

TEMPORAL REVERSAL ASYMMETRY: A PHYSICS-INSPIRED METRIC FOR EVALUATING WORLD MODELS

Kanpat Vesesook*

Department of Computer Science
Brown University
Providence, RI 02912, USA
kanpat_vesessook@brown.edu

Kevin Yang*

Department of Computer Science
Brown University
Providence, RI 02912, USA
kevin_c_yang@brown.edu

ABSTRACT

World models that understand physics should recognize the thermodynamic arrow of time: some processes look physically plausible when reversed, while others do not. We propose Temporal Reversal Asymmetry (TRA), comparing prediction loss on forward versus time-reversed videos. We evaluated four video models on 380 physics simulations and found that TRA for V-JEPA2 varies continuously with dissipation strength, peaking at intermediate values (restitution 0.5, damping 2.0) where energy loss remains visible throughout the video, while approaching zero for both genuinely reversible processes and motion that stops too quickly. This non-monotonic relationship suggests V-JEPA2 has learned to recognize irreversible processes by how they look as opposed to memorizing categories. In contrast, VideoMAE V2 shows inverted TRA (negative asymmetry), while MVD and HierA show near-zero TRA regardless of physics type, suggesting that latent prediction captures temporal causality more faithfully than pixel reconstruction or distillation. This work offers a training-free, physics-grounded probe that complements existing world model benchmarks.

1 INTRODUCTION

The second law of thermodynamics establishes an arrow of time: entropy tends to increase, making dissipative physical processes irreversible. A glass can shatter but spontaneously reassembling is vanishingly improbable. In contrast, idealized conservative systems like a frictionless pendulum are time-symmetric. A video of a pendulum played backwards remains physically plausible.

World models trained on natural video should, in principle, learn this distinction. If a model truly understands physics, it should find time-reversed dissipative processes surprising or difficult to predict, while treating low-dissipation processes symmetrically. This intuition motivates our proposed metric: Temporal Reversal Asymmetry (TRA).

Prior work has evaluated world models on intuitive physics benchmarks like IntPhys (Riochet et al., 2020) and PHYRE (Bakhtin et al., 2019), which test whether models distinguish possible from impossible scenarios. TRA offers a complementary approach: rather than asking whether a model knows something is impossible, we ask whether its prediction errors reflect thermodynamic asymmetry.

Our contributions are:

1. We introduce TRA, a training-free metric that probes whether world models respect the thermodynamic arrow of time (Section 2).
2. We show V-JEPA2 exhibits physically appropriate TRA while VideoMAE V2 shows an inverted pattern, and that TRA tracks visible irreversibility by peaking at intermediate physical parameters (Section 3).

*Equal contribution.

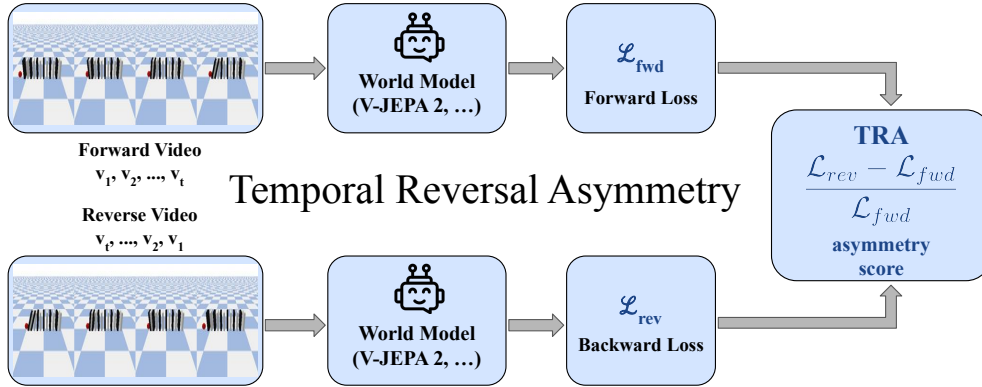


Figure 1: TRA computation: forward and time-reversed videos are passed through a world model; we compare their prediction losses.

