to the model, e.g.:

⁶During inference, all models receive both the meta-prompt and a starting frame. We use the shorthand T2V and Think2V (rather than TI2V and ThinkI2V) for simplicity.

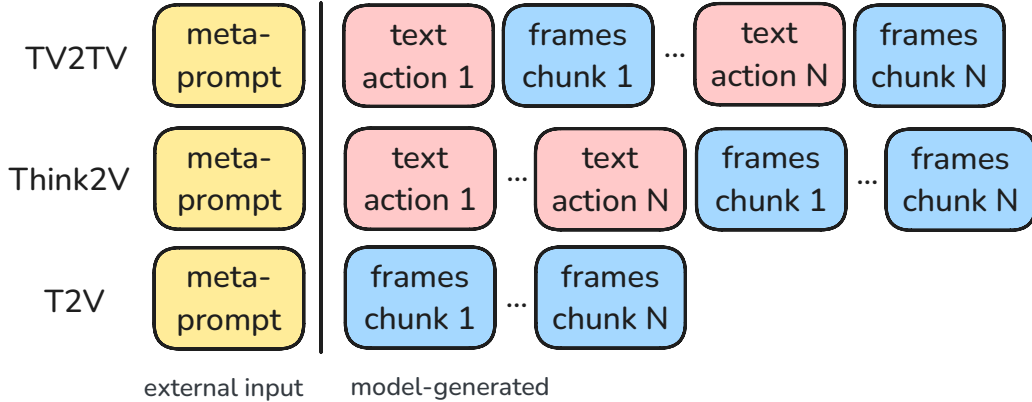


Figure 3 Illustration of sequence representations for TV2TV, T2V, and Think2V model variants. During inference, the meta-prompt (and a single conditioning frame) is provided and the rest of the sequence is model-generated.

$$\begin{array}{l}
 (w, d, shift). \ 10, \ 0, \ 0, \ 0. \\
 (d). \ 4, \ 0, \ 0, \ 0. \\
 (d). \ 4, \ 0, \ 0, \ 0. \\
 (d). \ 0, \ 0, \ 1, \ 0.
 \end{array}
 \quad \text{or} \quad
 \begin{array}{l}
 (d). \ -100, \ 10, \ 0, \ 0. \\
 (d). \ -100, \ 4, \ 0, \ 0. \\
 (). \ -60, \ 0, \ 0, \ 0. \\
 (). \ -4, \ 0, \ 0, \ 1.
 \end{array}$$

where the string includes: (keyboard inputs). horizontal mouse move, vertical mouse move. left mouse click, right mouse click. Example keyboard inputs are (w, a, s, d, space, ...) – walk forward, left, backward, right, jump, etc. See [Pearce and Zhu \(2022\)](#) for additional details on the CS:GO action space.

In the case of TV2TV, the meta-prompt text is inserted at the beginning of the sequence, and the controller text is inserted just before the associated block of frames. In the case of Think2V, all controller actions are inserted at the beginning of the sequence prior to any frames, along with the meta-prompt. In the case of T2V, no controller actions are provided. Again, refer to [Figure 3](#) for a visualization.

The CS:GO dataset has a resolution of 280×150 at 16 FPS, which we upsample to 320×192 for tokenizer compatibility. With a 2×2 U-Net patching and $4 \times 8 \times 8$ tokenizer, each chunk of 4 frames amounts to 240 tokens. The models’ context size fits 6.1 seconds (or 98 frames) of video per step per device during model training.

3.3 Evaluation set-up

We evaluate the overall video quality and user controllability of videos generated using the TV2TV approach with in-house blind human annotations.

Evaluating overall video quality To test the hypothesis that interleaved planning text improves video generation quality, we perform pairwise comparisons between videos generated by TV2TV and those from the non-interleaved baselines T2V and Think2V.

For all models, we generate 100 short-form and 40 long-form videos, yielding 280 pairwise annotation tasks. Short-form videos are approximately six seconds long, and long-form videos are 64 seconds long (see [§2.3](#) for details on how we generate long videos). For a given model, each generated video is conditioned on a single, unique frame from a held-out expert gameplay set and the prompt: “A brief gameplay clip from the iconic *Dust 2* map, showcasing classic *Counter-Strike* tactical combat in the legendary desert setting.” Videos evaluated side-by-side for visual quality share the same conditioning frames. Distortions and cloudiness in generated videos are penalized, along with physically implausible behaviors such as wall-clipping, abnormally slow gameplay, or spontaneous teleportation. Exact evaluation criteria can be found in [Appendix A.2](#).

Evaluating fine-grained controllability To test the hypothesis that interleaved planning text affords strong, fine-grained controllability relative to non-interleaved planning, we compare TV2TV’s response to

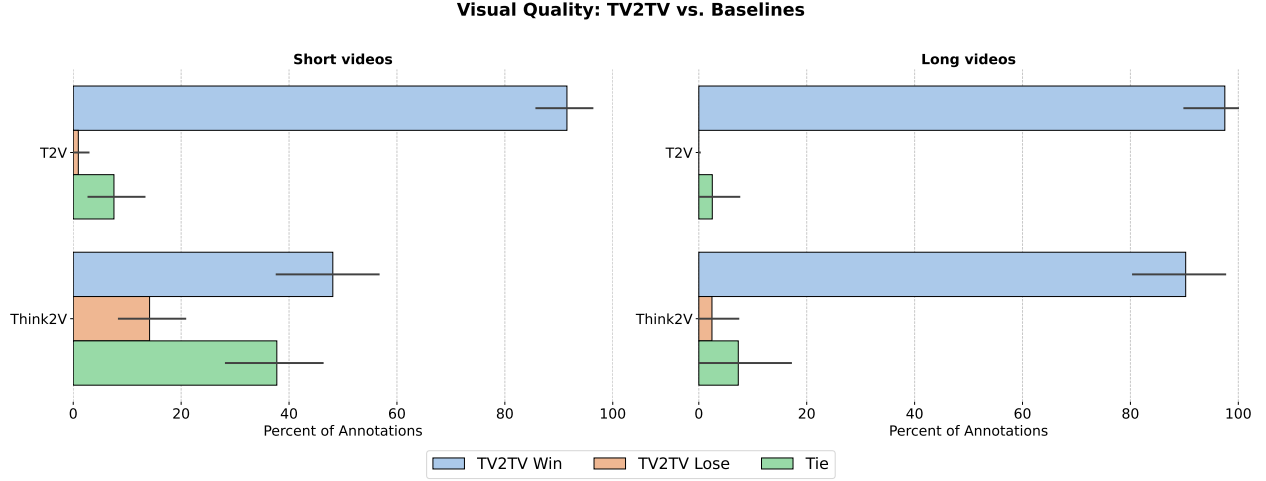


Figure 4 Visual quality human evaluation results. Pairwise visual quality comparison of videos generated by TV2TV, T2V, and Think2V. Results are shown for our standard 6-second context length (short videos) and extended 64-second videos (long videos) generated using sliding windows (see §2.3). 95% confidence intervals are shown. In both settings, TV2TV outperforms baselines.

intermediate text interventions against the Think2V baseline.

For TV2TV, we generate videos similar to those described for the overall video quality evaluation but with an additional invocation of a manual intervention (*e.g.* a “left-click” action, a “reloading” action, etc.) at an intermediate timestamp of the video. We experiment with interventions at $t = 0.8125s$ and $2.8125s$, which we refer to as 1s and 3s for simplicity. We compare against videos generated by the Think2V model, where the intervention is inserted into the model-generated text plan used as initial conditioning for video generation. To isolate the effect of these interventions, we include control videos from both models generated from the same conditioning frames and meta-prompt *without* any manual intervention.

In total, four user-controllable text interventions are evaluated: moving backwards, performing a “left-click” action that corresponds to shooting a weapon, initiating a weapon reload, and jumping. Control videos with no manual intervention are also included, denoted as “no-op” (no operation) interventions, to provide insight into how often an action is generated by the model even without manual intervention. For each model, we evaluate 150 generated videos with interventions, yielding 300 annotation tasks.

Human annotators provide point-wise evaluations of:

1. **Intervention Correctness**, i.e. how well the generated video reflects the user-specified intervention. Users select among three options: the video correctly reflects the intervention (+1), incorrectly reflects the intervention (0), or that they are unsure (+0.5). (Conversely, we mark *no-op* control videos as correct when the annotator marks that the action is not reflected in the video.) We average these results to obtain a score.
2. **Visual Quality**, using the same quality criteria defined above. Users have the option to rate the visual quality as strong (+3), moderate (+2), weak (+1), or none (0). We average these results and normalize to obtain a score between 0 and 1.

