

Motivation

Despite rapid progress in video generation, 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 Solution: MOTIVE

MOTION Training Influence 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

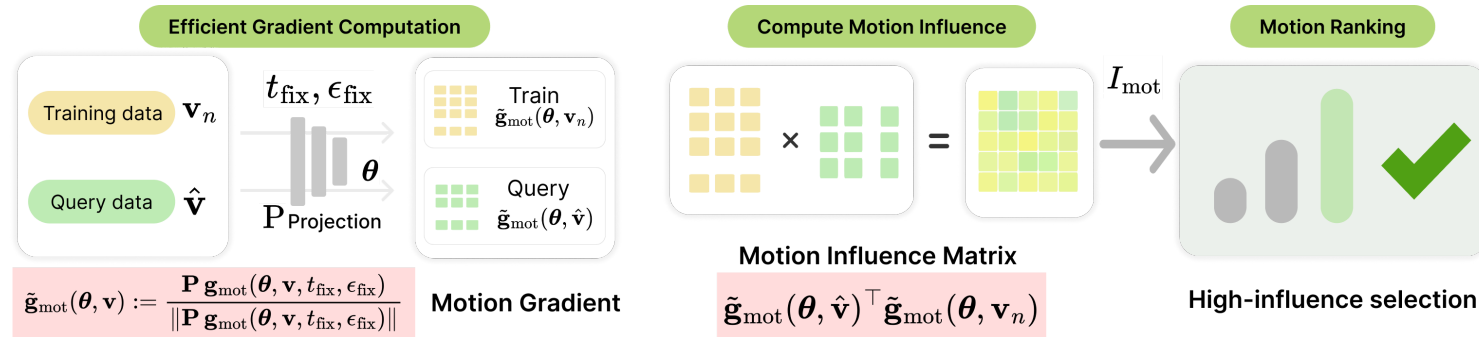
- Single-Sample Estimator
- Structured Projections (Fastfood)

2. Motion Attribution

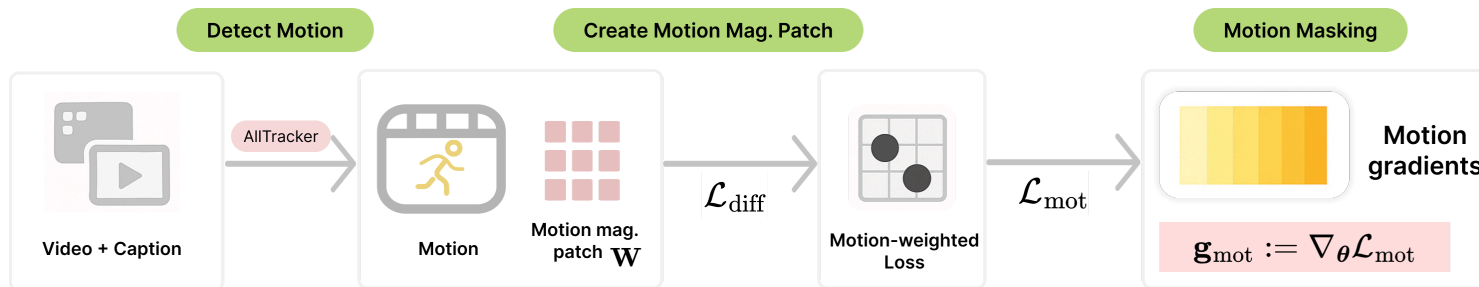
- Detect motion between frames w. AllTracker
- Create motion magnitude patches highlighting dynamic areas
- Apply motion-weighted loss to focus on moving regions and compute motion-specific gradients

Which training clips drive the motion in a video generation sample?

