

VideoPhy2: Challenging Action-Centric Physical Commonsense Evaluation of Video Generation

Hritik Bansal*, **Clark Peng***, Yonatan Bitton*,
Roman Goldenberg, Aditya Grover, Kai-Wei Chang

ICLR 2026



Video Generative Models



Gen-3 Alpha



COGVIDEOX

RAY2



Wan

Meta Movie Gen

Physical World Simulators



2025-3-28

GR00T N1: An Open Foundation Model for Generalist Humanoid Robots

NVIDIA¹



S

Learning Universal Policies via Text-Guided Video Generation

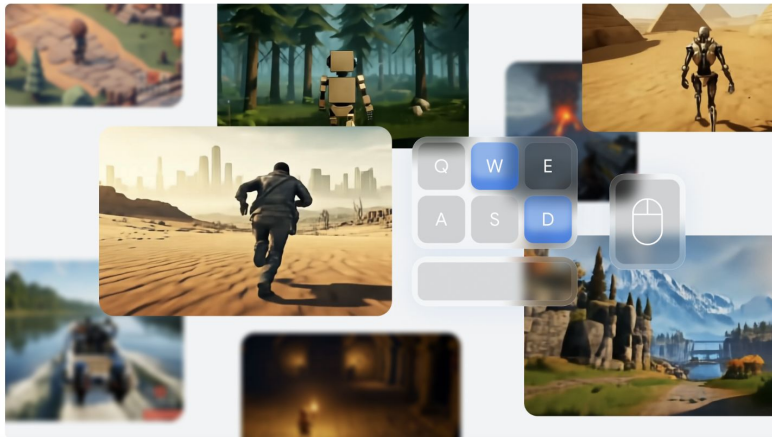
[§], ^{Bo Dai}[¶], ^{Hanjun Dai}[‡], ^{Ofir Nachum}[‡],
^{Dale Schuurmans}^{‡¶}, ^{Pieter Abbeel}[§]
^{erkeley}[§] ^{Georgia Tech}[¶] ^{University of Alberta}[‡]
universal-policy.github.io/

Genie 2: A large-scale foundation world model

4 DECEMBER 2024

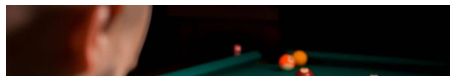
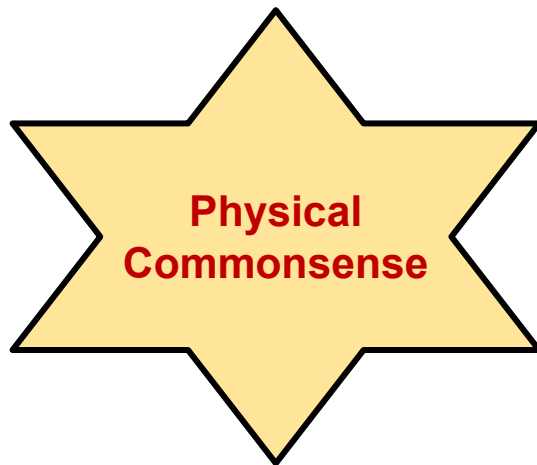
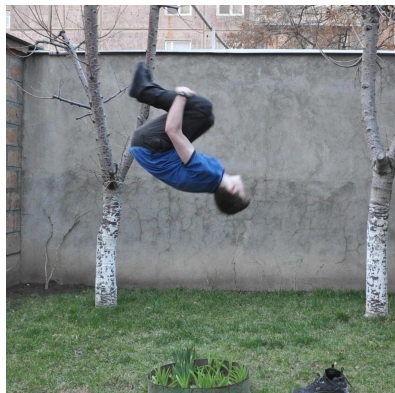
Jack Parker-Holder, Philip Ball, Jake Bruce, Vibhavari Dasagi, Kristian Holsheimer, Christos Kaplanis, Alexandre Moufarek, Guy Scully, Jeremy Shar, Jimmy Shi, Stephen Spencer, Jessica Yung, Michael Dennis, Sultan Kenjeyev, Shangbang Long, Vlad Mnih, Harris Chan, Maxime Gazeau, Bonnie Li, Fabio Pardo, Luyu Wang, Lei Zhang, Frederic Besse, Tim Harley, Anna Mitenkova, Jane Wang, Jeff Clune, Demis Hassabis, Raia Hadsell, Adrian Bolton, Satinder Singh, Tim Rocktäschel

[Share](#)



Humans Rely on Physical Commonsense In Everyday Activities

As humans, we develop an intuitive understanding of the object interactions through our experience with the real-world, without any formal education in physics.