2 METHOD

2.1 TEMPORAL REVERSAL ASYMMETRY

Given a video $\mathbf{v} = (v_1, v_2, \dots, v_T)$ and its time-reversed version $\tilde{\mathbf{v}} = (v_T, v_{T-1}, \dots, v_1)$, we define TRA as:

$$\text{TRA}(\mathbf{v}) = \frac{\mathcal{L}(\tilde{\mathbf{v}}) - \mathcal{L}(\mathbf{v})}{\mathcal{L}(\mathbf{v})} \quad (1)$$

where $\mathcal{L}(\cdot)$ is the model’s prediction loss. A TRA of zero indicates the model treats forward and reversed videos identically. Positive TRA means reversed videos are harder to predict; negative TRA means forward videos are harder.

For a world model with physically grounded temporal understanding, we expect $\text{TRA} \approx 0$ for low-dissipation processes (elastic collisions, pendulums) and $\text{TRA} > 0$ for dissipative processes (falling objects, toppling dominos). Figure 1 illustrates this computation.

2.2 MODELS

We evaluate four self-supervised video models: **V-JEPA2** (Assran et al., 2025), a Joint Embedding Predictive Architecture (ViT-Large, 303M params) that predicts latent representations with 3D rotary position embeddings; **VideoMAE V2** (Wang et al., 2023a), a masked autoencoder (ViT-Base, 94M params) that reconstructs pixel patches; **MVD** (Wang et al., 2023b), a teacher-student distillation model (ViT-Base, 94M params) that learns to match features from a pretrained teacher; and **Hiera** (Ryali et al., 2023), a hierarchical vision transformer (79M params) with multi-scale mask units. For all models, we compute prediction loss using sliding windows: given C context frames, the model predicts representations/pixels for remaining frames.

2.3 VIDEO DATASET

We generate 380 physics simulations using PyBullet (Coumans & Bai, 2016), shown in Figure 2: 160 videos across four scenarios (*bouncing ball*, *pendulum* = low-dissipation; *falling objects*, *dominos* = dissipative) plus 220 videos with varying restitution/damping for continuous probing. Each video: 48 frames at 256×256 . Bouncing ball uses restitution 0.9 (near-elastic); pendulum uses zero damping.

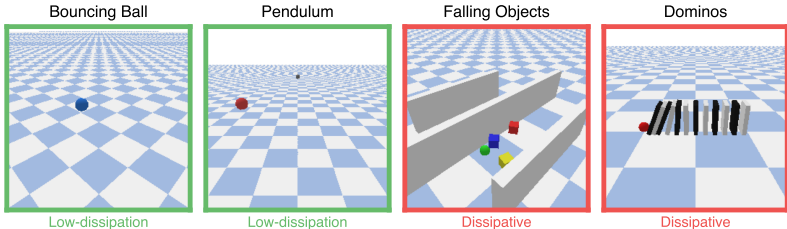


Figure 2: Physics simulation scenarios. Low-dissipation processes (green) look plausible when time-reversed; dissipative processes (red) do not; reversed falling objects would spontaneously jump upward.

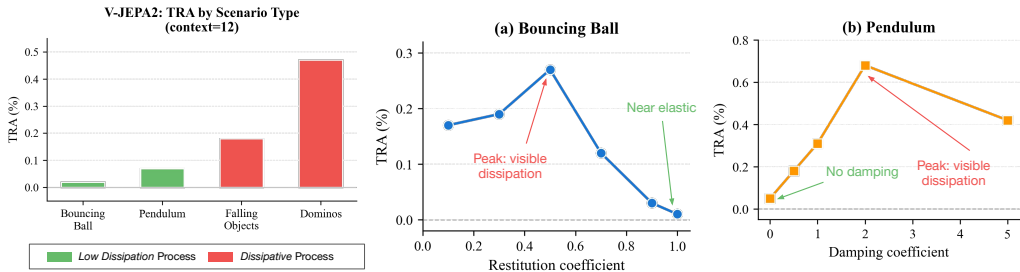


Figure 3: V-JEPA2 TRA reveals graded sensitivity to dissipation. **Middle + Right:** TRA varies continuously with physical parameters, peaking where energy loss remains visible throughout the video duration. **Left:** Discrete scenarios validate this pattern—low-dissipation processes (green) cluster near zero while dissipative ones (red) show elevated TRA.

