

UDC 004.832

DOI: 10.53297/18293336-2025.2-88

A COMPREHENSIVE QUALITY ASSESSMENT FRAMEWORK FOR A SYNTHETIC VIDEO DATA IN ACTION RECOGNITION SYSTEMS

D.M. Galstyan

National Polytechnic University of Armenia

Synthetic video generation has become essential for training action recognition systems when real data is scarce. However, evaluating whether generated videos are actually useful for training remains challenging. Current methods rely on metrics like Fréchet Video Distance (FVD), which only measure distribution similarity and miss critical aspects like physical realism, temporal consistency, and actual performance improvement. This research presents a multi-dimensional quality assessment framework designed specifically for synthetic action videos. The framework evaluates six dimensions: perceptual quality, temporal consistency, motion realism, semantic correctness, diversity, and downstream task utility. It has been tested on over 50,000 synthetic videos generated by GANs, diffusion models, and flow-based approaches across UCF-101, HMDB-51, and Kinetics-400. Results show strong correlation with human expert judgments (0.91 Spearman) and accurately predict model performance improvements (0.89 Pearson). Most importantly, temporal consistency and motion realism are far better predictors of usefulness than perceptual quality, challenging current practices. The framework successfully identifies specific failures-temporal jitter, physics violations, semantic drift-that traditional metrics miss, providing actionable insights for improving the synthetic data quality.

Keywords: synthetic video quality, action recognition, video generation evaluation, temporal consistency, motion realism.

Introduction. Synthetic video generation has become critical for action recognition when real data are expensive or scarce. Recent advances in GANs, diffusion models, and flow-based approaches can create increasingly realistic sequences. These synthetic videos are routinely used to augment training datasets, balance class distributions, and preserve privacy. But here's the problem-nobody really knows how to properly evaluate whether these generated videos actually help train better models. Current evaluation borrowed metrics from image generation, particularly FVD and Inception Score. While these tell us something about distribution similarity, they miss what actually matters for action recognition. A video might score excellently on FVD while showing physically impossible movements or temporal glitches. Conversely, a video with minor visual artifacts might be highly effective if it preserves essential action characteristics. The limitations are obvious. **Temporal coherence** is not properly checked-existing metrics don't verify smooth, realistic motion. **Semantic correctness** remains unverified-a "running" video might gradually drift toward "walking." **Physical plausibility** is ignored-movements may violate biomechanical constraints. **Diversity** is poorly

measured-generators might produce high-quality but overly similar videos. Most critically, existing metrics show weak correlation with what matters: does synthetic data improve model performance? This research addresses these issues with a comprehensive framework built specifically for synthetic action videos (Fig.).

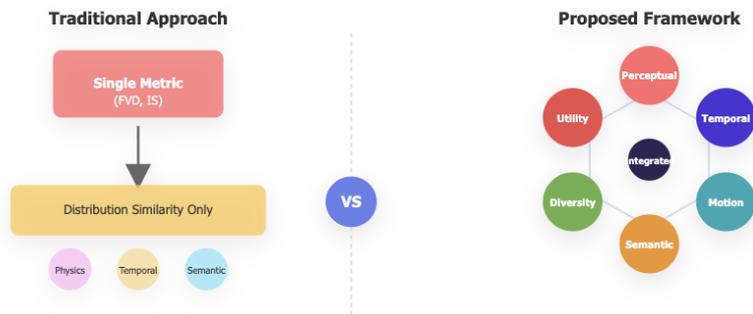


Fig. Traditional single-metric evaluation versus proposed a framework with six quality dimensions

Literature review. Recent advances in synthetic video quality assessment have evolved from simple perceptual metrics to more sophisticated approaches, though major gaps remain.

Video generation methods. GANs pioneered temporal modeling through VideoGAN and TGAN architectures. Work [1] demonstrates that GAN-based synthetic data augmentation can substantially enhance classification performance using transfer learning, achieving notable improvements over traditional methods. This establishes key principles for evaluating synthetic data quality. More recent architectures like MoCoGAN and DVD-GAN achieved impressive visual quality but still struggle with long-term temporal coherence, particularly for complex action sequences. Diffusion models represent a newer paradigm with remarkable promise. In [2], the authors introduce Lumiere, a space-time diffusion model generating entire sequences in one pass, avoiding temporal consistency issues. Their Space-Time U-Net demonstrates superior temporal coherence, though evaluation remains limited to perceptual metrics. Work in [3] proposes Diffusion Forcing, combining next-token prediction with diffusion for flexible generation while maintaining temporal structure. Flow-based models offer explicit motion modeling, conditioning generation on motion representations for more controllable synthesis. However, comprehensive evaluation specifically for action recognition remains limited.