Despite physical motivations, it is **unclear** how well do generated videos follow physical commonsense for **real-world** activities!

VideoPhy-2: Physical Commonsense for Real-World Actions



Categories

(e.g, sports, object interactions)

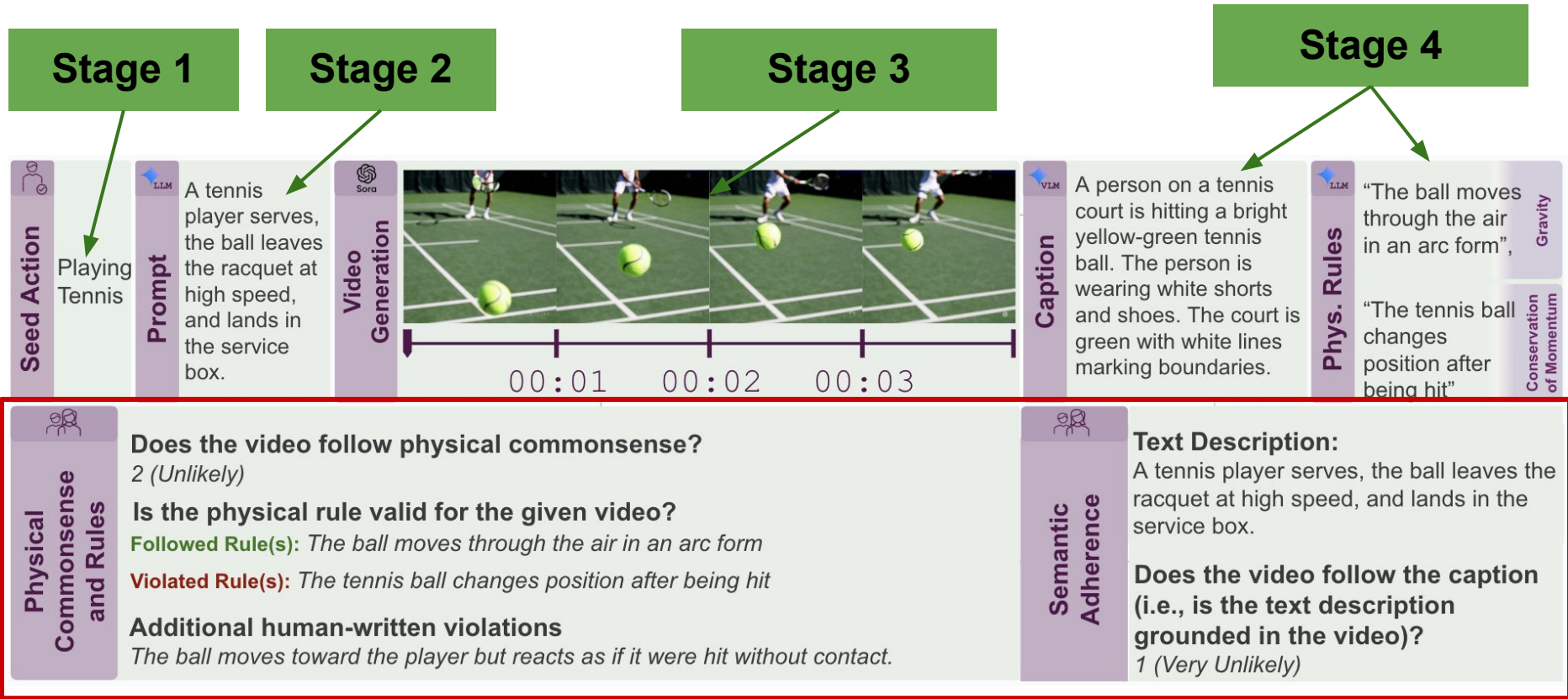
Diverse Physics Laws

(e.g, gravity, momentum)