3 RESULTS

3.1 TRA TRACKS VISIBLE IRREVERSIBILITY

To test whether TRA captures graded irreversibility, we varied physical parameters: restitution (0.1–1.0) for bouncing balls and damping (0.0–5.0) for pendulums. Figure 3 (right) reveals a non-monotonic pattern: TRA peaks at intermediate values. The bouncing ball shows maximum TRA at restitution 0.5 (+0.27%), approaching zero at both restitution 1.0 (elastic, +0.01%) and 0.1 (+0.17%). The pendulum peaks at damping 2.0 (+0.68%).

This reflects that TRA measures *visible* irreversibility: at high dissipation, motion stops quickly leaving a symmetric stationary scene; at low dissipation, the process is genuinely reversible. At restitution 0.5, the ball bounces throughout all 48 frames with visibly decreasing amplitude; time-reversing this creates an implausible video of a ball spontaneously gaining energy. V-JEPA2 responds to visual evidence of ongoing energy dissipation rather than abstract physical categories.

3.2 V-JEPA2 EXHIBITS PHYSICALLY APPROPRIATE TRA

These findings validate across discrete physics scenarios (Figure 3, left; Table 1). V-JEPA2 shows a clear distinction between low-dissipation and dissipative processes: low-dissipation processes have near-zero TRA (+0.04%, not significant), while dissipative processes show significantly elevated TRA (+0.32%, $p < 0.001$).

The bouncing ball (restitution 0.9, near the elastic extreme) shows essentially zero TRA (+0.02%, $p = 0.33$). In contrast, dominos shows the strongest signal (+0.47%), and falling objects shows robust positive TRA (+0.18%)— both using lower restitution values in the regime where TRA is elevated. A randomly initialized V-JEPA2 shows negligible TRA for all conditions ($< 0.01\%$), confirming that meaningful TRA sensitivity emerges from learned representations rather than model architecture. The effect is consistent across context lengths 6–12.

Table 1: Model comparison at context=8. V-JEPA2 shows physically appropriate TRA; VideoMAE V2 is inverted; MVD and Hiera show near-zero TRA; random baseline confirms learning is required.

Model	Objective	Low-dissip.	Dissipative	Diff.	<i>p</i> -value
V-JEPA2	Latent pred.	+0.03%	+0.22%	+0.20%	<0.001***
VideoMAE V2	Pixel recon.	-0.07%	-0.28%	-0.21%	<0.001***
MVD	Distillation	+0.00%	+0.00%	+0.00%	0.74 ns
Hiera	Hier. MAE	+0.19%	-0.09%	-0.29%	0.22 ns
Random	None	≈0%	≈0%	≈0%	0.56 ns

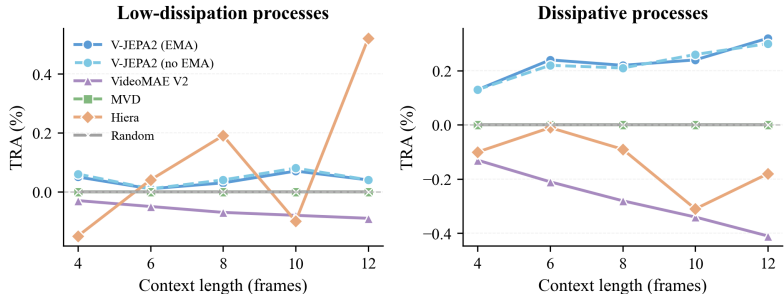


Figure 4: TRA across context lengths for all models. V-JEPA2 (both EMA and no-EMA variants) shows consistent positive TRA for dissipative processes; VideoMAE V2 shows inverted (negative) TRA; MVD and Random show near-zero TRA; Hiera shows inconsistent results.

