

---

# MOTION ATTRIBUTION FOR VIDEO GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Despite the rapid progress of video generation models, the role of data in influencing motion is poorly understood. We present **MOTIVE** (**MOTION attribution for Video gEneration**), a motion-centric, gradient-based data attribution framework that scales to modern, large, high-quality video datasets and models. We use this to study which fine-tuning clips improve or degrade temporal dynamics. **MOTIVE** isolates temporal dynamics from static appearance via motion-weighted loss masks, yielding efficient and scalable motion-specific influence computation. On text-to-video models, **MOTIVE** identifies clips that strongly affect motion and guides data curation that improves temporal consistency and physical plausibility. With **MOTIVE**-selected high-influence data, we improve both motion smoothness and dynamic degree on VBench, achieving a 74.1% human preference win rate compared with the pretrained base model. To our knowledge, this is the first framework to attribute motion rather than visual appearance in video generative models and to use it to curate fine-tuning data.

## 1 INTRODUCTION

Motion is the defining element of videos. Unlike image generation, which produces a single frame, video generative models must capture how objects move, interact, and obey physical constraints (Wiedemer et al., 2025; Kang et al., 2024). Yet a fundamental question remains: *which training data influence the motion in generated videos?*

**Why it matters.** Diffusion models are data-driven, and their progress tracks the scaling of data and compute (Saharia et al., 2022; Ho et al., 2022; Peebles & Xie, 2023). Prior work shows that training data shapes key generative properties, including visual quality (Rombach et al., 2022) and compositionality (Favero et al., 2025). Motion is no exception, if generated motion reflects the data distribution that shaped the model, then attributing motion to influential training clips provides a direct lens on why a model moves the way it does and enables targeted data selection for desired dynamics.

**Why existing methods are limited for motion.** Prior diffusion attribution focuses on images and static content. Extending these methods to videos naively collapses motion into appearance, missing temporal structure. Three challenges drive this gap: (i) localizing motion to focus attribution on dynamic regions, (ii) scaling to sequences where gradients integrate across time, and (iii) capturing temporal relations that single-frame attribution cannot measure.

**Our method.** We introduce **MOTIVE**, a motion attribution framework that isolates motion-specific influence via motion-aware gradient masking. Our key contributions are:

- A scalable gradient-based attribution for video generation that is computationally efficient at modern scales (§2).
- Addressing video-specific bias by correcting frame-length effects in gradient magnitudes.
- Motion-weighted attribution that emphasizes temporal dynamics over static appearance.
- Showing motion improvements on VBench and human evaluation, matching full-dataset fine-tuning with only 10% of data (Tab. 1).

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

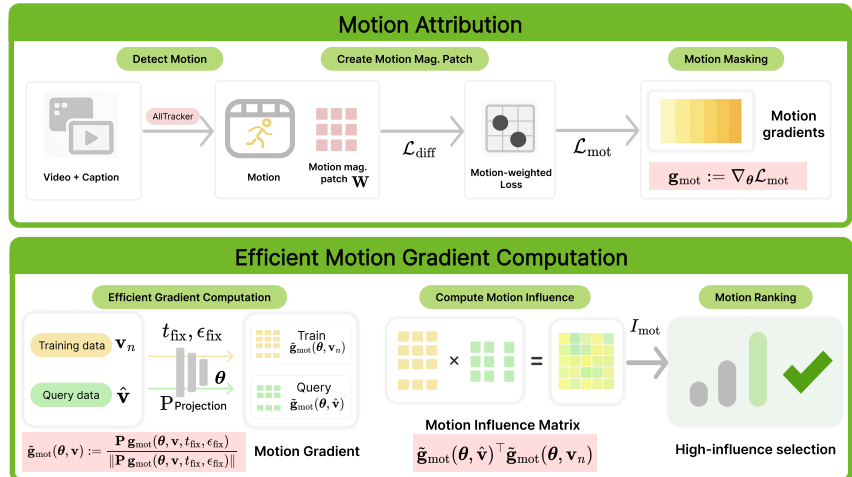


Figure 1: **MOTIVE**. *Top*. Motion-gradient computation (§2.3) has three steps: (1) detect motion with AllTracker; (2) compute motion-magnitude patches; (3) apply loss-space motion masks to focus gradients on dynamic regions. *Bottom*. Our method (§2.2) is made scalable via a single-sample variant with common randomness and a projection, computed for each pair of training and query data, aggregated (§D.2) for a final ranking, and eventually used to select fine-tuning subsets.

## 2 METHOD

We formalize the problem in §2.1 and develop a practical framework for motion attribution in video diffusion models with four components: scalable gradient computation (§2.2), frame-length bias fix (§D.1), motion-aware weighting (§2.3), and target data selection (§D.2). We provide efficiency analysis (§D.3) demonstrating scalability to billion-parameter models and large-scale video datasets.

### 2.1 PROBLEM FORMULATION

We study data attribution for motion in the fine-tuning setting. Let  $\mathcal{D}_{fit} = \{(\mathbf{v}_n, \mathbf{c}_n)\}_{n=1}^N$  be the fine-tuning corpus. Given a query video  $(\hat{\mathbf{v}}, \hat{\mathbf{c}})$ , we assign to each training clip  $(\mathbf{v}_n, \mathbf{c}_n)$  a motion-aware influence score  $I(\mathbf{v}_n, \hat{\mathbf{v}}; \theta)$  that explains how it contributes to the dynamics observed in  $\hat{\mathbf{v}}$ . The score should satisfy: (i) *predictivity*, rankings correlate with observed changes from fine-tuning on the most influential subsets; (ii) *efficiency*, scales to modern video generators, such as forgoing explicit Hessian inversion, expensive per-data integration, or prohibitive storage. To do this, we augment the influence target defined in Eq. 17 to be (a) lower variance for stable rankings with feasible levels of compute, (b) more scalable to store, and (c) motion-centric.

**Fine-tuning Subset Selection.** For a budget  $K \ll N$ , we get a motion-influential subset by ranking scores and taking the top- $K$  examples. For multiple query motions, we combine selections as described in §D.2. The resulting subsets serve as candidates for motion-centric fine-tuning.

### 2.2 SCALABLE GRADIENT-BASED ATTRIBUTION

We make attribution scalable for modern, large video datasets and models via inverse-Hessian approximations, lower-variance gradient-similarity estimators, low-cost single-sample estimators, and a Fastfood projection for tractable storage.

**Approximating the inverse-Hessian.** Computing exact inverse-Hessian-vector products is infeasible for modern neural networks. We estimate influence via gradient similarity, using an identity preconditioner for the inverse Hessian (Koh & Liang, 2017; Pruthi et al., 2020; Park et al., 2023).

**Common randomness for stable rankings.** To reduce variance without changing the target, we evaluate train and test gradients under the same  $(t, \epsilon)$  pairs and average over a small set  $\mathcal{T}$  (Xie et al., 2024; Lin et al., 2024). This paired averaging stabilizes rankings compared to independent draws:

$$I_{\text{diff}}^1(\mathbf{x}_n, \mathbf{x}_{\text{test}}) = \frac{1}{|\mathcal{T}|} \sum_{t, \epsilon \in \mathcal{T}} \frac{\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_{\text{test}}, t, \epsilon)}{\|\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_{\text{test}}, t, \epsilon)\|} \top \frac{\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_n, t, \epsilon)}{\|\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_n, t, \epsilon)\|}. \quad (1)$$

**Single-sample variant for reduced compute.** We fix a single  $t_{\text{fix}}$  and a shared draw  $\epsilon_{\text{fix}} \sim \mathcal{N}(0, \mathbf{I})$  for all train–test pairs at the final checkpoint. Sharing  $(t_{\text{fix}}, \epsilon_{\text{fix}})$  allows low enough variance for the low-cost single-sample estimator to maintain relative ordering (Xie et al., 2024; Lin et al., 2024). The estimator collapses to:

$$I_{\text{diff}}^2(\mathbf{x}_n, \mathbf{x}_{\text{test}}) = \frac{\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_{\text{test}}, t_{\text{fix}}, \epsilon_{\text{fix}})}{\|\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_{\text{test}}, t_{\text{fix}}, \epsilon_{\text{fix}})\|} \top \frac{\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_n, t_{\text{fix}}, \epsilon_{\text{fix}})}{\|\nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{x}_n, t_{\text{fix}}, \epsilon_{\text{fix}})\|}. \quad (2)$$

**Structured projection for reduced storage.** To operate at model scale, we apply a Johnson–Lindenstrauss projection via Fastfood (Le et al., 2014) and then normalize. Let

$$\mathbf{P} \in \mathbb{R}^{D' \times D} \text{ implemented as } \mathbf{P} := \frac{1}{\xi \sqrt{D'}} \mathbf{S} \mathbf{Q} \mathbf{G} \mathbf{\Pi} \mathbf{Q} \mathbf{B}, \quad (3)$$

where  $\mathbf{Q}$  is the Walsh–Hadamard matrix,  $\mathbf{B}$  is a diagonal Rademacher matrix,  $\mathbf{\Pi}$  is a random permutation,  $\mathbf{G}$  is a diagonal Gaussian scaling, and  $\mathbf{S}$  is a diagonal rescaling, and  $\xi$  normalizes the variance. The projected, normalized gradient is:

$$\tilde{\mathbf{g}}(\theta, \mathbf{x}) := \frac{\mathbf{P} \nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta, \mathbf{x}, t_{\text{fix}}, \epsilon_{\text{fix}})}{\|\mathbf{P} \nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta, \mathbf{x}, t_{\text{fix}}, \epsilon_{\text{fix}})\|}. \quad (4)$$

Then the influence score is the dot product of normalized projected gradients (i.e., cosine similarity) in  $\mathbb{R}^{D'}$ :

$$I_{\text{diff}}^3(\mathbf{x}_n, \mathbf{x}_{\text{test}}) = \tilde{\mathbf{g}}(\theta; \mathbf{x}_{\text{test}}) \top \tilde{\mathbf{g}}(\theta; \mathbf{x}_n). \quad (5)$$

This keeps compute  $\mathcal{O}(D' \log D')$  for projection and  $\mathcal{O}(D')$  per dot product, with storage  $\mathcal{O}(|\mathcal{D}| D')$ , while staying close to the ranking behavior of full-gradient cosine similarity (Park et al., 2023).

### 2.3 MOTION ATTRIBUTION

To move beyond whole-video influence, we introduce motion attribution, which isolates the contribution of training data to temporal dynamics. Unlike video-level attribution, which treats each clip as a single unit and conflates appearance with motion, motion attribution reweights per-location gradients using motion masks, assigning influence via dynamic behavior rather than static content.

**Motion Masking Attribution.** Motion is what distinguishes video diffusion from image diffusion. Our goal is to understand how training data shapes motion in video diffusion models. Prior work has emphasized architectural or algorithmic changes for motion modeling (Peebles & Xie, 2023; Blattmann et al., 2023; Guo et al., 2023), many of the largest generative gains have instead come from scaling and curating massive video corpora, which in turn enable impressive motion synthesis results in video diffusion models (Ho et al., 2022; Wan et al., 2025; Tan et al., 2024; Yang et al., 2024). Yet we lack tools that quantify how specific training clips shape particular motion patterns. We address this by attributing motion back to data via motion-weighted gradients, which yields actionable signals for targeted data selection, artifact diagnosis, and selective fine-tuning.