Multiple Events

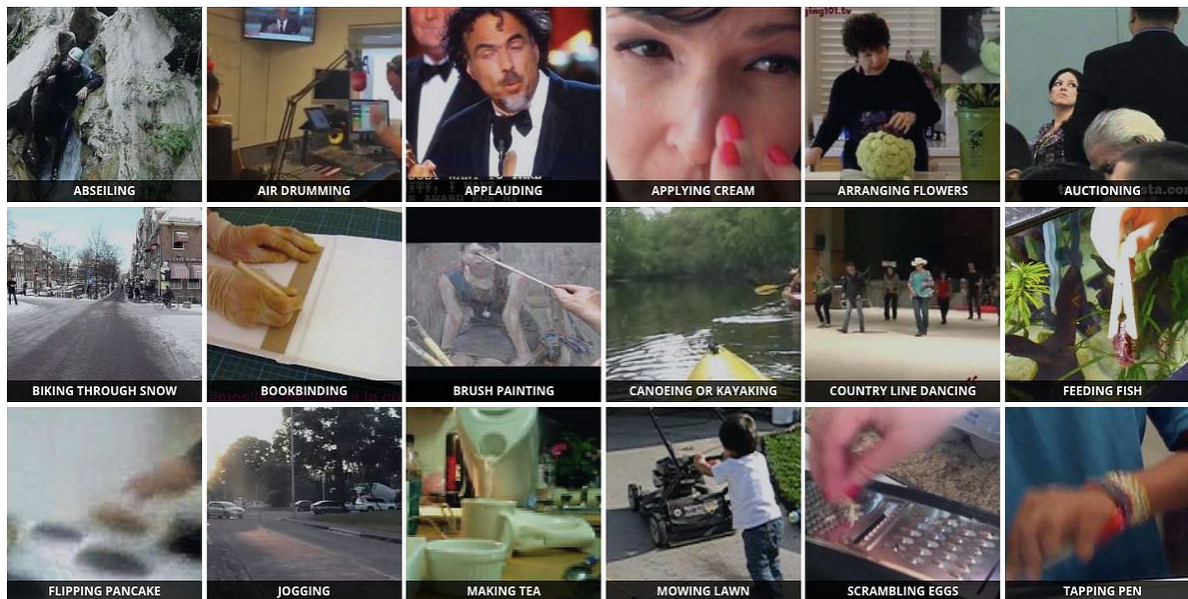
(e.g., moving and then stopping)

Pipeline



Stage 5: Human Eval

Stage 1: Seed Actions Are Taken from Video Datasets and Subsequently Filtered By Humans



200 actions are selected from 600 actions for physical commonsense evaluation!

Stage 3: Video Generation Conditioned on Prompt

Prompt:

A tennis player serves, the ball leaves the racquet at high speed, and lands in the service box.



Stage 4: Powerful Video-Language Model Captions And Lists Candidate Physical Rules



Video Caption:

A person on a tennis court is hitting a bright yellow-green tennis ball.



Candidate Physical Rules:

The ball moves through the air in an arc form (Gravity).

The tennis ball changes position after being hit (Conversation of Momentum).