3.3 VIDEOMAE V2 SHOWS INVERTED TRA

Remarkably, VideoMAE V2 exhibits the opposite pattern from V-JEPA2. All TRA values are negative, meaning the model finds forward videos harder to predict than reversed ones. At context length 12, low-dissipation processes show TRA of -0.09% while dissipative processes show -0.41% ($p < 0.001$), the inverse of what physics would predict. This pattern strengthens with longer context (from -0.10% at context 4 to -0.40% at context 14), suggesting a systematic bias rather than noise.

4 DISCUSSION

Why does VideoMAE V2 show inverted TRA? We hypothesize pixel reconstruction encourages learning low-level statistical regularities (motion blur, lighting gradients) that are easier to extrapolate in reverse, without corresponding to physical understanding. JEPA’s latent prediction may filter out such artifacts.

Why do MVD and Hiera show near-zero TRA? MVD’s distillation objective focuses on matching teacher features rather than predicting future states. Hiera’s hierarchical architecture may be better suited for spatial than temporal patterns. Both show TRA indistinguishable from random baseline.

Limitations. Our study uses synthetic physics simulations. On Something-Something V2 (real human-performed actions), human action timing confounds TRA; however, temporally cropping to isolate physics-only segments recovers the expected pattern (Appendix L).

5 CONCLUSION

We introduced Temporal Reversal Asymmetry (TRA) as a metric for probing whether world models respect the thermodynamic arrow of time. Among four video models, only V-JEPA2 exhibits the predicted pattern; VideoMAE V2 shows an inverted pattern, while MVD and Hiera show no significant TRA. This suggests *training objective matters more than architecture* for learning physics: V-JEPA2’s advantage is not simply scale, but rather its latent prediction objective. TRA offers a training-free, physics-grounded evaluation approach for world models.

REFERENCES

- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba Komeili, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, Sergio Arnaud, Abha Gejji, Ada Martin, Francois Robert Hogan, Daniel Dugas, Piotr Bojanowski, Vasil Khali-dov, Patrick Labatut, Francisco Massa, Marc Szafraniec, Kapil Krishnakumar, Yong Li, Xiaodong Ma, Sarath Chandar, Franziska Meier, Yann LeCun, Michael Rabbat, and Nicolas Ballas. V-jepa 2: Self-supervised video models enable understanding, prediction and planning, 2025.
- Anton Bakhtin, Laurens van der Maaten, Justin Johnson, Laura Gustafson, and Ross Girshick. Phyre: A new benchmark for physical reasoning, 2019.
- Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning, 2016.
- Ronan Riochet, Mario Ynocente Castro, Mathieu Bernard, Adam Lerer, Rob Fergus, Véronique Izard, and Emmanuel Dupoux. Intphys: A framework and benchmark for visual intuitive physics reasoning, 2020.
- Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, Jitendra Malik, Yanghao Li, and Christoph Feichtenhofer. Hiera: A hierarchical vision transformer without the bells-and-whistles, 2023.
- Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking, 2023a.
- Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Lu Yuan, and Yu-Gang Jiang. Masked video distillation: Rethinking masked feature modeling for self-supervised video representation learning, 2023b.

APPENDIX

A IMPLEMENTATION DETAILS

V-JEPA2. We use the official ViT-Large checkpoint (303M parameters) with 3D rotary position embeddings (RoPE). Prediction loss is computed using the EMA target encoder with layer normalization applied to targets. We use a sliding window approach: window size of 16 frames with stride 2. For each window, we provide C context frames and compute MSE loss between predicted and target latent representations for the remaining frames.

VideoMAE V2. We use the ViT-Base checkpoint (94M parameters) pretrained on Kinetics-400 with decoder depth 4. We compute pixel reconstruction loss using the same sliding window approach. The model reconstructs masked patches given visible context patches.

MVD. We use the ViT-Base checkpoint (94M parameters) pretrained on Kinetics-400 with ViT-Base teacher. MVD uses teacher-student distillation where the student learns to predict teacher features for masked regions. We use the same sliding window approach as other models.

Hiera. We use the Hiera-Base checkpoint (79M parameters) pretrained with MAE on Kinetics-400. Hiera uses a hierarchical architecture with multi-scale mask units. We compute reconstruction loss using the same sliding window approach.

Random baseline. We initialize V-JEPA2 with random weights (no pretraining) to verify that TRA sensitivity requires learned representations.