The exact evaluation instructions are provided in the Appendix §A.2.

3.4 Results

Results for overall visual quality and controllability experiments are shown in Figure 4 and Figure 5, respectively. We make the following observations:

- **Generating interleaved planning text substantially improves overall video quality.** As shown in Figure 4, TV2TV achieves substantial improvements in overall visual quality compared to both T2V and Think2V

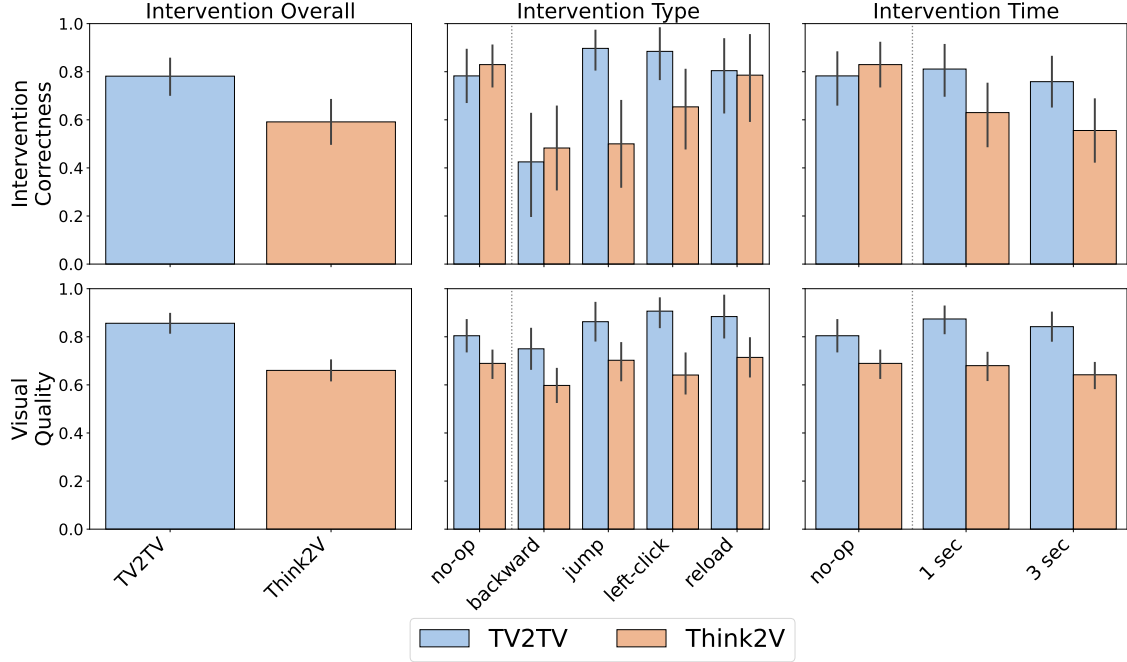


Figure 5 Fine-grained controllability human evaluation results. Pointwise controllability (Intervention Correctness) and quality (Visual Quality) comparison of videos generated by TV2TV and Think2V with intermediate action interventions. Results show that TV2TV demonstrates significant controllability advantages over Think2V with particularly strong control for *jump* and *left-click* and similar control for *reload* and *backward*. Moreover, TV2TV produces significantly higher visual quality generations with these interventions. Both scores range from 0 to 1. We also evaluate the *no-op* (no operation) baseline for reference – in the case of *no-op*, the video is marked as correct if the action is *not shown* in the video. We find that, when not manually intervened on, models generate the action only 17-22% of the time, showing that the intervention correctness during manual intervention is significantly above random. Details on metrics are provided in §3.3.

across short and long videos. The performance gains are more substantial relative to T2V (91% *vs.* 1%, with 8% ties) than to Think2V (48% *vs.* 14%, with 38% ties). Additionally, improvements are more pronounced for long videos compared to short videos across both baseline comparisons, which underscores TV2TV’s effectiveness for long-form video generation.

- **Interleaved text affords strong, fine-grained controllability over video generation.** Figure 5 shows intervention correctness and visual quality results for four types of intermediate interventions: firing (*left-click*), jumping (*jump*), reloading (*reload*), and moving backward (*backward*), each applied at timestamps $t = 0.8125s$ and $2.8125s$, denoted 1s and 3s for simplicity. TV2TV demonstrates significant controllability advantages over Think2V (78% *vs.* 59%), with particularly strong control for *jump* and *left-click* interventions. Controllability is similar between the two models for *reload* and *backward*. TV2TV shows stronger intervention controllability than Think2V for both intervention timestamps. For reference, we also evaluate the *no-op* (no operation) baseline, where no manual intervention is applied – in the case of *no-op*, the video is marked as correct if the action is *not shown* in the video. We find that, when not manually intervened on, models generate the action only 17-22% of the time. This highlights that the intervention correctness during manual intervention is significantly above random. Additionally, TV2TV produces significantly higher visual quality generations during manual interventions than Think2V.

4 Scaling TV2TV to Real World Data

While video games offer a convenient source of interleaved text-and-video data, most real-world video data lacks clean, temporally-aligned (timed), interleaving captions, presenting a challenge for extending TV2TV to open domain settings. In this section, following recent work on dense, differential video captioning (Chen

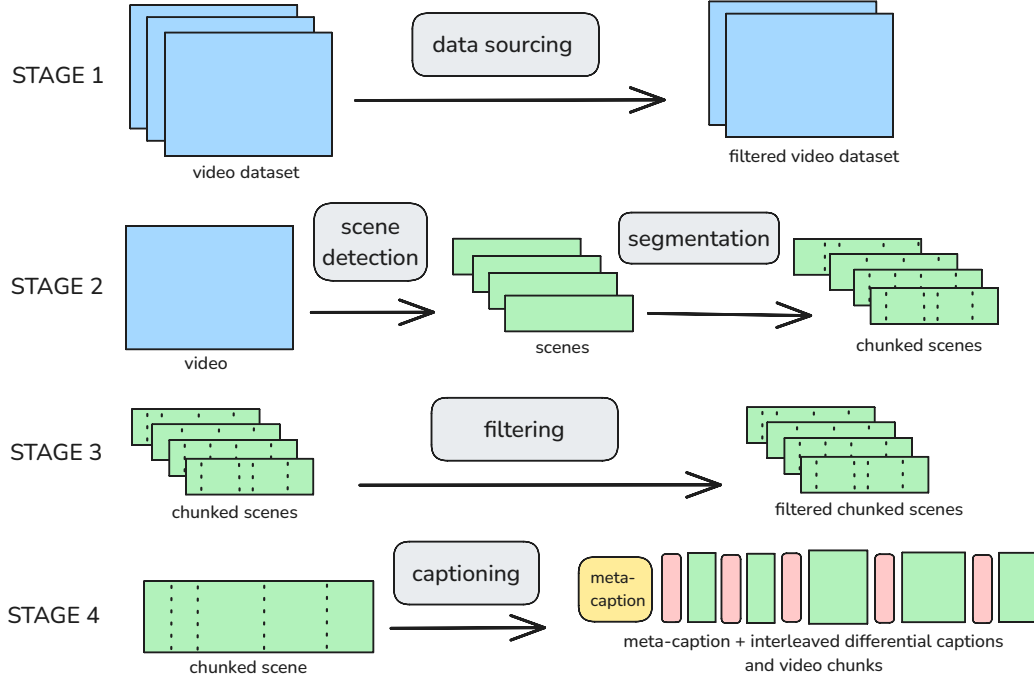


Figure 6 Interleaved text and video data pipeline. We present our methodology for using constructing an interleaved text and video training dataset from real world videos. The pipeline consists of several stages, including (1) data sourcing, (2) scene detection and segmentation, (3) filtering, and finally, (4) interleaved captioning with a VLM.

et al., 2024b,a; Peng et al., 2025a; Chen et al., 2025b; Wang et al., 2023; Xue et al., 2025b), we present our methodology for developing a data augmentation pipeline that synthetically augments real-world videos with interleaved natural language action descriptions using vision-language models (VLMs). This enables us to extend and generalize TV2TV to large-scale real-world data. We apply this pipeline to *sports videos* – chosen for their dynamic motion and rich action content – to create a dataset of 8K hours of interleaved text-and-video data. Finally, we evaluate TV2TV trained on this augmented data from scratch against established video generation models⁷ as well as controlled baselines.