Stage 5: Human Evaluation Via Amazon Mechanical Turk

Semantic Adherence

Rate: 1-5

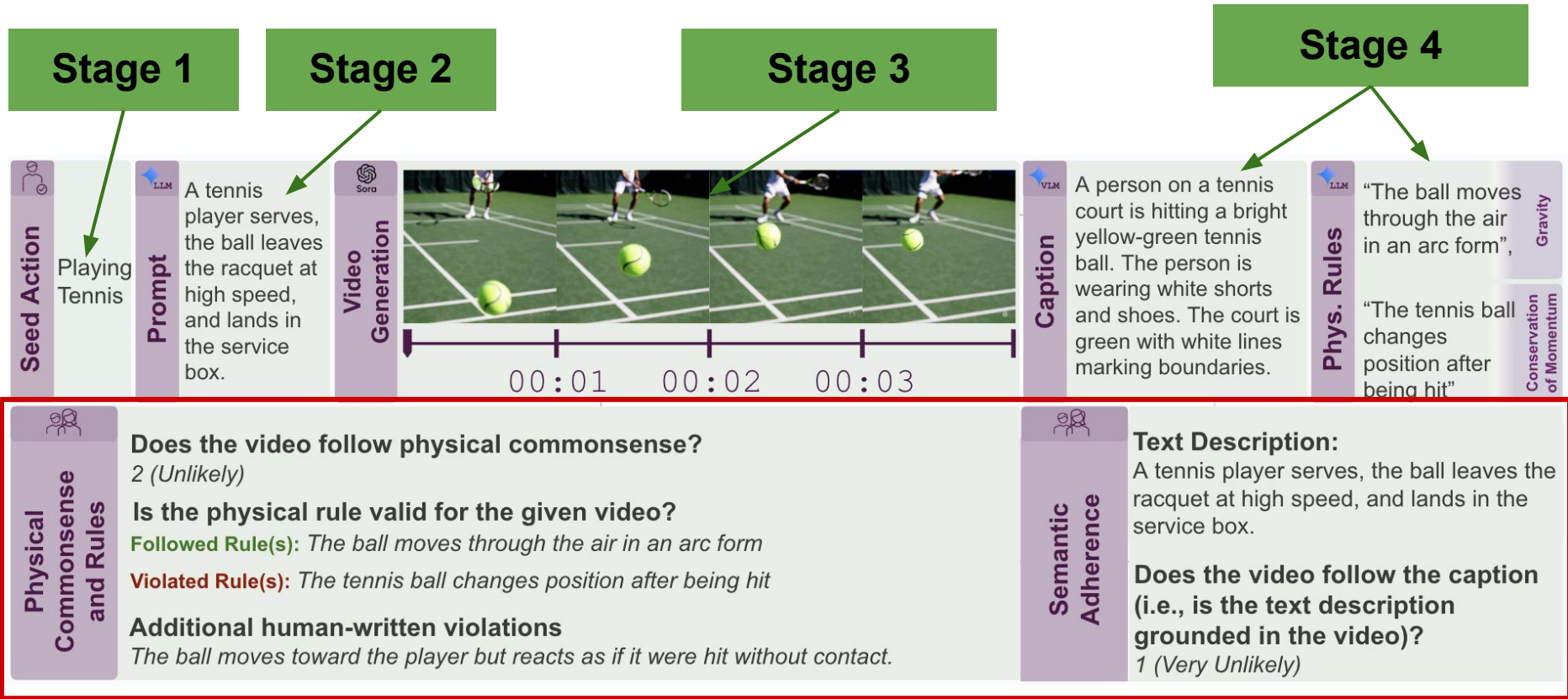
Physical Commonsense

Rate: 1-5

Physical Rule Classification

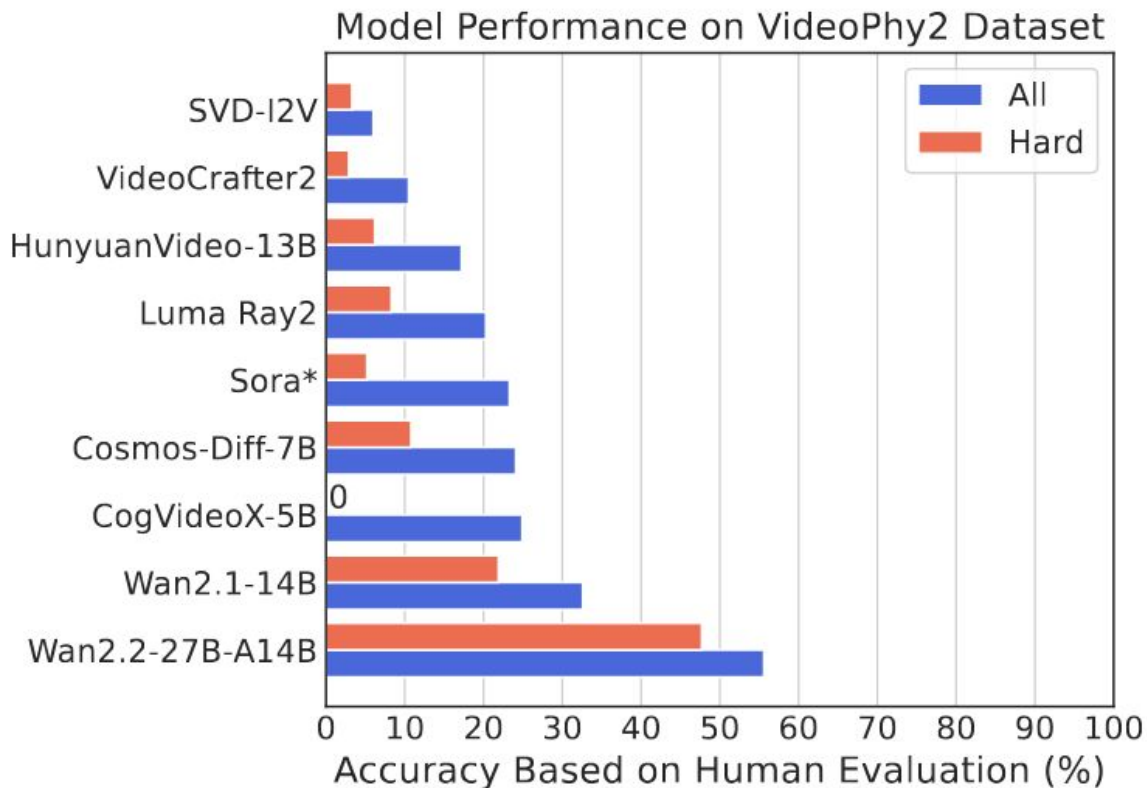
