

# Wan2.1 I2V-14B with LoRA Adaptation Performance on Out-of-Domain Cinematic Scenes

Assignee Research

May 31, 2026

## Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does Wan2.1 I2V-14B with LoRA adaptation perform on out-of-domain cinematic scenes (e.g., sci-fi) compared to its performance on historical scenes, as evaluated by CLIP-based metrics like FID. 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. 17 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. 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: How does Wan2.1 I2V-14B with LoRA adaptation perform on out-of-domain cinematic scenes (e.g., sci-fi) compared to its performance on historical scenes, as evaluated by CLIP-based metrics like FID (Fréchet Inception Distance) and KID (Kernel Inception Distance)?.

## 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 4.5/10.

### **3 Results**

13 papers retrieved. 17 claims extracted; 2 independently verified. Quality review score: 4.5/10.

### **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 fine-tuning pipeline uses a LoRA rank of 8 and an alpha value of 16.	×	0.08
The learning rate used for training is $3 \times 10^{-5}$ with a cosine schedule and 5% warm-up.	×	0.03
The optimizer used is AdamW with $\beta_1=0.9$ , $\beta_2=0.999$ , and weight decay=0.01.	×	0.01
The effective batch size is 2, calculated as 1 video multiplied by a gradient accumulation of 4.	×	0.02
The model was trained for 4000 steps using bf16 precision.	×	0.04
Activation checkpointing is enabled to reduce the VRAM footprint.	×	0.00
The framework used is PyTorch combined with DeepSpeed, utilizing Fully Sharded Data Parallelism (FSDP).	×	0.02
Training employs early stopping based on the LPIPS plateau.	×	0.03
On a single A100-80GB GPU, the configuration time is 187 seconds.	×	0.02
The single A100-80GB configuration serves as the $1.0 \times$ speedup baseline.	×	0.01
The pipeline expands inputs into coherent 720p video sequences.	×	0.09
Evaluations were conducted using FVD, CLIP-SIM, and LPIPS metrics.	✓	0.17
A small expert user study was conducted to support the quantitative evaluations.	×	0.15
The proposed method demonstrates measurable improvements in cinematic fidelity and temporal stability over the base mode	✓	0.20
Diffusion transformers have evolved to produce coherent multi-second videos from textual descriptions.	×	0.05
VideoCrafter, ModelScope, and Wan2.x are open-source video generation efforts.	×	0.11
Runway Gen-2, Pika, and Sora are commercial video generation systems.	×	0.03

## References

- <http://arxiv.org/abs/2510.27364v1>
- <http://arxiv.org/abs/1707.09465v5>
- <http://arxiv.org/abs/2512.04830v1>