

KeyVID: Keyframe-Aware Video Diffusion for Audio-Synchronized Visual Animation

Xingrui Wang^{1,2} Jiang Liu¹ Ze Wang¹ Xiaodong Yu¹ Jialian Wu¹
Ximeng Sun¹ Yusheng Su¹ Alan Yuille² Zicheng Liu¹ Emad Barsoum¹
¹Advanced Micro Devices ²Johns Hopkins University

Abstract

Generating video from conditions, such as text, image, and audio, enables both spatial and temporal control, leading to high-quality generation results. Videos with dramatic motions often require a higher frame rate to ensure smooth motion. Currently, most audio-to-visual animation models use uniformly sampled frames from video clips. However, these uniformly sampled frames fail to capture significant key moments in dramatic motions at low frame rates and require significantly more memory when increasing the number of frames directly. In this paper, we propose KeyVID, a keyframe-aware audio-to-visual animation framework that significantly improves the generation quality for key moments in audio signals. Given an image and an audio input, we first localize keyframe time steps from the audio. Then, we use a keyframe generator to generate the corresponding visual keyframes. Finally, we generate all intermediate frames using the motion interpolator. Through extensive experiments, we demonstrate that KeyVID significantly improves audio-video synchronization and video quality across multiple datasets, particularly for highly dynamic motions. The code and demo will be released after acceptance.

1. Introduction

Recent years have witnessed remarkable progress in video generation, driven by advancements in diffusion-based models [1, 3, 24]. These frameworks typically condition the generation process on text prompts and/or image inputs, where the text provides semantic guidance, while the image specifies spatial composition. Despite their success, these methods largely focus on aligning visual outputs with static text or image, leaving dynamic, time-sensitive modalities such as audio underexplored.

Audio-Synchronized Visual Animation (ASVA) [26] aims to animate objects in a static image into a video with motion synchronized with the input audio. To achieve precise synchronization, it is crucial to align key visual actions

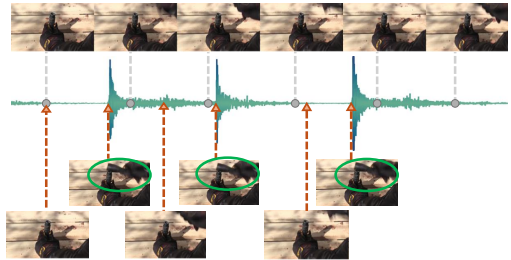


Figure 1. *Top*: Uniformly sampled sparse frames fail to capture the key moments in audio. *Bottom*: Key frames precisely aligned with the hammer strike, matching the critical moments.

with their corresponding audio signals. For example, given hammering sounds, the hammer should strike at the exact moment the impact sound occurs. However, this synchronization is constrained by frame rates—AVSyncD [26] operates at 6 FPS, while audio carries fine-grained temporal information, causing key moments to be lost in sparse low frame rate videos (see Fig. 1).

Directly training on high frame rate videos incurs substantial computational costs in GPU memory and training time. A common solution adopts a two-stage strategy that generates low frame rate videos then applies frame interpolation [1, 18]. However, this approach struggles in highly dynamic sequences, where critical events may be lost due to the sparsity of initial uniform frames. To ensure accurate audio-visual synchronization while maintaining computation efficiency, we propose **KeyVID**, a **Keyframe-aware Video Diffusion** framework. We first develop a keyframe selection strategy that identifies critical moments based on optical flow-based motion score. A *Keyframe Localizer* predicts keyframe positions directly from audio. Instead of uniform downsampling, we train a *Keyframe Generator* that explicitly captures crucial moments without requiring excessive frames. A specialized *Motion Interpolator* then synthesizes intermediate frames between the uneven keyframes. This approach mimics animation workflows in the animation industry where a “Key Animator” establishes crucial moments and a “In-betweener” fills gaps.

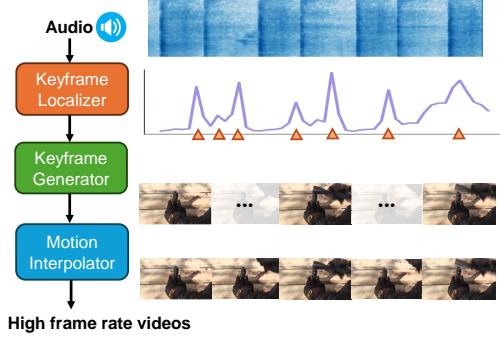


Figure 2. KeyVID video generation pipeline.

We demonstrate that our keyframe-aware approach outperforms state-of-the-art methods in video quality and audio-video synchronization in ASVA task. The main contributions are:

- A keyframe-aware framework that localizes keyframes from audio and generates them as video diffusion model.
- A motion interpolation network enabling high frame rate generation while maintaining efficiency.
- Superior performance in audio-synchronized video generation, particularly for dynamic scenes.

