

Comic2Video: From Panels to Motion

Jet Wu (jetw) & Ananya Sharma (ananyash) & Rithwick Sethi (rithwics)
10-423/623 Generative AI Course Project

December 9, 2025

1 Introduction

The goal of this project is to determine whether a modern video diffusion model can be fine-tuned to produce context-consistent anime-style animation. We focus on *Attack on Titan*, a long-running anime with distinctive art style and complex narratives. Our objective is to finetune a latent video model to learn the spatial style, character appearance and motion dynamics from multiple episodes of the series, then evaluate its ability to recreate scenes from another episode. The problem is challenging because anime episodes contain long scene cuts and semantic jumps; generating a sequence of frames within a single scene is comparatively easy, while crossing cuts requires maintaining character identity and narrative continuity.

We adapt WAN 2.1 14 B first-last-frame, a large video diffusion model, and train a smaller variant to fit our compute budget. We treat our training episode as a dense sequence of key frames and extract sparse panels for conditioning. During training, the model sees only the first and last frame of each clip and must reconstruct the intervening frames; at test time it must hallucinate frames between panels extracted from a different episode. We evaluate with perceptual metrics (LPIPS and SSIM) and a temporal metric (Fréchet Video Distance, FVD). Human-subject evaluation assesses qualitative aspects such as style consistency and narrative plausibility.

Our final executive summary brings together the baseline LoRA tuning pipeline, the completed scene aware extensions and also the updated 13-frame/30-step experiments. We ran LoRA ranks 8, 16, 32 with the rank-16 checkpoint showing a modest FVD improvement over the frozen base model.

With the completed runs, rank 32 now yields the strongest improvement, reducing FVD by nearly 470 points while maintaining perceptual quality.

2 Dataset and Task

Our dataset comprises the first publicly available 30-minute episode of *Attack on Titan* downloaded from

YouTube. Using our `create_data.py` script we pre-processed this video. `create_data.py` resamples the video to 16 fps and resizes each frame to 832×480 pixels using bicubic interpolation. It then generates fixed-length clips of 81 frames ($4k + 1$) with a stride of 16 frames, basically sampling about one clip per second of video. Each clip has the first and last frame along with the entire frame sequence. Frames are saved as PNGs in a Hugging Face dataset at `attack-on-genai/video-frames`. The dataset is around 35GB after processing and has 1,238 training clips. We hold out the next chronological episode from the series as unseen data for testing. This test episode shares the same characters and art style but contains different scenes and plot points which makes it a stringent test of generalization.

Given a sparse set of conditioning frames $\{F_{t_0}, F_{t_N}\}$ from a clip of length N frames, our model is tasked with predicting intermediate frames $F_{t_1 \dots t_{N-1}}$ which (a) preserve the style and character appearance of the input frames, (b) exhibit smooth and physically plausible motion, and (c) remain consistent with the narrative context. We consider two regimes: within-scene interpolation (no shot boundaries) and across-scene interpolation (conditioning frames straddle a shot boundary or narrative jump). For both cases we report FVD, SSIM, PSNR and LPIPS, and Section 5 details the updated 13-frame, 30-step evaluation setup we adopted to stretch our limited GPU hours.

To provide conditioning text for the diffusion model, `create_metadata.py` loads the pre-processed clips, encodes the first and last frames to base64 and uses an OpenAI GPT-5.1 model to generate positive prompts describing the scene, characters and motion and negative prompts describing unwanted artifacts such as low resolution or jitter. The outcome is written to a `metadata.csv` file inside the dataset. In parallel we prototyped a lighter Moondream2-based prompt generator (Listing ??) that summarizes the first, last and strided middle frames for every clip. Although we ran out of compute before fully integrating those prompts into the training loop, the text assets already exist for future experiments.

84 3 Related Work

85 **Stable Video Diffusion (SVD)** uses a high-resolution
86 latent diffusion model with temporal convolutions
87 and cross-attention that is trained using text-to-image
88 and video pretraining and fine-tuning. Its multi-
89 conditioning design (text, image, frame rate) accom-
90 plishes strong temporal coherence and content fidelity
91 and has consequences for the effective scaling of latent
92 diffusion to video.