**Motion Detection and Latent Space Mapping.** Given a video  $\mathbf{v} \in \mathbb{R}^{F \times H \times W \times 3}$  with  $F$  frames of resolution  $H \times W$ , we first encode it into the VAE latent space as  $\mathbf{h} = E(\mathbf{v}) \in \mathbb{R}^{F \times H/s \times W/s \times C}$ , with downsampling factor  $s = 8$  and  $C = 16$  following the **wan2.1** backbone used in our experiments.

For motion computation, we use AllTracker (Harley et al., 2025) to extract motion information in pixel space:  $A = \mathcal{A}(\mathbf{v}) \in \mathbb{R}^{F \times H \times W \times 4}$ , where the first two channels contain optical flow maps  $A_{:, :, :, 0:2}$  indicating pixel displacement between frames, and the remaining channels  $A_{:, :, :, 2:4}$  encode visibility and confidence scores. We extract displacement vectors at each pixel location as:

$$\mathbf{D}_f(h, w) = (A_{f, h, w, 0}, A_{f, h, w, 1}) = (dw, dh). \quad (6)$$

We then bilinearly downsample motion quantities from  $(H, W)$  to the latent grid  $(\frac{H}{s}, \frac{W}{s})$  so that our masking lives where gradients are computed.

**Motion-Weighted Gradient Computation.** We define the motion magnitude at each location as:  $M_f(h, w) = \|\mathbf{D}_f(h, w)\|_2$ . To obtain comparable motion weights across frames and pixels, we min–max normalize over all frames and pixels, ensuring values lie in  $[0, 1]$ :

$$\mathbf{W}(f, h, w) = \frac{M_f(h, w) - m_{\min}}{m_{\max} - m_{\min} + \zeta}, \quad (7)$$

Table 1: **VBench Evaluation.** Performance comparison on VBench (Huang et al., 2024) across different baselines (all values in %, higher is better). All selection methods use 10% of training data; our method uses majority vote aggregation (§D.2) across motion queries. MM: motion masking. Subj.: Subject Consistency, Bg.: Background Consistency, Mot.: Motion Smoothness, Dyn.: Dynamic Degree, Aesth.: Aesthetic Quality, Img.: Imaging Quality.

Method	Wan2.1-T2V-1.3B						Wan2.2-TI2V-5B					
	Subj.	Bg.	Mot.	Dyn.	Aesth.	Img.	Subj.	Bg.	Mot.	Dyn.	Aesth.	Img.
Base	95.3	96.4	96.3	39.6	45.3	<b>65.7</b>	94.9	96.4	97.5	42.0	44.4	65.5
Full FT	95.9	<b>96.6</b>	96.3	42.0	45.0	63.9	<b>95.3</b>	96.5	97.5	45.3	44.8	<b>66.2</b>
Random	95.3	96.6	96.3	41.3	45.7	65.1	94.7	96.2	97.3	41.6	44.6	65.2
Motion mag.	95.6	96.2	95.7	40.1	45.1	63.2	95.0	96.3	97.4	44.9	45.0	65.1
V-JEPA	95.7	96.0	95.6	41.6	44.9	62.7	95.2	96.4	97.3	45.6	44.9	64.8
Ours w/o MM	95.4	96.1	96.3	43.8	45.7	63.2	94.9	96.5	97.4	43.8	45.2	64.8
<b>Ours (MOTIVE)</b>	<b>96.3</b>	96.1	96.3	<b>47.6</b>	<b>46.0</b>	64.6	95.1	<b>96.6</b>	<b>97.6</b>	<b>48.3</b>	<b>45.6</b>	65.5

where  $m_{\min} = \min_{f', h', w'} M_f(h', w')$ ,  $m_{\max} = \max_{f', h', w'} M_f(h', w')$ , and  $\zeta = 10^{-6}$  ensures a positive denominator. This normalization mitigates bias from absolute motion scale, yielding weights that emphasize relative motion saliency rather than raw magnitude, following prior practice in video saliency detection (Fang et al., 2013). Let  $(\tilde{h}, \tilde{w})$  index the latent grid. We obtain latent-aligned weights by bilinear downsampling:

$$\tilde{\mathbf{W}}(f, \tilde{h}, \tilde{w}) = \text{Bilinear}(\mathbf{W}(\cdot, \cdot, \cdot), F, \frac{H}{s}, \frac{W}{s}). \quad (8)$$

We compute per-location squared error at fixed  $(t_{\text{fix}}, \epsilon_{\text{fix}})$  at each frame  $f$  and ‘‘latent pixel’’  $(\tilde{h}, \tilde{w})$ :

$$\tilde{\mathcal{L}}_{\theta, \mathbf{v}, \mathbf{c}}(f, \tilde{h}, \tilde{w}) = \left( [\epsilon_{\theta}(\mathbf{z}(\mathbf{v}, t_{\text{fix}}, \epsilon_{\text{fix}}), t_{\text{fix}}, \mathbf{c})]_{f, \tilde{h}, \tilde{w}} - [\epsilon_{\text{target}}(t_{\text{fix}}, \epsilon_{\text{fix}})]_{f, \tilde{h}, \tilde{w}} \right)^2, \quad (9)$$

and define the motion-weighted loss by averaging over frames and latent spatial locations:

$$\mathcal{L}_{\text{mot}}(\theta; \mathbf{v}, \mathbf{c}) = \frac{1}{F_{\mathbf{v}}} \text{mean}_{f, \tilde{h}, \tilde{w}} \left[ \tilde{\mathbf{W}}_{\mathbf{v}, \mathbf{c}}(f, \tilde{h}, \tilde{w}) \cdot \tilde{\mathcal{L}}_{\theta, \mathbf{v}, \mathbf{c}}(f, \tilde{h}, \tilde{w}) \right]. \quad (10)$$

Notably, when  $\tilde{\mathbf{W}}$  is all ones, this recovers the standard objective with no motion emphasis. The  $1/F_{\mathbf{v}}$  factor corrects for frame-length bias and  $F_{\mathbf{v}}$  signifies how the number of frames may be video-dependent. The corresponding motion-weighted gradient for attribution is:

$$I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}}) = \tilde{\mathbf{g}}_{\text{mot}}(\theta, \hat{\mathbf{v}})^{\top} \tilde{\mathbf{g}}_{\text{mot}}(\theta, \mathbf{v}_n), \quad (11)$$

$$\text{where } \mathbf{g}_{\text{mot}} := \nabla_{\theta} \mathcal{L}_{\text{mot}}, \tilde{\mathbf{g}}_{\text{mot}}(\theta, \mathbf{v}) := \frac{\mathbf{P} \mathbf{g}_{\text{mot}}(\theta, \mathbf{v}, t_{\text{fix}}, \epsilon_{\text{fix}})}{\|\mathbf{P} \mathbf{g}_{\text{mot}}(\theta, \mathbf{v}, t_{\text{fix}}, \epsilon_{\text{fix}})\|}. \quad (12)$$

Loss-space masking leaves forward noising and generation unchanged and reweights only attribution, avoiding interactions between motion weighting and noise injection. In contrast, our motion-aware attribution emphasizes dynamic regions and de-emphasizes static backgrounds, so rankings identify training clips that most strongly shape the model’s motion rather than appearance.

The most influential fine-tuning subset selection are detailed in App. §D.2 and §D.3.

### 3 EXPERIMENTS

More experiment setups details are in App. §E and we provide evaluation results on VBench in Tab. 1.

### 4 CONCLUSION

We introduced **MOTIVE**, a motion-aware data attribution framework for video diffusion. By tracing generated motion to influential training clips, our method enables targeted data curation that improves motion quality with a fraction of the data. As video models scale, such data-level understanding will be key to diagnosing failure modes and building more controllable generative systems.

---

216 REFERENCES  
217  
218 Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chat-  
219 topadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform  
220 for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.

221 Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying  
222 framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.

223  
224 Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar  
225 Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhohus, et al. V-jepa 2: Self-supervised video models  
226 enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.

227 Juhan Bae, Wu Lin, Jonathan Lorraine, and Roger B Grosse. Training data attribution via approximate  
228 unrolling. *Advances in Neural Information Processing Systems*, 37:66647–66686, 2024.

229  
230 David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel.  
231 The mechanics of n-player differentiable games. In *International Conference on Machine Learning*,  
232 pp. 354–363. PMLR, 2018.

233  
234 Yuxiang Bao, Di Qiu, Guoliang Kang, Baochang Zhang, Bo Jin, Kaiye Wang, and Pengfei Yan.  
235 Latentwarp: Consistent diffusion latents for zero-shot video-to-video translation. *arXiv preprint*  
236 *arXiv:2311.00353*, 2023.

237 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
238 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
239 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.

240  
241 Jonathan Brokman, Omer Hofman, Roman Vainshtein, Amit Giloni, Toshiya Shimizu, Inderjeet  
242 Singh, Oren Rachmil, Alon Zolfi, Asaf Shabtai, Yuki Unno, et al. Montrage: Monitoring training  
243 for attribution of generative diffusion models. In *European Conference on Computer Vision*, pp.  
244 1–17. Springer, 2024.

245 Weifeng Chen, Yatai Ji, Jie Wu, Hefeng Wu, Pan Xie, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang  
246 Lin. Control-a-video: Controllable text-to-video diffusion models with motion prior and reward  
247 feedback learning. *arXiv preprint arXiv:2305.13840*, 2023.

248  
249 Zhaoxi Chen, Tianqi Liu, Long Zhuo, Jiawei Ren, Zeng Tao, He Zhu, Fangzhou Hong, Liang  
250 Pan, and Ziwei Liu. 4dnex: Feed-forward 4d generative modeling made easy. *arXiv preprint*  
251 *arXiv:2508.13154*, 2025.

252 RDWS Cook et al. Residuals and influence in regression. 1982.

253  
254 Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Stable audio  
255 open. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal*  
256 *Processing (ICASSP)*, pp. 1–5. IEEE, 2025.

257  
258 Yuming Fang, Weisi Lin, Zhenzhong Chen, Chia-Ming Tsai, and Chia-Wen Lin. A video saliency  
259 detection model in compressed domain. *IEEE transactions on circuits and systems for video*  
260 *technology*, 24(1):27–38, 2013.

261  
262 Alessandro Favero, Antonio Sclocchi, Francesco Cagnetta, Pascal Frossard, and Matthieu Wyart.  
263 How compositional generalization and creativity improve as diffusion models are trained. *arXiv*  
264 *preprint arXiv:2502.12089*, 2025.

265  
266 Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long  
267 tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891,  
268 2020.

269  
270 Kristian Georgiev, Joshua Vendrow, Hadi Salman, Sung Min Park, and Aleksander Madry. The  
271 journey, not the destination: How data guides diffusion models. *arXiv preprint arXiv:2312.06205*,  
272 2023.

