

# Inferring Dynamic Physical Properties from Video Foundation Models

Anonymous CVPR submission

Paper ID 4

## Abstract

001 *We study the task of predicting dynamic physical proper-*  
002 *ties from videos. More specifically, we consider physical*  
003 *properties that require temporal information to be inferred:*  
004 *elasticity of a bouncing object, viscosity of a flowing liq-*  
005 *uid, and dynamic friction of an object sliding on a surface.*  
006 *To this end, we make the following contributions: (i) We*  
007 *collect a new video dataset for each physical property, con-*  
008 *sisting of synthetic training and testing splits, as well as a*  
009 *real split for real world evaluation. (ii) We explore three*  
010 *ways to infer the physical property from videos: (a) an or-*  
011 *acle method where we supply the visual cues that intrin-*  
012 *sically reflect the property using classical computer vision*  
013 *techniques; (b) a simple read out mechanism using a vi-*  
014 *sual prompt and trainable prompt vector for cross-attention*  
015 *on pre-trained video generative and self-supervised mod-*  
016 *els; and (c) prompt strategies for Multi-modal Large Lan-*  
017 *guage Models (MLLMs). (iii) We show that video founda-*  
018 *tion models trained in a generative or self-supervised man-*  
019 *ner achieve a similar performance, though behind that of*  
020 *the oracle, and MLLMs are currently inferior to the other*  
021 *models, though their performance can be improved through*  
022 *suitable prompting. The datasets will be publicly released.*

## 023 1. Introduction

024 Humans are remarkably adept at intuitively estimating  
025 physical properties from visual observations. Without di-  
026 rect interaction, people can often estimate how bouncy a  
027 ball is, how thick a liquid seems, or how slippery a surface  
028 might be—simply by watching how objects move. While  
029 these estimations are not precise in a scientific sense, they  
030 are sufficiently accurate for guiding perception, prediction,  
031 and action. Bringing this capability to machines is an im-  
032 portant step towards building more general and physically  
033 grounded artificial intelligence. In particular, visual systems  
034 that can infer dynamic physical properties from raw video  
035 could enhance robotic manipulation, embodied agents, and  
036 video understanding tasks in ways that go beyond the tra-  
037 ditional perception tasks of recognition, detection, and seg-

mentation.

Recent progress in video foundation models, including  
generative models [26, 52], self-supervised models [3, 5]  
and multi-modal large language models (MLLMs) [9, 18,  
19], have shown impressive capability in synthesizing re-  
alistic dynamics, learning general-purpose video represen-  
tations, and tackling semantic understanding tasks, for ex-  
ample, video question answering. However, a question that  
remains underexplored is: **do these models acquire an un-  
derstanding of dynamic physical properties from videos**  
?

In this paper, we address this question by focusing  
on several representative physical properties that are not  
directly observable in static frames but instead emerge  
through temporal dynamics: the elasticity of a bouncing  
object, the viscosity of a flowing liquid, and the dynamic  
friction between a surface and a sliding object. These prop-  
erties are especially compelling because their inference re-  
quires temporal reasoning and sensitivity to subtle visual  
cues—such as deformation, deceleration, spreading, or os-  
cillation. By examining how well current video foundation  
models capture these dynamic attributes, we aim to assess  
their physical understanding beyond static appearance.

To support this investigation, we introduce a new  
dataset, *PhysVid*, specifically designed to evaluate the dy-  
namic physical properties from video. Existing datasets  
lack ground-truth annotations for such properties, so  
we construct *PhysVid* using a combination of synthetic  
videos—rendered via a physics simulator—and real-world  
videos sourced from the internet or captured in-house. Each  
video is annotated with physical property values, either de-  
rived from simulation parameters or estimated manually.  
The dataset is designed to facilitate the study of out-of-  
domain generalization, both within the synthetic domain  
and from synthetic to real-world data. To establish an upper  
bound on what is inferable from visual input alone, we im-  
plement an oracle method for each property. These oracles  
leverage privileged access to the visual cues that directly  
reflect the corresponding property.

We evaluate three categories of video foundation mod-  
els: generative models, self-supervised models, and multi-

079 modal large language models (MLLMs). For the genera-  
080 tive and self-supervised models, we propose a simple yet  
081 effective readout mechanism that extracts dynamic physi-  
082 cal properties from pre-trained, frozen representations. Our  
083 method introduces a learnable query vector that attends  
084 to internal representation tokens via cross-attention, en-  
085 abling the selective extraction of relevant information. This  
086 approach is both lightweight and training-efficient. For  
087 MLLMs, we explore various prompting strategies to elicit  
088 predictions of dynamic physical properties directly from  
089 video input. These strategies include few-shot prompt-  
090 ing to provide task context, as well as procedural prompt-  
091 ing that guides the model through the oracle estimation  
092 steps—helping it focus on the intrinsic visual cues that re-  
093 veal the target properties.

## 094 2. Related Work

095 **Physics Prediction from Images and Videos.** Inferring  
096 physical properties from visual observations remains a core  
097 challenge in computer vision. Early methods estimate la-  
098 tent physical parameters (e.g., mass, friction, stiffness) via  
099 differentiable physics engines or learning-based simula-  
100 tors [10, 20, 24, 43, 44, 50], while later works infer salient  
101 attributes like viscosity or elasticity from task-specific vi-  
102 sual cues [2, 21, 22, 28–31], yet both rely heavily on simu-  
103 lation supervision, domain priors, or handcrafted heuristics.  
104 More recently, unsupervised learning of intuitive physics  
105 has emerged via next-frame prediction from large-scale ev-  
106 eryday physical scenes [1, 4, 11, 12, 15, 16, 27, 42], captur-  
107 ing latent dynamics without explicit physical supervision.  
108 However, the resulting representations are usually implicit  
109 and lack interpretability in terms of concrete physical quan-  
110 tities. In contrast, we infer physical properties by directly  
111 prompting pre-trained video foundation models, enabling  
112 explicit estimation without reliance on task-specific heuris-  
113 tics, or end-to-end prediction pipelines from scratch.