B DATASET GENERATION

All videos are generated using PyBullet physics engine at 256×256 resolution with 48 frames per video.

Bouncing ball. A sphere is dropped onto a plane. Default restitution coefficient: 0.9 (near-elastic). For continuous probing, we vary restitution from 0.1 to 1.0.

Pendulum. A rigid body swings from a fixed pivot point. Default damping: 0.0 (no energy loss). For continuous probing, we vary damping from 0.0 to 5.0.

Falling objects. Multiple rigid bodies fall and collide, settling into a pile. Restitution: 0.3.

Dominos. A row of rectangular blocks topples sequentially after the first is pushed. Restitution: 0.2.

C V-JEPA2 RESULTS ACROSS CONTEXT LENGTHS

Table 2 shows V-JEPA2 TRA across all tested context lengths. The effect is robust from context 6–12, with context 12 showing the strongest separation.

Table 2: V-JEPA2 (EMA) TRA by context length.

Context	Low-dissip.	Dissipative	Difference	p -value
4	+0.05%	+0.13%	+0.08%	0.037*
6	+0.01%	+0.24%	+0.23%	<0.001***
8	+0.03%	+0.22%	+0.20%	<0.001***
10	+0.07%	+0.24%	+0.17%	<0.001***
12	+0.04%	+0.32%	+0.28%	<0.001***
14	+0.09%	-0.11%	-0.21%	<0.001***

D V-JEPA2 NO-EMA RESULTS

Table 3 shows results using the online encoder instead of the EMA target encoder. The pattern is consistent, confirming robustness.

Table 3: V-JEPA2 (no EMA) TRA by context length.

Context	Low-dissip.	Dissipative	Difference	<i>p</i> -value
4	+0.06%	+0.13%	+0.07%	0.048*
6	+0.01%	+0.22%	+0.21%	<0.001***
8	+0.04%	+0.21%	+0.17%	<0.001***
10	+0.08%	+0.26%	+0.18%	<0.001***
12	+0.04%	+0.30%	+0.26%	<0.001***
14	+0.09%	-0.11%	-0.20%	<0.001***

E VIDEOMAE V2 FULL RESULTS

Table 4 shows VideoMAE V2 TRA across all context lengths. All values are negative, and the inverted pattern strengthens with longer context.

Table 4: VideoMAE V2 TRA by context length. All values negative (inverted pattern).

Context	Low-dissip.	Dissipative	Difference	<i>p</i> -value
4	-0.03%	-0.13%	-0.10%	<0.001***
6	-0.05%	-0.21%	-0.16%	<0.001***
8	-0.07%	-0.28%	-0.21%	<0.001***
10	-0.08%	-0.34%	-0.26%	<0.001***
12	-0.09%	-0.41%	-0.32%	<0.001***
14	-0.10%	-0.50%	-0.40%	<0.001***

F RANDOM BASELINE RESULTS

Table 5 shows results for randomly initialized V-JEPA2 (no pretraining). TRA values are negligible for all conditions (<0.01%). While some *p*-values reach nominal significance at longer context lengths due to the large sample size ($n = 80$ per group), the effect sizes are negligible (Cohen’s $d < 0.05$), confirming that meaningful temporal asymmetry sensitivity requires learned representations.

Table 5: Random baseline (untrained V-JEPA2). TRA values are negligible (<0.01%) for all conditions. Some *p*-values reach significance due to large sample size, but effect sizes are negligible.

Context	Low-dissip.	Dissipative	Difference	<i>p</i> -value
4	+0.00%	+0.00%	+0.00%	0.02*
6	+0.00%	+0.00%	+0.00%	0.05*
8	+0.00%	+0.00%	+0.00%	0.56 ns
10	+0.00%	+0.00%	+0.00%	0.02*
12	+0.00%	+0.00%	+0.00%	<0.001***
14	+0.00%	+0.00%	+0.00%	<0.001***

G MVD FULL RESULTS

Table 6 shows MVD (Masked Video Distillation) TRA across all context lengths. TRA values are essentially zero, indicating that teacher-student distillation does not encode temporal asymmetry.

