



Project page

Motivation

Despite rapid progress in video world models, how data shapes motion quality remains **poorly understood**.

Key Goals

Focus on Motion

Separate motion from static appearance

Scale Efficiently

Modern, large-scale models & datasets

Guide Curation

Identify clips that improve motion quality

Our Method: MOTIVE

MOTION attribution for Video gENERation

Problem Formulation

Given a query video and finetuning dataset, assign each training clip a **motion-aware influence score** to quantify its contribution to target generation.

Method Components

1. Efficient Motion Gradient Computation

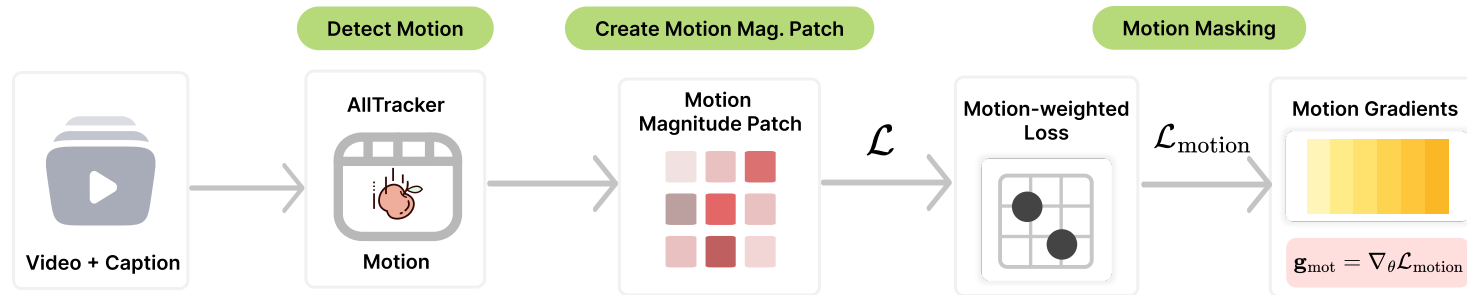
- Single-Sample Estimator
- Structured Projections (Fastfood)

2. Motion Attribution

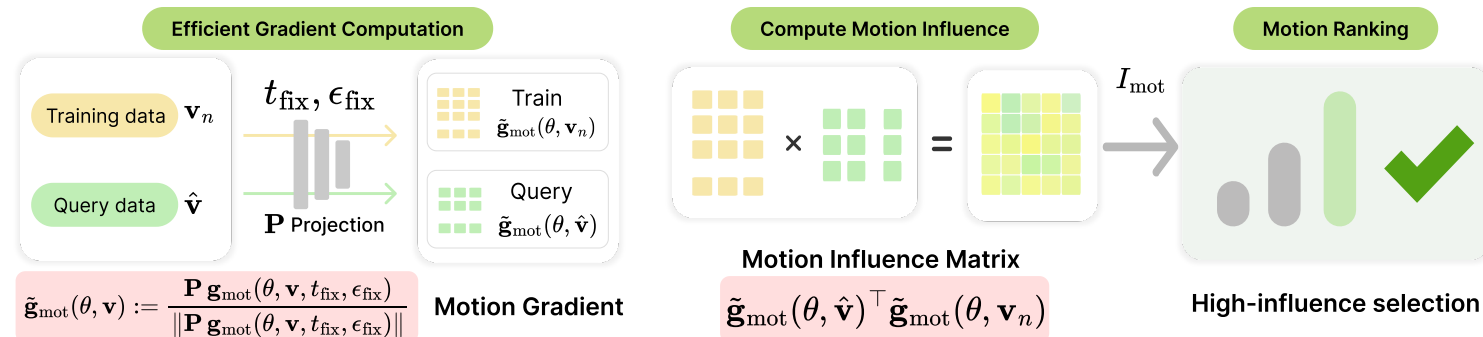
- Detect motion between frames w. AllTracker
- Create motion magnitude patches highlighting dynamic areas
- Apply motion-weighted loss to focus on moving regions and compute motion-specific gradients

Which training data influence the motion in generated videos?

Motion Attribution

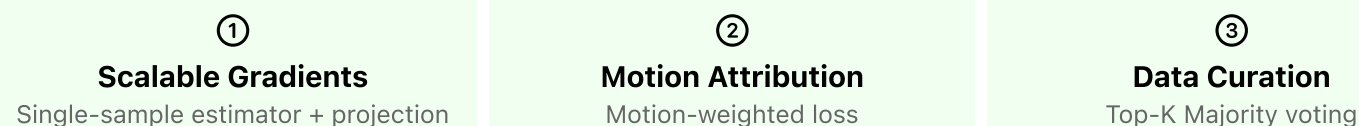


Efficient Motion Gradient Computation



MOTIVE: A scalable, gradient-based, motion-centric data attribution framework for video generation models

Three Key Components



Motion Attribution Samples

Qualitative Results

Generated Videos

Quantitative Results

VBench Evaluation *Wan2.1-T2V-1.3B*

Method	Subj.	Bg.	Mot.	Dyn.	Aesth.	Img.
Base	95.3	96.4	96.3	39.6	45.3	65.7
Full FT	95.9	96.6	96.3	42.0	45.0	63.9
Random selection	95.3	96.6	96.3	41.3	45.7	65.1
Motion magnitude	95.6	96.2	95.7	40.1	45.1	63.2
V-JEPA	95.7	96.0	95.6	41.6	44.9	62.7
Ours w/o MM	95.4	96.1	96.3	43.8	45.7	63.2
Ours (MOTIVE)	96.3	96.1	96.3	47.6	46.0	64.6

✓ Maintains smoothness, improves dynamics with only 10% data. All selection methods use 10% of training data.

Why Motion Masking? *Dynamic Degree*

Without: 43.8% **With:** 47.6% (+3.8%)

Human Evaluation

vs. Base	Ours wins	Tie	Baseline wins
vs. Base	74.1%	12.3%	13.6%
vs. Random	58.9%	12.1%	29.0%
vs. Full FT	53.1%	14.8%	32.1%
vs. Ours w/o MM	46.9%	20.0%	33.1%

Ablation Findings

Single Timestep: t=751 achieves **66%** agreement.

Projection: D'=512 reaches **74.7%** Spearman ρ.

Conclusion

First motion-centric attribution framework for video generation

Scalable via **single-sample estimator & structured projection**

74.1% human preference vs. baseline with **10% data**; Motion masking: **+3.8%** Dynamic Degree