4.1 Interleaved data augmentation pipeline

We describe how we construct interleaved text-and-video data from videos. The pipeline consists of several stages:

1. **Data sourcing:** Collecting sports domain data from the large-scale YT-Temporal-1B (Zellers et al., 2022) dataset, using keyword-based and other filtering criteria.
2. **Scene detection and segmentation:** Dividing videos into 6 to 16-second scenes, and further sub-dividing those scenes into segments based on detected key frames and other heuristics.
3. **Filtering:** Selecting high-quality scenes using model-based and other filters that assess motion characteristics, semantic content, and overall quality.
4. **Interleaved captioning:** Generating captions at multiple levels using a VLM: an overall meta-caption for each scene and fine-grained, interleaved captions describing changes between short segments.

See Figure 6 for a visualization of this pipeline.

⁷Sports highlights contain fast motion and relatively complex semantics. This makes them a challenging test case for existing pretrained video generation models, where sports content is present but not well represented in the training distribution.

Data sourcing We focus on building a dataset of sports content by filtering the YT-Temporal-1B dataset (Zellers et al., 2022) using keyword-based filters (e.g. “*game highlights*”). We chose the sports domain for its high action density, which provides a strong testbed for interleaved reasoning capabilities. This yields 38K total hours of data.

Scene detection and segmentation We segment the data into scenes using TransNetV2 (Souček and Lokoč, 2020), a shot boundary detection model based on 3D convolutional networks. To identify clips containing interesting content, we employ a two-step approach. First, we use the Perception Encoder (Bolya et al., 2025) to embed video frames and compute the cosine distance between consecutive frames, producing a time-series of semantic change. Peaks in this time-series may indicate moments where significant action likely occurs, so we refer to the frames associated with these peaks as *key frames*. We then apply an 8-second sliding window to each scene, extracting clips that are between 6 and 16 seconds long and contain the highest number of peaks. On average, these clips are 8.2 seconds in length.

Finally, we use a combination of key frames, hierarchical clustering of Perception Encoder embeddings (Chen et al., 2025b), and heuristics to further segment the clips into chunks of frames suitable for interleaved captioning. Each clip is divided into an average of 4.3 chunks, though the number of chunks can range from as few as 2 to as many as 10 depending on the length of the video and the number of key frames detected. We impose a minimum chunk length of 1-second; the average chunk length is 1.9 seconds.

Filtering We apply several scene-level filters to further refine our selection:

- **VLM-based quality classifier:** To select high quality scenes we prompt Gemma-3-12B-Instruct (Team, 2025) to select semantically-relevant content. For each scene, we sample consecutive frames from the start, center, and end of the video and ask the model to provide a score of 1-10 based on the perceived quality and relevance; see §B.1 for the full prompt.
- **Face bounding box filter:** We observed that a considerable portion of videos consisted of people talking directly to the camera without meaningful action or motion in the foreground. To remove such videos, we use RetinaFace (Deng et al., 2019) to obtain face bounding boxes and analyze both their coverage and temporal stability throughout the video. Clips with large, stable face bounding boxes are filtered out.
- **Motion filter:** We compute the optical flow for each clip (Farnebäck, 2003) and calculate its average magnitude across frames as a motion score. Clips with low motion scores, indicating static or minimal movement, are filtered out.

After filtering, our final dataset comprises 8K hours of sports video data.

Interleaved captioning Finally, following recent work on differential video captioning (Chen et al., 2024b,a; Peng et al., 2025a; Chen et al., 2025b; Wang et al., 2023; Xue et al., 2025b), we use Qwen3-VL-30B-A3B-Instruct (Bai et al., 2025a) to generate (1) an overall *meta-caption* for the video and (2) differential captions describing action changes across subsequent frame chunks. Detailed prompts passed to the VLMs are provided in the Appendix §B.2. An example interleaved document produced by this pipeline is provided in Table 1.

4.2 Experiments and analysis

4.2.1 Modeling and data details

We adopt an 8B-MoT backbone with modality-specific 8B-parameter text and video towers (§2). The text tower is initialized with Llama-3.1-8B (Dubey et al., 2024). We train for 250K steps utilizing a batch size of 512 and a cosine learning rate scheduler with a maximum learning rate of $3e-4$. Detailed model and training configurations are provided in Table 7 in §B.3.

We downsample the video data to 320×192 resolution at 16 FPS and randomly sample 6.1 seconds (or 98 frames) of video per step per device for model training.

⁸Original video source: <https://www.youtube.com/watch?v=v6htOcLa7KM>


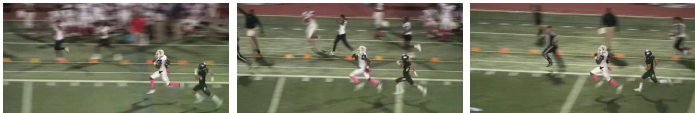
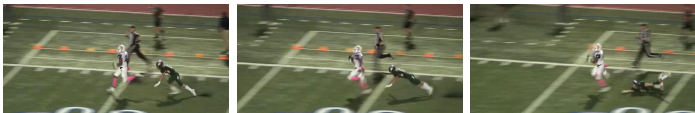


Interleaved caption	Subsampled frames
[0s - 2.2s] A player in a white uniform with pink socks runs with the ball, evading a defender in a black uniform. The player in white sprints towards the sideline, maintaining possession of the ball. The defender in black trails behind, attempting to catch up. Spectators are visible in the background.	
[2.2s - 4.3s] The player in the white uniform with pink socks continues running with the ball, moving further downfield. The defender in the black uniform remains in pursuit, closing the distance. Additional players in black uniforms join the chase, running towards the player with the ball. The player in white maintains possession and speed, evading the approaching defenders.	
[4.3s - 6.0s] The defender in the black uniform attempts a tackle but misses, falling to the ground. The player in the white uniform continues running unopposed towards the end zone.	
[6.0s - 8.6s] The player in the white uniform continues running towards the end zone, approaching the goal line.	
[8.6s - 10.1s] The player in the white uniform crosses the goal line, scoring a touchdown. The ball is now on the ground near the end zone.	
Overall meta-caption: A player in a white uniform with pink socks runs with the ball, evading defenders, and scores a touchdown.	

Table 1 Example interleaved training document. The source data is from YT-Temporal-1B (Zellers et al., 2022).⁸ Interleaved captions and the overall meta-caption are generated by Qwen3-VL-30B-A3B-Instruct (Bai et al., 2025a).

4.2.2 Evaluation set-up

We conduct pairwise evaluations comparing videos generated with TV2TV and (1) established, external video generation models and (2) controlled T2V and Think2V baselines. Pairs are evaluated by a pool of professional external annotators via the Turing platform for increased robustness.

We curate a held-out sports evaluation set consisting of major sports highlights and use meta-prompts captioned by VLMs as text prompts. The first frame of each video is used as the initial image condition. For each pairwise comparison, annotators assess prompt alignment, visual quality, real-world fidelity, and overall preference. Evaluation instructions used by annotators can be found in the Appendix §B.3.

External models We compare videos generated with TV2TV to those produced by several established video generation models: Cosmos-Predict2-Video2World 2B and 14B variants (Agarwal et al., 2025), MAGI-1 4.5 and 24B variants (Teng et al., 2025), and WAN-2.2 T12V 5B (Wan et al., 2025). All of these models include a T12V mode, which conditions video generation on a text prompt and an initial frame. We note that the external models we compare to were not specifically tuned for the sports domain and as such we do not expect them to perform well; however, we include them as out-of-domain baselines to obtain insights into TV2TV’s visual quality and prompt alignment relative to high-performing general-purpose models in this challenging domain.