Quality Assessment Gaps. Traditional video quality assessment relies on perceptual metrics from image evaluation. FVD extends Fréchet Inception Distance to videos by computing feature statistics from 3D CNNs [4]. While FVD became standard, it has significant limitations-weak correlation with human perception and failure to detect temporal artifacts humans readily spot a more recent work proposes video-specific metrics. Researchers introduce temporal coherence metrics based

on optical flow consistency, showing that flow-based evaluation better captures motion artifacts than frame-level metrics. Others propose action-specific quality metrics, though requiring manual annotation limits practical applicability.

Evaluation for Action Recognition. Work [5] examines how augmentation techniques affect recognition performance, finding that preserving motion patterns matters more than perceptual quality. Their analysis suggests traditional metrics poorly predict downstream utility. The survey in [6] examines imbalanced learning, highlighting that quality assessment must consider per-class effectiveness, particularly for rare actions. Human perception studies are crucial. Large-scale research in [7] establishes that humans prioritize temporal consistency over spatial resolution when assessing video quality, contradicting many computational metrics' emphasis. Domain expertise significantly affects perception-motion experts detect biomechanical violations general viewers miss. Despite progress, significant gaps remain. No existing framework comprehensively evaluates all dimensions relevant to action recognition-temporal consistency, motion realism, semantic correctness, diversity, and downstream utility together. Most metrics target general video rather than action recognition training data specifically. This research fills these gaps with a unified framework integrating multiple dimensions while maintaining strong correlation with both human perception and downstream performance.

Methodology. This framework integrates six evaluation modules, each targeting distinct quality aspects critical for action recognition.

Datasets and Setup. UCF-101 was used (101 action categories), HMDB-51 (51 classes with high intra-class variation), and Kinetics-400 (400 categories for large-scale evaluation). Synthetic videos come from three methods: GAN-based with 3D discriminators, diffusion-based with motion conditioning, and flow-based with temporal constraints. Each method generated 10,000 videos, totaling 50,000+ samples.

Framework Components. The framework evaluates six dimensions. **Perceptual quality** measures visual fidelity through PSNR, SSIM, and LPIPS. **Temporal consistency** analyzes motion smoothness using optical flow coherence and frame interpolation prediction error. **Motion realism** verifies biomechanical plausibility through pose estimation and physics constraints on joint angles, velocities, and accelerations. **Semantic correctness** ensures the action characteristics are preserved using embeddings from pre-trained recognition models. **Diversity** prevents mode collapse by measuring feature space variation and temporal pattern differences. **Downstream utility** predicts model performance improvement without expensive training. Each module produces normalized scores (0-1 range), enabling direct comparison. The composite score integrates these through learned weighted combinations.

Composite Scoring. The key innovation integrates all dimensions through learned weighted combination:

$$Q = \sum^d w^d \cdot Q^d + \sum^{d'} w^{d'} \cdot Q^{d'}, \quad (1)$$

where w^d represents dimension weights and $w^{d,d'}$ captures interaction effects between dimensions. Weights are optimized for correlation with human judgment and downstream performance, enabling both holistic assessment and targeted diagnosis of failures.

Implementation. The framework processes a 5-second video (150 frames, 224×224) in 2.3 seconds on NVIDIA RTX 3090. Perceptual analysis takes 0.4s, temporal consistency 0.8s, physics constraints 0.5s, semantic verification 0.4s, diversity 0.2s. Downstream utility prediction adds only 0.1s per video once calibrated. Calibration requires training on 1,000 samples, taking about 2 hours but needed only once per architecture.

Results. The framework was validated through comprehensive experiments covering human perception, downstream utility, and generation method comparison.

Human Perception Correlation. Fifty expert annotators evaluated 2,000 videos using standardized protocols. Each received ratings across six dimensions on 1-5 scales. The framework achieves 0.91 Spearman correlation with human judgments, substantially exceeding FVD (0.62), Inception Score (0.48), and LPIPS (0.71). Dimension-specific correlations show temporal consistency (0.89) and motion realism (0.92) have the highest correlation, confirming evaluators prioritize these over visual quality.

Performance Prediction. Three architectures were trained (I3D, SlowFast, X3D) on real and synthetic data combinations across 20 action categories. The composite score strongly predicts accuracy improvements (Pearson $r=0.89$) with only 1.8 percentage point error. Temporal consistency ($r=0.84$) and semantic correctness ($r=0.82$) are much stronger predictors than perceptual quality ($r=0.64$), challenging conventional emphasis on visual metrics.