Rate: Yes-No

Pipeline



Stage 5: Human Eval

Results

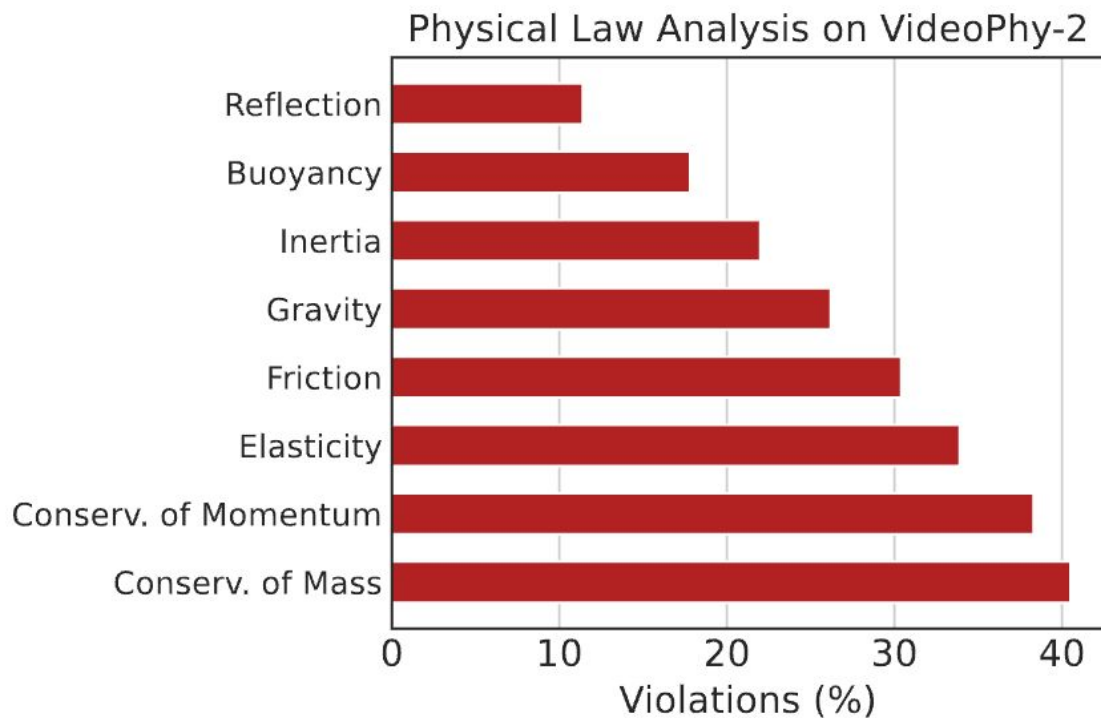


Joint score
1 if (SA \geq 4 and PC \geq 4)

Modern video generative models (closed and open) are far from simulating the real-world actions!