---

270 Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features  
271 for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023.  
272

273 Google DeepMind. Veo-3: Advancing controllable and physically plausible video generation.  
274 Technical report, Google DeepMind, 2025. URL [https://storage.googleapis.com/  
275 deepmind-media/veo/Veo-3-Tech-Report.pdf](https://storage.googleapis.com/deepmind-media/veo/Veo-3-Tech-Report.pdf).

276 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala,  
277 Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models  
278 without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.  
279

280 Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson,  
281 Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, Poriya Panet, Sapir Weissbuch, Victor Kulikov,  
282 Yaki Bitterman, Zeev Melumian, and Ofir Bibi. Ltx-video: Realtime video latent diffusion. *arXiv  
283 preprint arXiv:2501.00103*, 2024.

284 Zayd Hammoudeh and Daniel Lowd. Training data influence analysis and estimation: A survey.  
285 *Machine Learning*, 113(5):2351–2403, 2024.  
286

287 Adam W Harley, Yang You, Xinglong Sun, Yang Zheng, Nikhil Raghuraman, Yunqi Gu, Sheldon  
288 Liang, Wen-Hsuan Chu, Achal Dave, Pavel Tokmakov, et al. Alltracker: Efficient dense point  
289 tracking at high resolution. *arXiv preprint arXiv:2506.07310*, 2025.

290 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in  
291 neural information processing systems*, 33:6840–6851, 2020.

292 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J  
293 Fleet. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646,  
294 2022.  
295

296 Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):  
297 185–203, 1981.

298 Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing  
299 Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video  
300 generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
301 Recognition*, pp. 21807–21818, 2024.  
302

303 Ruoxi Jia, Fan Wu, Xuehui Sun, Jiachen Xu, David Dao, Bhavya Kailkhura, Ce Zhang, Bo Li, and  
304 Dawn Song. Scalability vs. utility: Do we have to sacrifice one for the other in data importance  
305 quantification? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
306 Recognition*, 2021.

307 Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, Kaixin Wang, Gao Huang, and Jiashi  
308 Feng. How far is video generation from world model: A physical law perspective. *arXiv preprint  
309 arXiv:2411.02385*, 2024.

310 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In  
311 *International Conference on Machine Learning*, pp. 1885–1894. PMLR, 2017.  
312

313 Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. Datainf: Efficiently estimating data influence  
314 in lora-tuned llms and diffusion models. *arXiv preprint arXiv:2310.00902*, 2023.

315 Quoc Viet Le, Tamás Sarlós, and Alexander Johannes Smola. Fastfood: Approximate kernel  
316 expansions in loglinear time. *arXiv preprint arXiv:1408.3060*, 2014.  
317

318 Jinxu Lin, Linwei Tao, Minjing Dong, and Chang Xu. Diffusion attribution score: Evaluating training  
319 data influence in diffusion models. *arXiv preprint arXiv:2410.18639*, 2024.

320 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
321 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.  
322

323 Jonathan Lorraine and David Duvenaud. Stochastic hyperparameter optimization through hypernet-  
works. *arXiv preprint arXiv:1802.09419*, 2018.

---

324 Jonathan Lorraine and Safwan Hossain. Jacnet: Learning functions with structured jacobians. *arXiv*  
325 *preprint arXiv:2408.13237*, 2024.  
326

327 Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by  
328 implicit differentiation. In *International conference on artificial intelligence and statistics*, pp.  
329 1540–1552. PMLR, 2020.

330 Jonathan Lorraine, Paul Vicol, Jack Parker-Holder, Tal Kachman, Luke Metz, and Jakob Foerster.  
331 Lyapunov exponents for diversity in differentiable games. *arXiv preprint arXiv:2112.14570*, 2021.  
332

333 Jonathan Lorraine, Nihesh Anderson, Chansoo Lee, Quentin De Laroussilhe, and Mehadi Hassen.  
334 Task selection for automl system evaluation. *arXiv preprint arXiv:2208.12754*, 2022a.

335 Jonathan P Lorraine, David Acuna, Paul Vicol, and David Duvenaud. Complex momentum for  
336 optimization in games. In *International Conference on Artificial Intelligence and Statistics*, pp.  
337 7742–7765. PMLR, 2022b.

338 Jonathan Peter Lorraine. *Scalable nested optimization for deep learning*. PhD thesis, University of  
339 Toronto (Canada), 2024.  
340

341 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
342 *arXiv:1711.05101*, 2017.  
343

344 Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to  
345 stereo vision. In *IJCAI’81: 7th international joint conference on Artificial intelligence*, volume 2,  
346 pp. 674–679, 1981.

347 Nikhil Mehta, Jonathan Lorraine, Steve Masson, Ramanathan Arunachalam, Zaid Pervaiz Bhat, James  
348 Lucas, and Arun George Zachariah. Improving hyperparameter optimization with checkpointed  
349 model weights. In *European Conference on Computer Vision*, pp. 75–96. Springer, 2024.

350 Bruno Mlodozieniec, Runa Eschenhagen, Juhan Bae, Alexander Immer, David Krueger, and Richard  
351 Turner. Influence functions for scalable data attribution in diffusion models. *arXiv preprint*  
352 *arXiv:2410.13850*, 2024.  
353

354 Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak:  
355 Attributing model behavior at scale. *arXiv preprint arXiv:2303.14186*, 2023.

356 Yonghyun Park, Chieh-Hsin Lai, Satoshi Hayakawa, Yuhta Takida, Naoki Murata, Wei-Hsiang  
357 Liao, Woosung Choi, Kin Wai Cheuk, Junghyun Koo, and Yuki Mitsufuji. Concept-trak: Under-  
358 standing how diffusion models learn concepts through concept-level attribution. *arXiv preprint*  
359 *arXiv:2507.06547*, 2025.  
360

361 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
362 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.

363 Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data  
364 influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33,  
365 2020.

366 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
367 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
368 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.  
369

370 Aniruddh Raghu, Jonathan Lorraine, Simon Kornblith, Matthew McDermott, and David K Duvenaud.  
371 Meta-learning to improve pre-training. *Advances in Neural Information Processing Systems*, 34:  
372 23231–23244, 2021.

373 Jessie Richter-Powell, Antonio Torralba, and Jonathan Lorraine. Score distillation sampling for audio:  
374 Source separation, synthesis, and beyond. *arXiv preprint arXiv:2505.04621*, 2025.  
375

376 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
377 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*  
*ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, 2022.

---

378 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
379 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
380 text-to-image diffusion models with deep language understanding. *Advances in neural information*  
381 *processing systems*, 35, 2022.

382  
383 Qingyu Shi, Jianzong Wu, Jinbin Bai, Jiangning Zhang, Lu Qi, Yunhai Tong, and Xiangtai Li.  
384 Decouple and track: Benchmarking and improving video diffusion transformers for motion transfer.  
385 *arXiv preprint arXiv:2503.17350*, 2025.

386  
387 Yanke Song, Jonathan Lorraine, Weili Nie, Karsten Kreis, and James Lucas. Multi-student diffusion  
388 distillation for better one-step generators. *arXiv preprint arXiv:2410.23274*, 2024.

389  
390 Zhiyu Tan, Xiaomeng Yang, Luozheng Qin, and Hao Li. Vidgen-1m: A large-scale dataset for  
391 text-to-video generation. *arXiv preprint arXiv:2408.02629*, 2024.

392  
393 Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European*  
394 *conference on computer vision*, pp. 402–419. Springer, 2020.

395  
396 Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion  
397 and content for video generation. In *Proceedings of the IEEE conference on computer vision and*  
398 *pattern recognition*, pp. 1526–1535, 2018.

399  
400 Paul Vicol, Jonathan P Lorraine, Fabian Pedregosa, David Duvenaud, and Roger B Grosse. On  
401 implicit bias in overparameterized bilevel optimization. In *International Conference on Machine*  
402 *Learning*, pp. 22234–22259. PMLR, 2022.

403  
404 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,  
405 Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai  
406 Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi  
407 Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang,  
408 Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng  
409 Zhou, Wenteng Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan  
410 Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You  
411 Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen  
412 Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models.  
413 *arXiv preprint arXiv:2503.20314*, 2025.

414  
415 Jiachen T Wang, Prateek Mittal, Dawn Song, and Ruoxi Jia. Data shapley in one training run. *arXiv*  
416 *preprint arXiv:2406.11011*, 2024a.

417  
418 Jiachen T Wang, Tianji Yang, James Zou, Yongchan Kwon, and Ruoxi Jia. Rethinking data shapley  
419 for data selection tasks: Misleads and merits. *arXiv preprint arXiv:2405.03875*, 2024b.

420  
421 Jiangshan Wang, Yue Ma, Jiayi Guo, Yicheng Xiao, Gao Huang, and Xiu Li. Cove: Unleashing the  
422 diffusion feature correspondence for consistent video editing. *Advances in Neural Information*  
423 *Processing Systems*, 37:96541–96565, 2024c.

424  
425 Sheng-Yu Wang, Alexei A Efros, Jun-Yan Zhu, and Richard Zhang. Evaluating data attribution for  
426 text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer*  
427 *Vision*, 2023a.

428  
429 Sheng-Yu Wang, Aaron Hertzmann, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Fast data  
430 attribution for text-to-image models. *arXiv preprint arXiv:2511.10721*, 2025.

431  
432 Xiang Wang, Shiwei Zhang, Han Zhang, Yu Liu, Yingya Zhang, Changxin Gao, and Nong Sang.  
433 Videolcm: Video latent consistency model. *arXiv preprint arXiv:2312.09109*, 2023b.

434  
435 Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky,  
436 Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners.  
437 *arXiv preprint arXiv:2509.20328*, 2025.

---

432 Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu,  
433 Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion  
434 models for text-to-video generation. In *Proceedings of the IEEE/CVF international conference on*  
435 *computer vision*, 2023.

436 Xindi Wu, Mengzhou Xia, Rulin Shao, Zhiwei Deng, Pang Wei Koh, and Olga Russakovsky. Icons:  
437 Influence consensus for vision-language data selection. *arXiv preprint arXiv:2501.00654*, 2024.  
438

439 Tong Xie, Haoyu Li, Andrew Bai, and Cho-Jui Hsieh. Data attribution for diffusion models: Timestep-  
440 induced bias in influence estimation. *arXiv preprint arXiv:2401.09031*, 2024.  
441

442 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
443 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
444 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

445 Tianjun Zhang, Yi Zhang, Vibhav Vineet, Neel Joshi, and Xin Wang. Controllable text-to-image  
446 generation with gpt-4. *arXiv preprint arXiv:2305.18583*, 2023.  
447

448 Xiaosen Zheng, Tianyu Pang, Chao Du, Jing Jiang, and Min Lin. Intriguing properties of data  
449 attribution on diffusion models. *arXiv preprint arXiv:2311.00500*, 2023.