2. Related Work

Video Diffusion Models. Recent video diffusion models [1, 3, 8, 24] generate high-quality videos by learning to denoise Gaussian noise through a reversed diffusion process. For efficiency, latent video generation [2, 24] encodes video into latent space. These models incorporate text [3, 9] and image [6, 24] to guide generation, but typically use uniform frame sampling, limiting their ability under critical motion moments at low frame rates.

Audio-to-Video Generation. While early audio-conditioned video synthesis focused on domain-specific tasks [13, 19], recent work leverage pre-trained audio encoders [4, 5] for general video generation. Some approaches use audio as global features [10, 21], while others [11, 12, 17] consider temporal alignment. AVSyncD [26] introduced time-dependent audio features for finer temporal control, but remains limited by low frame rates (6 FPS) for dynamic motions. Directly increasing frame density requires prohibitive computational resources.

3. Methods

We present KeyVID, a keyframe-aware video generation framework. Given input audio and first frame, we follow a three-stage process to get the video output (see Fig. 2): (1) the *Keyframe Localizer* first predicts keyframe locations from audio; (2) the *Keyframe Generator* produces the keyframes conditioned on audio and image for each positions; (3) *Motion Interpolator* synthesize all the intermediate frames to obtains a smooth high frame rate videos.

3.1. Keyframe Localization from Audio

We train a *Keyframe Localizer* to infer keyframe locations from audio by exploiting correlations between acoustic events and motion changes. These motion changes are defined by analyzing the motion score from training videos and serve as pseudo labels. Using RAFT [22], we compute optical flow \mathbf{OF}_t between consecutive frames and calculate the motion score as: $M(t) = \sum_{i,j} (|u_t(i,j)| + |v_t(i,j)|)$, where u_t, v_t are horizontal/vertical flow components. The keyframe localizer takes audio spectrograms as input, which consists of pretrained ImageBind [5] as feature extractor and fully connected layers to predict motion scores. The training process is guided by L2 loss.

3.2. Audio-conditioned Keyframe Generation

Our keyframe generator produces T_K keyframes from a T -frame video, conditioned on audio, first frame, and text.

3.2.1. Keyframe Data Selection

Rather than the uniformly sampled T_K frames [24] to train the video diffusion model, we select $T_K \ll T$ keyframes from the peaks and valleys of motion score which represents the most crucial moments of motions in a video clip. We choose the first frame, randomly select up to $\frac{T_K}{2} - 1$ peaks, add valleys between consecutive peaks, then evenly sample remaining frames. This ensures coverage of critical moments while approximating uniform sampling for smooth sequences. The details of the selection algorithm are in Appendix B. The selected keyframe indices $\{t_1, \dots, t_{T_K}\}$ serve as additional conditions to the following step.

3.2.2. Keyframe Generator Diffusion Model

The keyframe generator introduces two key enhancements: (1) A **Frame index embedding** encodes each frame’s absolute position, ensuring coherence when generating non-uniformly distributed frames; (2) **Multi-modal condition features** consist of global text features, and audio and image features extracted from corresponding keyframe timesteps.

As shown in Fig. 3(b), we build upon latent diffusion models [1, 2] with pretrained encoder and decoder from Xing et al. [24]. The latent features are represented as $\mathbf{z} \in \mathbb{R}^{B \times T_k \times C \times h \times w}$ where B denotes the batch size, T_k the number of frames at each denoising step, h and w the spatial dimensions, and C is the feature channels.

Frame Index Embedding. We introduce an embedding layer to encode the absolute index of each keyframe within the original video sequence $\{i_t\}_{t=1}^{T_K}$. The frame index embedding $\mathbf{f}_{\text{emb}} \in \mathbb{R}^{B \times T_K \times C}$ is added up with the latent features \mathbf{z} before passing them into the denoising U-Net, ensuring explicit positional information is provided to the network to enable generation of non-uniformly spaced frames.