For each model, we generate five videos for each of 30 unique conditioning images and prompts, yielding a total of 750 pairwise comparison tasks between TV2TV and external models. For the external models, we

Preferences: TV2TV vs. Sports Baselines

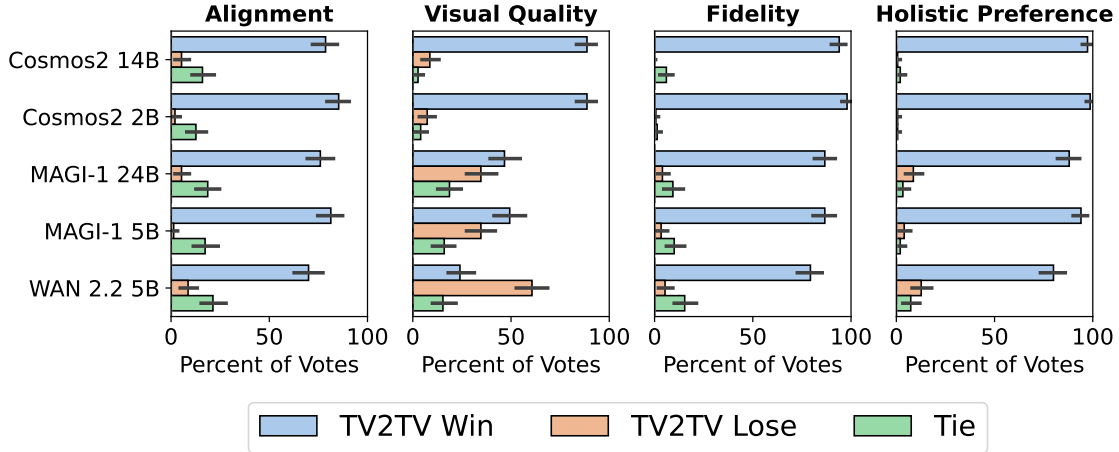


Figure 7 Evaluation of TV2TV trained on real world interleaved sports data vs. external video models We note that the external models we compare to were not specifically tuned for the sports domain and as such we do not expect them to perform well; however, we include them as out-of-domain baselines to obtain insights into TV2TV’s visual quality and prompt alignment relative to high-performing general-purpose models in this challenging domain. In blind pairwise human evaluations by external annotators, TV2TV outperforms all models on sports prompts in prompt alignment, real-world fidelity, and overall holistic preference. For visual quality, we find that TV2TV surpasses Cosmos2 variants, has similar performance as MAGI-1 variants, and underperforms compared to WAN 2.2 5B. 95% confidence intervals are shown.

generate videos in their native resolutions. To allow for a balanced comparison, prior to human evaluation stage we downsample the videos generated with external models to the shared lowest resolution across all models (320×192) and FPS (16) of TV2TV.

Controlled baselines We compare videos generated with TV2TV to those generated by T2V and Think2V baselines trained under the same settings (but varying the data representation for each framework). Similar to §3.1, T2V is trained without any interleaved text, conditioning directly on a concise meta-prompt. Think2V is trained on a concatenation of the meta-prompt and an extended, detailed prompt. Compared to the meta-prompt, the extended prompt explicitly instructs VLMs (Qwen3-VL-30B-A3B-Instruct) to provide more thorough and detailed descriptions of objects, actions, event progression, and so on, based on more densely sampled video frames. During inference, Think2V first self-generates (thinks) such detailed descriptions conditioned on the meta-prompt, and then generates video frames in a non-interleaved manner. This setup is similar to the use of highly descriptive synthetic captions in Betker et al. (2023).

For each model, we generate five videos for each of 30 unique conditioning images and prompts, yielding a total of 300 pairwise comparison tasks between TV2TV and these baselines.

4.2.3 Results

Results for the pairwise comparison with external models and controlled baselines are shown in Figure 7 and Figure 8, respectively.

External models As shown in Figure 7, we find that TV2TV outperforms all external models on this sports data in prompt alignment, real-world fidelity, and overall holistic preference when evaluated by external annotators. For visual quality, we find that TV2TV surpasses Cosmos2 variants, has similar performance as MAGI-1 variants, and is worse than WAN 2.2 5B. For Cosmos2, this is expected, as this model was largely tuned for performance on domains like Robotics and Autonomous Driving. For more general domain pretrained models like MAGI-1 and WAN 2.2, sports provides a challenging test case.

Preferences: TV2TV vs. Sports Baselines

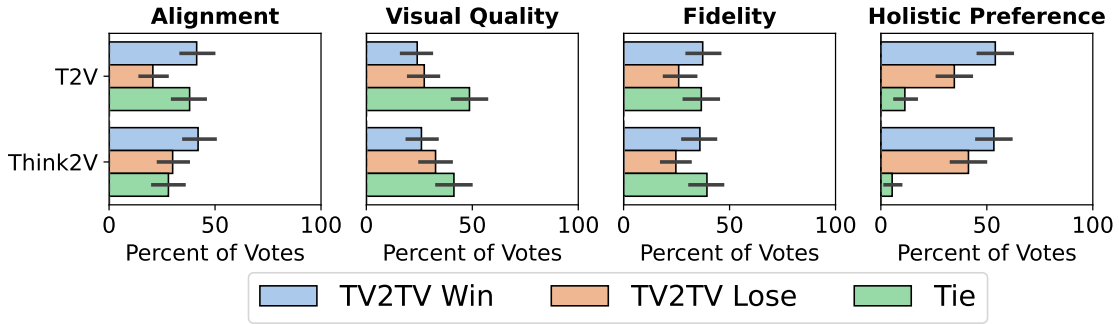


Figure 8 Evaluation of TV2TV trained on real world interleaved sports data vs. T2V and Think2V in a controlled setup. Compared to T2V, TV2TV shows stronger alignment and holistic preference, with similar visual quality and fidelity. When compared to Think2V, TV2TV shows a similar pattern of improved alignment and overall preference, though it is not statistically significant.

Controlled baselines In Figure 8, the TV2TV model shows improvements on holistic preference against both baselines: TV2TV videos have a win-rate 19 points higher than T2V (54% *vs.* 35%, with 11% ties) and 12 points higher than Think2V (53% *vs.* 41%, with 6% ties), although the latter is not statistically significant. Furthermore, for alignment, TV2TV has a win-rate 20 points higher than T2V and 12 points higher than Think2V (although again, the latter is not statistically significant). TV2TV shows similar real world fidelity and visual quality compared to T2V and Think2V.

Together, these results demonstrate that strong generation performance in real-world, dynamic contexts can be achieved with the TV2TV method.

4.2.4 Qualitative analysis

We provide qualitative examples of generations from TV2TV conditioned on a meta-prompt and one frame. Table 2, Table 3, and Table 4 illustrate the evolution of the scene with model-generated interleaved textual planning. Each generated segment of text helps plan for the subsequent set of frames.

A key advantage of TV2TV is its controllability through interleaved text: users can inspect, edit, or steer the intermediate textual plan to modify video generation at any timestep. We demonstrate this capability using the 8B-MoT TV2TV model trained on sports data, providing qualitative examples of how interleaved captions dynamically alter generation trajectories. See Table 5.

4.3 Discussion

We demonstrated that co-generating text and video in an interleaved sequence provides advantages over non-interleaved approaches such as T2V and Think2V. The gains in gaming (CS:GO) are more pronounced than in sports (§3, Figure 4 *vs.* §4, Figure 8). Two factors potentially contribute to this gap. First, the density of interleaved text: CS:GO provides frame-level textual signals, while sports rely on synthetic captions inserted every 1.9 seconds on average. Second, the quality of the interleaved text: CS:GO uses ground-truth actions, whereas sports captions are generated by VLMs and often contain hallucinations. Despite these challenges, TV2TV scales effectively to real-world data and still outperforms non-interleaved baselines.

In addition, TV2TV naturally supports user interventions at any point during generation, and its interleaved autoregressive setup allows simple extension to longer videos through sliding windows. Overall, interleaved text helps the model learn video generation more effectively and provides flexible control at inference time. Future work could focus on improving the granularity and accuracy of the interleaved text in the training data across more video domains.

Meta-prompt: *A player in a white jersey dribbles past defenders, kicks the ball towards the goal, and scores, celebrating as the goalkeeper remains on the ground.*

Generated video (conditioned on 1 frame):



Timestamp	Generated interleaved text during video generation
t=0.0	A soccer player in a white jersey dribbles the ball while being closely defended by two players in red jerseys. The player in white maintains possession, maneuvering around the defenders. Another player in a red jersey is positioned further back, observing the play. The action takes place on a grass field with a crowd of spectators in the background.
t=0.56	The player in the white jersey kicks the ball forward, and a teammate in a white jersey runs towards it. The defenders in red jerseys chase after the ball, with one attempting to intercept. The player in the white jersey gains possession and continues dribbling.
t=2.56	The player in the white jersey kicks the ball towards the goal. The goalkeeper dives to the left in an attempt to save it.
t=4.63	The player in the white jersey celebrates the goal by running towards the right side of the frame. The goalkeeper remains on the ground.