450 Shangchen Zhou, Peiqing Yang, Jianyi Wang, Yihang Luo, and Chen Change Loy. Upscale-a-video:  
451 Temporal-consistent diffusion model for real-world video super-resolution. In *Proceedings of the*  
452 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2535–2545, 2024.  
453

454 Yixuan Zhu, Jiaqi Feng, Wenzhao Zheng, Yuan Gao, Xin Tao, Pengfei Wan, Jie Zhou, and Jiwen  
455 Lu. Astra: General interactive world model with autoregressive denoising. *arXiv preprint*  
456 *arXiv:2512.08931*, 2025.  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

## 486 A NOTATION

488 Table 2: Glossary and notation.

491 Symbol	Description
492 <i>Acronyms and Basic Notation</i>	
493 VAE	Variational Autoencoder
494 DiT	Diffusion Transformer backbone
495 $\mathbf{I}$	Identity matrix
496 <hr/>	
497 <i>Video Generation</i>	
498 $p_{\theta}(\mathbf{v} \mid \mathbf{c})$	Conditional video generator with parameters $\theta$
499 $\mathbf{v} \in \mathbb{R}^{F \times H \times W \times 3}$	Video clip with frames $F$ , height $H$ , width $W$
500 $\mathbf{c}$	Conditioning signal such as text or multimodal metadata
501 $\theta$	Trainable model parameters
502 $f \in \{1, \dots, F\}$	Frame index
503 $h \in \{1, \dots, H\}, w \in \{1, \dots, W\}$	Spatial indices for height and width respectively
504 $\tilde{h}, \tilde{w}$	Latent grid indices
505 <hr/>	
506 <i>Datasets</i>	
507 $\mathcal{D} = \{(\mathbf{v}_n, \mathbf{c}_n)\}_{n=1}^N$	Training corpus with size $N$ and index $n$
508 $\mathcal{D}_{\text{fit}} \subseteq \mathcal{D}$	Fine-tuning dataset
509 $\mathcal{S} \subseteq \mathcal{D}$	Selected influential subset
510 $k \in \{1, \dots, K\}$	The selected subset size
511 $Q$	Number of query data
512 $q \in \{1, \dots, Q\}$	Query index
513 $\mathbf{x}$	Generic input data pair
514 $\mathbf{x}_{\text{test}}, \mathbf{x}_n$	Test/query pair and training pair
515 $\hat{\mathbf{v}}, \hat{\mathbf{c}}$	Query video and its conditioning
516 <hr/>	
517 <i>Latent Space and Diffusion Components</i>	
518 $E$	VAE encoder
519 $\mathbf{h} = E(\mathbf{v}) \in \mathbb{R}^{F \times (H/s) \times (W/s) \times C}$	Latent video with spatial factor $s$ and channels $C$
520 $\mathbf{z}$	Noisy latent variable used in diffusion or flow matching
521 $\epsilon \sim \mathcal{N}(0, \mathbf{I})$	Gaussian noise
522 $\epsilon_{\theta}(\mathbf{z}, \mathbf{c}, t)$	Predicted noise network in diffusion training
523 $\mathbf{f}_{\theta}(\mathbf{z}, \mathbf{c}, t)$	Time-indexed vector field in flow matching
524 $\dot{\mathbf{z}}$	Time derivative of the latent trajectory
525 $\alpha_t, \sigma_t$	Scheduler signal and noise scales at timestep $t$
526 $\epsilon_{\text{target}}$	Target noise or velocity used for supervision
527 $t \in \{1, \dots, T\}$	Diffusion or flow-matching timestep, with total timesteps $T$
528 $t_{\text{fix}}, \epsilon_{\text{fix}}$	Fixed timestep and shared noise draw used for low-variance gradients
529 <hr/>	

## 532 B BACKGROUND

533 **Diffusion and flow matching in latent space.** Let  $p_{\theta}(\mathbf{v} \mid \mathbf{c})$  be a conditional generator with  
534 parameters  $\theta$ , where  $\mathbf{v} \in \mathbb{R}^{F \times H \times W \times 3}$  is a clip of height  $H$ , width  $W$ , and  $F$  frames, and  $\mathbf{c}$  denotes  
535 conditioning such as text or other multimodal metadata (e.g., fps, depth, pose). We operate in VAE  
536 latents:  $\mathbf{h} = E(\mathbf{v})$  and train a denoiser or velocity field on noisy latents. A noise scheduler supplies  
537 time-dependent coefficients  $(\alpha_t, \sigma_t)$  controlling signal and noise scales, and the forward noising is:  
538

$$539 \mathbf{z}(t, \epsilon) = \alpha_t \mathbf{h} + \sigma_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}), \quad t \in \{1, \dots, T\}. \quad (13)$$

Table 3: Glossary and notation (continued).

Symbol	Description
<i>Attribution and Influence</i>	
$I(\mathbf{v}_n, \hat{\mathbf{v}}; \boldsymbol{\theta})$	Influence score between a train clip and a query clip
$I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}}; \boldsymbol{\theta})$	Motion-aware influence score
$\text{TopK}(\cdot)$	Top- $K$ operator for selecting highest scores
$\text{MajVote}(\cdot)$	Majority-vote aggregation across queries
$\tau$	Percentile cutoff for voting
$\rho$	Spearman correlation coefficient
<i>Loss Functions</i>	
$\mathcal{L}$	Generic loss
$\mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{v}, \mathbf{c}), \mathcal{L}_{\text{flow}}(\boldsymbol{\theta}; \mathbf{v}, \mathbf{c})$	Diffusion and flow-matching objective
$\mathcal{L}_{\text{mot}}(\boldsymbol{\theta}; \mathbf{v}, \mathbf{c})$	Motion-weighted objective used for attribution
$\tilde{\mathcal{L}}$	Per-location squared error in latent space
$\mathbf{g}, \tilde{\mathbf{g}}$	Gradient and its projected version
$\mathbf{g}_{\text{mot}}, \tilde{\mathbf{g}}_{\text{mot}}$	Motion-weighted gradient and its projection
$\mathbf{H}_{\boldsymbol{\theta}}$	Hessian with respect to $\boldsymbol{\theta}$
<i>Motion Representations</i>	
$\mathcal{A}(\mathbf{v}) = A$	AllTracker motion extraction
$A \in \mathbb{R}^{F \times H \times W \times 4}$	Motion tensor containing flow, visibility, and confidence
$\mathbf{D}_f(h, w)$	Displacement vector at frame $f$ and location $(h, w)$
$M_f(h, w)$	Motion magnitude at a location, computed from the displacement
$\mathbf{W}(f, h, w) \in [0, 1]$	Normalized motion weights used to mask per-location losses
$\zeta$	Small numerical bias for stability (e.g., $10^{-6}$ )
<i>Projections and Computational Details</i>	
$D, D'$	Full and projected gradient dimensions
$\mathbf{P} \in \mathbb{R}^{D' \times D}$	Projection matrix used for Fastfood-style JL projection
$\xi$	Variance normalization constant for projection
$\mathcal{T}$	Set of sampled $(t, \epsilon)$ pairs for gradient estimation
$B$	Unit compute cost used in complexity accounting

*Denosing diffusion* (Ho et al., 2020) trains a network  $\epsilon_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{c}, t)$  to predict the injected noise:

$$\mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{v}, \mathbf{c}) = \mathbb{E}_{t, \epsilon} \left[ \|\epsilon_{\boldsymbol{\theta}}(\mathbf{z}(t, \epsilon), \mathbf{c}, t) - \epsilon\|_2^2 \right]. \quad (14)$$

*Flow matching* (Lipman et al., 2022; Albergo et al., 2023) learns a time-dependent vector field  $\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{z}_t, \mathbf{c}, t)$  that matches the instantaneous velocity  $\dot{\mathbf{z}} = \frac{d}{dt}\mathbf{z}$  induced by a chosen interpolant:

$$\mathcal{L}_{\text{flow}}(\boldsymbol{\theta}; \mathbf{v}, \mathbf{c}) = \mathbb{E}_{t, \epsilon} \left[ \|\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{z}(t, \epsilon), \mathbf{c}, t) - \dot{\mathbf{z}}(t, \epsilon)\|_2^2 \right]. \quad (15)$$

Both objectives train time-indexed predictors over the latent space by integrating over  $t$  and  $\epsilon$ , thus gradient-based methods like attribution share similar challenges.

**From images to video for generation.** Adding a temporal axis materially changes modeling and training. Generation must capture spatial appearance and temporal dynamics such as object and camera motion, deformations, and interactions. Modern systems extend image backbones with temporal capacity, for example, 3D U-Nets or 2D U-Nets augmented with temporal attention, causal or sliding-window context, and factorized space-time blocks, often trained in a latent-video VAE that compresses frames while preserving temporal cues. Training departs from images along several axes, which we address in §2: (i) *Compute and storage*. Longer sequences multiply the cost of sampling timesteps, noise draws, and frames, motivating fixed-timestep or small-subset estimators that reduce variance without prohibitive cost (§2.2). (ii) *Variable horizon*. Clips vary in  $F$  and frame rate (§D.1). (iii) *Time-specific failure modes*. Typical artifacts include inconsistent trajectories, temporal flicker, identity drift, and physically implausible dynamics despite sharp individual frames (§2.3).

**Motion representations in videos.** We denote our video as  $\mathbf{v} = [\mathbf{f}_f]_{f=1}^F$  with  $\mathbf{f}_f \in \mathbb{R}^{H \times W \times 3}$  being the  $f$ -th frame. We represent motion via optical flow between consecutive frames:  $\mathbf{F}_f : \{1, \dots, H\} \times \{1, \dots, W\} \rightarrow \mathbb{R}^2$ , where each flow vector in  $\mathbb{R}^2$  encodes the horizontal displacement  $dw$  and vertical displacement  $dh$  of a pixel. The motion magnitude is  $M_f(h, w) = \|\mathbf{F}_f(h, w)\|_2$ . The  $M_f$  over frames  $f$  and pixels  $h, w$  summarizes the amount and spatial layout of motion in a clip, which we use to provide masks in our motion-weighted loss in §2.

## B.1 DATA ATTRIBUTION

Data attribution measures how individual training samples affect a model’s predictions (Bae et al., 2024). A classic approach to data attribution is to use influence functions (Koh & Liang, 2017). Intuitively, the influence of a training sample measures: if we upweight this training example, how much would the model’s prediction on a test datum change? Consider a loss function  $\mathcal{L}(\boldsymbol{\theta}; \mathbf{x})$  and a test sample  $\mathbf{x}_{\text{test}}$ , the influence of a training point  $\mathbf{x}_n$  can be quantified as:

$$I(\mathbf{x}_n, \mathbf{x}_{\text{test}}) = -\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_{\text{test}})^\top \mathbf{H}_{\boldsymbol{\theta}}^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_n),$$

$$\mathbf{H}_{\boldsymbol{\theta}} = \frac{1}{N} \sum_{n=1}^N \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_n), \quad (16)$$

where the inverse Hessian captures the curvature of the loss landscape, yet computing or storing it is infeasible at modern model and dataset scales. Thus, practical methods (e.g., TracIn (Pruthi et al., 2020) and TRAK (Park et al., 2023)) approximate influence via gradient inner products or gradient feature projections.