Audio Feature Conditions. We extract audio features using pretrained ImageBind [5], which encodes spectrograms into global and local tokens capturing semantic and tempo-

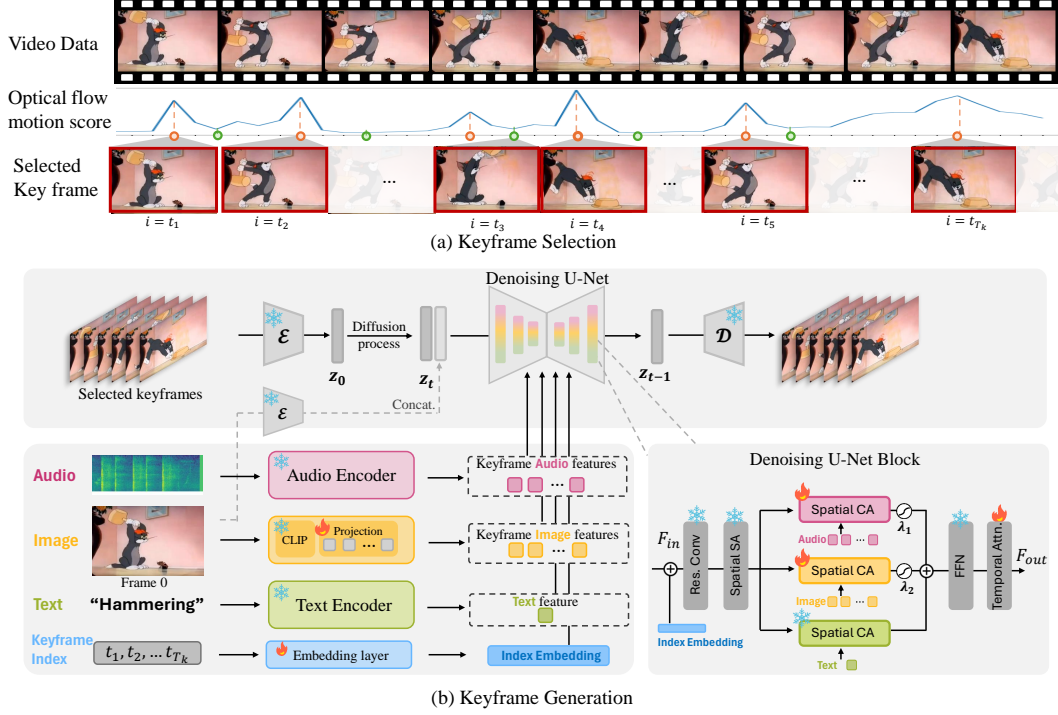


Figure 3. (a) Keyframe selection from motion score curve. (b) Diffusion model for keyframe generation with audio, image, text conditions and frame index embedding.

ral information. We segment the features into T timesteps matching the final video length, then select T_K features at keyframe indices $\{i_t\}_{t=1}^{T_K}$ for cross-attention fusion in the U-Net, ensuring audio-visual synchronization.

Image/Text Feature Conditions. We extract first-frame features using CLIP [16] and project them to T frames via learned tokens, yielding $\mathbf{f}_{\text{img}} \in \mathbb{R}^{B \times T \times C \times H \times W}$. We select T_K features at keyframe indices $\{i_t\}_{t=1}^{T_K}$ for the cross attention during denoising. The text descriptions are encoded using CLIP [16] and repeated for all T_K keyframes.

Feature Fusion. Each conditioning feature (audio, image, and text) is processed separately through spatial cross-attention layers during the denoising step. Given input latent features \mathbf{F}_{in} , one denoising step computes:

$$\mathbf{F}_{\text{out}} = \mathbf{F}_{\text{in}} + \text{SpatialAttn}(\mathbf{F}_{\text{in}}, \mathbf{f}_{\text{txt}}) + \lambda_1 \cdot \text{SpatialAttn}(\mathbf{F}_{\text{in}}, \mathbf{f}_{\text{aud}}) + \lambda_2 \cdot \text{SpatialAttn}(\mathbf{F}_{\text{in}}, \mathbf{f}_{\text{img}}),$$

where λ_1, λ_2 are learnable fusion weights.

3.3. Motion Interpolator

After generating T_K keyframes, we synthesize intermediate frames using a motion interpolator. Unlike uniform interpolation [1, 24] that predicts between first/last frames, we adapt our keyframe generator to condition on all generated keyframes via masked frame conditioning. We use FreeNoise [15] to generate all T frames in a single pass. Details are in Appendix C.

4. Experiments

4.1. Implementation Details

Datasets. We evaluate on AVSync15 [26], Landscapes [12], and TheGreatestHits [14]. AVSync15 contains 15 activity classes with synchronized audio-video. Landscapes features natural scenes with ambient sounds. TheGreatestHits contains percussive hitting sounds aligned with motions. We use 2-second clips at 24fps (48 frames) resized to 320x512, with $T_K = 12$ keyframes.

Baselines. We compare with: (1) **T+A**: TPoS [10], TempoToken [25]; (2) **I+T**: DynamiCrafter [24], I2VD [26]; (3) **I+T+A**: CoDi [21], AADiff [11], AVSyncD [26].

Metrics. We use the Frechet Image Distance (FID) [7] and Frechet Video Distance (FVD) [23] to evaluate the visual quality and temporal coherence of synthesized videos. We evaluate the audio synchronization with the generated video by **RelSync** and **AlignSync** proposed by Zhang et al. [26].