114 **Physics Datasets and Benchmarks.** An increasing num-  
115 ber of physics-related datasets have been collected in recent  
116 years to provide ground truth annotations for different physi-  
117 cal properties, including material [14, 36], shadow [45, 46],  
118 support relations [39], occlusion [53, 54], mass and vol-  
119 ume [51]. Another line of work [6–8, 35, 37, 41] proposes  
120 broad benchmarks with video-image-text QA tasks to as-  
121 sess physical understanding in vision-language models, but  
122 the questions are typically qualitative and categorical. More  
123 recently, Zhang et al. [56] introduces a benchmark consist-  
124 ing of 130 real-world videos capturing physical phenomena  
125 guided by conservation laws, to evaluate the physics plau-  
126 sibility of video generative models by assessing the trajec-  
127 tory of objects in their generated videos. In contrast, our  
128 datasets consist of both *synthetic* and *real-world* videos an-  
129 notated with the *quantitative value* for the associated physi-  
130 cal parameter of the coefficient of friction, elasticity, and

viscosity.

## 3. Problem Scenario and The *PhysVid* Datasets

In this paper, we address the problem of estimating physi-  
cal properties from videos. Specifically, we focus on  
three properties: **elasticity** of a bouncing object, **viscos-  
ity** of a flowing liquid, and the **dynamic friction coeffi-  
cient** between a surface and a sliding object. Given a video  
 $v \in \mathbb{R}^{T \times H \times W \times 3}$ , we consider two formulations, the first  
is **absolute value prediction**, where the input is a single  
video and the model is tasked with predicting the numer-  
ical value of the physical property, *i.e.*,  $y_{\text{abs}} = \Phi(v; \theta_1)$ .  
The second is **relative value comparison**, where the in-  
put is a pair of videos captured from the same viewpoint,  
and the model must determine whether the first video ex-  
hibits a higher physical property value than the second, *i.e.*,  
 $y_{\text{rel}} = \Phi(v_1, v_2; \theta_2)$ , and  $y_{\text{rel}}$  is binary.

Each scenario is parameterized by a set of variables, in-  
cluding the value of the target *physical property* (e.g., elas-  
ticity, viscosity, or friction), and a set of *nuisance param-  
eters* (including camera viewpoint, object appearance, light-  
ing, *etc.*). While the model must be sensitive to changes in  
the physical property, it should be robust (ideally invariant)  
to variations in nuisance parameters.

To assess generalization, we define two domains of nui-  
sance parameters, denoted as  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , which differ in  
their distributions. For instance,  $\mathcal{A}_2$  may have different  
camera viewpoints or different lighting conditions to  $\mathcal{A}_1$   
(full details of these differences are given in Supplemen-  
tary Section B). We generate a dataset using a physics-based  
simulator, consisting of one training split and two test splits.  
The models are only trained on the training split from the  
simulator for all the evaluations. The training and `test-1`  
splits are sampled from  $\mathcal{A}_1$ , sharing the same nuisance dis-  
tribution; `test-2` is drawn from  $\mathcal{A}_2$ , introducing a distri-  
bution shift. The target property values are sampled from  
a shared range across all splits to ensure consistency. Fi-  
nally, `test-3` consists of real-world videos, used to eval-  
uate generalization beyond simulation.

### 3.1. The *PhysVid* Datasets

To study the dynamic physical properties of elasticity, vis-  
cosity, and friction, we construct a dataset for each, con-  
taining both synthetic and real-world videos. Synthetic  
ones are generated with the Genesis simulator [58], and  
real ones are captured with an iPhone in slow-motion mode  
or downloaded from the Internet. For each property we  
have: 10,000 videos for `train`; 1000 videos for each of  
`test-1` and `test-2`; and 100 videos for `test-3`. Sam-  
ple frames are shown in Figure 1. In the following we de-  
scribe how each property is realized in the video. Please  
refer to Supplementary Section B for more details of the  
datasets.

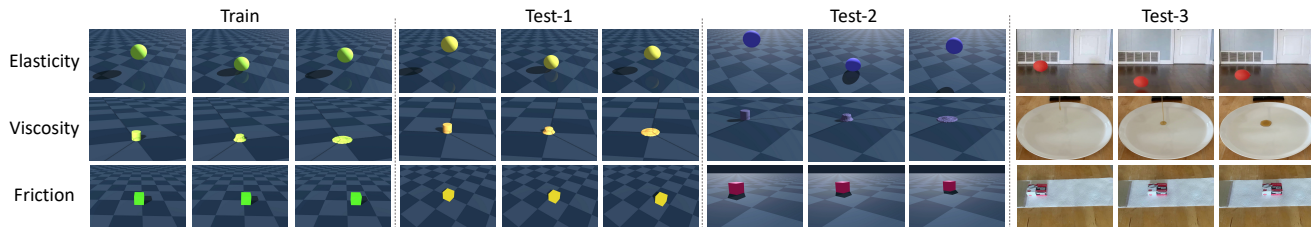


Figure 1. **Examples of the PhysVid dataset.** Each row shows a different property, and each column shows three frames from video samples in the synthetic sets (train, test-1, and test-2) and the real test-3 set. The train and test-1 sets are from the same distribution. In test-2 parameters, such as lighting, viewpoint and color, differ from those in test-1.

### 182 3.1.1. Elasticity

183 We study an object’s elasticity by analyzing the motion of  
184 a ball dropped onto the ground and its subsequent bounces.  
185 In physics, elasticity  $e$  is quantified as the ratio of the re-  
186 bound velocity  $v_{\text{after impact}}$  to the impact velocity  $v_{\text{before impact}}$ ,  
187 and also equals  $\sqrt{h_{\text{bounce}}/h_{\text{drop}}}$ , where  $h_{\text{drop}}$  is the drop-  
188 ping height and  $h_{\text{bounce}}$  is the bouncing height. Here and  
189 for the following properties, please refer to Supplementary  
190 Section C for the detailed derivations. These expressions  
191 are used for the oracle estimation in Section 4.1.