Wan2.1-14B achieves the highest joint performance.

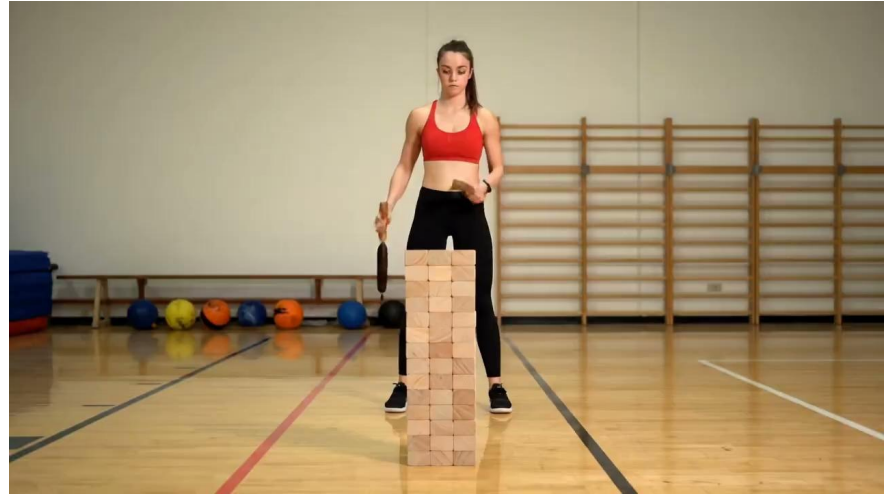
Our Analysis Suggests Conservation Laws are the Most Violated



Examples (Veo3)



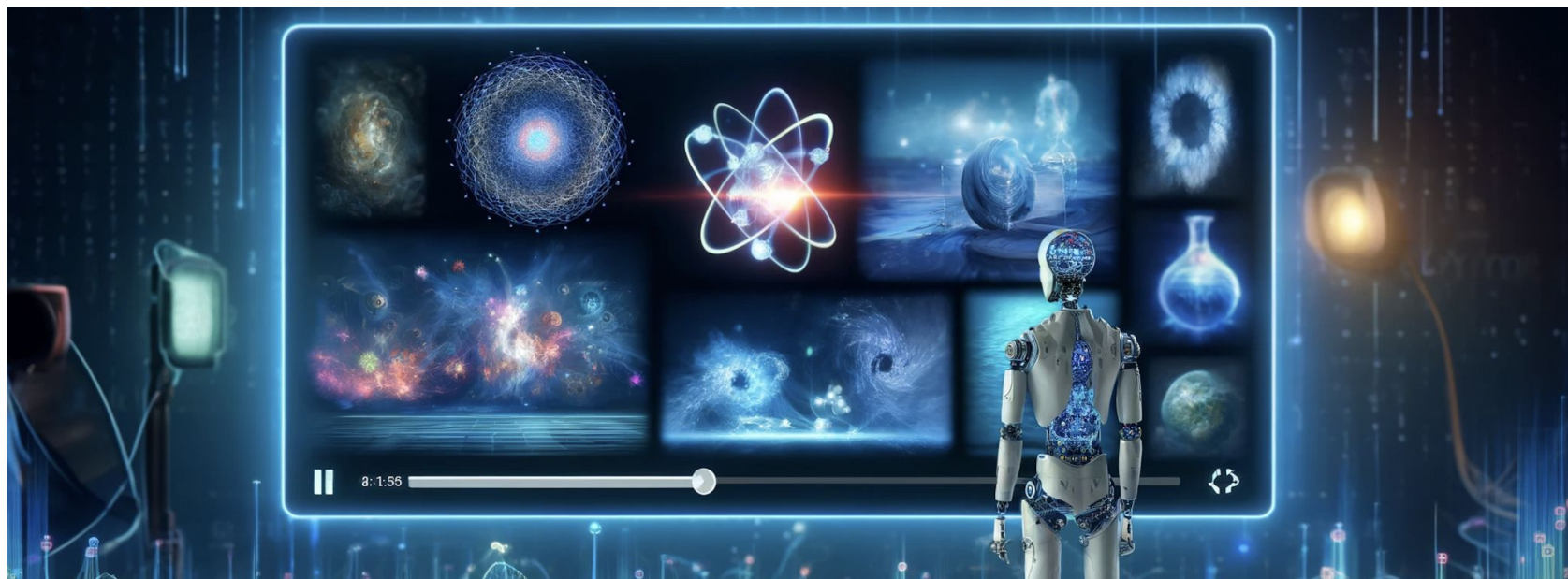
an archer draws a recurve bow, the string stretching taut, then releases the arrow, which hits a target's center.



a person uses nunchucks to break a stack of wooden blocks, the blocks scattering on impact.