93 **VideoCrafter 1** presents open-source text-to-video
94 and image-to-video models, generating 1024x576
95 videos that outperform the other open-sourced systems
96 while maintaining the structure of reference images,
97 proving the feasibility and usefulness of community-
98 maintained video diffusion models.

99 **Make-A-Video (Meta)** transfers text-to-image dif-
100 fusion to video by pretraining on paired image-text and
101 unlabeled video, using temporal U-Nets and spatial-
102 temporal diffusion to produce high resolution videos
103 without paired text-video data, motivating our use of
104 strong image priors plus motion learning.

105 **Imagen Video (Google)** extends diffusion to high-
106 definition videos using a cascade: a base generator
107 plus spatial and temporal super-resolution stages, with
108 progressive distillation and classifier-free guidance that
109 enables fast sampling. This informs our training strat-
110 egy.

111 **MCVD** trains a diffusion model on sequences with
112 randomly masked past/future frames and allows for all
113 three tasks of prediction, unconditional generation, and
114 interpolation in one framework. Its block-conditioning
115 and autoregressive block generation inspire our first-
116 last-frame conditioning.

117 **ToonCrafter** adapts diffusion priors to cartoons
118 through toon-rectification learning, a dual-reference
119 3D decoder, and a sketch encoder to achieve convinc-
120 ing animation with support for large nonlinear motion.

121 **W.A.L.T.** uses a transformer-based latent diffu-
122 sion model; it adopts a causal encoder with win-
123 dowed spatial/spatiotemporal attention, which allows
124 for memory-efficient training on images and videos
125 and can generate high-resolution outputs, 512x896,
126 and demonstrates that transformers scale well for video
127 generation.

128 **Tune-A-Video** fine-tunes a text-to-image diffusion
129 model into a text-to-video generator with a single
130 text-video pair and a tailored spatiotemporal attention
131 mechanism. Its ability to adapt large models with min-
132 imal video data directly motivates our attack-on-titan
133 fine-tuning.

4 Methods

4.1 Baseline: First-Last Frame Fine-Tuning

134 Our baseline model is built around the WAN 2.1
136 FLF2V-14B architecture, a video diffusion trans-
137 former(DiT) paired with a 3D causal variational au-
138 toencoder(VAE) that conditions on two reference im-
139 ages (first and last frames) to generate intermediate
140 frames. We have implemented a complete LoRA fin-
141 tunning pipeline to learn a small delta on this 14B trans-
142 former while freezing its components (VAE, CLIP im-
143 age encoder and text encoder).

144 **Data loading and dynamic 4k+1 sam-
145 pling:** WanFLF2VDatasetFromHF in
146 train.py scans the Hugging Face dataset, while
147 WanFLF2VDatasetFromLocal supports cached
148 local-splits for repeated runs and its evaluation. Both
149 the functions provide the first, last and the intermediate
150 frames, along with the prompts. In our final pipeline,
151 clips are no longer assumed to be fixed-length 81-
152 frame sequences. Instead, we subsample 4k+1 frames.
153 For example, with num_frames=9, we select evenly
154 spaced frames from the 81 available (indices 0, 10, 20,
155 ..., 80), effectively taking k = 2. This is done via the
156 training configuration’s valid 4k+1 values. Because
157 clips vary in length, we keep effective batch size at 1
158 for stability and memory predictability.

159 **LoRA injection and rank sweep:** We inject
160 LoRA adapters into the spatial and temporal at-
161 tention modules of the transformer (to_q, to_k,
162 to_v, to_out.0). To address ablation require-
163 ments and quantify parameter efficiency trade offs,
164 we run a lora_rank sweep over {8, 16, 32} with
165 lora_alpha=8 and dropout 0 under consistent train-
166 ing/evaluation settings. We train only LoRA parame-
167 ters with AdamW, using gradient checkpointing, op-
168 tional mixed precision (bf16), gradient clipping and
169 support for checkpoint resume (with optional quanti-
170 zation).