Table 2 TV2TV interleaved text and video generation rollout. 16 subsampled frames illustrate the evolution of the scene with model-generated interleaved planning. Each generated text helps plan for the video segment following its timestamp.

5 Related Work

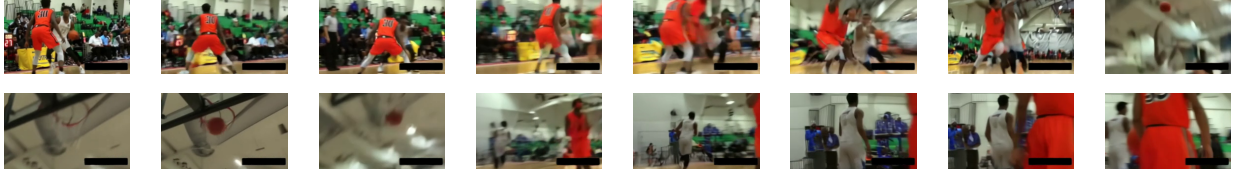
Unified multimodal architectures. The field has rapidly progressed toward unified foundation models capable of understanding and generating content across multiple modalities. Early unified models like Flamingo (Alayrac et al., 2022) pioneered this area, integrating vision and language components via cross-attention. Emu2 (Sun et al., 2023) advanced this paradigm by employing a unified autoregressive (AR) objective for both modalities. More recently, Chameleon (Chameleon Team, 2024) adopted an early-fusion, purely AR architecture that unifies images and text as discrete tokens, demonstrating flexible multi-modal reasoning.

Subsequent work has focused on improving scalability using hybrid autoregressive and diffusion approaches (Zhao et al., 2024; Zhou et al., 2024; Xie et al., 2024; Tong et al., 2024; Ma et al., 2025; Shi et al., 2024). For instance, Transfusion (Zhou et al., 2024) improved upon Chameleon by leveraging modality-specific losses, using AR for text and diffusion for images, which allows for the use of continuous visual tokens and scales better. Similarly, the Janus series collectively introduced decoupled visual pathways for encoding and generation (Ma et al., 2025; Chen et al., 2025c) to improve scalability while preserving shared semantics. BAGEL (Deng et al., 2025) scales multimodal pretraining using a Mixture-of-Transformers (MoT) (Liang et al., 2024). Manzano (Li et al., 2025) advances this line by unifying image, text, and video understanding under a coherent generative framework, producing seamless multimodal interleaving. Here, we propose a scalable recipe to extend these works to the generation of video and interleaved multimodal sequences.

Action-conditioned video generation. Video generation models are increasingly capable of simulating diverse environments (Bruce et al., 2024; Brooks et al., 2024; Xiang et al., 2024; Teng et al., 2025; Ye et al., 2025). The first successes of action-conditioned video generation have been in computer games (Alonso et al., 2024; Yu et al., 2025; HunyuanWorld, 2025); for example, GameNGen successfully simulates Doom (Valevski et al., 2024). However, models focusing on computer games remain limited to discrete, task-specific actions. Other models extend beyond games to visually diverse environments with physical action spaces, such as navigation (Bar et al., 2025; Lu et al.) or full-body control (Bai et al., 2025b). Building a truly general world model requires a flexible and expressive action space. Genie (Bruce et al., 2024) addresses this by

Meta-prompt: *A player in a white jersey dribbles past a defender in an orange jersey, shoots, and scores. The players then walk away from the hoop.*

Generated video (conditioned on 1 frame):



Timestamp	Generated interleaved text during video generation
t=0.0	A player in a white jersey dribbles the basketball while being guarded by a player in an orange jersey. The player in white maneuvers around the defender, maintaining control of the ball as they move across the court.
t=2.56	The player in the white jersey jumps and shoots the basketball toward the hoop. The ball arcs through the air and goes through the net, scoring a basket. The camera follows the ball's trajectory.
t=4.88	Players in white and orange jerseys walk across the court. The player in the white jersey with the number 1 walks toward the sideline, while the player in the orange jersey with the number 33 walks in the opposite direction.

Table 3 TV2TV interleaved text and video generation rollout. 16 subsampled frames illustrate the evolution of the scene with model-generated interleaved planning. Each generated text helps plan for the video segment following its timestamp.

learning latent actions, but interpreting and self-generating such actions remains challenging. In contrast, we aim to build models which can both generate *and* condition video generation on open-ended actions with a vocabulary grounded in text, spanning diverse environments including computer games and sports videos.

Autoregressive video generation. Modeling long, sequential multimodal inputs, particularly video, remains a critical challenge. Several approaches take an autoregressive approach to iteratively predict the next frame or block of frames conditioned on previous frames (Yin et al., 2024; Chen et al., 2025a; Agarwal et al., 2025; HaCohen et al., 2024; Lin et al., 2025; Cheng et al., 2025; Yuan et al., 2025), with notable work like MAGI-1 (Teng et al., 2025) focusing on high-fidelity AR video generation. However, these works, including MAGI-1, Cosmos (Agarwal et al., 2025), and Chen et al. (2025a), primarily focus on video generation and cannot generate interleaved text. Furthermore, Cosmos does not condition on the entire history, limiting long-term coherence. Meanwhile, token-based autoregressive generation—popularized by large language models (LLMs)—is utilized by models like VideoPoet (Kondratyuk et al., 2024), which treats video generation as a next-token prediction task and can straightforwardly leverage KV caching to improve generation efficiency.

To address the computational intensity and mitigate exposure bias during long-sequence generation, optimization techniques like diffusion forcing (Chen et al., 2025a) and self-forcing (Huang et al., 2025) are employed. Specifically, diffusion forcing introduced progressive noise to parallelize denoising across frames. To alleviate exposure bias in sequential training, a popular strategy is to inject a small amount of noise into clean data (Teng et al., 2025).

Full-sequence diffusion models. Another line of work (Zheng et al., 2024; Peng et al., 2025b; Wan et al., 2025) follows the full-sequence paradigm where the video is treated as one large tensor and all frames are denoised simultaneously. However, this results in prohibitively large attention computations, making development of such methods prohibitive for long form generation.

Multi-prompt video generation. A complementary line of work involves extending video generation length using multiple prompts (Villegas et al., 2023; Oh et al., 2024; Bansal et al., 2024; Cai et al., 2025). The interleaved design of TV2TV allows for a natural extension of video length by simply using a sliding window

Meta-prompt: *A weightlifter in a blue uniform lifts a barbell from the ground to his shoulders, then raises it above his head, holding it steady. An observer in a blue jacket watches closely.*

Generated video (conditioned on 1 frame):



Timestamp	Generated interleaved text during video generation
t=0.0	A weightlifter in a blue uniform lifts a barbell from the ground to his shoulders in a swift motion, while another man in a blue jacket watches from the side.
t=1.56	The weightlifter holds the barbell at shoulder height, standing upright with a focused expression, while the observer remains stationary, watching the lifter’s form.
t=3.13	The weightlifter raises the barbell above his head, fully extending his arms, maintaining a stable stance with legs slightly apart, as the observer continues to watch.
t=5.13	The weightlifter holds the barbell overhead with arms fully extended, standing still with a steady posture, while the observer remains in place, observing the lifter’s position.

Table 4 TV2TV interleaved text and video generation rollout. 16 subsampled frames illustrate the evolution of the scene with model-generated interleaved planning. Each generated text helps plan for the video segment following its timestamp.

approach: additional prompts for continuing the generation can be generated by the model (or inserted by the user) at any timestep.

Dense captioning pipelines. Recent work has developed dense and temporally aligned video captioning approaches that generate hierarchically structured descriptions over long videos, often combining coarse global summaries with fine-grained local captions (Chen et al., 2024a; Peng et al., 2025a; Chen et al., 2025b). Our pipeline follows this general recipe but emphasize feature-based adaptive segmentation compared to existing pipelines such as InternVid (Wang et al., 2023). On the data selection side, our motion-, face-, and VLM-based filters are conceptually related to large-scale video curation efforts such as OpenVid (Nan et al., 2024), UltraVideo (Xue et al., 2025a), VideoUFO (Wang and Yang, 2025), and Koala-36M (Wang et al., 2025), but are tailored to the sports domain and explicitly optimized for producing high-quality *interleaved* text–video training corpora. Generating action change descriptions is relevant to differential captioning works such as ShareGPT4Video (Chen et al., 2024b), ProgressCaptioner (Xue et al., 2025b), and CI-VID (Ju et al., 2025). Finally, the Think2V baseline described in §4 closely resembles the descriptive caption “upsampling” approach described in Betker et al. (2023), which was used to improve DALLE-3 prompt following.