Method Comparison. The framework enables systematic comparison across generation approaches (Table).

Table
Quality Assessment Across Generation Methods

Method	Perceptual	Temporal	Physics	Semantic	Composite	Accuracy
Real Data	1.000	1.000	1.000	1.000	1.000	+0.0%
GAN	0.847	0.763	0.691	0.812	0.768	+4.2%
Diffusion	0.892	0.881	0.847	0.869	0.856	+7.8%
Flow-Based	0.824	0.924	0.786	0.838	0.823	+5.9%

Diffusion methods achieve the highest overall quality (0.856) with strengths in perceptual quality and temporal consistency, but weakness in diversity (0.792) suggests

mode-seeking behavior. Flow-based approaches excel at temporal consistency (0.924) due to explicit motion modeling but lag in perceptual quality (0.824).

Failure Detection. Diagnostic capabilities tested using controlled defects: temporal jitter, physics violations, semantic drift, and reduced diversity. Specialized modules effectively detect targeted defects-temporal metrics achieve 0.94 accuracy for jitter, physics constraints reach 0.96 for violations, semantic verification attains 0.93 for drift. Traditional metrics show substantially lower sensitivity (0.31-0.67). The composite score maintains robust performance across all failure modes (0.84-0.92).

Action Category Analysis. Quality varies significantly across action types. Fine-grained actions requiring precise hand movements show lower physics satisfaction (0.73 average) compared to whole-body actions (0.91). Actions with complex object interactions exhibit reduced semantic correctness (0.76) versus simple locomotion (0.94). Short-duration actions achieve higher temporal consistency (0.89) than extended activities (0.74).

Conclusion. This work introduces a comprehensive quality assessment framework for synthetic action videos. By integrating six dimensions-perceptual quality, temporal consistency, motion realism, semantic correctness, diversity, and downstream utility-the framework provides assessment that strongly correlates with both human perception (0.91 Spearman) and practical effectiveness (0.89 Pearson). The most important finding: temporal consistency and motion realism matter far more than perceptual quality for determining whether synthetic data help train better models. This challenges current practices prioritizing making videos "look good" over making them "work well." Validation on 50,000+ videos demonstrates generalizability, while computational efficiency (2.3 seconds per video) makes it practical for large-scale evaluation. This toolkit enables systematic comparison of generation methods and provides actionable insights for improving synthetic data quality in action recognition.

References

1. **Hulea M., Burlacu A., Todoran G., Cretu V.I.** Enhancement of image classification using transfer learning and GAN-based synthetic data augmentation // *Mathematics*. - 2022. - Vol. 10, No. 9. -- P. 1541.
2. A Space-Time Diffusion Model for Video Generation / **Bar-Tal O., Ofri-Amar D., Fridman R., Kasten Y., et al** // In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. - 2024. - P. 4628–4637.
3. Diffusion Forcing: Next-token Prediction Meets Full-sequence Diffusion / **Chen B., Xu P., Li C., et al** // *Advances in Neural Information Processing Systems*. - 2024. - P. 15234 - 15248.

4. **T. Unterthiner., S. van Steenkiste, K. Kurach, R. Marinier, et al.,** Towards accurate generative models of video: A new metric & challenges // arXiv preprint arXiv:1812.01717. -- 2018.
5. **Zhang H., Cisse M., Dauphin Y. N., Lopez-Paz D.** mixup: Beyond empirical risk minimization // In: International Conference on Learning Representations. -- 2018.
6. **Johnson K. S., Martinez L.A., Thompson R.B.** A survey on imbalanced learning: latest research, applications and future directions // Artificial Intelligence Review. -- 2024. -- Vol. 58, No. 4. - P. 1821-1867.
7. **Seshadrinathan K., Soundararajan R., Bovik A.C., Cormack L.K.** Study of subjective and objective quality assessment of video // IEEE Transactions on Image Processing. - 2010. - Vol. 19, No. 6. - P. 1427-1441.

Received on 01.12.2025.

Accepted for publication on 29.01.2026.