**Attribution in diffusion models.** Diffusion training aggregates gradients over timesteps  $t$  and noise draws  $\epsilon$ , where gradient norms vary systematically with  $t$ , producing a timestep bias where examples aligned with large-norm timesteps appear spuriously influential. Diffusion-ReTrac (Xie et al., 2024) reduces this bias by normalizing gradients and sub-sampling  $t$  and  $\epsilon$  for influence. Let  $\mathcal{L}_{\text{diff}}$  denote the diffusion loss, and with the sampled-timestep-and-noise set  $\mathcal{T}$ , we compute a cosine-style score for normalized test and train gradients:

$$I_{\text{diff}}(\mathbf{x}_n, \mathbf{x}_{\text{test}}) = \frac{1}{|\mathcal{T}_{\text{test}}|} \sum_{t, \epsilon \in \mathcal{T}_{\text{test}}} \frac{\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{x}_{\text{test}}, t, \epsilon)^\top}{\|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{x}_{\text{test}}, t, \epsilon)\|}$$

$$\frac{1}{|\mathcal{T}_n|} \sum_{t, \epsilon \in \mathcal{T}_n} \frac{\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{x}_n, t, \epsilon)}{\|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{diff}}(\boldsymbol{\theta}; \mathbf{x}_n, t, \epsilon)\|}. \quad (17)$$

Averaging gradients over  $(t, \epsilon)$  stabilizes estimates, and normalization mitigates timestep-induced scale effects. Attribution quality is also sensitive to the measurement function used to score examples, such as denoising loss versus likelihood proxies (Zheng et al., 2023).

**Why vanilla attribution is insufficient for videos.** Naïvely applying gradient-based attribution to video diffusion risks treating appearance and motion alike, overemphasizing appearance (objects, textures, backgrounds) while overlooking dynamics (Park et al., 2025; Tulyakov et al., 2018). Its cost grows with clip length, sampled timesteps, noise draws, and gradient dimensionality, making naïve methods impractical at modern video scales. To improve motion, we need attribution that suppresses static appearance, emphasizes motion-specific signals, and remains efficient (§2). Motion is distributed across frames and entangled with static cues, thus influence cannot be assigned frame-independently.

---

## C EXTENDED RELATED WORK

### C.1 DATA ATTRIBUTION

Understanding how individual training examples shape model behavior has been a long-standing goal. Modern data attribution methods fall into two main groups (Hammoudeh & Lowd, 2024): retraining-based methods (e.g., leave-one-out (Cook et al., 1982; Jia et al., 2021), downsampling (also known as subsampling or counterfactual influence) (Feldman & Zhang, 2020), Shapley-value (Wang et al., 2024a;b)) and gradient-based methods (influence-function family, including Influence Functions (Koh & Liang, 2017; Lorraine & Hossain, 2024), TracIn (Pruthi et al., 2020), and TRAK (Park et al., 2023)). Influence functions provide a principled framework by approximating the effect of removing a training point. TracIn (Pruthi et al., 2020) and TRAK (Park et al., 2023) make attribution feasible at scale. While effective for classification, these assume a mapping between training gradients and predictions, which becomes more complex in generative models.

Data attribution traces how individual training examples (or subsets) influence a model’s predictions or behavior. Formally, it assigns an attribution score to each training sample, estimating the extent to which that sample contributes (positively or negatively) to the model’s output on a given test query or behavior. Influence data attribution is an example of nested optimization (Vicol et al., 2022; Lorraine, 2024) with other examples including differentiable games (Balduzzi et al., 2018; Lorraine et al., 2021; 2022b) and hyperparameter optimization (Raghu et al., 2021; Lorraine & Duvenaud, 2018; Mehta et al., 2024; Lorraine et al., 2020). Before diffusion models, attribution methods were applied to supervised learning tasks such as classification and regression, where influence functions (Koh & Liang, 2017; Lorraine et al., 2022a) and scalable approximations such as TracIn (Pruthi et al., 2020), TRAK (Park et al., 2023), and TDA (Bae et al., 2024) quantified the impact of training examples on downstream predictions. Recent work adapted data attribution to diffusion models (Georgiev et al., 2023; Zheng et al., 2023; Wang et al., 2025; 2024a; Lin et al., 2024; Brokman et al., 2024; Kwon et al., 2023), where iterative denoising introduces timestep-dependent bias. Mlodozienec et al. (2024) propose scalable approximations, while Xie et al. (2024) identify timestep-induced artifacts and normalization schemes. Concept-TRAK (Park et al., 2025) extends attribution to concepts by reweighting gradients with concept-specific rewards, enabling attribution to semantic factors. Wang et al. (2023a) instead design a customization-based benchmark for text-to-image models, where models are fine-tuned on exemplar images with novel tokens and attribution is evaluated by whether it can recover the responsible exemplars. However, these are limited to image diffusion, which captures static appearance but not temporal dynamics.

### C.2 MOTION IN VIDEO GENERATION

Video diffusion extends image generation to time, requiring coherent motion across frames (Ho et al., 2022; Blattmann et al., 2023; Peebles & Xie, 2023; Wan et al., 2025; Agarwal et al., 2025). A large body of work builds temporal structure via attention layers (Wu et al., 2023), control signals (Chen et al., 2023; Zhang et al., 2023), feature correspondences (Geyer et al., 2023; Bao et al., 2023; Wang et al., 2024c), or consistency distillation (Wang et al., 2023b; Zhou et al., 2024). Recent work has highlighted the challenge of decoupling motion from appearance in video diffusion transformers, where spatial and temporal information become entangled in the model’s representations (Shi et al., 2025). However, understanding which training clips influence specific motion patterns in generated videos remains an open challenge.

In parallel, motion has long been studied using optical flow and correspondence, from classical formulations (Horn & Schunck, 1981; Lucas & Kanade, 1981) to modern approaches such as RAFT (Teed & Deng, 2020), which improve accuracy and generalization. These priors are often repurposed during generation to guide dynamics, but they do not explain which training examples shaped a model’s motion behavior. Our work addresses both gaps by introducing a motion-aware data attribution framework specifically designed for video diffusion. We use motion-weighted gradients that disentangle temporal dynamics from static appearance, enabling us to trace generated motion patterns back to the most influential training clips.

---

## 702 D ADDITIONAL DETAILS OF METHOD

### 703 D.1 VIDEO-SPECIFIC FRAME-LENGTH BIAS FIX

704 Raw gradient magnitudes depend on the number of frames  $F$  in the video  $\mathbf{v}$ , thereby biasing  
705 scores toward longer videos. We correct this by normalizing for frame count before the projec-  
706 tion-normalization step:

$$707 \nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{v}, t_{\text{fix}}, \epsilon_{\text{fix}}) \leftarrow \frac{1}{F} \nabla_{\theta} \mathcal{L}_{\text{diff}}(\theta; \mathbf{v}, t_{\text{fix}}, \epsilon_{\text{fix}}). \quad (18)$$

708 We still apply  $\ell_2$  normalization in Eq. 5, further stabilizing scales across examples. Together, single-  
709 timestep, common randomness, projection, and frame-length correction form a compact, scalable  
710 estimator that we use throughout. However, naïve video-level attribution conflates appearance with  
711 motion, often ranking clips high due to shared backgrounds or objects, while offering little insight  
712 into dynamics.

### 713 D.2 MOST INFLUENTIAL FINE-TUNING SUBSET SELECTION

714 **Goal.** Given a query clip  $(\hat{\mathbf{v}}, \hat{\mathbf{c}})$ , we compute a motion-aware attribution value for each candidate  
715 fine-tuning sample  $(\mathbf{v}_n, \mathbf{c}_n) \in \mathcal{D}_{\text{ft}}$  using  $I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}})$  (Eq. 11). We construct a fine-tuning set  $\mathcal{S}$  for  
716 one or many query videos  $\hat{\mathbf{v}}$ .

717 **Single-query-point fine-tuning selection.** For a budget of  $K$  data points, we select the  $K$  highest-  
718 scoring examples. In practice,  $K$  is chosen as a percentile of the dataset size (e.g., top 1–10%),  
719 ensuring the subset scales consistently across datasets.

720 **Multi-query-point fine-tuning selection: aggregating attribution scores.** For  $Q$  queries, we adopt  
721 the majority voting approach from ICONS (Wu et al., 2024) and aggregate motion-aware influence  
722 scores across queries by percentile thresholding and voting. A sample receives a vote if the score  
723 is above the percentile cutoff  $\tau$  for that query. The consensus score of a candidate  $\mathbf{v}_n$  is the total  
724 number of queries that vote for it. We then rank all training samples by  $\text{MajVote}(\mathbf{v}_n)$  and select  
725 the top- $K$  to form the fine-tuning subset. This formulation emphasizes samples that are consistently  
726 influential across multiple queries, without requiring cross-query calibration of raw scores:

$$727 \text{MajVote}_n = \sum_{q=1}^Q \mathbb{I}[I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}}_q) > \tau], \quad (19)$$
$$728 \mathcal{S}_{\text{vote}}(K) = \{\mathbf{v}_n | \mathbf{v}_n \text{ in top-}K \text{ by MajVote}\}.$$

### 729 D.3 COMPUTATIONAL EFFICIENCY ANALYSIS

730 **Gradient Compute.** Naïvely averaging over timesteps and noise for every example costs  
731  $\mathcal{O}(|\mathcal{D}| |\mathcal{T}| B)$ , where  $B$  is a single forward+backward cost and  $|\mathcal{T}|$  is the number of sampled  $t$ ,  $\epsilon$  per  
732 data. Using a single sample reduces this to  $\mathcal{O}(|\mathcal{D}| B)$ , which is key to keeping the cost reasonable for  
733 modern video datasets and models, while reusing a sample across data allows low enough variance  
734 for stable rankings. Projection adds  $\mathcal{O}(D' \log D')$  per example using Fastfood (Le et al., 2014),  
735 negligible relative to a backward pass.

736 **Gradient Storage.** Storing full gradients is  $\mathcal{O}(|\mathcal{D}| D)$ . We instead store only projected vectors,  
737  $\mathcal{O}(|\mathcal{D}| D')$ , plus the structured Fastfood state,  $\mathcal{O}(D)$ . Since  $D'$  is typically orders of magnitude  
738 smaller than  $D$ , this transformation makes storage tractable for billion-parameter models.

739 **Data Ranking Compute.** Influence computation in Eq. 5 is an inner product in  $\mathbb{R}^{D'}$ , so evaluating  
740 all train examples against a query is  $\mathcal{O}(|\mathcal{D}| D')$ , and sorting is  $\mathcal{O}(|\mathcal{D}| \log |\mathcal{D}|)$ .

741 **Additional Motion-Emphasis Compute.** Motion-specific overhead primarily stems from AllTracker  
742 mask extraction with complexity  $\mathcal{O}(|\mathcal{D}| \cdot H \cdot W \cdot F)$  for clip length  $F$  and frame resolution  $H \times W$ .  
743 Masks are extracted once, cached, and negligible relative to gradient cost. We provide a detailed  
744 runtime breakdown in App. I.1.