192 **Synthetic Dataset.** All synthetic videos are generated using  
193 Genesis [58], with object’s elasticity as the target property.  
194 Nuisance factors include drop height, camera view-  
195 point, object appearance, and lighting conditions. The object  
196 is of the same size in all videos. Note, here and for  
197 the following properties, the ground truth property value is  
198 obtained directly from the simulator.

199 **Real-World Dataset.** The real-world videos are collected  
200 from YouTube using the search term “ball bouncing experi-  
201 ments”. Each clip is manually trimmed to include the drop-  
202 and-bounce sequence of a single ball. The dataset includes a  
203 wide range of materials (*e.g.*, rubber balls, tennis balls, bas-  
204 ketballs, balloons, *etc.*), resulting in diverse elasticity values.  
205 The ground truth elasticity values for the real sequences  
206 are estimated by computing  $\sqrt{h_{\text{bounce}}/h_{\text{drop}}}$ : the videos are  
207 chosen such that the balls bounce in a fronto-parallel plane,  
208 which means that ratios of image heights (differences in  $y$ -  
209 coordinates) are approximately equal to the ratio of heights  
210 in 3D. These image differences are obtained by manual an-  
211 notation.

### 212 3.1.2. Viscosity

213 We study the viscosity by observing a liquid column drop-  
214 ping and spreading on the ground. The viscosity can be re-  
215 flected by the growth rate of the liquid area on the ground.  
216 The viscosity  $\mu$  is negatively correlated to the liquid area  
217 growth rate  $\frac{d(A(t))}{dt}$ , given the controlled liquid density  $\rho$ ,  
218 controlled liquid column diameter  $D$ , and controlled drop-  
219 ping velocity  $v$  of the liquid column when it reaches the  
220 ground.

221 **Synthetic Dataset.** The synthetic videos are generated us-

ing Genesis [58], where the target property is the viscosity  
222 of liquid. Nuisance factors include camera viewpoint, ob-  
223 ject appearance, and lighting conditions. The liquid column  
224 is of the same size in all videos. 225

226 **Real-World Dataset.** Since it is challenging to find real-  
227 world videos online that provide ground-truth viscosity val-  
228 ues while controlling for other relevant physical paramet-  
229 ers—such as  $\rho$ ,  $D$  and  $v$ , we collected real videos under  
230 controlled conditions. We use a funnel with a fixed nozzle  
231 diameter to produce a consistent liquid column. A funnel  
232 holder allows us to fix the height from which the liquid is  
233 poured, thereby controlling the initial velocity  $v$ . Ground-  
234 truth viscosity values for each liquid are obtained from stan-  
235 dard physics reference tables. The selected liquids span  
236 a wide range of viscosities, from 1.2 (*e.g.*, coffee) to 225  
237 (*e.g.*, maple syrup), allowing for a diverse and comprehen-  
238 sive evaluation.

### 239 3.1.3. Friction

240 We study friction between an object and a surface by ob-  
241 serving how the object slows down as it slides with an initial  
242 velocity. The dynamic friction coefficient  $\mu_k$  is proportional  
243 to the (negative) acceleration of the object  $a$ . 244

245 **Synthetic Dataset.** The synthetic videos are generated using  
246 Genesis [58], where the target property is the dynamic  
247 friction coefficient at the contacting surface of the object  
248 and the ground. Nuisance factors include initial location  
249 and initial velocity of the object, camera viewpoint, object  
250 appearance, and lighting conditions. The object is of the  
251 same size in all videos. 252

253 **Real-World Dataset.** While many online videos depict ob-  
254 jects sliding on surfaces, they lack ground-truth annotations  
255 for friction coefficients. We therefore collect a real video  
256 dataset featuring 5 different objects and 6 surface materi-  
257 als, spanning a wide range of dynamic friction values. Each  
258 object is given an initial velocity by sliding it down from  
259 a slope and it then slides on a horizontal plane. To obtain  
260 ground-truth friction coefficients, we use a spring dyna-  
261 mometer to measure the friction force  $F$  for each object-  
surface pair (by dragging the object at constant speed), and  
record the object’s weight  $G$ . The dynamic friction coeffi-

cient is then computed as:  $\mu_k = F/G$ .

## 4. Inferring Physical Properties

This section presents the three different ways for inferring dynamic physical properties: an oracle method via classical computer vision techniques (Section 4.1); a visual prompt mechanism for video generative and self-supervised models (Section 4.2); and prompts for MLLMs (Section 4.3).

### 4.1. Oracle Estimation

#### 4.1.1. Elasticity

We aim to estimate elasticity from both synthetic and real-world videos. The key visual cue is the relative height of the ball during its drop and subsequent bounce, observed in 3D. As noted earlier, the ratio in 3D can be approximated from their corresponding image-space measurements. This approximation is exact when the motion occurs in a fronto-parallel plane, and remains reasonably accurate otherwise—since the ratio of lengths between parallel line segments is invariant under affine transformations [17]. Given that perspective effects are minimal in our videos, the affine approximation provides a reliable estimate for elasticity. To automate this process, we extract the ball’s trajectory  $y(t)$  from the video and input the sequence of ratios into a GRU network to regress the elasticity. In detail, we segment the ball in each frame and use their centroids as the  $y$ -coordinate. From this trajectory, we identify key points: the initial drop position, the first ground contact, and the peak of the first bounce. The resulting trajectory is normalized to the range  $[0, 1]$ , by subtracting the  $y$ -coordinate of the first ground contact and dividing by the initial drop height. This normalization not only ensures invariance to viewpoint and scale, but also simplifies learning for the GRU by standardizing the input distribution. We train a GRU, as it is noisy to directly obtain  $h_{\text{drop}}$  and  $h_{\text{bounce}}$  using heuristics (*e.g.*, determining the maximum and minimum points), and in practice a GRU provides a good estimate. The full pipeline is illustrated in Figure 2 (top row). For the **absolute prediction**, the normalized trajectory is fed into a GRU network, which directly regresses the elasticity value. For the **relative comparison**, the binary decision score between two videos  $v_1$  and  $v_2$  is calculated as:

$$\text{score} = \sigma\left(\log\left(\frac{e_1}{e_2}\right)\right), \quad (1)$$

where  $e_1$  and  $e_2$  are the estimated elasticities based on height ratios, and  $\sigma(\cdot)$  denotes the sigmoid function.