4.2. Quantitative results

Table 1 presents results on three datasets. On AVSync15, our KeyVID achieves the highest AlignSync (24.09) and RelSync (48.30) scores, demonstrating the effectiveness of our keyframe selection strategy in capturing crucial dynamic moments. Our method also achieves competitive visual quality (FID=11.0, FVD=262.3), outperforming AVSyncD. On Landscapes, which has less dynamics and is used for evaluating visual quality, our method achieves the lowest FVD score (391.09). On TheGreatestHits, featuring distinct percussive audio events, our approach achieves the best performance across all metrics, with notable improvements over AVSyncD.

Input	Model	AVSync15				Landscapes				The Greatest Hit			
		FID↓	FVD↓	AlignSync↑	RelSync↑	FID↓	FVD↓	AlignSync↑	RelSync↑	FID↓	FVD↓	AlignSync↑	RelSync↑
T+A	TPoS [10]	13.5	2671.0	19.52	42.50	16.5	2081.3	23.12	48.15	33.85	3327.90	21.48	44.90
	TempoToken [25]	12.2	4466.4	19.74	44.05	16.4	2480.0	24.21	48.65	25.90	3300.53	21.56	45.38
I+T	I2VD [26]	12.1	398.2	21.80	43.92	16.7	539.5	24.74	49.89	9.1	425.0	22.05	44.58
	DynamiCrafter [24]	11.7	400.7	21.76	43.68	23.51	445.8	24.17	49.63	12.4	337.71	22.82	45.85
I+T+A	CoDi [20]	14.5	1522.6	19.54	41.51	20.5	982.9	22.63	45.48	21.78	1336.00	22.30	45.35
	TPoS [10]	11.9	1227.8	19.67	39.62	16.2	789.6	23.51	47.05	28.43	1370.57	22.04	45.55
	AVSyncD [26]	11.7	349.1	22.62	45.52	16.2	415.2	24.82	49.93	8.7	249.3	22.83	45.95
	KeyVID (Ours)	11.00	262.3	24.08	48.33	23.28	391.0	24.35	49.95	12.1	202.1	22.91	46.03
	Static Groundtruth	-	1220.4	21.83	43.66	-	1177.5	25.79	51.59	-	348.9	24.36	48.73
		-	-	25.04	50.00	-	-	25.01	50.00	-	-	25.02	50.00

Table 1. Performance on AVSync15 and Landscapes.

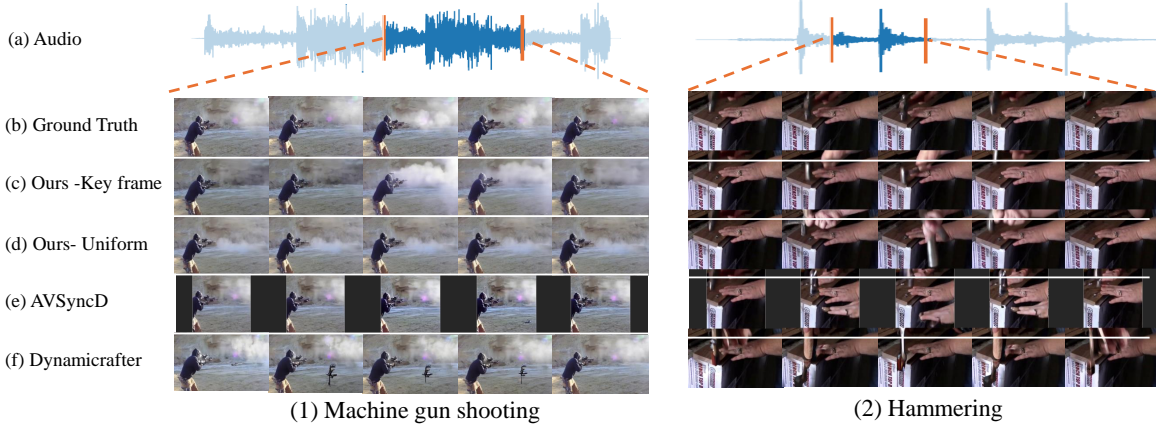


Figure 4. Qualitative comparison of KeyVID and baseline methods. We crop a key motions on audio waveform in (a) and the corresponding ground truth video in (b) as references and compare the generated video clip between models from (c) to (f). Compared with other models, our KeyVID with keyframe awareness in (c) have better alignment of the motion peaks with audio signals, for example, hitting the hammering, smoke in gun shooting.