## 745 E ADDITIONAL EXPERIMENTS

746 **Setup.** We evaluate on **VIDGEN-1M** (Tan et al., 2024) and **4DNeX-10M** (Chen et al., 2025) using  
747 10k videos. Query videos span 10 motion categories (compress, bounce, roll, spin, slide, free

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

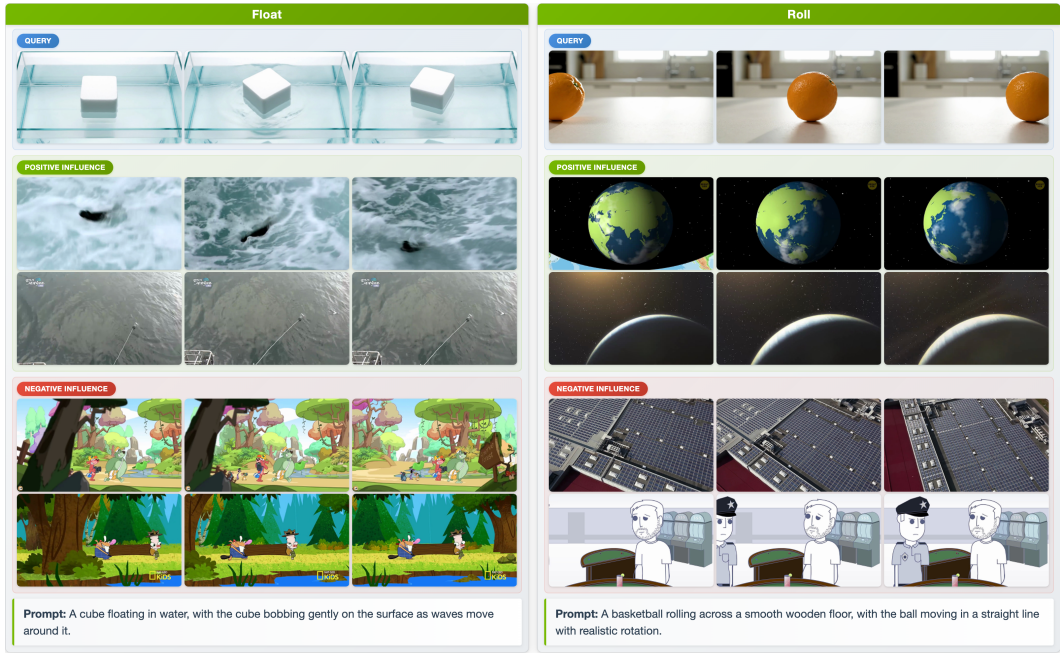


Figure 2: **Motion Attribution Samples** with Wan2.1-T2V-1.3B. *Top*: Query clips showing float (*left*) and roll (*right*) motions. *Middle*: Top-ranked positive training samples identified by **MOTIVE** with high influence scores. *Bottom*: Negative influence samples with minimal, camera-only motion, or cartoon-style content that conflict with target motions. Videos are included in the supplementary material.

fall, float, inflate, deflate, oscillate) with 5 videos each. We fine-tune **Wan2.1-T2V-1.3B** and **Wan2.2-TI2V-5B** models, updating only the DiT backbone with LoRA (rank 128). Baselines include random selection, motion magnitude filtering, V-JEPA embeddings (Assran et al., 2025), and whole-video attribution without motion masking. We evaluate using VBench (Huang et al., 2024) metrics: motion smoothness, dynamic degree, aesthetic quality, and imaging quality.

**Quantitative Results.** Tab. 1 shows **MOTIVE** achieves the highest dynamic degree (47.6% and 48.3%) across both models, significantly outperforming random selection (+6.3% / +6.7%) and whole-video attribution without motion masking (+3.8% / +4.5%). Notably, using only 10% of training data, we surpass full fine-tuning on dynamic degree by +5.6% / +3.0%, demonstrating superior performance of motion-specific attribution for targeted fine-tuning. Our method also maintains competitive motion smoothness (96.3% / 97.6%) and improves aesthetic quality (46.0% / 45.6%).

**Human Evaluation.** We conducted a human study with 17 annotators comparing 50 generated videos across methods. Results show 74.1% preference for **MOTIVE** over the base model and 53.1% over full fine-tuning, confirming that motion-attributed data selection produces perceptually better motion. Additional results on **LTX-2B** are in the Appendix.

**Qualitative Analysis.** Fig. 3 shows example attributions. Our method ranks clips with coherent physical dynamics highly, objects floating, rolling, and spinning with realistic physics, and downranks static footage, camera-only motion, and cartoon-style content with non-transferable kinematics. Positive-influence samples share motion semantics with queries (e.g., buoyant objects for “float”, rotating objects for “roll”), while negative-influence samples contain minimal motion or incompatible dynamics. This pattern holds across all motion categories, aligning with the quantitative gains.

**Ablations.** We validate key design choices: (1) Single-timestep attribution ( $t = 751$ ) achieves 66% Spearman correlation with multi-timestep ground truth while reducing compute by  $10\times$ . (2) Projection to 512 dimensions preserves 74.7% ranking correlation with full gradients.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

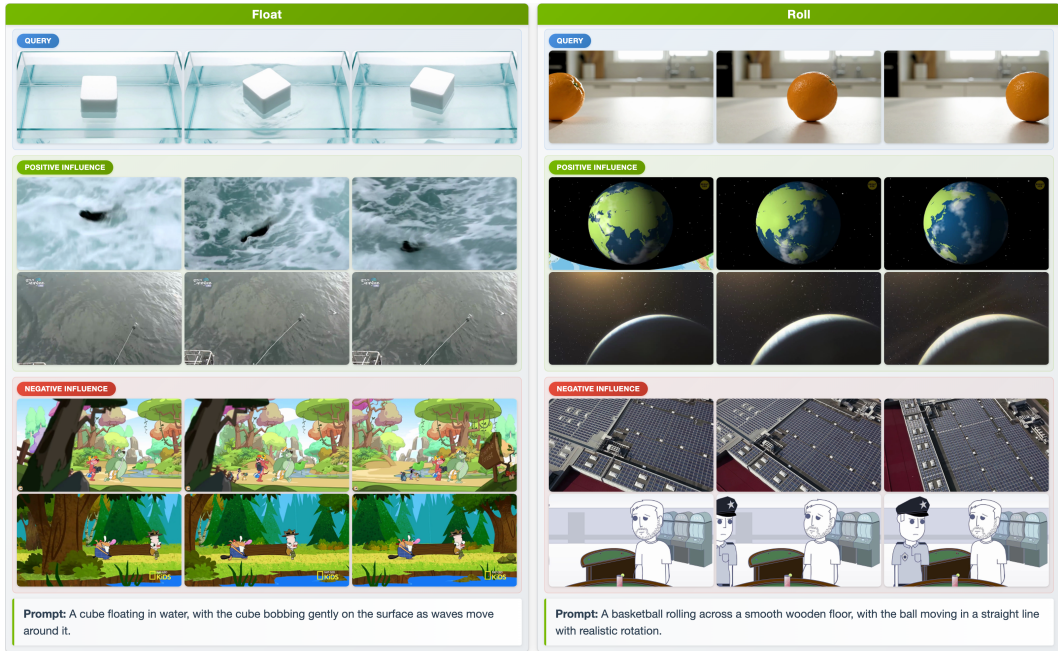


Figure 3: **Motion Attribution Examples.** *Top:* Query clips showing float (left) and roll (right) motions. *Middle:* Top-ranked positive training samples with high influence scores showing similar motion patterns. *Bottom:* Negative influence samples with static content or cartoon-style motion that conflicts with target dynamics.

Table 4: **VBench Evaluation on LTX-2B.** Performance comparison on VBench (Huang et al., 2024) across different baselines (all values in %, higher is better). All selection methods use 10% of training data; MM: motion masking.

Method	LTX-2B Video					
	Subject Consist.	Background Consist.	Motion Smooth.	Dynamic Degree	Aesthetic Qual.	Imaging Qual.
Base	93.8	95.2	94.6	36.5	32.1	48.2
Full FT	94.4	<b>95.7</b>	94.8	38.9	32.6	<b>49.1</b>
Random	93.7	95.4	94.7	37.8	<b>32.8</b>	48.4
Motion mag.	94.1	95.3	95.2	39.6	32.4	47.9
V-JEPA	94.2	95.5	95.1	40.2	31.9	47.5
Ours w/o MM	94.3	95.6	95.3	41.8	<b>32.8</b>	48.7
<b>Ours (MOTIVE)</b>	<b>94.8</b>	95.6	<b>95.5</b>	<b>45.1</b>	32.7	48.6

### E.1 RESULTS ON ADDITIONAL VIDEO GENERATION MODELS

We further test our framework on additional video generation architectures beyond the **Wan** models (Tab. 1). We evaluate **MOTIVE** on **LTX-2B** HaCohen et al. (2024), which represents a different architectural design and training paradigm. The results in Tab. 4 demonstrate that **MOTIVE** works effectively across different model architectures and scales.

## F ANALYSIS

### F.1 MOTION DISTRIBUTION ANALYSIS

**MOTIVE is not simply selecting “motion-rich” clips:** The key distinction is that our influence scores are computed via gradients, and training videos are considered influential only when they directly allow the model to lower the loss, improving the model’s ability to generate the target motion dynamics, not because they contain more motion overall.

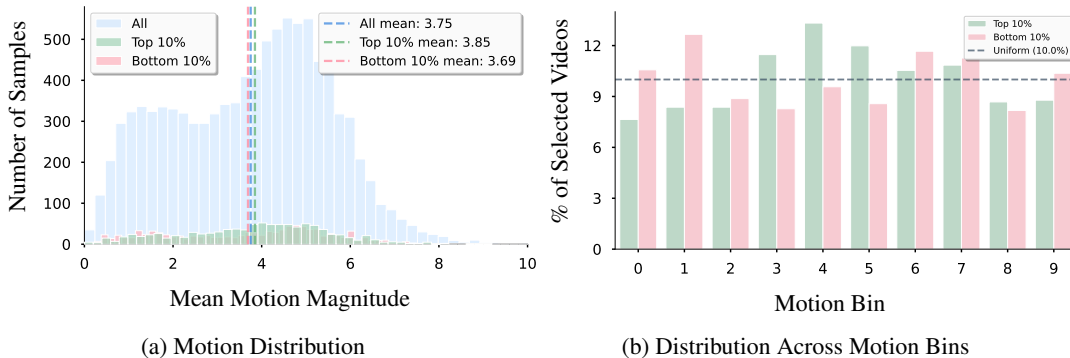


Figure 4: **MOTIVE is not simply selecting “motion-rich” clips.** Our influence scores are computed via gradients, and training videos are considered influential only when they directly improve the model’s ability to generate the target motion dynamics, not because they contain more motion overall.

To empirically validate this, we further analyze the distribution of motion magnitudes in our selected data. We compute the mean motion magnitude for 10k videos in the **VIDGEN** dataset and compare the distributions of the top 10% (highest influence scores) and the bottom 10% (lowest influence scores).

As shown in Fig. 4, the top 10% selected videos have a mean motion magnitude of 3.85, which is only 4.3% higher than the bottom 10% (3.69), despite representing opposite extremes of influence scores. The analysis also shows that within the moderate-motion range (bins 3, 4, and 5), the top 10% of positive-influence samples outnumber the bottom 10% of negative-influence samples. Yet, both groups also appear in low-motion bins (0-2) and high-motion bins (6-9).