#### 4.1.2. Viscosity

The key visual cue for estimating viscosity is the rate at which the liquid spreads on the ground-plane, measured as an area ratio normalized by the initial area of the liquid column. As with elasticity, we approximate perspective using

an affine transformation – here of the ground-plane. Since area ratios are invariant under affine transformations [17], the liquid’s normalized image-space area growth approximates its true normalized ground-plane expansion (in our setup the liquid spreads only within a limited area around the release point, and the camera is distant; consequently an affine viewing approximation is adequate). Specifically, we extract segmentation masks for each frame and compute the liquid’s area over time. This area sequence is normalized by the area in the first frame where the liquid contacts the surface, ensuring invariance to viewpoint and scale. The process is illustrated in Figure 2 (middle row). For **absolute prediction**, we calculate the slope  $k$  of  $A(t)$  and use  $1/k$  to represent the viscosity value; For **relative comparison**, the binary decision score between two videos  $v_1$  and  $v_2$  is calculated as in Equation 1, where  $e_1$  and  $e_2$  are the estimated viscosities based on area growth rate.

#### 4.1.3. Friction

The key visual cue for estimating dynamic friction is the acceleration of the sliding object—*i.e.*, how quickly its velocity decreases due to friction—which can be inferred from its position over time. Since the object moves significantly in the video, we do not use an affine approximation, but instead take account of the projective geometry by mapping the object’s motion to a bird’s-eye view, allowing for consistent trajectory analysis. This is achieved by estimating a homography between the image and bird’s eye view (normal to the plane) from the four corners of the object’s top surface (see Figure 2, bottom row). We fit a parabola  $x = at^2 + \beta t + c$  to the transformed top surface trajectory to estimate the acceleration  $a$  from the coefficient  $\alpha$ , and the coefficient of friction  $\mu_k = 2\alpha/g$ . For **absolute prediction**, we use the estimated  $\mu_k$  to represent the friction coefficient value; For **relative comparison**, the binary decision score between two videos  $v_1$  and  $v_2$  is calculated as in Equation 1, where  $e_1$  and  $e_2$  are the estimated friction coefficients based on the transformed object trajectory.

## 4.2. Video Generative and Self-Supervised Models

### 4.2.1. Video Feature Extraction

Given a video  $v \in \mathbb{R}^{T \times H \times W \times 3}$ , we extract features with a pre-trained video backbone, that can either be generative or self-supervised, resulting into spatiotemporal feature representations, *i.e.*,  $r = \psi(v) \in \mathbb{R}^{t \times h \times w \times c}$ , which can be detailed as follows.

**Generative Model as Backbone.** We adopt a pre-trained video diffusion model (Figure 3, left), namely DynamicCrafter [52], to compute the visual features. Specifically, given an input video, we add noise to the latent representations after the pre-trained VAE encoder, and replace the text prompt with a learnable embedding. We extract multi-scale features from all U-Net layers at diffusion time step 50,

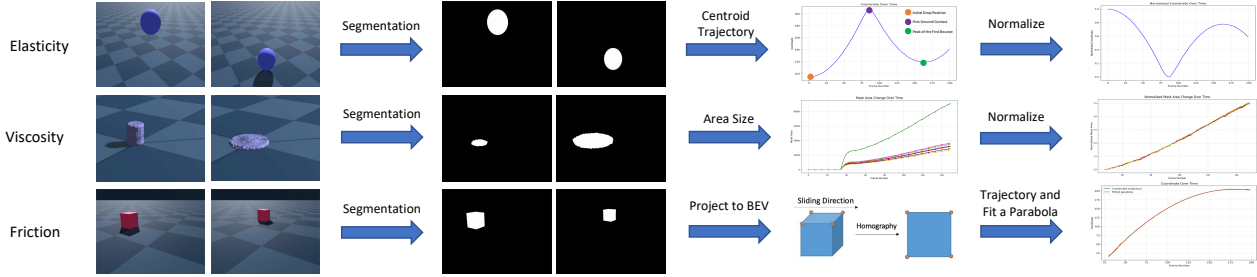


Figure 2. **Oracle methods for physical properties.** The objective in each case is to extract a measurement from the sequence that can directly be used to predict the property. For elasticity, we extract the centroid trajectory from segmentation masks, and then normalize the  $y$ -coordinates into 0-1; the ratio of bouncing to dropping height over the sequence indicates the elasticity. For viscosity, we calculate the area size in the image via segmentation masks, and then normalize the area sizes by the area in the frame when the liquid first touches the ground; the slope of the normalized area size sequence reflects the viscosity. For friction, we transform to a bird’s eye view (using a homography transformation based on 4 corner points of the top surface of the sliding object), and fit a parabola  $x = \alpha t^2 + \beta t + c$  to the transformed trajectory; the parabola coefficient  $\alpha$  predicts the friction coefficient. For each video, we show the segmentation for two frames (left  $\rightarrow$  right).