4.3. Visualization

Figure 4 presents qualitative comparisons between our method and baseline approaches. Our keyframe-aware approach more accurately captures motion peaks that align with audio events, such as the exact moment of impact in hammering or the smoke in gun shooting. This demonstrates the effectiveness of keyframe-aware training across both high- and low-intensity motion scenarios.

4.4. Ablation Study

We conduct ablation studies to validate the effectiveness of keyframe awareness, as shown in Tab. 2. Specially, we train a variant of the KeyVID model using uniformly sampled 12 frames. Since our method generates high-frame-rate videos (48 frames/2s), we evaluate under two settings: (1) downsample our output to 12 frames to compare with the baseline’s 12 uniform frames; (2) interpolate the baseline’s 12 frames to 48 using the same method from Sec. 3.3 and evaluate on 48 frames. KeyVID consistently outperforms uniform sampling in both settings, with notable improvements in synchronization metrics (AlignSync and RelSync). These results support our hypothesis that selecting keyframes based on audio and motion cues enhances temporal alignment between audio events and visual dynamics.

Method	FID↓	FVD↓	AlignSync↑	RelSync↑
<i>Evaluate on 12 frames</i>				
KeyVID	11.00	262.34	24.08	48.33
Uni. Frame	11.01	273.40	23.53	47.23
<i>Evaluate on 48 frames</i>				
KeyVID	4.83	335.68	24.08	48.37
Uni. Frame	4.90	337.10	23.96	48.09

Table 2. Ablation study comparing keyframe-based generation with uniform sampling. KeyVID achieves better performance in both audio synchronization and visual quality with keyframes.

5. Conclusion

We introduce a keyframe-aware, audio-synchronized visual animation model that improves video quality and audio alignment, especially under dynamic motion. Our approach first detects keyframes locations from audio input, then generates them with video diffusion model, and then interpolates intermediate frames for smooth, high-frame-rate output with low memory cost. Experiments on multiple datasets show significant gains in both visual quality and synchronization.

References

- [1] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 1, 2, 3
- [2] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575, 2023. 2
- [3] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. 1, 2
- [4] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 2
- [5] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023. 2
- [6] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. In *ICLR*, 2024. 2
- [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 3
- [8] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 2
- [9] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 2
- [10] Yujin Jeong, Wonjeong Ryoo, Seunghyun Lee, Dabin Seo, Wonmin Byeon, Sangpil Kim, and Jinkyu Kim. The power of sound (tpos): Audio reactive video generation with stable diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7822–7832, 2023. 2, 3, 4
- [11] Seungwoo Lee, Chaerin Kong, Donghyeon Jeon, and Nojun Kwak. Aadiff: Audio-aligned video synthesis with text-to-image diffusion. *arXiv preprint arXiv:2305.04001*, 2023. 2, 3
- [12] Seung Hyun Lee, Gyeongrok Oh, Wonmin Byeon, Chanyoung Kim, Won Jeong Ryoo, Sang Ho Yoon, Hyunjun Cho, Jihyun Bae, Jinkyu Kim, and Sangpil Kim. Sound-guided semantic video generation. In *European Conference on Computer Vision*, pages 34–50. Springer, 2022. 2, 3
- [13] Evonne Ng, Javier Romero, Timur Bagautdinov, Shaojie Bai, Trevor Darrell, Angjoo Kanazawa, and Alexander Richard. From audio to photoreal embodiment: Synthesizing humans in conversations. In *ArXiv*, 2024. 2
- [14] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2405–2413, 2016. 3
- [15] Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. Freenoise: Tuning-free longer video diffusion via noise rescheduling. *arXiv preprint arXiv:2310.15169*, 2023. 3
- [16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. Pmlr, 2021. 3
- [17] Ludan Ruan, Yunzhi Ma, Hongjie Yang, Haoxian He, Bing Liu, Jianlong Fu, Nenghai Yuan, Qin Jin, and Bing Guo. MM-Diffusion: Learning multi-modal diffusion models for joint audio and video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14012–14021, 2023. 2
- [18] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations*, 2023. 1
- [19] Xusen Sun, Longhao Zhang, Hao Zhu, Peng Zhang, Bang Zhang, Xinya Ji, Kangneng Zhou, Daiheng Gao, Liefeng Bo, and Xun Cao. Vividtalk: One-shot audio-driven talking head generation based on 3d hybrid prior. *arXiv preprint arXiv:2312.01841*, 2023. 2
- [20] Zhaohan Tang, Zhilin Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023. 4
- [21] Zhaohan Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. CoDi: Any-to-any generation via composable diffusion. *arXiv preprint arXiv:2305.11846*, 2023. 2, 3
- [22] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 2
- [23] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 3