This distribution pattern shows that high-influence videos selected by **MOTIVE** span the entire motion spectrum, not just high-motion regions. Many high-motion videos receive low influence scores, while numerous influential videos exhibit modest or even low motion magnitude. These findings show that our motion attribution approach captures training influence, focusing on motion rather than simply acting as a motion-saliency filter.

## F.2 CROSS-MOTION INFLUENCE PATTERNS

To analyze cross-motion influence patterns, we examine the percentage overlap of top-100 influential training data across different motion categories in both 4DNEX and VIDGEN datasets. We use 5 query samples to identify the top-100 most influential training videos, aggregating results across queries. As shown in Fig. 5, both datasets exhibit remarkably similar patterns with mean overlaps of 24.0% and 24.3%, respectively, indicating moderate sharing of influential data across motion categories.

Both datasets consistently identify the same high-overlap pairs: bounce-float (44.4%/46.3%), compress-float (40.1%/34.0%), and compress-spin (36.9%/39.6%), suggesting these motions share fundamental characteristics that the model learns from similar training examples. Conversely, low-overlap pairs such as free fall-stretch (12.8%/12.7%) and float-slide (14.0%/10.9%) indicate more specialized influential data for mechanically dissimilar motions. The influence matrices are asymmetric because the number of unique influential samples shared across the 5 query samples differs across motion categories. The similar cross-motion influence patterns observed across both the 4DNEX and VIDGEN datasets demonstrate that these relationships are generalizable across different video datasets and reflect dynamic similarity.

## G ADDITIONAL METHOD DETAILS

**Tracker-agnostic scope.** We treat the motion estimator as a pluggable source of saliency rather than a training dependency. Given displacement magnitudes, we construct latent-space weights via bilinear mapping and normalization. Our implementation supports alternative estimators (such as dense optical flow or point tracking) with identical interfaces, enabling users to swap AllTracker without modifying the attribution code.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

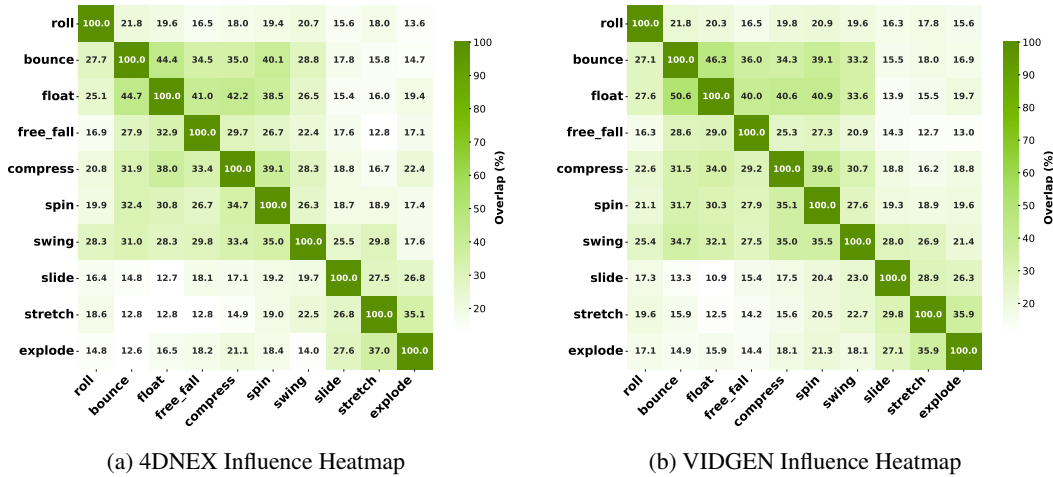


Figure 5: **Cross-motion influence overlap across datasets.** Heatmaps showing the percentage overlap of top-100 influential training samples across motion categories for (a) 4DNEX and (b) VIDGEN datasets. Each cell  $(i, j)$  shows the percentage of motion category  $i$ 's influential data (aggregated from 5 queries per category) that also appears in motion category  $j$ 's top-100 influential samples. The asymmetric nature of the matrices (e.g., bounce→float  $\neq$  float→bounce) arises because different motion categories have different numbers of unique influential videos, leading to directional overlap percentages. Consistent high-overlap pairs (e.g., bounce–float: 44.4%/46.3%) and low-overlap pairs (e.g., free fall–stretch: 12.8%/12.7%) across datasets validate that these influence patterns reflect fundamental aspects of motion representation in video generation models.

---

**Algorithm 1** MOTIVE: Motion-Aware Data Attribution Framework

---

**Require:** fine-tuning corpus  $\mathcal{D}_{\text{fit}} = \{(\mathbf{v}_n, \mathbf{c}_n)\}_{n=1}^N$ , query video  $(\hat{\mathbf{v}}, \hat{\mathbf{c}})$ , fixed  $(t_{\text{fix}}, \epsilon_{\text{fix}})$ , projection matrix  $\mathbf{P}$

**Ensure:** Motion-aware influence scores  $\{I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}})\}_{n=1}^N$

- 1: **for**  $(\mathbf{v}_n, \mathbf{c}_n) \in \mathcal{D}_{\text{fit}}$  **do**
- 2:  $A_n \leftarrow \text{ALLTRACKER}(\mathbf{v}_n)$  {Extract per-pixel flow displacements  $\mathbf{D}_f$  (Eq. 6)}
- 3: Downsample and normalize to latent-space motion mask  $\mathbf{W}_n$  (Eqs. 7–9)
- 4: Evaluate motion-weighted loss  $\mathcal{L}_{\text{mot}}(\theta; \mathbf{v}_n, \mathbf{c}_n)$  (Eq. 10)
- 5: Compute motion gradient  $\mathbf{g}_{\text{mot}}(\theta, \mathbf{v}_n, t_{\text{fix}}, \epsilon_{\text{fix}}) = \nabla_{\theta} \mathcal{L}_{\text{mot}}(\theta; \mathbf{v}_n, \mathbf{c}_n, t_{\text{fix}}, \epsilon_{\text{fix}})$
- 6: Normalize by frame length:  $\mathbf{g}_{\text{mot}} \leftarrow \mathbf{g}_{\text{mot}}/F$
- 7: Project motion gradient:  $\tilde{\mathbf{g}}_{\text{mot}}(\theta, \mathbf{v}_n) := \frac{\mathbf{P}\mathbf{g}_{\text{mot}}(\theta, \mathbf{v}_n, t_{\text{fix}}, \epsilon_{\text{fix}})}{\|\mathbf{P}\mathbf{g}_{\text{mot}}(\theta, \mathbf{v}_n, t_{\text{fix}}, \epsilon_{\text{fix}})\|}$  (Eq. 11)
- 8: **end for**
- 9: Compute query gradient:  $\tilde{\mathbf{g}}_{\text{mot}}(\theta, \hat{\mathbf{v}})$  using the same procedure for  $(\hat{\mathbf{v}}, \hat{\mathbf{c}})$
- 10: **for**  $n = 1, \dots, N$  **do**
- 11:  $I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}}) = \tilde{\mathbf{g}}_{\text{mot}}(\theta, \hat{\mathbf{v}})^{\top} \tilde{\mathbf{g}}_{\text{mot}}(\theta, \mathbf{v}_n)$  (Eq. 11)
- 12: **end for**
- 13: Rank all training clips by  $I_{\text{mot}}(\mathbf{v}_n, \hat{\mathbf{v}})$  and select top- $K$  influential samples using majority vote aggregation (Eq. 19):
- 14:  $\mathcal{S} = \mathcal{S}_{\text{vote}}(K) = \{\mathbf{v}_n \mid \mathbf{v}_n \text{ in top-}K \text{ by MajVote}\}$
- 15: **return**  $\mathcal{S}$

---

**Model-agnostic scope.** Our attribution only requires per-example gradients under matched  $(t_{\text{fix}}, \epsilon_{\text{fix}})$ , and therefore applies to both diffusion and flow-matching objectives. The score reduces to a gradient inner product under a fixed preconditioner; the generator architecture affects gradient statistics but not the definition of influence. In practice, replacing the denoiser or velocity field leaves the weighting and aggregation unchanged.

**Algorithm Summary.** For completeness, Algorithm 1 summarizes the full MOTIVE pipeline, detailing the computation of motion-weighted gradients, projection into low-dimensional space, and the subsequent influence-based ranking and selection of training clips.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

## H ADDITIONAL EXPERIMENT DETAILS

### H.1 HYPERPARAMETER SETTINGS

For reproducibility, we document the hyperparameters used throughout attribution, subset selection, and fine-tuning. Where values were not explicitly tuned, we adopted defaults from DiffSynth-Studio and the official Wan repo for Wan models, and Diffusion-Pipe for LTX models.

**Attribution.** Motion-aware influence estimation is computed at a single fixed timestep corresponding to the midpoint of the denoising trajectory, which strongly correlates with multi-timestep averaging. A shared Gaussian draw  $\epsilon_{\text{fix}} \sim \mathcal{N}(0, \mathbf{I})$  is used across all training–query pairs to reduce stochastic variance. Gradients are projected from dimension  $D = 1418996800$  to  $D' = 512$  using a Fastfood Johnson–Lindenstrauss projection  $\mathbf{P}$  selected via the search to balance performance and storage. Motion weights  $\mathbf{W}$  are computed from AllTracker flow magnitudes  $M_f$ , min–max normalized to  $[0, 1]$  with a small bias  $\zeta = 10^{-6}$ . All computations use bfloat16 precision for memory efficiency.

**Subset Selection & fine-tuning.** For any number of query points, we select top-10% data of the datasets. We finetune the backbone while freezing both the text encoder (Raffel et al., 2020) and the VAE. The input resolution is fixed to  $480 \times 832$  pixels. We use a learning rate of  $1 \times 10^{-5}$  and the AdamW optimizer (Loshchilov & Hutter, 2017) following the DiffSynth-Studio defaults. We train the models for 1 epoch, repeating the dataset 50 times.

**Evaluation.** The test set consists of the same 10 motion categories, but with different visual appearances, compared with the query set. We provide the prompt samples below.

#### Samples from Query Set

We illustrate representative prompts from our query set used to generate query videos with Veo-3.

[leftmargin=\*, itemsep=2pt, parsep=0pt]

- **compress:** “A slice of sandwich bread flattened by a flat metal plate, steady camera, soft studio lighting, plain backdrop; emphasize air pockets collapsing.”
- **bounce:** “A ping-pong ball bouncing on a white table, steady side camera, neutral light, seamless backdrop; emphasize consistent bounce height and timing.”
- **roll:** “A spool of thread rolling from left to right, close-up static camera, bright studio light; highlight axle rotation and smooth travel.”
- **explode:** “A single balloon bursting into fragments, captured in high-speed slow motion with a fixed camera, bright even lighting, seamless background; emphasize outward debris and air release.”
- **float:** “A foam cube floating on the surface of water, static overhead camera, bright light, clean tank; emphasize buoyancy and slight rocking.”

#### Samples from Test Set

We illustrate representative prompts from our test set that our fine-tuned models use to generate test videos.