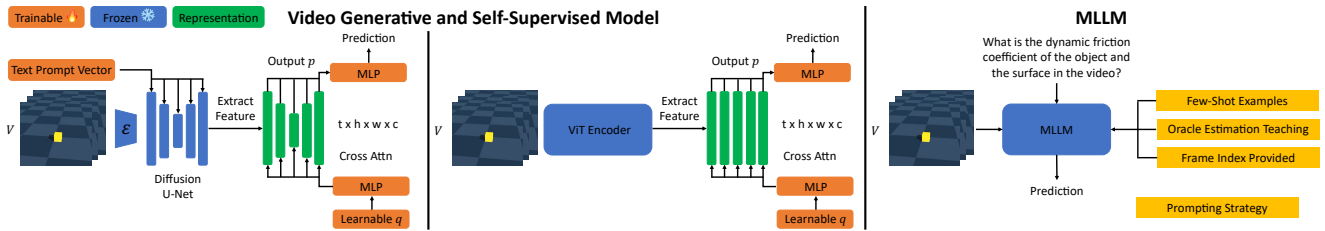


Figure 3. **Architectures for dynamic physical property prediction.** **Left:** video generative model as backbone; **Middle:** video self-supervised model as backbone; **Right:** multimodal large language model (MLLM). For the pre-trained video diffusion model (U-Net, left) and the pre-trained self-supervised model (ViT, middle), the representations are kept frozen, and a ‘visual prompt’ learns to infer the physical properties. For the MLLMs, the physical properties are inferred using a language prompt (right).

361 which was shown to be effective for capturing 3D physics 381  
 362 in prior work [40, 55]. To aggregate the features, we 382  
 363 introduce a learnable query vector  $q$ , which is first mapped 383  
 364 to the different dimensions of the multi-scale features (see 384  
 365 Supplementary Section A.2 for details), and then attends 385  
 366 to the diffusion tokens ( $r_i$ ) via cross-attention: 386

$$367 \quad p = \sum_{i=1}^{t \times h \times w} \text{softmax}(q \cdot r_i) \cdot r_i \quad (2) \quad 387$$

368 The resulting vectors  $p$  from different layers are then 390  
 369 mapped by another MLP network to a common dimension 391  
 370 and average pooled to generate the final video feature 392  
 371 representation  $P$ . To predict the physical properties, we train 393  
 372 the text token of the generative model, together with the 394  
 373 ‘visual prompt’ architecture that includes the query  $q$  and the 395  
 374 MLPs. 396

375 **Self-Supervised Model as Backbone.** Here, we adopt 397  
 376 a pre-trained self-supervised model (Figure 3, middle), 398  
 377 namely V-JEPA-2 [3], as the visual backbone. The input 399  
 378 video is passed through the model, and we extract feature 400  
 379 tokens from all layers of the ViT encoder. Similar to the  
 380 generative setting, we introduce a learnable query vector  $q$

to extract the video feature representation  $P$  from the ViT 381  
 tokens via attentive pooling. Although the feature dimension 382  
 at each ViT layer is the same, we still use a MLP network 383  
 to map  $q$  to generate the query vector of each layer 384  
 (keeping it similar to the generative setting in terms of MLP 385  
 network architecture), and use another MLP network to map 386  
 the output vectors  $p$  to a same dimension as the generative 387  
 setting before average pooling them to get  $P$ . Please see 388  
 Supplementary Section A.2 for more details. 389

#### 4.2.2. Physical Property Prediction 390

Given the computed feature  $P$  from video foundation mod- 391  
 els, we train a MLP network to predict the physical prop- 392  
 erties using the synthetic video dataset training split. The 393  
 network for each property is trained separately. 394

**Absolute Value Prediction.** Given the resulting video fea- 395  
 ture ( $P$ ), we pass it through a MLP network  $\gamma$  to predict the 396  
 absolute value  $\chi$  of the physical property: 397

$$398 \quad \chi = \gamma(P) \quad (3) \quad 398$$

For elasticity and friction, the absolute value prediction is 399  
 supervised with L1 loss with the ground truth value; For 400

401 viscosity, as the ground truth values may have very different  
402 scales, *i.e.*, from  $1e^{-5}$  to  $1e^{-2}$ , the absolute value predic-  
403 tion is trained with Log L1 loss, which calculates L1 loss  
404 between the log of the predicted value and the log of the  
405 ground truth value.

406 **Relative Value Prediction.** Given the resulting features  
407 for a pair of videos,  $P_1$  and  $P_2$ , we concatenate them and  
408 formulate a binary classification problem, indicating which  
409 video has a larger physical property value via a MLP net-  
410 work  $\gamma$ :

$$411 \quad \xi = \gamma([P_1, P_2]) \quad (4)$$

412 The binary prediction for all three tasks is trained with bi-  
413 nary cross entropy loss with the binary ground truth.

414 **Bridging the Sim2real Gap.** Since our models are trained  
415 on synthetic datasets, they may not generalize well to real-  
416 world test videos due to the domain gap. To mitigate this  
417 sim-to-real gap, for both synthetic training and real test, we  
418 draw a red circle on each video frame, enclosing the full  
419 trajectory of the target object or liquid, as illustrated in Fig-  
420 ure 4 (middle). The red circle is obtained automatically as  
421 a bounding ellipse enclosing the merged masks of the target  
422 object or liquid across all frames. This visual cue directs the  
423 model’s attention to the relevant region [38], effectively sig-  
424 naling which object to focus on for physical reasoning. The  
425 red circle serves as a lightweight yet effective form of weak  
426 annotation that helps the model localize and interpret the  
427 dynamics of interest. Please refer to Supplementary Sec-  
428 tion G for the quantitative results demonstrating the effec-  
429 tiveness of drawing such red circles to mitigate the sim-to-  
430 real gap.

### 431 4.3. Multimodal Large Language Models

432 This section studies off-the-shelf multimodal large lan-  
433 guage models (MLLMs) for understanding dynamic phys-  
434 ical properties from video. We explore various prompting  
435 strategies on state-of-the-art MLLMs, including Qwen2.5-  
436 VL-Max [18], GPT-4o [19], and Gemini 2.5 Pro [9], as il-  
437 lustrated in Figure 3 (right). Examples of the prompting  
438 strategies are provided in Supplementary Section E.

#### 439 4.3.1. Preliminary

440 The MLLM receives video frames as visual input. The text  
441 prompt includes (1) a brief description of the target prop-  
442 erty—for example: “we are studying the viscosity of the  
443 liquid, where water is 1.0 and honey is 5000.0.” This is fol-  
444 lowed by (2) a query, such as: “what is the viscosity value of  
445 the liquid in the video?” (absolute) or “which video shows a  
446 liquid with higher viscosity? please output a decision score  
447 between 0 and 1, indicating the likelihood that the first video  
448 exhibits a higher property value.” (relative). All the follow-  
449 ing prompt strategies provide (1) and (2) by default.