6 Conclusion

This work introduces TV2TV, a unified modeling framework that decomposes video generation into an interleaved text and video generation process. By generating both text and video in an interleaved manner, this approach offers two key advantages: (1) it offloads much of the semantic complexity of video generation to the model’s text generation components, and (2) it enables more flexible and effective user control during generation. Through controlled experiments on video game data, we demonstrate that TV2TV outperforms baseline models in both video generation quality and controllability. We also show how this paradigm scales to real-world data by building a data augmentation pipeline that uses VLMs to enrich video data with interleaved action descriptions and comparing with established and controlled baselines. We believe TV2TV represents a promising step toward unifying advances in language model planning and reasoning with highly controllable video generation systems into a single generative framework.













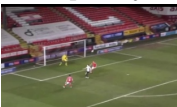
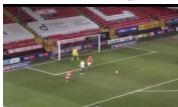










Subsampled frames					
t=0.0	t=1.19	t=2.38	t=3.62	t=4.81	t=6.06
Intervention 1 at t=1.56: “The man completes his golf swing, raising the club above his shoulder as he follows through. The golf ball is no longer visible, having been struck and sent flying forward. His body turns slightly to the left, maintaining balance after the powerful swing.”					
					
Intervention 2 at t=1.56: “The man completes his golf swing, raising the club above his shoulder as he follows through. The camera pans to track the ball as it soars through the air.”					
					
t=0.0	t=0.75	t=1.5	t=2.31	t=3.06	t=3.88
Intervention 1 at t=1.56: “The player in the white jersey takes control of the ball and runs towards the goal. He kicks the ball powerfully towards the net. The goalkeeper in the yellow jersey leaps to the left, extending his arms in an attempt to block the shot. The ball moves swiftly towards the goalpost as the goalkeeper’s jump reaches its peak.”					
					
Intervention 2 at t=1.56: “The player in the red jersey intercepts the ball near the center of the field and starts dribbling towards the right side of the frame, moving away from the goal. He evades an approaching defender in a white jersey by maneuvering the ball skillfully. As he advances, other players adjust their positions, preparing for the next phase of play. The goalkeeper remains near the goal, observing the developing action.”					
					

Table 5 Comparison of video rollouts steered by different interleaved captions. We alter the interleaved caption at second 1.56s and observe how the generation trajectory is altered.

7 Acknowledgements

We would like to thank Tariq Berrada, Rohit Girdhar, Sachin Mehta, Hritik Bansal, Mike Lewis, and Adriana Romero Soriano for helpful discussions throughout this project.

References

- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan.

- Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736, 2022.
- Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos J Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural Information Processing Systems*, 37:58757–58791, 2024.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuezhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025a.
- Yutong Bai, Danny Tran, Amir Bar, Yann LeCun, Trevor Darrell, and Jitendra Malik. Whole-body conditioned egocentric video prediction. *arXiv preprint arXiv:2506.21552*, 2025b.
- Hritik Bansal, Yonatan Bitton, Michal Yarom, Idan Szepes, Aditya Grover, and Kai-Wei Chang. TALC: time-aligned captions for multi-scene text-to-video generation. *CoRR*, abs/2405.04682, 2024. doi: 10.48550/ARXIV.2405.04682. <https://doi.org/10.48550/arXiv.2405.04682>.
- Amir Bar, Gaoyue Zhou, Danny Tran, Trevor Darrell, and Yann LeCun. Navigation world models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15791–15801, 2025.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, Junke Wang, Marco Monteiro, Hu Xu, Shiyu Dong, Nikhila Ravi, Daniel Li, Piotr Dollár, and Christoph Feichtenhofer. Perception encoder: The best visual embeddings are not at the output of the network. *CoRR*, abs/2504.13181, 2025. doi: 10.48550/ARXIV.2504.13181. <https://doi.org/10.48550/arXiv.2504.13181>.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. <https://openai.com/research/video-generation-models-as-world-simulators>.
- Jake Bruce, Michael D. Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle, Feryal M. P. Behbahani, Stephanie C. Y. Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero, Sherjil Ozair, Scott E. Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando de Freitas, Satinder Singh, and Tim Rocktäschel. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. <https://openreview.net/forum?id=bJbSbJskOS>.
- Minghong Cai, Xiaodong Cun, Xiaoyu Li, Wenzhe Liu, Zhaoyang Zhang, Yong Zhang, Ying Shan, and Xiangyu Yue. Ditctrl: Exploring attention control in multi-modal diffusion transformer for tuning-free multi-prompt longer video generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 7763–7772. Computer Vision Foundation / IEEE, 2025. doi: 10.1109/CVPR52734.2025.00727. https://openaccess.thecvf.com/content/CVPR2025/html/Cai_DiTCtrl_Exploring_Attention_Control_in_Multi-Modal_Diffusion_Transformer_for_Tuning-Free_CVPR_2025_paper.html.
- Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models, 2024. <https://arxiv.org/abs/2405.09818>.
- Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. In *NeurIPS*, 2025a.
- Delong Chen, Samuel Cahyawijaya, Etsuko Ishii, Ho Shu Chan, Yejin Bang, and Pascale Fung. What makes for good image captions? *arXiv preprint arXiv:2405.00485*, 2024a.
- Delong Chen, Théo Moutakanni, Willy Chung, Yejin Bang, Ziwei Ji, Allen Bolourchi, and Pascale Fung. Planning with reasoning using vision language world model. *CoRR*, abs/2509.02722, 2025b. doi: 10.48550/ARXIV.2509.02722. <https://doi.org/10.48550/arXiv.2509.02722>.
- Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Lin Bin, Zhenyu Tang, Li Yuan, Yu Qiao, Dahua Lin, Feng Zhao, and Jiaqi Wang. Sharegpt4video: Improving video understanding and generation with better captions. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information*

- Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b. http://papers.nips.cc/paper_files/paper/2024/hash/22a7476e4fd36818777c47e666f61a41-Abstract-Datasets_and_Benchmarks_Track.html.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025c.
- Xinle Cheng, Tianyu He, Jiayi Xu, Junliang Guo, Di He, and Jiang Bian. Playing with transformer at 30+ fps via next-frame diffusion. *arXiv preprint arXiv:2506.01380*, 2025.
- Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Shi Guang, and Haoqi Fan. Emerging properties in unified multimodal pretraining. *CoRR*, abs/2505.14683, 2025. doi: 10.48550/ARXIV.2505.14683. <https://doi.org/10.48550/arXiv.2505.14683>.
- Jiankang Deng, Jia Guo, Yuxiang Zhou, Jinke Yu, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-stage dense face localisation in the wild. *CoRR*, abs/1905.00641, 2019. <http://arxiv.org/abs/1905.00641>.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- Gunnar Farnebäck. Two-frame motion estimation based on polynomial expansion. In Josef Bigün and Tomas Gustavsson, editors, *Image Analysis, 13th Scandinavian Conference, SCIA 2003, Halmstad, Sweden, June 29 - July 2, 2003, Proceedings*, volume 2749 of *Lecture Notes in Computer Science*, pages 363–370. Springer, 2003. doi: 10.1007/3-540-45103-X_50. https://doi.org/10.1007/3-540-45103-X_50.
- Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024.
- Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the train-test gap in autoregressive video diffusion, 2025. <https://arxiv.org/abs/2506.08009>.
- Team HunyuanWorld. Hunyuanworld 1.0: Generating immersive, explorable, and interactive 3d worlds from words or pixels. *arXiv preprint*, 2025.
- Yiming Ju, Jijin Hu, Zhengxiong Luo, Haoge Deng, hanyu Zhao, Li Du, Chengwei Wu, Donglin Hao, Xinlong Wang, and Tengfei Pan. Ci-vid: A coherent interleaved text-video dataset, 2025. <https://arxiv.org/abs/2507.01938>.
- Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, Krishna Somandepalli, Hassan Akbari, Yair Alon, Yong Cheng, Josh Dillon, Agrim Gupta, Meera Hahn, Anja Hauth, David Hendon, Alonso Martinez, David Minnen, Mikhail Sirotenko, Kihyuk Sohn, Xuan Yang, Hartwig Adam, Ming-Hsuan Yang, Irfan Essa, Huisheng Wang, David A. Ross, Bryan Seybold, and Lu Jiang. Videopoet: A large language model for zero-shot video generation, 2024. <https://arxiv.org/abs/2312.14125>.
- Yanghao Li, Rui Qian, Bowen Pan, Haotian Zhang, Haoshuo Huang, Bowen Zhang, Jialing Tong, Haoxuan You, Xianzhi Du, Zhe Gan, Hyunjik Kim, Chao Jia, Zhenbang Wang, Yinfei Yang, Mingfei Gao, Zi-Yi Dou, Wenze Hu, Chang Gao, Dongxu Li, Philipp Dufter, Zirui Wang, Guoli Yin, Zhengdong Zhang, Chen Chen, Yang Zhao, Ruoming Pang, and Zhifeng Chen. MANZANO: A simple and scalable unified multimodal model with a hybrid vision tokenizer. *CoRR*, abs/2509.16197, 2025. doi: 10.48550/ARXIV.2509.16197. <https://doi.org/10.48550/arXiv.2509.16197>.
- Weixin Liang, Lili Yu, Liang Luo, Srinivasan Iyer, Ning Dong, Chunting Zhou, Gargi Ghosh, Mike Lewis, Wen-tau Yih, Luke Zettlemoyer, et al. Mixture-of-transformers: A sparse and scalable architecture for multi-modal foundation models. *arXiv preprint arXiv:2411.04996*, 2024.
- 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. *arXiv preprint arXiv:2506.09350*, 2025.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.