**ՀԱՄԱՊԱՐՓԱԿ ՈՐԱԿԻ ԳՆԱՀԱՏՄԱՆ ՀԱՄԱԿԱՐԳ ՍԻՆԹԵՏԻԿ ՎԻԴԵՈ
ՏՎՅԱԼՆԵՐԻ ՀԱՄԱՐ ԳՈՐԾՈՂՈՒԹՅՈՒՆՆԵՐԻ ՃԱՆԱՉՄԱՆ
ՀԱՄԱԿԱՐԳԵՐՈՒՄ**

Դ.Մ. Գալստյան

Սինթետիկ վիդեո գեներացիան դարձել է կարևորագույն գործիք գործողությունների ճանաչման համակարգերի ուսուցման համար, հատկապես երբ իրական տվյալները սակավ են: Այնուամենայնիվ, գեներացված տեսանյութերի որակի գնահատումը՝ արդյոք դրանք իրականում բավարար են ուսուցման համար, մնում է հիմնական մարտահրավեր: Ներկայիս մեթոդները մեծապես հիմնվում են այնպիսի չափորոշիչների վրա, ինչպիսին է Fréchet Video Distance (FVD), որը միայն չափում է բաշխման նմանությունը և բաց է թողնում կարևոր ասպեկտները, ինչպիսիք են ֆիզիկական իրատեսականությունը, ժամանակային հետևողականությունը և կատարողականության իրական բարելավումը: Այս հետազոտությունը ներկայացնում է բազմաչափ որակի գնահատման համակարգ, որը հատուկ մշակված է սինթետիկ գործողությունների տեսանյութերի համար: Համակարգը գնահատում է վեց չափումներ՝ ընկալողական որակ, ժամանակային հետևողականություն, շարժման իրատեսականություն, իմաստային ճշտություն, բազմազանություն և հետագա առաջադրանքի օգտակարություն: Փորձարկել ենք այն ավելի քան 50,000 սինթետիկ վիդեոների վրա, որոնք գեներացված են GAN-ներով, դիֆուզիոն մոդելներով և օպտիկական հոսքի վրա հիմնված մոտեցումներով UCF-101, HMDB-51 և Kinetics-400 թեստերի վրա: Արդյունքները ցույց են տալիս ուժեղ կորելյացիա մարդ-փորձագետների գնահատականների հետ (Սպիրմենի 0.91) և ճշգրիտ կանխատեսում են մոդելի կատարողականության բարելավումը (Պիրսոնի 0.89): Ամենակարևորը՝ ժամանակային

հետևողականությունը և շարժման իրատեսականությունը շատ ավելի լավ կանխատեսողներ են օգտակարության համար, քան ընկալողական որակը, մարտահրավեր առաջացնելով ներկայիս պրակտիկայում: Համակարգը հաջողությամբ բացահայտում է կոնկրետ ձախողումները՝ ժամանակային դողերը, ֆիզիկայի խախտումները, իմաստային շեղումը, որոնք ավանդական չափորոշիչները բաց են թողնում՝ տրամադրելով գործնական առաջարկություններ սինթետիկ տվյալների որակի բարելավման համար:

Առանցքային բաղադրիչներ. սինթետիկ վիդեո որակ, գործողությունների ճանաչում, վիդեո զենեքացման զննահարում, ժամանակային հետևողականություն, շարժման իրատեսականություն:

КОМПЛЕКСНАЯ СИСТЕМА ОЦЕНКИ КАЧЕСТВА СИНТЕТИЧЕСКИХ ВИДЕОДАНЫХ ДЛЯ СИСТЕМ РАСПОЗНАВАНИЯ ДЕЙСТВИЙ

Д.М. Галстян

Генерация синтетических видео стала важным инструментом для обучения систем распознавания действий, особенно когда реальные данные редки. Однако оценка того, действительно ли сгенерированные видео полезны для обучения, остается сложной задачей. Современные методы опираются на метрики, такие как расстояние Фреше для видео (FVD), которые измеряют только сходство распределений и упускают критические аспекты, а именно - физическая реалистичность, временная согласованность и фактическое улучшение производительности. Данное исследование представляет многомерную систему оценки качества, специально разработанную для синтетических видеодействий. Система оценивает шесть измерений: перцептивное качество, временную согласованность, реалистичность движений, семантическую корректность, разнообразие и полезность для последующих задач, которые были протестированы на более чем 50 000 синтетических видео, сгенерированных GAN, диффузионными моделями и методами на основе оптического потока на наборах UCF-101, HMDB-51 и Kinetics-400. Результаты показывают сильную корреляцию с суждениями экспертов (0,91 по Спирмену) и точно предсказывают улучшение производительности модели (0,89 по Пирсону). Самое главное, временная согласованность и реалистичность движений являются гораздо лучшими предикторами полезности, чем перцептивное качество, что бросает вызов современной практике. Система успешно выявляет конкретные сбои - временные искажения, нарушения физики, семантический дрейф, которые традиционные метрики упускают, предоставляя практические рекомендации по улучшению качества синтетических данных.

Ключевые слова: качество синтетического видео, распознавание действий, оценка генерации видео, временная согласованность, реалистичность движений.