#### 450 4.3.2. Baseline Prompt

451 For *relative* tasks, we specify that the first  $n$  frames belong  
452 to the first video and the last  $n$  to the second. For *absolute*  
453 tasks, the default prompt is used. Supplementary Figure 8  
454 and Figure 13 provide an example of *baseline prompt* for  
455 the absolute formulation and the relative formulation, re-  
456 spectively.

#### 457 4.3.3. Oracle Estimation Teaching

458 For both *relative* and *absolute* settings, we provide the key  
459 cue to concentrate on from the Section 4.1 description to  
460 teach the MLLM how to estimate the properties step by  
461 step. Supplementary Figure 9 and Figure 14 provide an ex-  
462 ample of *oracle estimation teaching* for the absolute formu-  
463 lation and the relative formulation, respectively.

#### 464 4.3.4. Few-Shot Examples

465 For both *relative* and *absolute* settings, we provide sev-  
466 eral examples, including the video input and desired ground  
467 truth. For fair comparison with visual prompting, we use  
468 examples in the synthetic training split. Supplementary Fig-  
469 ure 10 and Figure 15 provide an example of *few-shot exam-  
470 ples* for the absolute formulation and the relative formu-  
471 lation, respectively.

#### 472 4.3.5. Frame Index Provided

473 For both *relative* and *absolute* settings, we input the text  
474 of the index of each frame along with the frames. In this  
475 way the MLLMs may have a better understanding about the  
476 temporal relations between the input video frames. Sup-  
477plementary Figure 11 and Figure 16 provide an example of  
478 *frame index provided* for the absolute formulation and the  
479 relative formulation, respectively.

#### 480 4.3.6. Black Frames in Between

481 This strategy is only used for the *relative* setting. We in-  
482 sert black frames between the two video segments to clearly  
483 separate them. In the prompt, we refer to the videos as the  
484 frames before and after the black frames, rather than as the  
485 first and last  $n$  frames. Supplementary Figure 17 provides  
486 an example of *black frames in between* for the relative for-  
487 mulation.

## 488 5. Experiments

489 **Implementation Details.** During oracle estimation, we  
490 train the GRU network with a learning rate of  $1e^{-3}$  and the  
491 batch size is 128. For the generative and self-supervised  
492 video models, the backbones are frozen, the trainable pa-  
493 rameters are optimised with a learning rate of  $1e^{-5}$  and the  
494 batch size 16. For MLLMs, we perform prompt selection,  
495 and use the best strategy that we find for each of the ab-  
496 solute and relative settings for the experiments. *Few-shot  
497 examples* and *oracle estimation teaching* work best for the  
498 absolute and relative settings, respectively, as they directly

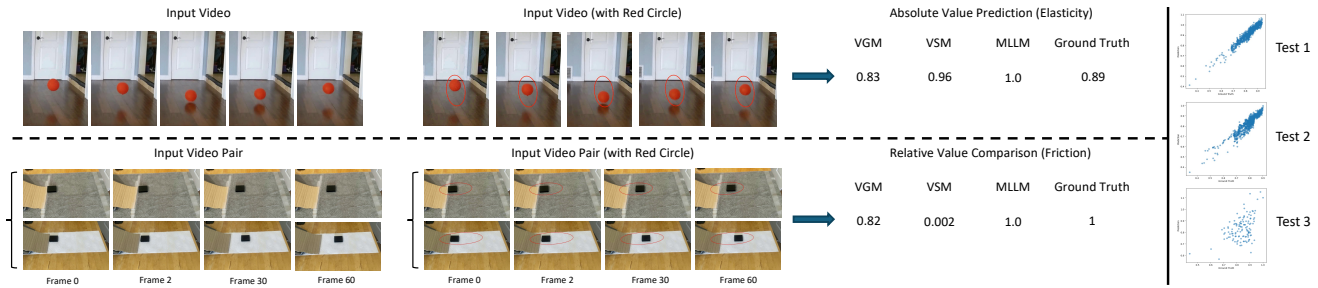


Figure 4. **Qualitative results.** **Top Left:** An example for elasticity absolute value prediction; **Bottom Left:** An example for friction relative value comparison. For each example, the original input video is shown on the left. A static red circle is overlaid in the center to highlight the full trajectory of the object on every frame, shown in the middle. Model predictions are shown on the right, including results from the Video Generative Model (VGM), Video Self-Supervised Model (VSM), and a MLLM (Gemini in this case). For the relative formulation, the ground truth value of ‘1’ indicates that the first (top) video has larger dynamic friction coefficient than the second video. In this example, the initial velocity of the lego brick in the two videos is similar (note the same displacement from frame 0 to 2), but the velocity reduces to 0 at frame 30 in the first video, while the object is still moving in frame 30 to 60 in the second video. **Right:** Scatter plots of prediction vs ground truth for the elasticity property from the V-JEPA-2 model.

499 provide the model with more context information about the  
 500 properties. Please refer to Supplementary Section D  
 501 for the comparison results and analysis. All models are trained  
 502 on H100/A6000/A40 GPUs. Please refer to Supplementary  
 503 Section A for more implementation details.

504 **Evaluation Metrics.** For *relative value comparison*, we  
 505 report the ROC AUC score; for *absolute value prediction*, we  
 506 use the Pearson Correlation Coefficient between the predic-  
 507 tion and ground truth as this automatically calibrates the  
 508 predictions to the scale of the ground truth. Please refer  
 509 to Supplementary Section A.4 for more details and motiva-  
 510 tions on the evaluation metrics.