- [24] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. Dynamicrafter: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer, 2024. [1](#), [2](#), [3](#), [4](#)
- [25] Guy Yariv, Itai Gat, Sagie Benaim, Lior Wolf, Idan Schwartz, and Yossi Adi. Diverse and aligned audio-to-video generation via text-to-video model adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6639–6647, 2024. [3](#), [4](#)
- [26] Lin Zhang, Shentong Mo, Yijing Zhang, and Pedro Morgado. Audio-synchronized visual animation. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024. [1](#), [2](#), [3](#), [4](#)

KeyVID: Keyframe-Aware Video Diffusion for Audio-Synchronized Visual Animation

Supplementary Material

A. Details of Keyframe Localization Network from Audio

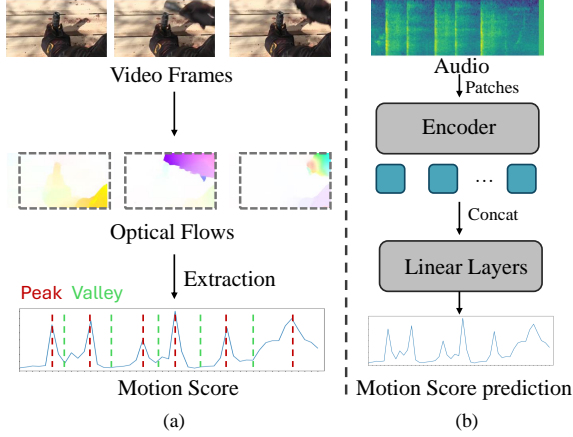


Figure 5. (a) We first calculate the optical flow and then take the average across all pixels for each frame to form a curve of motion score. The peaks (red) and valleys (green) indicate key frames. (b) key frame prediction network from audio, described in Sec. 3.1.

In the Sec. 3.1 of main paper, we introduce that we need to know the position of key frame at the begin of inference by predicting optical motion scores. Here is the detailed structure of this network. The network processes raw audio by converting it into a spectrogram $\mathbf{A} \in \mathbb{R}^{C_A \times T_A}$, where C_A denotes the number of frequency channels and T_A represents the temporal length. The original ImageBind preprocessing pipeline applies a CNN with a kernel stride of (10, 10) to patchify the input spectrogram, producing feature embeddings that are then processed by a transformer-based encoder $f_{\text{audio}} \in \mathbb{R}^{B \times T \times C}$. However, this results in T (e.g. $T=19$) being misaligned with the temporal resolution of the dense motion curve sequence (e.g. 48).

To address this, we modify the CNN stride to (10, 4), increasing the temporal resolution of extracted features (e.g. increase to 46). The transformer encoder then processes the updated feature sequence:

$$\mathbf{F}_{\text{audio}} = f_{\text{audio}}(\mathbf{A}), \quad \mathbf{F}_{\text{audio}} \in \mathbb{R}^{B \times T' \times C}, \quad (1)$$

where $T' > T$ reflects the increased temporal resolution. Since the transformer relies on positional embeddings, we interpolate the pretrained positional embeddings to match the new sequence length T'_A and keep them frozen during training.

The extracted features are passed through fully connected layers to predict a sequence of confidence scores $\mathbf{s} \in \mathbb{R}^{B \times T'}$, where each s_t represents the likelihood of a keyframe occurring at time step t :

$$\mathbf{s} = \sigma(\mathbf{W}\mathbf{F}_{\text{audio}} + \mathbf{b}), \quad (2)$$

where $\mathbf{W} \in \mathbb{R}^{C \times 1}$ and $\mathbf{b} \in \mathbb{R}^{T'_A}$ are learnable parameters, and $\sigma(\cdot)$ is the sigmoid activation function. The model is trained using an L1 loss:

$$\mathcal{L} = \|\mathbf{s} - \hat{\mathbf{s}}\|_1, \quad (3)$$

where $\hat{\mathbf{s}}$ represents the ground-truth keyframe labels derived from optical flow analysis.

B. Details of Keyframe Selection

B.1. Detect Peak and Valley

To identify the local maxima (*peaks*) and minima (*valleys*) from a one-dimensional motion score $\{M(t)\}_{t=1}^T$, we perform the following steps:

1. **Smoothing:** Convolve the raw score $M(t)$ with a short averaging filter with a window size 5, producing a smoothed label $\tilde{M}(t)$. This helps reduce noise and minor fluctuations.
2. **Peak Detection:** Finds all local maxima by simple comparison of neighboring values for $\tilde{M}(t)$. We force a minimum distance of 5 frames between any two detected peaks and requiring a prominence (height relative to its surroundings) of at least 0.1. This returns the indices of the local maxima.
3. **Valley Detection:** Repeat the same peak-finding procedure on the negative of the smoothed signal.