Efficient Motion Gradient Computation



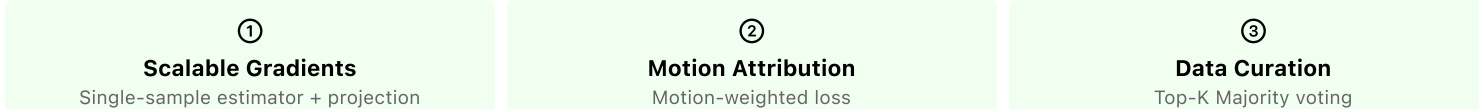
Motion Attribution



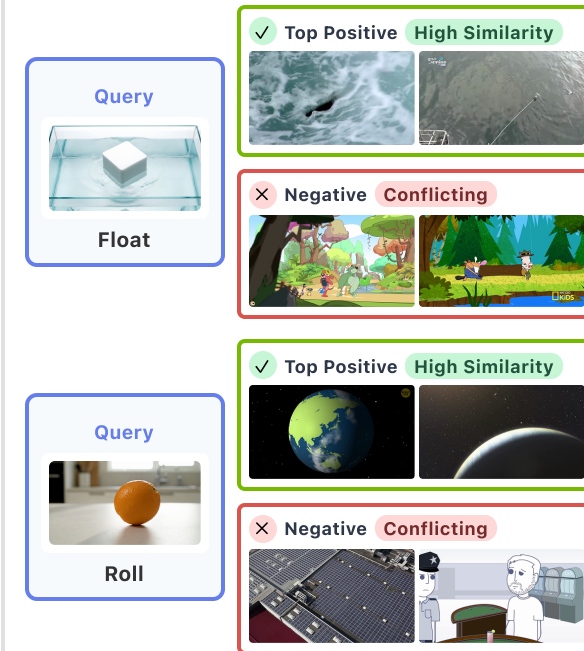
Motion-gradient computation has three steps: (1) detect motion with AllTracker; (2) compute motion-magnitude patches; (3) apply loss-space motion masks to focus gradients on dynamic regions.

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

Three Key Components

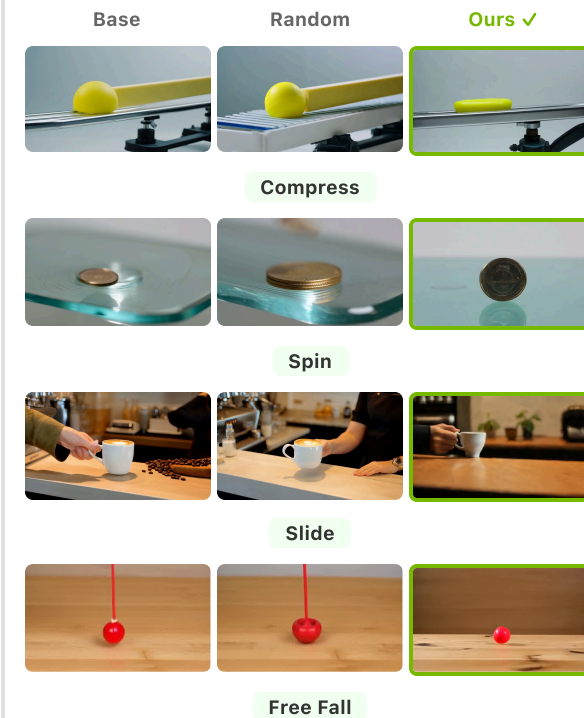


Motion Attribution Samples



Qualitative Results

Generated Videos



Quantitative Results

VBench Evaluation

Method	Motion Smooth.	Dynamic Deg.
Base	96.3	82.3
Full FT	96.3	84.7
Random 10%	96.3	81.6
Ours w/o mask	96.3	85.3
MOTIVE	96.3	89.4

✓ Maintains smoothness, improves dynamics with only 10% data

Why Motion Masking?

Without: 85.3%	With: 89.4% (+4.1%)
--------------------------	-------------------------------

Human Evaluation

vs. Base:	76.7% win
vs. Random:	66.7% win
vs. Full FT:	57.5% win

Ablation Findings

Single Timestep: t=500 achieves **68%** agreement.
Projection: D'=512 reaches **74.7%** Spearman ρ .

Conclusion

First motion-centric attribution framework for video generation

Scalable via **projection & majority voting**

76.7% human preference vs. baseline with **10% data**; Motion masking: **+4.1%** Dynamic Degree