## 511 5.1. Results for Relative Value Comparison

512 Table 1 (left) shows relative value comparison results across  
 513 physical properties and model types. The oracle estimator  
 514 performs nearly perfectly on *test-1* and *test-2*, and  
 515 strongly on *test-3*, indicating that the task is largely solv-  
 516 able using visual cues, geometry, and physics. Both genera-  
 517 tive and self-supervised video models achieve strong results  
 518 on synthetic splits (*test-1* and *test-2*). Notably, they  
 519 can also generalize well to the real-world split (*test-3*)  
 520 for viscosity and elasticity, which rely on simple height ra-  
 521 tios and expansion. However, friction proves more chal-  
 522 lenging. Models trained on synthetic data struggle to gen-  
 523 eralize, likely due to the fact that reliance on visual refer-  
 524 ences (e.g., ground plane grids) is absent in real videos, and  
 525 due to friction’s inherent complexity involving higher-order  
 526 motion and projective geometry of the viewpoint. To fur-  
 527 ther confirm, we introduce an additional real-world train-  
 528 ing split for friction videos with disjoint objects and sur-  
 529 faces from the test set (see Supplementary Section B.2 for  
 530 more details). Fine-tuning the visual prompting architec-  
 531 ture on this data improves performance on the real test split,  
 532 as shown by the \* values in Table 1. Multimodal large  
 533 language models (MLLMs), though not working very well

with *Baseline Prompt* (see Supplementary Section D), when  
 534 prompted properly, also perform well, especially on real  
 535 videos, which are more *in-distribution* for them – while on  
 536 synthetic splits, their performance drops significantly. This  
 537 is likely due to the fact that the models tend to leverage se-  
 538 mantic cues, e.g., the type of liquid or the category of object  
 539 and surface, rather than visual motion. 540

## 541 5.2. Results for Absolute Value Prediction

542 Table 1 (right) shows results for absolute value predic-  
 543 tion across physical properties and methods. This task  
 544 is more challenging than relative comparison, as models  
 545 must regress quantitative physical values rather than com-  
 546 pare video pairs from the same viewpoint. Similar to the  
 547 relative setting, the oracle estimator achieves near-perfect  
 548 performance on *test-1* and *test-2*, and strong perform-  
 549 ance on *test-3*, confirming that the task is largely solv-  
 550 able through visual cues, multi-view geometry, and physi-  
 551 cal laws. We highlight several key observations: (i) **com-**  
 552 **parable performance across backbones.** Despite being  
 553 trained for generative tasks, video generative models per-  
 554 form on par with self-supervised models when predicting  
 555 dynamic physical properties. (ii) **friction remains chal-**  
 556 **lenging.** Similar to the relative setting, both generative and  
 557 self-supervised models struggle with friction estimation.  
 558 Performance again improves with domain adaptation. (iii)  
 559 **MLLMs better on real test split than synthetic.** MLLMs  
 560 continue to perform better on the real test split than syn-  
 561 thetic test splits, benefiting from their familiarity with real-  
 562 world visual semantics. (iv) **greater gap from oracle.** The  
 563 performance gap between video foundation models and the  
 564 oracle is more pronounced here than in the relative setting,  
 565 indicating that accurate physical value regression remains a  
 566 significant challenge for current video models. (v) **more**  
 567 **difficult to generalise to real situations.** Compared to  
 568 the relative setting, video generative model and video self-

Table 1. **Results for relative value comparison and absolute value prediction.** Left: ROC AUC scores for relative comparisons (range  $[0, 1]$ ). Right: Pearson correlation coefficients for absolute predictions (range  $[-1, 1]$ ). \* indicates results after domain adaptation using a disjoint real training set. `test-1` is the synthetic in-distribution test split; `test-2` is the synthetic out-of-distribution test split; `test-3` is the real-world test split.

Property	Method	Relative – ROC AUC			Absolute – Pearson Corr.		
		Test-1	Test-2	Test-3	Test-1	Test-2	Test-3
Elasticity	<b>Oracle</b>	1.00	1.00	1.00	0.99	0.98	0.87
	Video Generative Model	1.00	0.98	0.84	0.92	0.82	0.07
	Video Self-Supervised Model	0.89	0.96	0.77	0.96	0.93	0.47
	Qwen2.5VL-max	0.59	0.50	0.54	-0.05	0.11	0.16
	GPT-4o	0.51	0.66	0.62	0.19	0.11	0.30
	Gemini-2.5-pro	0.64	0.80	0.47	0.04	0.15	0.24
Viscosity	<b>Oracle</b>	0.99	1.00	1.00	0.99	0.98	0.80
	Video Generative Model	1.00	1.00	1.00	0.99	0.95	0.76
	Video Self-Supervised Model	1.00	1.00	0.99	1.00	0.97	0.79
	Qwen2.5VL-max	0.64	0.61	0.86	0.16	0.06	0.02
	GPT-4o	0.63	0.59	0.99	0.18	0.08	0.55
	Gemini-2.5-pro	0.48	0.69	0.95	-0.06	-0.05	0.60
Friction	<b>Oracle</b>	1.00	1.00	0.87	0.99	1.00	0.83
	Video Generative Model + Domain Adaptation	0.98	0.89	0.47	0.95	0.78	0.21
	Video Self-Supervised Model + Domain Adaptation	1.00	0.97	0.58	0.71	0.58	0.28
	Qwen2.5VL-max	0.50	0.62	0.80	0.03	0.14	0.06
	GPT-4o	0.34	0.42	0.67	-0.10	0.03	0.38
	Gemini-2.5-pro	0.54	0.59	0.97	-0.03	-0.05	0.12

569 supervised model exhibit a significantly larger performance  
570 drop from synthetic to real test splits (especially for elastic-  
571 ity and viscosity). This demonstrates the absolute setting is  
572 challenging.