VideoPhy-AutoEval: Automatic Evaluator For Scalable Evls!

A subset of human annotations finetune a video language model



VideoPhy2-AutoEval Outperforms Gemini-2.0-Flash!

Table 4: **Auto-rater evaluation results.** We present the pearson’s correlation ($\times 100$) between the predicted scores and ground-truth scores (1-5) on the unseen prompts and unseen video models.

	Unseen prompts			Unseen video models		
	Avg.	SA	PC	Avg.	SA	PC
VideoCon-Physics [5]	28.5	32.0	25.0	26.5	27.0	26.0
VideoCon [3]	12.5	23.0	2.0	8.9	17.0	0.8
VideoLlava [34]	16.0	30.0	2.0	19.0	33.0	5.0
VideoScore [20]	13.5	17.0	10.0	9.0	5.0	13.0
Gemini-2.0-Flash-Exp	18.5	26.0	11.0	21.0	31.0	11.0
VIDEOPHY-2-AUTOEVAL	42.0	47.0	37.0	41.0	45.0	37.0
<i>Rel. to Best (%)</i>	+47.4	+46.9	+48.0	+49.0	+36.4	+61.5
<i>Rel. to Gemini (%)</i>	+127.0	+80.8	+236.4	+107.1	+45.2	+281.8

Table 6: **Auto-rater evaluation on physical rule classification.** We present the accuracy results for VIDEOPHY-2-AUTOEVAL and other video-language models on the rule classification tasks.

	Unseen prompts	Unseen video models
Random	34.5	31.2
VideoLlava [34]	38.1	38.7
Gemini-2.0-Flash-Exp	59.2	57.1
VIDEOPHY-2-AUTOEVAL	78.7	72.9
<i>Rel. to Best (%)</i>	+32.9	+27.7

Community Adoption of VideoPhy Series

Benchmark

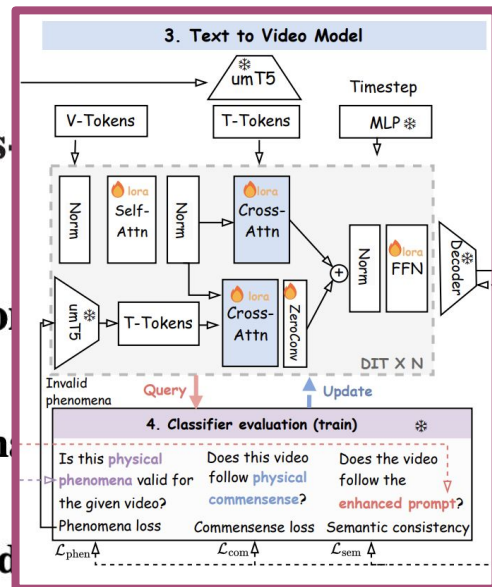
Reward Model

Think Before You Diffuse: LLMs-Guided Physics-Aware Video Generation

Articulated Kinematics Distillation from Video Diffusion

PhysMotion: Physics-Grounded Dynamics from a Single Image

WISA: World Simulator Assistant for Physics-Aware Text-to-Video



TL;DR

- **VideoPhy2** is a high-quality dataset for physical commonsense evaluation of the open and closed video generative models.
- **VideoPhy2-AutoEval** is an automatic evaluator for assessing physical commonsense for the novel video models.

Paper, code, and dataset are publicly available

VIDEOPHY 2

Challenging Action-Centric Physical Commonsense Evaluation
of Video Generation

Hritik Bansal^{*1}, Clark Peng^{*1}, Yonatan Bitton^{*2}, Roman Goldenberg², Aditya Grover¹, Kai-Wei Chang¹,

(* Equal Contribution)

¹University of California, Los Angeles

²Google Research