4.2 Prompt Metadata Generation and Condi- 172 tioning

173 The final pipeline incorporates prompt conditioning
174 systematically: training and inference accept posi-
175 tive and negative prompts per sample and text em-
176 beddings are computed via a helper that clamps
177 tokenizer max length to avoid pathological set-
178 tings. To reduce reliance on external APIs, we cre-
179 ated moondream_clip_minimal_pipeline.py,
180 which uses Moondream2 to summarize first/last
181 frames and build a motion synopsis from strided
182 intermediate frames, writing per-clip files such as
183 clip_X_transition_prompt.txt. Although we
184

185 collected prompts, we deferred training with them because each additional LoRA run would consume ~ 200
 186 GPU units on Colab. At inference time, prompts are ex-
 187 posed as CLI arguments with defaults emphasizing motion
 188 smoothness and character consistency, plus negative
 189 prompts to suppress blur/distortion.
 190

191 4.3 Loss Functions

192 Our final objective augments the standard FLF2V dif-
 193 fusion loss with an optical-flow-based auxiliary term.
 194 The baseline follows the Wan latent routine (VAE
 195 encode, concatenate first+video+last, channel projec-
 196 tion, noise injection, transformer forward) and com-
 197 putes MSE versus target noise. We implement an
 198 optimized loss, `compute_wan_flf2v_loss`, which
 199 stages VAE and CLIP encoders on GPU sequentially
 200 to reduce memory footprint and GPU–CPU transfers
 201 while adhering to the official Wan loop. To explicitly
 202 supervise motion, `OpticalFlowLoss` uses RAFT
 203 to compute flow discrepancies between generated and
 204 ground-truth frames, with optional downsampling and
 205 a configurable weighting, enabling and tuning this term
 206 is controlled via training flags.

207 4.4 Training Regime and Practical Con- 208 straints

209 All models are trained under single-GPU (A100) con-
 210 straints with bf16 when available. We limit training
 211 to 50 clips (4350 frames) for rapid iteration and run
 212 500 steps with learning rate 5×10^{-5} . The flow-loss
 213 integration increases compute cost but provides mo-
 214 tion supervision without requiring full-length training
 215 on every clip. We log losses and checkpoints, LoRA
 216 weights are stored in `wan_flf2v_lora` and loaded
 217 by `inference.py`.

218 5 Experiments

219 5.1 Experimental Setup

220 Our code implements the full baseline pipeline:
 221 data processing (`create_data.py`), prompt
 222 generation (`prompt_gen.py`), dataset caching
 223 (`train_setup.py`), LoRA training (`train.py`),
 224 evaluation (`evaluation.py`) and infer-
 225 ence (`inference.py`), all tested in Colab
 226 (`colab_script.ipynb`). The WAN model
 227 can be fine-tuned and used to generate videos from
 228 first/last-frame conditioning.

229 We evaluate the frozen WAN base and LoRA vari-
 230 ants at ranks 8, 16 and 32 using `num_frames`
 231 = 13, `num_inference_steps` = 30,

LoRA	Model	FVD \downarrow	SSIM \uparrow	PSNR \uparrow	LPIPS \downarrow
8	Base	1374.21	0.98235	37.41	0.00879
8	Fine-tuned	1299.20	0.98200	36.68	0.00889
16	Base	1316.51	0.98230	37.59	0.00879
16	Fine-tuned	1129.60	0.98207	36.63	0.00883
32	Base	1365.40	0.98220	37.46	0.00870
32	Fine-tuned	898.36	0.98310	38.75	0.00880

Table 1: Metrics on the withheld episode for 13-frame, 30-step evaluation.

guidance_scale = 5.0, and only_final = True to fit within a 12 GB budget while extending beyond the 9-frame training sequence. Inference lets you adjust the usual params (resolution, 4k+1 frame sampling, guidance, steps, flow shift, dtype, device). Metrics include FVD, SSIM, PSNR and LPIPS.

Training uses a single Colab A100 (bf16). A 500-step run converges in ~ 2.5 GPU-hours (0.55 it/s). Evaluation on 10 held-out clips adds ~ 125 minutes due to repeated I3D embedding computation. Results are logged to W&B.

243 5.2 Quantitative Results

Table 1 compares the base model and LoRA ranks 8/16/32 on the 13-frame evaluation. Fine-tuning reduces FVD by up to 476 points while keeping SSIM (≈ 0.98) and LPIPS (≈ 0.0088) stable this indicates strong reconstruction with modest temporal blur.