[leftmargin=\*, itemsep=2pt, parsep=0pt]

- **compress:** “A rubber ball being compressed under a flat press, filmed with a stationary camera. Bright, shadow-free lighting and a clean background emphasize the deformation as it flattens.”
- **bounce:** “A basketball bouncing vertically on a wooden court plank, unmoving camera, balanced indoor lighting, plain wall background; clearly show deformation at impact.”
- **roll:** “A bike tire rolling freely on a stand, static side camera, indoor neutral light; show uniform rotation without wobble.”
- **explode:** “A fragile glass ornament breaking apart mid-air, fixed camera, bright controlled lighting, plain backdrop; capture shards and reflections crisply.”
- **float:** “A green leaf floating gently on perfectly still water in a transparent tank, fixed top-down camera, bright even lighting; emphasize surface tension ripples.”

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079



Figure 6: **Illustration of motion query set.** We generate near-realistic video queries with Veo-3 across ten motion categories. Each category contains five query videos synthesized with controlled prompts and manually screened for clarity and physical plausibility.

## H.2 DETAILS ON MOTION QUERY DATA

A small, controlled set of query videos is constructed to isolate specific motion primitives while minimizing confounding factors (e.g., textured backgrounds, uncontrolled camera motion). Such clean and consistent clips are challenging to obtain from natural data sources. To address this, we synthesize the query set using Veo-3 (Google DeepMind, 2025) and apply a strict post-generation screening for physical plausibility and generation realism. We target ten motion types: *compress*, *bounce*, *roll*, *explode*, *float*, *free fall*, *slide*, *spin*, *stretch*, *swing*. For each category, we retain 5 query samples, yielding a total of 50 queries. This scale provides adequate coverage of the motion taxonomy used in our evaluations while maintaining tractable attribution computation. We further provide a few examples of the generation prompts and the generated video query set in Fig. 6.

**Rationale for synthetic queries.** The query set is not used as training data; instead, it specifies targets for attribution and for multi-query aggregation. Synthetic generation offers controllability that is difficult to achieve at scale with web videos. This design yields near-realistic yet standardized stimuli aligned with our goal of probing motion-specific influence.

# I DISCUSSION

## I.1 RUNTIME

All training runs are conducted on 4-8 NVIDIA A100 GPUs. We provide a detailed runtime breakdown from our experiments on 10k training samples with **Wan2.1-T2V-1.3B** model to address scalability concerns. The key insight is that the dominant cost of our pipeline is computing per-training-sample gradients, which is performed once and then can be reused for all subsequent queries. Each training clip’s gradient is projected into a compact 512-dimensional vector, and adding a new query requires only (i) a single backward pass to obtain its own projected 512-dimensional gradient and (ii) computing cosine similarity between the query vector and stored training vectors, which is exceptionally lightweight (on the order of seconds). Thus, the computational burden does not scale with the number of queries but only with the size of the training set. As shown in Tab. 5, the dominant cost is the one-time training data gradient computation for 10k samples (~ 150 hours on 1 A100), which is amortized across all queries. Once computed, adding a new query requires only ~ 54 seconds (gradient computation) + 46ms (influence computation) = ~54 seconds total. The training data gradient computation is embarrassingly parallel and can be reduced to ~2.3 hours with 64 GPUs.

Table 5: **Runtime Breakdown.** Detailed computational complexity and runtime for each component of our motion attribution framework on 10k training samples with **Wan2.1-T2V-1.3B** model.

Component	Complexity	Runtime	Notes
Gradient computation	$\mathcal{O}(B)$ per sample	Query: $\sim 54$ seconds Training: $\sim 150$ hours	1 A100 GPU; Single forward+backward pass; training is dominant cost but amortized over all queries; embarrassingly parallel (with 64 GPUs, $\sim 2.3$ hours)
Projection	$\mathcal{O}( \mathcal{D}  \cdot D' \log D')$	$\sim 1.97$ seconds per sample	$D' = 512$
Influence computation	$\mathcal{O}( \mathcal{D}  \cdot D')$	$\sim 46$ milliseconds per query	1 query $\times$ 10k training samples
Majority-vote aggregation	$\mathcal{O}( \mathcal{D}  \cdot q)$	$\sim 139$ milliseconds	50 queries $\times$ 10k samples

Table 6: **Runtime Comparison with Baselines.** Total computational time required for processing 10k training samples with **Wan2.1-T2V-1.3B** model on a single GPU across different data selection methods.

Method	Random	Motion Magnitude	Optical Flow	V-JEPA	Ours
Total for 10k (1 GPU)	$< 1$ second	$\sim 5.5$ hours	$\sim 5.7$ hours	$\sim 3$ hours	$\sim 150$ hours

**Runtime comparison with baselines.** We compare the computational cost of our method with the baseline approaches for processing 10k training samples on a single GPU (Table 6). While our method requires more upfront computation than the baseline approaches, this cost is amortized across all queries, and the computed gradients can be reused for multiple selection queries, making it practical for large-scale data curation scenarios.

## I.2 LIMITATIONS

Gradient-based attribution is computationally expensive, requiring high upfront cost for per-sample gradient computation (see §I.1), though this cost is amortized across queries. Our analysis treats each video as a whole unit, thereby avoiding collapsing motion into frame-level appearance, but it risks overlooking the fact that only certain segments may carry motion-relevant information. Highly informative intervals can be diluted when averaged with static or redundant portions of the same clip. This suggests an open direction toward finer-grained attribution at the motion segment or motion event level, which could reveal more precise insights into how different phases of a trajectory shape motion learning. Another limitation is that our motion masks may overemphasize camera-only motion; we detect this by spatial uniformity of  $\mathbf{W}$  and down-weight such clips, but a full disentanglement of ego and object motion remains future work.

Additionally, our framework does not explicitly account for classifier-free guidance (CFG), which is widely used in practice to steer video generation but introduces discrepancies between training-time attribution and inference-time dynamics. As a result, our influence estimates may not fully capture how guidance alters motion behavior. In addition, while attribution-driven fine-tuning improves targeted motion quality, it may introduce trade-offs with base model capabilities. This raises the need for future work on balancing targeted motion adaptation with the preservation of broader generative capabilities.

## I.3 FUTURE DIRECTIONS

**Improving the attribution pipeline.** Several directions can strengthen our motion attribution methodology. *Tracker-robust motion saliency:* replace or ensemble AllTracker with alternative estimators and use its confidence/visibility channels to weight masks. *Self-generated video queries:* use model-generated videos as queries to trace problematic motion patterns (e.g., unrealistic physics) back to training data, enabling iterative diagnostics and targeted motion improvement.

---

1134 **Scaling data curation.** *Closed-loop data curation:* move from one-shot ranking to active selection:  
1135 iteratively attribute, finetune, and re-attribute, or replace simple majority voting with learned query  
1136 weights. *Sophisticated finetuning:* move to more sophisticated finetuning setups, such as multi-student  
1137 distillation (Song et al., 2024).

1138 **Extending to new domains.** *Other modalities:* extend our methodology to other modalities, in-  
1139 cluding world models (Zhu et al., 2025), audio (Evans et al., 2025; Richter-Powell et al., 2025), or  
1140 video+audio (Wiedemer et al., 2025).

1141 **Safety and governance.** Use negative-influence filtering to suppress undesirable or unsafe dynamics,  
1142 document curator choices, and audit motion behaviors exposed by our framework.  
1143

1144

## 1145 J VISUALIZATION

1146

### 1147 J.1 MOTION VISUALIZATION

1148

1149 To provide intuition for the behavior of our motion-weighted loss, we visualize the motion magnitude  
1150 as an overlay. We compute per-pixel motion magnitude using optical flow and apply motion-based  
1151 weighting that preserves the appearance of dynamic regions while attenuating static backgrounds.  
1152 This motion overlay directly illustrates the spatial weighting applied by our motion loss during  
1153 training.

1154 Fig. 7 presents representative frames from our dataset, comparing original frames with their cor-  
1155 responding motion overlays for seven distinct video samples. These visualizations show that the  
1156 motion-weighted loss preferentially emphasizes dynamic content while down-weighting static scene  
1157 elements.

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

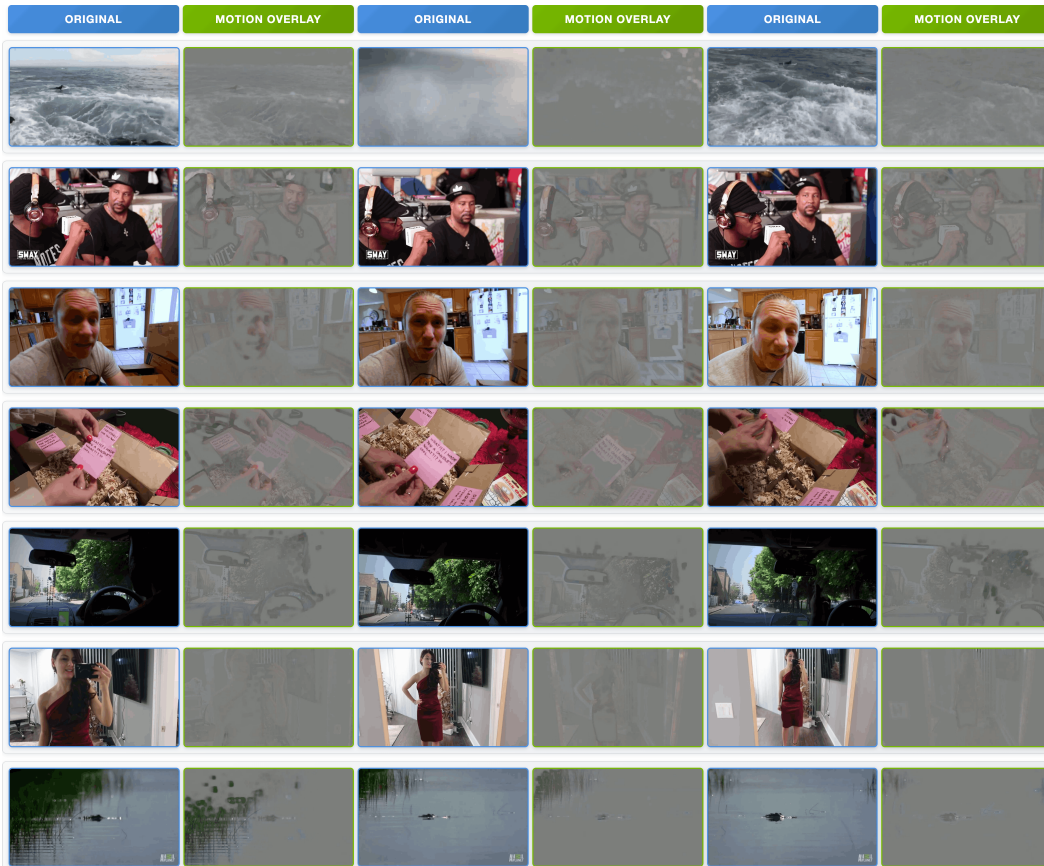


Figure 7: **Motion Overlay Visualization.** Comparison of original frames and motion overlays for seven video samples across three time points (early, middle, late). The motion overlay demonstrates the spatial weighting of our motion loss: dynamic regions remain visible, while static backgrounds are attenuated to neutral gray. *Takeaway:* This provides heuristic intuition into what information our motion attribution focuses on: the information in grayer regions, which lack motion, is down-weighted by our method.