B.2. Sample keyframes

In the main text, we discuss the process of selecting $T_K \ll T$ keyframes based on the motion score $M(t)$ for each frame. Specifically, we first pick the initial frame, then select up to $\frac{T_K}{2} - 1$ peaks among all detected ones (or all peaks if fewer are found). Next, we include a valley between each consecutive pair of selected peaks. Finally, we sample any remaining frames by an evenly distributed (proportional) strategy, which approximates uniform downsampling if few peaks and valleys are present. This approach ensures that smooth motion or weak audio signals, producing limited peaks and valleys, do not degrade the consistency of training for video diffusion models.

Algorithm 1 is the detailed pseudo-code for the full procedure, including both peak and valley selection and the final proportional allocation of remaining key frames.

C. Structure of Motion Interpolation

As shown in Fig. 6, we present the pipeline of motion interpolation network as introduced in Sec. 3.3. After generating T_K keyframes, we use a *motion interpolator* to generate the missing frames back to the full video sequence of length T . Interpolation has been widely used in uniform frame generation [1, 24], where a model predicts a fixed number of intermediate frames given the

Algorithm 1: Keyframe Selection Algorithm

Input: Motion scores $\{M(t)\}_{t=1}^T$, desired keyframe count $T_K \ll T$.

Output: A set of T_K keyframes.

- 1 **Step 1: Detect peaks and valleys** based on $M(t)$.
 - 2 **Step 2: Initialize keyframe list:**
Keyframes $\leftarrow \{\text{first_frame}\}$.
 - 3 **Step 3: Randomly select peaks**
Choose up to $\lfloor \frac{T_K}{2} - 1 \rfloor$
from the detected peaks and add to Keyframes.
 - 4 **Step 4: Insert valleys**
for each pair of consecutive peaks in
Keyframes do
 Select one valley in between and add it to
 Keyframes.
 - 5 **Step 5: Compute how many more keyframes are needed:**
 $R \leftarrow T_K - |\text{Keyframes}|$.
 - 6 **if** $R > 0$ **then**
 - 7 Define a list of N remaining frames (unselected)
 with some weights $\{w_1, \dots, w_N\}$.
 - 8 $W \leftarrow \sum_{i=1}^N w_i$
 - 9 **for** $i \leftarrow 1$ **to** N **do**
 - ideal_share $_i \leftarrow R \cdot \frac{w_i}{W}$;
 - allocated $_i \leftarrow \lfloor \text{ideal_share}_i \rfloor$;
 - 10 $r \leftarrow R - \sum_{i=1}^N \text{allocated}_i$; // Remainder
 after flooring
 - 11 **if** $r > 0$ **then**
 - for** $i \leftarrow 1$ **to** N **do**
 - frac $_i \leftarrow \text{ideal_share}_i - \text{allocated}_i$;
 - Sort frames by frac $_i$ in descending order.
 - for** $j \leftarrow 1$ **to** r **do**
 - $i^* \leftarrow$ index of the j -th largest frac $_i$;
 - allocated $_{i^*} \leftarrow \text{allocated}_{i^*} + 1$;
 - 12 **for** $i \leftarrow 1$ **to** N **do**
 - if** allocated $_i > 0$ **then**
 - Keyframes $\leftarrow \text{Keyframes} \cup \{\text{frame}_i\}$;
 - 13 **return** Keyframes
-

first and last frame. However, for keyframe-based generation, the positions of missing and available frames vary, introducing additional challenges.

To address this, we adapt our *keyframe generator* diffusion model into a *motion interpolator* model that generates T_K frames at once using masked frame conditioning. The overall architecture remains nearly unchanged, with the primary difference lying in how image conditions are incorporated. Rather than conditioning solely on the first frame, the model utilizes the features of generated keyframes as conditions, thereby learning to synthesize the missing frames in between. This approach facilitates interpolation

between non-uniformly distributed keyframes while maintaining temporal consistency. A pipeline can be found in Appendix C.

To generate a full video with T frames in a single pass, we incorporate FreeNoise [15] to increase the number of output frames during inference. This allows the interpolation model to take all generated keyframes as conditioning inputs and predict all missing frames in one single step.

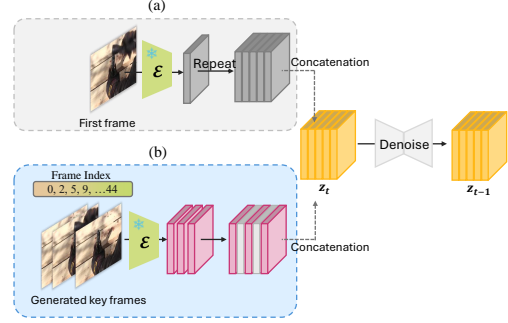


Figure 6. The frame interpolation models has the same structure as the original keyframe generation model, but has different image features for concatenation. (a) For the keyframe generation in Sec 3.2.2, the first frame features are repeated to match the frame length of the latent vector; (b) In frame interpolation, the condition feature from keyframes are padded with zero tensor between keyframe locations to match the frame length.