Relative to earlier midway runs (without optical flow loss and with a smaller training setup), the updated pipeline results in larger motion-aware improvements, especially at rank 32. Because frame-level metrics remain tightly grouped, the gains mainly reflect improved temporal realism and distribution-level video quality, consistent with flow loss acting as a motion-focused regularizer.

257 5.3 Optical Flow Loss Ablation

We compare LoRA training with and without optical flow supervision (controlled via `use_flow_loss` and `flow_loss_weight`), keeping all other hyperparameters fixed. As expected, flow loss mainly improves temporal smoothness and reduces jitter, reflected in lower FVD and greater perceptual stability.

Figure 1 shows diffusion MSE, total loss and flow loss for LoRA ranks 8, 16 and 32. Motion supervision introduces a consistent auxiliary signal without destabilizing diffusion training.

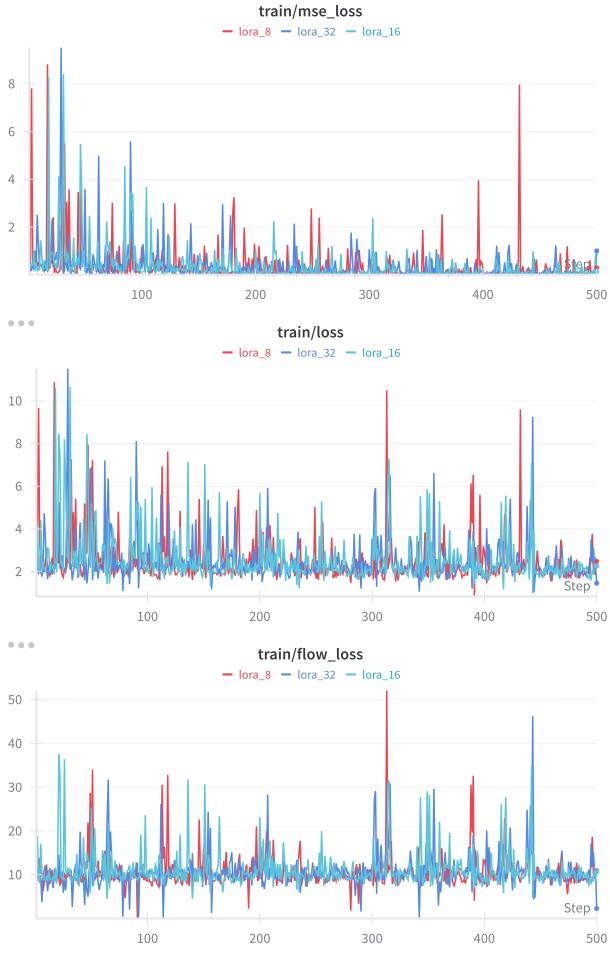


Figure 1: **Optical flow ablation training curves.**



(a) Example 1



(b) Example 2



(c) Example 3

Figure 2: **Qualitative comparison across base and LoRA ranks.** We show three held-out clips under the 13-frame evaluation setting. For each example, the top row is ground truth, followed by the base model and LoRA adapters at ranks 8, 16, and 32. Frames are sampled at consistent timestamps between the conditioning first and last frames. Higher-rank adapters, especially LoRA-32, more reliably preserve character structure and reduce motion artifacts, qualitatively aligning with the larger FVD improvements reported in Table 1.

283 6 Code Overview

284 Our main contributions span dataset loading and dy-
285 namic frame sampling, training configuration and val-
286 idation, custom first/last-frame conditioning, an op-
287 tional optical-flow auxiliary loss and an optimized
288 FLF2V loss designed for large video models under
289 strict VRAM limits.

290 6.1 Dataset Loading and Frame Sampling

291 We provide two Dataset classes.
292 WanFLF2VDatasetFromHF scans the HF-style
293 repository, builds clip-level samples and supports
294 *dynamic 4k+1* frame sampling, ensuring first/last-
295 frame inclusion and evenly spaced intermediates.
296 This enables training across variable temporal spans
297 with a controlled budget (Lines 38–270, `train.py`).
298 WanFLF2VDatasetFromLocal is a fast-path
299 loader for locally downloaded datasets, using a pickle
300 metadata cache and the same dynamic/fixed sampling
301 logic without HF I/O (Lines 271–397).