H HIERA FULL RESULTS

Table 7 shows Hiera (Hierarchical MAE) TRA across all context lengths. Results are inconsistent across context lengths and mostly not statistically significant, suggesting Hiera does not reliably encode temporal asymmetry.

Table 6: MVD TRA by context length. All values near zero.

Context	Low-dissip.	Dissipative	Difference	<i>p</i> -value
4	+0.00%	+0.00%	+0.00%	0.006**
6	+0.00%	+0.00%	+0.00%	0.047*
8	+0.00%	+0.00%	+0.00%	0.74 ns
10	+0.00%	-0.00%	-0.00%	0.07 ns
12	+0.00%	-0.00%	-0.00%	0.015*
14	+0.00%	+0.00%	-0.00%	0.97 ns

Table 7: Hiera TRA by context length. Results inconsistent and mostly not significant.

Context	Low-dissip.	Dissipative	Difference	<i>p</i> -value
4	-0.15%	-0.10%	+0.05%	0.80 ns
6	+0.04%	-0.01%	-0.06%	0.77 ns
8	+0.19%	-0.09%	-0.29%	0.22 ns
10	-0.10%	-0.31%*	-0.21%	0.32 ns
12	+0.52%**	-0.18%	-0.70%	0.019*
14	-0.01%	+0.06%	+0.07%	0.84 ns

I PER-SCENARIO BREAKDOWN

Table 8 shows detailed per-scenario results for V-JEPA2 (EMA) at context length 12.

Table 8: V-JEPA2 TRA by individual scenario (context=12).

Scenario	Type	TRA	<i>t</i> -stat	<i>p</i> -value
Bouncing ball	Low-dissip.	+0.02%	0.99	0.33 ns
Pendulum	Low-dissip.	+0.07%	3.29	0.002**
Falling objects	Dissipative	+0.18%	4.10	<0.001***
Dominos	Dissipative	+0.47%	∞	<0.001***

J CONTINUOUS PROBE DATA

Table 9 shows TRA as a function of physical parameters. Both scenarios exhibit non-monotonic patterns: TRA peaks at intermediate values where energy dissipation is visually apparent throughout the video, approaching zero at both extremes (perfectly elastic/undamped, or motion stops too quickly to observe dissipation).

K CONTEXT LENGTH 14 ANOMALY

At context length 14 (out of 16 total frames in each window), V-JEPA2’s TRA pattern inverts: dissipative processes show *negative* TRA (-0.11%) while low-dissipation processes remain slightly positive (+0.09%). We hypothesize that this anomaly occurs because:

1. With only 2 frames left to predict, the model’s task at hand changes significantly
2. Short prediction horizons may not provide enough temporal extent to distinguish dissipative dynamics
3. The prediction task becomes more about single-frame extrapolation than sequence modeling

We exclude context=14 from our main analysis but report it here for completeness. The consistent pattern across context lengths 6–12 suggests this anomaly reflects a boundary condition rather than a fundamental limitation of TRA.

Table 9: TRA vs. physical parameters (context=12). Peak values in bold.

Bouncing Ball			Pendulum		
Restitution	TRA	Note	Damping	TRA	Note
0.1	+0.17%	Stops quickly	0.0	+0.05%	No damping
0.3	+0.19%		0.5	+0.18%	
0.5	+0.27%	Peak	1.0	+0.31%	
0.7	+0.12%		2.0	+0.68%	Peak
0.9	+0.03%		5.0	+0.42%	Stops quickly
1.0	+0.01%	Elastic			

L REAL-WORLD VIDEO: SOMETHING-SOMETHING V2

To test generalization beyond synthetic physics, we evaluated TRA on Something-Something V2 (SSv2), a dataset of human-performed actions. We selected 250 “reversible” actions (spinning, rolling) and 250 “irreversible” actions (dropping, pouring, tearing).