D. Motion score prediction evaluation

Quantitative result. We evaluate the keypoint detected from the predicted motion score with the ground truth score. We calculate the average precision with a distance threshold t . In this way, for each keypoint in ground truth motion score curve, if it can match with a predicted keypoint with distance lower than t , it will be consider as a successful match. The average precision means the the average of $N_{match}/N(\text{total})$ across all instance, denoted as $AP@t$. We achieve the $AP@3 = 60.57\%$ and $AP@5 = 77.92\%$.

Visualization

E. More Qualitative Results of Video Generation

As the generation result need to be watch with audio for the best experience, we have put more visualization result into the supplementary as mp4 files.

F. Experimental Details

For the experiments of KeyVID on three dataset AVSyncD, Landscape, and TheGreatestHit, we all train on resolution 320×512 as Dynamicrafter [24]. During the inference time, we use ddim sampling with step 90. The temporal length of both key frame generation and interpolation model are all 12. As our interpolation module use freenoise[15] technique to obtain the final 48 frames in one run. we change the windows size 12 and the stride 6 to fit our temporal length.

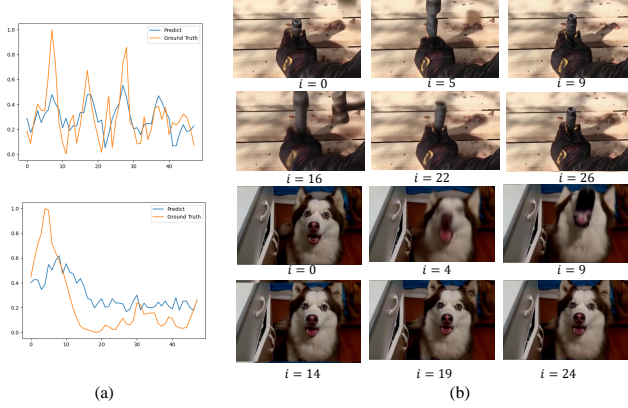


Figure 7. Visualization of (a) Predicted motion score from audio with the ground truth caluate from video data; and (b) the generated video keyframe by diffusion network described in Sec. 3.2.2 before interpolations.

G. Multimodal Classifier Free Guidance

Similar to Xing et al. [24], we introduce three guidance scales s_{img} , s_{txt} , and s_{aud} to extend video generation with additional audio control. These scales allow balancing the influence of different conditioning modalities in video generation. The modified noise estimation function is defined as:

$$\begin{aligned} \hat{\epsilon}_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \mathbf{c}_{\text{txt}}, \mathbf{c}_{\text{aud}}) &= \epsilon_{\theta}(\mathbf{z}_t, \emptyset, \emptyset, \emptyset) \\ &+ s_{\text{img}} (\epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \emptyset, \emptyset) - \epsilon_{\theta}(\mathbf{z}_t, \emptyset, \emptyset, \emptyset)) \\ &+ s_{\text{txt}} (\epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \mathbf{c}_{\text{txt}}, \emptyset) - \epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \emptyset, \emptyset)) \\ &+ s_{\text{aud}} (\epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \mathbf{c}_{\text{txt}}, \mathbf{c}_{\text{aud}}) - \epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}_{\text{img}}, \mathbf{c}_{\text{txt}}, \emptyset)). \end{aligned} \quad (4)$$

Here, \mathbf{c}_{img} , \mathbf{c}_{txt} , and \mathbf{c}_{aud} represent image, text, and audio conditioning, respectively. The newly introduced audio guidance scale s_{aud} enables the model to integrate temporal audio cues, ensuring synchronized motion generation in audio-reactive video synthesis. By adjusting these guidance parameters, we can control the relative impact of each modality in the final video output.

In experiment, we choose the audio guidance scale to 7.5 and image guidance scale to 2, for both keyframe generation network and frame interpolation network. As we add the audio guidance as a new feature, we compare the result from different audio guidance from 4.0 to 11.0 as list in Tab. 3. Although the higher audio guidance obtains a better audio synchronization score (RelSync and AlignSync) we finally choose the one with the best visual quality (FVD and FID) but still ahiveve compatible audio synchronization score.

s_{aud}	FID↓	FVD↓	AlignSync↑	RelSync↑
4.0	11.4	270.5	48.18	24.14
7.5	11.0	262.3	48.33	24.08
9.0	11.1	277.2	48.55	24.16
11.0	11.1	278.6	48.66	24.22

Table 3. Performance metrics for different guidance values.