302 6.2 Training Configuration and Safety 303 Checks

304 A `TrainConfig` dataclass consolidates model/-
305 data/optimization settings, including LoRA rank/al-
306 pha/dropout, mixed precision and the optical-flow tog-
307 gle. Its `__post_init__` validates 4k+1 constraints, de-
308 fines the allowed frame-count set for the RNG sampler
309 and prevents silent misconfiguration (Lines 398–460).

310 6.3 Text Conditioning Utility

311 `encode_text` produces stable prompt em-
312 beddings while guarding against pathological
313 `model_max_length` values. It clamps lengths,
314 normalizes input formatting and ensures compatibility
315 with CLIP-style encoders in the Wan pipeline (Lines
316 461–503).

317 6.4 First/Last Frame Image Conditioning

318 `encode_first_last_images` constructs the two-
319 frame conditioning tensor in the $(2B, S, D)$ for-
320 mat expected by the transformer. We eliminate PIL
321 round-trips, apply CLIP preprocessing manually and
322 temporarily move the image encoder to the target de-
323 vice to reduce VRAM use when training large models
324 (Lines 629–699).

325 6.5 Optional Optical Flow Auxiliary Loss

326 `OpticalFlowLoss` adds motion supervision via
327 RAFT-based flow comparisons between generated and

ground-truth frames. RAFT inference is forced to
328 float32 for stability, supports spatial downsampling and
329 includes an OOM-safe fallback that skips the flow term
330 when needed. This keeps long runs stable while en-
331 hancing temporal coherence (Lines 504–628).
332

333 6.6 Optimized FLF2V Loss Computation

compute_wan_flf2v_loss handles the full FLF2V
334 objective with explicit staging: the VAE is moved to
335 GPU only during encoding and then returned to CPU
336 to free VRAM—critical for 14B-scale training. When
337 optical flow is enabled, predicted latents are decoded
338 back to RGB to compute the auxiliary loss (Lines 700–
339 887).
340

341 6.7 Main Training Orchestration

The main function aggregates CLI overrides, de-
342 vice selection, W&B logging, model/component load-
343 ing, LoRA injection, gradient checkpointing, and
344 DataLoader setup. Its data-loading priority order
345 (`train_data_dir` → checkpoint restore → HF)
346 makes local training the default while preserving repro-
347 ducibility. The loop manages mixed precision, grad ac-
348 cumulation/clipping, and logs key diagnostics includ-
349 ing dynamic frame-count usage (Lines 911–1496).
350

351 7 Research Log

We had initially underestimated the GPU usage in the
352 WAN pipeline. After dataset preparation and metadata
353 generation, our first Colab sessions were used to val-
354 idate LoRA injection on 9-frame subsamples. Early
355 runs exposed bugs in the dynamic frame sampler. Fix-
356 ing these took a day but stabilized the dataloader for all
357 later experiments.
358

The main early challenge was to make the baseline
359 reliable under tight compute and memory limits. The
360 WAN FLF2V model is big enough that small imple-
361 mentation choices (frame indexing, tensor shapes, de-
362 vice placement, mixed precision) would trigger either
363 silent failures or OOMs. Loading full 81-frame se-
364 quences was too slow and heavy in memory, so effi-
365 cient subsampling became a necessity rather than an
366 optimization. We therefore put our focus on strict 4k+1
367 sampling, light local caching (`train_setup.py`),
368 and a reduced 10-clip prototype regime for rapid de-
369 bugging.
370

Once the baseline converged, we tried the scene-
371 aware extension, adding shot-boundary metadata, but
372 these additional tensors exceeded the 80 GB A100 limit
373 for 81-frame clips, which motivated our compromise
374 of a 13-frame evaluation: short enough to be stable
375

376 for inference, yet long enough to test behavior beyond
377 the 9-frame training horizon. This also made apparent
378 that long clips would routinely be unrealistic to evaluate
379 under the single-GPU budget and shifted our plan
380 toward focusing on improvements that could be validated
381 on short clips-loss design, sampling robustness,
382 conditioning-prior to scaling