- TaiMing Lu, Tianmin Shu, Alan Yuille, Daniel Khashabi, and Jieneng Chen. Genex: Generating an explorable world. In *The Thirteenth International Conference on Learning Representations*.
- Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan, Zhenda Xie, Haowei Zhang, Xingkai Yu, et al. Janusflow: Harmonizing autoregression and rectified flow for unified multimodal understanding and generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7739–7751, 2025.
- Kepan Nan, Rui Xie, Penghao Zhou, Tiehan Fan, Zhenheng Yang, Zhijie Chen, Xiang Li, Jian Yang, and Ying Tai. Openvid-1m: A large-scale high-quality dataset for text-to-video generation. *arXiv preprint arXiv:2407.02371*, 2024.
- Gyeongrok Oh, Jaehwan Jeong, Sieun Kim, Wonmin Byeon, Jinkyu Kim, Sungwoong Kim, and Sangpil Kim. MEVG: multi-event video generation with text-to-video models. In Ales Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol, editors, *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XLIII*, volume 15101 of *Lecture Notes in Computer Science*, pages 401–418. Springer, 2024. doi: 10.1007/978-3-031-72775-7_23. https://doi.org/10.1007/978-3-031-72775-7_23.
- Tim Pearce and Jun Zhu. Counter-strike deathmatch with large-scale behavioural cloning. In *IEEE Conference on Games, CoG 2022, Beijing, China, August 21-24, 2022*, pages 104–111. IEEE, 2022. doi: 10.1109/COG51982.2022.9893617. <https://doi.org/10.1109/CoG51982.2022.9893617>.
- Ruotian Peng, Haiying He, Yake Wei, Yandong Wen, and Di Hu. Patch matters: Training-free fine-grained image caption enhancement via local perception. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3963–3973, 2025a.
- Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu, Hongxin Liu, Mingyan Jiang, Wenjun Li, Yuhui Wang, Anbang Ye, Gang Ren, Qianran Ma, Wanying Liang, Xiang Lian, Xiwen Wu, Yuting Zhong, Zhuangyan Li, Chaoyu Gong, Guojun Lei, Leijun Cheng, Limin Zhang, Minghao Li, Ruijie Zhang, Silan Hu, Shijie Huang, Xiaokang Wang, Yuanheng Zhao, Yuqi Wang, Ziang Wei, and Yang You. Open-sora 2.0: Training a commercial-level video generation model in \$200k. 2503.09642, 2025b.
- Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton, Bethanie Brownfield, Gavin Buttimore, Max Cant, Sarah Chakera, et al. Scaling instructable agents across many simulated worlds. *arXiv preprint arXiv:2404.10179*, 2024.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III* 18, pages 234–241. Springer, 2015.
- Weijia Shi, Xiaochuang Han, Chunting Zhou, Weixin Liang, Xi Victoria Lin, Luke Zettlemoyer, and Lili Yu. Lmfusion: Adapting pretrained language models for multimodal generation. *arXiv preprint arXiv:2412.15188*, 2024.
- Tomáš Souček and Jakub Lokoč. Transnet v2: An effective deep network architecture for fast shot transition detection. *arXiv preprint arXiv:2008.04838*, 2020.
- Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality. *arXiv preprint arXiv:2307.05222*, 2023.
- Gemma Team. Gemma 3 technical report. *CoRR*, abs/2503.19786, 2025. doi: 10.48550/ARXIV.2503.19786. <https://doi.org/10.48550/arXiv.2503.19786>.
- Hansi Teng, Hongyu Jia, Lei Sun, Lingzhi Li, Maolin Li, Mingqiu Tang, Shuai Han, Tianning Zhang, W. Q. Zhang, Weifeng Luo, Xiaoyang Kang, Yuchen Sun, Yue Cao, Yunpeng Huang, Yutong Lin, Yuxin Fang, Zewei Tao, Zheng Zhang, Zhongshu Wang, Zixun Liu, Dai Shi, Guoli Su, Hanwen Sun, Hong Pan, Jie Wang, Jiexin Sheng, Min Cui, Min Hu, Ming Yan, Shucheng Yin, Siran Zhang, Tingting Liu, Xianping Yin, Xiaoyu Yang, Xin Song, Xuan Hu, Yankai Zhang, and Yuqiao Li. MAGI-1: autoregressive video generation at scale. *CoRR*, abs/2505.13211, 2025. doi: 10.48550/ARXIV.2505.13211. <https://doi.org/10.48550/arXiv.2505.13211>.
- Shengbang Tong, David Fan, Jiachen Zhu, Yunyang Xiong, Xinlei Chen, Koustuv Sinha, Michael Rabbat, Yann LeCun, Saining Xie, and Zhuang Liu. Metamorph: Multimodal understanding and generation via instruction tuning. *arXiv preprint arXiv:2412.14164*, 2024.
- Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
- Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open

- domain textual descriptions. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. <https://openreview.net/forum?id=vOEXS39nOF>.
- Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Fei Wu 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, Wenten 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.
- Qiuhe Wang, Yukai Shi, Jiarong Ou, Rui Chen, Ke Lin, Jiahao Wang, Boyuan Jiang, Haotian Yang, Mingwu Zheng, Xin Tao, et al. Koala-36m: A large-scale video dataset improving consistency between fine-grained conditions and video content. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8428–8437, 2025.
- Wenhao Wang and Yi Yang. Videoufo: A million-scale user-focused dataset for text-to-video generation. *arXiv preprint arXiv:2503.01739*, 2025.
- Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2307.06942*, 2023.
- Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, and Ping Luo. Janus: Decoupling visual encoding for unified multimodal understanding and generation, 2024. <https://arxiv.org/abs/2410.13848>.
- Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin Shi, et al. Pandora: Towards general world model with natural language actions and video states. *arXiv preprint arXiv:2406.09455*, 2024.
- Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.
- Zhucun Xue, Jiangning Zhang, Teng Hu, Haoyang He, Yanan Chen, Yuxuan Cai, Yabiao Wang, Chengjie Wang, Yong Liu, Xiangtai Li, et al. Ultravideo: High-quality uhd video dataset with comprehensive captions. *arXiv preprint arXiv:2506.13691*, 2025a.
- Zihui Xue, Jounghbin An, Xitong Yang, and Kristen Grauman. Progress-aware video frame captioning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13639–13650, 2025b.
- Deheng Ye, Fangyun Zhou, Jiacheng Lv, Jianqi Ma, Jun Zhang, Junyan Lv, Junyou Li, Minwen Deng, Mingyu Yang, Qiang Fu, et al. Yan: Foundational interactive video generation. *arXiv preprint arXiv:2508.08601*, 2025.
- Tianwei Yin, Qiang Zhang, Richard Zhang, William T. Freeman, Frédo Durand, Eli Shechtman, and Xun Huang. From slow bidirectional to fast causal video generators. *CoRR*, abs/2412.07772, 2024. doi: 10.48550/ARXIV.2412.07772. <https://doi.org/10.48550/ARXIV.2412.07772>.
- Jiwen Yu, Yiran Qin, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Gamefactory: Creating new games with generative interactive videos. *arXiv preprint arXiv:2501.08325*, 2025.
- Hangjie Yuan, Weihua Chen, Jun Cen, Hu Yu, Jingyun Liang, Shuning Chang, Zhihui Lin, Tao Feng, Pengwei Liu, Jiazheng Xing, et al. Lumos-1: On autoregressive video generation from a unified model perspective. *arXiv preprint arXiv:2507.08801*, 2025.
- Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. Merlot reserve: Multimodal neural script knowledge through vision and language and sound. In *CVPR*, 2022.
- Chuyang Zhao, Yuxing Song, Wenhao Wang, Haocheng Feng, Errui Ding, Yifan Sun, Xinyan Xiao, and Jingdong Wang. Monoformer: One transformer for both diffusion and autoregression. *arXiv preprint arXiv:2409.16280*, 2024.
- Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-Sora: Democratizing efficient video production for all. 2412.20404, 2024.

Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024.

A Additional details: Controlled Experiments with Video Game Data

A.1 Model configuration details

See Table 6 for key training and inference hyperparameters.

Layers	28
Model Dimension	3072
FFN Dimension	8192
Attention Heads	24
Key/Value Heads	8
Activation Function	SwiGLU
Vocabulary Size	128K
Positional Embeddings – Interleaved Sequence	1D RoPE
Positional Embeddings – Video Only	2D APE
Training Steps	50K
Batch Size	128
Learning rate	3e-4
Max Context Length	15360
Tokens per Frame Chunk	240
Timestep t	$\text{logistic}(\mathcal{N}(0, 1.96))$
Text Dropout Rate $p_{\text{txt-drop}}$	0
Clean Video Flip Rate $p_{\text{clean-vid-flip}}$	0.5
Text Sampling Temperature	0.7
ODE Sampler	Euler
ODE Sampling Steps	50

Table 6 Model configuration details for TV2TV and baselines for experiments on video game data in §3. All model variants adopt a 3B-MoT Transfusion architecture.

A.2 Human Evaluation Task

Pairwise comparisons Annotators evaluate overall video quality within pairings according to the question:

Which video has better visual quality? Examples of poor visual quality include cloudy/flickering generation quality, jumping through walls, static player movement (e.g. moving very slowly in one direction), and random jumps to elsewhere in the map. Do not penalize ghost-like or translucent characters.

Answer choices:

- Left has significantly better visual quality.
- Left has marginally more visual quality.
- Unsure or both seem equally good/bad.
- Right has marginally more visual quality.
- Right has significantly better visual quality.

Pointwise comparisons Annotators evaluate alignment between the generated video and intervention prompt according to the following instruction. In this instruction, ‘Caption A’ refers to the user-controllable intervention prompt included with the generated video, and ‘START’ and ‘STOP’ are assistive visual indicators added posthoc to the generated video corresponding to the intervention timestamp:

Please watch the video. Does Caption A correctly reflect the clip shown between ‘START’ and ‘STOP’? Consider primarily the period shown between START/STOP, although if the caption refers to jumping, reloading, or moving backwards you may also consider the period *immediately* following ‘STOP’. Please use the rest of the video only for general context.

Answer choices:

- The caption correctly reflects the video.
- The caption does not correctly reflect the video.
- Unsure - the player is not actively playing (taken down by enemy and can’t move).
- Unsure - other (please specify in comments).

Additionally, visual quality is evaluated with the question:

How is the overall visual quality of the video? Examples of poor visual quality include cloudy/flickering generation quality, jumping through walls, static player movement (e.g. moving very slowly in one direction), and random jumps to elsewhere in the map. Do not penalize ghost-like or translucent characters.

Answer choices:

- The video has strong visual quality.
- The video has moderate visual quality.
- The video has weak visual quality.
- The video has no visual quality.

B Additional details: Scaling TV2TV to Real World Data

B.1 VLM-based quality filtering prompt

You are a capable model that can determine if the video is high quality or low quality, here are criteria to determine the quality:

- 1) Video is low quality if it has person talking to camera without any other motion, face in corners is OK.
- 2) Video is low quality if there is jittery motion due to camera movement.
- 3) Video is low quality if there is no motion with blank or static screen or image with just zoom in and zoom out.
- 4) Video is high quality if it has meaningful sports content like highlights of a game being played.

Now please rate video with a score between 1 and 10, where 1 is low quality and 10 is high quality, return the score in json format e.g.: {‘quality_score’: <predicted score>}.

Here are frames from new video sampled from start, middle and end of video:

B.2 VLM-based differential captioning prompt

You are a cautious video describer.
You will be shown multiple video segments from the same source video, shown in chronological order.
Describe what happens in EACH segment separately.
DO NOT reference ‘the video’ or ‘the segment’ in your descriptions.
DO NOT describe any text in the video.
Most important: Describe the actions or movements of the main characters or objects in each segment.
DO NOT anticipate future actions; only describe actions that are clearly visible in the current segment.
DO NOT repeat descriptions from previous segments.
If nothing meaningfully changes in the current segment compared to previous segments, use an empty string "" for the description of the current segment.
Keep each description concise (under 50 words). DO NOT hallucinate.

Format your output EXACTLY as follows:
Description of segment 1: [your description here]
Description of segment 2: [your description here]
...
Description of segment <N>: [your description here]
Here are the video segments in order: <VIDEO SEGMENTS>

B.3 Model configuration details

See [Table 7](#) for key training and inference hyperparameters.

B.4 Human Evaluation Task

Video Comparison

Input Instruction
Female runners compete in a track race, with the leading runner in yellow maintaining her position while others adjust their pace. The runner in red lags further behind the main group.

Your Evaluation
Video ID: 0

Your Ratings
Please rate each video on the following criteria:
Which video is most aligned with the text language instruction?
Which video has better overall quality?
Which video has better fidelity to the race event?
Which video has better visual quality?
Which video do you prefer (subjectively)?
Additional Comments (Optional):

Submit Ratings

[Continue to Next Video](#)

Figure 9 Example UI for evaluating alignment, fidelity, quality, and overall preference in TV2TV and comparison models.

For the human evaluation, all pairs are evaluated by a pool of professional external annotators via the Turing platform for increased robustness. A similar user interface as the one used by annotators is shown in Figure 9.

We include the evaluation questions answer by the annotators for the results discussed in §4:

Layers	32
Model Dimension	4096
FFN Dimension	14336
Attention Heads	32
Key/Value Heads	8
Activation Function	SwiGLU
Vocabulary Size	128K
Positional Embeddings – Interleaved Sequence	1D RoPE
Positional Embeddings – Video Only	2D APE
Training Steps	250K
Batch Size	512
Learning Rate	3e-4
Max Context Length	13056
Tokens per Frame Chunk	240
Timestep t	$\text{logistic}(\mathcal{N}(0, 1.96))$
Text Dropout Rate $p_{\text{txt-drop}}$	0.05
Clean Video Flip Rate $p_{\text{clean-vid-flip}}$	0.2
Text Sampling Temperature	0.7
ODE Sampler	Euler
ODE Sampling Steps	50

Table 7 Model configuration details for TV2TV and baselines for experiments on real world sports data in §4. All model variants adopt an 8B-MoT Transfusion architecture.

- **Prompt alignment:** Which video is more aligned with the input language instruction?
 - Left is significantly more aligned with the text instruction.
 - Left is marginally more aligned with the text instruction.
 - Unsure or both seem equally good/bad.
 - Right is marginally more aligned with the text instruction.
 - Right is significantly more aligned with the text instruction.
- **Real world fidelity:** Which video has better fidelity to the real world?
 - Left has significantly better fidelity to the real world.
 - Left has marginally better fidelity to the real world.
 - Unsure or both seem equally good/bad.
 - Right has marginally better fidelity to the real world.
 - Right has significantly better fidelity to the real world.
- **Visual quality:** Which video has better visual quality?
 - Left has significantly better visual quality.
 - Left has marginally better visual quality.
 - Unsure or both seem equally good/bad.

- Right has marginally better visual quality.
- Right has significantly better visual quality.
- **Holistic preference:** Which video do you prefer holistically?
 - I strongly prefer Left.
 - I somewhat prefer Left.
 - Unsure or both seem equally good/bad.
 - I somewhat prefer Right.
 - I strongly prefer Right.