

LoRA Trade-offs in Video Diffusion Models: Latency and Temporal Consistency Across Hardware

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What is the trade-off between inference latency and temporal consistency (measured by TSSIM or LPRO) when applying LoRA to video diffusion models like Make-A-Video or AnimateDiff across different hardware configurations? We present a practical pipeline for fine-tuning open-source video diffusion transformers to synthesize cinematic scenes for television and film production from small datasets. The proposed two-stage process decouples visual style learning from motion generation. 9 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.6/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Fine-Tuning Open Video Generators for Cinematic Scene Synthesis: A Small-Data Pipeline with LoRA and Wan2.1 I2V. Research question: What is the trade-off between inference latency and temporal consistency (measured by TSSIM or LPRO) when applying LoRA to video diffusion models like Make-A-Video or AnimateDiff across different hardware configurations?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.6/10.

3 Results

13 papers retrieved. 9 claims extracted; 4 independently verified. Quality review score: 6.6/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The hyperparameters used for fine-tuning include LoRA rank / α of 8 / 16, a learning rate of 3×10^{-5} , AdamW optimizer w	×	0.04
The configuration time for a single A100-80GB is 187 seconds with a speedup of $1.0 \times$.	×	0.02
The methodology involves applying costume, lighting, and color grading, which are then temporally expanded into coherent	✓	0.23
Lightweight parallelization and sequence partitioning strategies are applied to accelerate inference without quality deg	✓	0.20
Quantitative and qualitative evaluations using FVD, CLIP-SIM, and LPIPS metrics demonstrate measurable improvements in c	✓	0.32
The complete training and inference pipeline is released to support reproducibility and adaptation across cinematic doma	✓	0.27
Diffusion transformers have evolved into powerful spatio-temporal generators capable of producing coherent multi-second	×	0.06
Open-source efforts such as VideoCrafter, ModelScope, and Wan2.x have narrowed the gap with commercial systems like Runw	×	0.04
Cinematic generation—the ability to reproduce film-like motion, controlled lighting, lens depth, and storytelling rhythm	×	0.06

References

- <http://arxiv.org/abs/2510.27364v1>
- <http://arxiv.org/abs/2307.04725v2>
- <http://arxiv.org/abs/2503.11495v1>