383 Prompt metadata introduced another bottleneck. Our
384 proposal assumed rich GPT-based prompts per clip, but
385 repeated API calls were costly and slow to iterate on.
386 We tried using Moondream2 as a cheaper alternative;
387 encoding strided frames for all 1,238 clips would have
388 consumed \sim 200 additional Colab GPU units, so we
389 stopped after 64 clips. Thus, we implemented the full
390 prompt-conditioning path but held off on large-scale
391 prompt regeneration to fit the LoRA+ loss ablations
392 within our compute budget.

393 We then migrated to AWS to avoid Colab quota re-
394 sets. We faced EC2 capacity and vCPU limits when
395 requesting GPU instances. Even after getting a suit-
396 able A100-class machine, storage costs, reconfigura-
397 tion time, and environment rebuilds imposed further
398 overhead. AWS increased flexibility but did not elim-
399 inate resource scarcity, credits still capped total train-
400 ing and evaluation time. These constraints dictated the
401 methodological pivots necessary in the final system:

- **Dynamic clip handling:** Dataset logic updated to handle variable-length clips, while maintaining valid 4k+1 subsamples, eliminating brittleness and better reflecting real comic-style transitions.
- **Optical flow loss:** Added flow-based supervision that penalizes motion inconsistencies. This required careful VRAM control, downsampling options, and OOM-safe fallbacks.
- **Rank sweep reframing:** Each LoRA rank was extremely computationally expensive. We first focused on stabilizing rank 16 and then expanded to rank 32 when GPU time allowed; rank 8 is still planned for future runs.

415 The final week focused on evaluation: collecting
416 the metrics in Table
417 reftab:metrics13, generating qualitative samples
418 and scripting the rank sweep for automatic con-
419 tinuation when more compute becomes available.
420 The evaluation pipeline also highlighted that FVD
421 dominates runtime since embeddings must be re-
422 computed per frame count and per model, rein-
423 forcing our choice of the 13-frame, 30-step evalua-
424 tion as the only feasible standardized setting.

425 Several limitations remain: Optical flow loss is
426 promising yet lacks a full ablation across ranks
427 and longer frame count; the prompt-generation

428 pipeline, while functional, is incomplete at scale,
429 scene-aware extension remains blocked by mem-
430 ory limits and most likely requires multi-GPU
431 training or even more aggressive latent compres-
432 sion. Lastly, the dataset is still comparably small
433 with respect to model size. adding more episodes
434 with diverse transitions would improve general-
435 ization. The project scope changed from “fine-
436 tune a big model” to making large-video diffusion
437 emphractable under student-level compute. But
438 the most important outcome is a complete, mod-
439 ular pipeline, dynamic datasets, prompt hooks,
440 LoRA configurability, and motion-aware losses,
441 now prepared for any future scaling once more re-
442 sources become available.

8 Conclusion

443 Our completed LoRA 8/16/32 results show that
444 fine-tuning meaningfully improves temporal real-
445 ism even with small adapters. Across all ranks,
446 perceptual scores (SSIM, LPIPS) remain nearly
447 unchanged from the base model which means that
448 our changes keep frame level fidelity. The main
449 gains appear in FVD, where rank 16 reduces error
450 by \sim 187 points and rank 32 achieves the strongest
451 improvement, lowering FVD by nearly 470 points
452 while slightly boosting SSIM and PSNR. These
453 trends confirm that motion-aware supervision and
454 LoRA capacity together have a substantial impact
455 on distribution-level video quality.

456 Evaluating under the 13-frame, 30-step setting
457 provides a more realistic test of interpolation be-
458 yond the 9-frame training window and the frame-
459 work now supports reproducible comparisons
460 across ranks and loss configurations. The partially
461 implemented scene-aware and style-conditioned
462 extensions, as well as the prompt-generation path,
463 remain promising next steps but require additional
464 compute to run at scale. Future work should (1)
465 complete controlled ablations of optical-flow su-
466 pervision, (2) integrate scene labels and prompt
467 conditioning directly into the loss, (3) revisit long-
468 clip evaluation once multi-GPU resources are
469 available and (4) expand the dataset to strengthen
470 generalization across diverse transitions.