SSv2 Example: "Spinning a pen so it continues spinning"

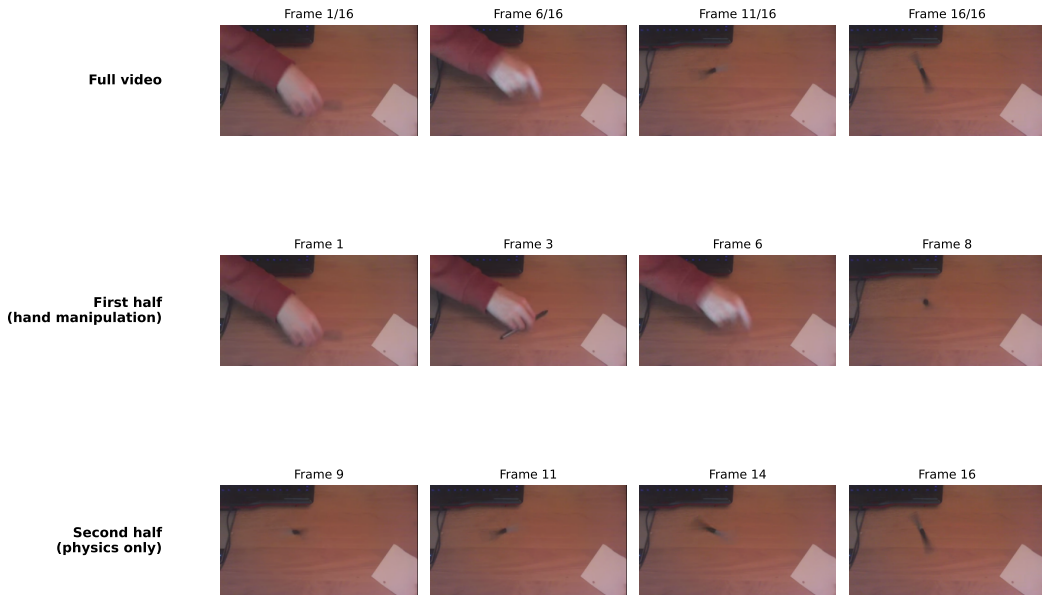


Figure 5: SSv2 temporal cropping. The first half of videos (red border) shows human hand manipulation; the second half (green border) shows physics-dominated motion after release. Cropping to the second half isolates the physics signal from human action timing.

Initial results. V-JEPA2 shows *negative* TRA for *all* SSv2 actions (-0.6% to -0.7%), with no significant difference between reversible and irreversible categories ($p > 0.5$). This contradicts the pattern observed on synthetic physics.

Root cause: human action initiation. We hypothesized that human action initiation (hands entering the frame to manipulate objects) creates a temporal asymmetry unrelated to physics. In forward videos, predicting *when* a hand will appear and initiate action is difficult. In reversed videos, the action “undoes” itself predictably. This human-induced asymmetry may dominate over the smaller physics-based asymmetry.

Temporal cropping experiment. To isolate physics from human agency, we temporally cropped videos using a simple heuristic: the **second half** of SSv2 videos typically contains physics-

dominated motion after the hand releases the object (Figure 5). This is an approximate method; the exact transition from manipulation to physics varies across videos. However, it provides a proof-of-concept that the physics signal exists when human action is reduced.

Results are shown in Table 10. When using only the second half, a significant difference emerges: irreversible actions show near-zero TRA (-0.04% , ns) while reversible actions remain negative (-0.29% , $p < 0.001$). The difference is significant ($p = 0.0008$).

Table 10: SSv2 cropped experiment (V-JEPA2, context=8). Cropping to second half reveals physics signal.

Crop	Reversible	Irreversible	Diff.	p -value
First half (hands)	$-0.30\%^{***}$	$-0.35\%^{***}$	-0.05%	0.51 ns
Second half (physics)	$-0.29\%^{***}$	-0.04% ns	$+0.25\%$	0.0008 ^{***}
Middle 50%	$-0.25\%^{***}$	$-0.16\%^{***}$	$+0.09\%$	0.21 ns

This confirms that TRA measures physics understanding when physical dynamics dominate, but human action timing confounds the signal in videos with visible agents. Future work could develop more sophisticated methods to segment human manipulation from physics-dominated phases, such as hand detection or motion-based segmentation. Alternatively, datasets with minimal human presence (surveillance footage, nature videos, robotic manipulation) may provide cleaner signals for evaluating physics understanding.