9 Thought-Experiment on Com- 472 pute

473 We used around 500 Colab GPU units which is
474 equivalent to about 120 A100-80GB GPU-hours

(roughly \$230 at current Colab Pro+ pricing). Only one teammate consistently received Colab Priority access so we were unable to parallelize runs. Also only one teammate received the AWS credits but even with this our AWS account could not provision A100s on short notice. Due to this, we skipped the prompt-conditioned fine-tuning run because it would have required an additional 200 GPU units for data preparation plus training. If we had an extra \$450 in cloud credits, we would reserve two on-demand A100-80GB instances for a week: one to finish the LoRA rank 8/32 sweeps with full-scene clips and another to jointly train the scene-aware loss with real-time Moondream2 prompts so that we can quantify how textual guidance changes motion quality. The extra budget would also go to funding the remaining user-study renders so that evaluation can go beyond the single 13-frame metric slice presented here.

10 Bibliography

References

- [1] Pan, Z., Zhu, Z., Wu, T., Xiao, J., Huang, J., & Wang, J. (2023). *Sakuga-42M: A large-scale anime video generation dataset and benchmark*. Retrieved from https://zhenglinpan.github.io/sakuga_dataset_webpage/
- [2] Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., & Fleet, D. (2022). *Video Diffusion Models*. arXiv:2204.03458. <https://arxiv.org/abs/2204.03458>
- [3] Wu, H., Fan, H., Xiong, B., Mangalam, K., Li, Y., Mahajan, D., Yan, Z., & Feichtenhofer, C. (2022). *VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training*. arXiv:2203.12602. <https://arxiv.org/abs/2203.12602>
- [4] Jiang, H., Sun, D., Jampani, V., Yang, M.-H., Learned-Miller, E., & Kautz, J. (2018). *Super SloMo: High quality estimation of multiple intermediate frames for video interpolation*. arXiv:1712.03087. <https://arxiv.org/abs/1712.03087>
- [5] Guo, Y.-C., Zhang, Y., Chen, Z., & Song, L. (2023). *AnimateDiff: Animate your personalized text-to-image diffusion models without*

specific tuning. arXiv:2307.04760. <https://arxiv.org/abs/2307.04760>

- [6] Stability AI. (2023). *Stable Video Diffusion*. Retrieved from <https://www.emergentmind.com/topics/stable-video-diffusion-svd>
- [7] VideoCrafter Team. (2023). *VideoCrafter: Open Diffusion Models for High-Quality Video Generation*. arXiv:2310.19512. <https://arxiv.org/abs/2310.19512>
- [8] Singer, A., Shetty, R., Mahajan, D., Chechik, G., & Vedaldi, A. (2022). *Make-A-Video: Text-to-Video Generation Without Text-Video Data*. arXiv:2209.14792. <https://arxiv.org/abs/2209.14792>
- [9] Ho, J., Saharia, C., Chan, W., Fleet, D., Norouzi, M., & Salimans, T. (2022). *Imagen Video: High Definition Video Generation with Diffusion Models*. arXiv:2210.02303. <https://arxiv.org/abs/2210.02303>
- [10] Voleti, V., Jolicoeur-Martineau, A., & Pal, C. (2022). *Masked Conditional Video Diffusion Models*. arXiv:2205.09853. <https://arxiv.org/abs/2205.09853>
- [11] Li, Y., Xu, J., Yang, L., & Fang, L. (2024). *ToonCrafter: Adaptation of Motion Priors to Cartoon Animation*. arXiv:2405.17933. <https://arxiv.org/abs/2405.17933>
- [12] Kong, C., Zhang, Z., Chai, L., & Wang, Y. (2024). *W.A.L.T: Window Attention Latent Transformer for Video Generation*. ECCV 2024. https://www.ecva.net/papers/eccv_2024/papers_ECCV/papers/10270.pdf
- [13] Wu, J., Xia, W., Chen, X., Duan, Y., & Wu, B. (2022). *Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation*. arXiv:2212.11565. <https://arxiv.org/abs/2212.11565>