### 573 5.3. Qualitative Results

574 Figure 4 (left) shows qualitative examples comparing model  
575 predictions across different tasks. In the **first row**, we il-  
576 lustrate an example from the elasticity absolute value pre-  
577 diction task. The video generative model, self-supervised  
578 model, and MLLMs predict values of 0.83, 0.96, and 1.0,  
579 respectively—all reasonably close to the ground-truth value  
580 of 0.89. In the **second row**, we present a friction relative  
581 value comparison task. The input consists of two videos,  
582 where the first exhibits a higher dynamic friction coefficient  
583 than the second. Both the video generative model and the  
584 MLLM correctly assign high likelihoods to this relationship  
585 (0.82 and 1.0, respectively), aligning with the ground truth.  
586 In contrast, the self-supervised model incorrectly predicts  
587 the reverse and does so with high confidence. Figure 4  
588 (right) shows examples of the scatter plots for the abso-  
589 lute value prediction. More specifically, we show the scat-  
590 ter plots of video self-supervised model on the three test  
591 splits. It can be observed that the performance degrades  
592 from `test-1` to `test-3`, as `test-1` is of the same dis-  
593 tribution as the synthetic training split, while `test-2` is  
594 out-of-distribution synthetic test and `test-3` is for real

evaluation. We provide more scatter plots in Supplemen-  
tary Section F.

## 597 6. Conclusion

598 We investigate the task of inferring dynamic physical  
599 properties—elasticity, viscosity, and friction—from videos.  
600 To support this, we introduce a benchmark dataset with  
601 ground-truth annotations and evaluate a range of video  
602 foundation models under both absolute prediction and rela-  
603 tive comparison settings. We adopt a simple architecture  
604 to extract physical cues from off-the-shelf generative and  
605 self-supervised video models, and explore prompting strate-  
606 gies to elicit predictions from MLLMs. Experiments show  
607 that generative and self-supervised models have similar and  
608 reasonable performance. This confirms the findings in pre-  
609 vious works, where the prior of video generative models  
610 can be used to estimate the material field of objects [57],  
611 and the video self-supervised models can be used for robots  
612 to interact with the physical world [3]. MLLMs perform  
613 worse overall but improve with more informative prompt-  
614 ing, especially on real-world data. The worse performance  
615 of MLLMs is consistent with previous work [13], where it  
616 is observed that the visual information is not properly fused  
617 in the language model. However, all models fall short of  
618 the oracle, particularly in absolute value prediction. These  
619 results highlight the need to enhance physical reasoning in  
620 video models—a key direction for future research.

621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676**References**

- [1] Pulkit Agrawal, Ashvin V Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to poke by poking: Experiential learning of intuitive physics. In *Advances on Neural Information Processing Systems (NeurIPS)*, 2016. 2
- [2] Jan Assen, Pascal Barla, and Roland Fleming. Visual features in the perception of liquids. *Current Biology*, 2018. 2
- [3] Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhohus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025. 1, 5, 8
- [4] Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction. *arXiv preprint arXiv:2106.13195*, 2021. 2
- [5] Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mido Assran, and Nicolas Ballas. V-jepa: Latent video prediction for visual representation learning. *OpenReview*, 2023. 1
- [6] Daniel M Bear, Elias Wang, Damian Mrowca, Felix J Binder, Hsiao-Yu Fish Tung, RT Pramod, Cameron Holdaway, Sirui Tao, Kevin Smith, Fan-Yun Sun, et al. Physion: Evaluating physical prediction from vision in humans and machines. *arXiv preprint arXiv:2106.08261*, 2021. 2
- [7] Florian Bordes, Quentin Garrido, Justine T Kao, Adina Williams, Michael Rabbat, and Emmanuel Dupoux. Intphys 2: Benchmarking intuitive physics understanding in complex synthetic environments. *arXiv preprint arXiv:2506.09849*, 2025.
- [8] Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Benchmarking and enhancing vision-language models for physical world understanding. *International Conference on Learning Representations (ICLR)*, 2025. 2
- [9] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 6
- [10] Mingyu Ding, Zhenfang Chen, Tao Du, Ping Luo, Josh Tenenbaum, and Chuang Gan. Dynamic visual reasoning by learning differentiable physics models from video and language. *Advances In Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [11] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 2
- [12] Katerina Fragkiadaki, Pulkit Agrawal, Sergey Levine, and Jitendra Malik. Learning visual predictive models of physics for playing billiards. In *International Conference on Learning Representations (ICLR)*, 2016. 2
- [13] Stephanie Fu, Tyler Bonnen, Devin Guillory, and Trevor Darrell. Hidden in plain sight: Vlms overlook their visual representations. *arXiv preprint arXiv:2506.08008*, 2025. 8
- [14] Jensen Gao, Bidipta Sarkar, Fei Xia, Ted Xiao, Jiajun Wu, Brian Ichter, Anirudha Majumdar, and Dorsa Sadigh. Physically grounded vision-language models for robotic manipulation. In *International Conference on Robotics and Automation (ICRA)*, 2024. 2
- [15] Alejandro Castañeda Garcia, Jan Warchocki, Jan van Gemert, Daan Brinks, and Nergis Tomen. Learning physics from video: Unsupervised physical parameter estimation for continuous dynamical systems. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, 2025. 2
- [16] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning (ICML)*, 2019. 2
- [17] Richard I. Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004. 4
- [18] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024. 1, 6
- [19] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1, 6
- [20] Krishna Murthy Jatavallabhula, Miles Macklin, Florian Golemo, Vikram Voleti, Linda Petrini, Martin Weiss, Brendan Considine, Jerome Parent-Levesque, Kevin Xie, Kenny Erleben, Liam Paull, Florian Shkurti, Derek Nowrouzezahrai, and Sanja Fidler. gradsim: Differentiable simulation for system identification and visuomotor control. In *International Conference on Learning Representations (ICLR)*, 2021. 2
- [21] Takahiro Kawabe and Shin'ya Nishida. Seeing jelly: Judging elasticity of a transparent object. In *Proceedings of the ACM Symposium on Applied Perception*, 2016. 2
- [22] Takahiro Kawabe, Kazushi Maruya, Roland Fleming, and Shin'ya Nishida. Seeing liquids from visual motion. *Vision Research*, 2014. 2
- [23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023. 12
- [24] Yunzhu Li, Toru Lin, Kexin Yi, Daniel Bear, Daniel L.K. Yamins, Jiajun Wu, Joshua B. Tenenbaum, and Antonio Torralba. Visual grounding of learned physical models. In *International Conference on Machine Learning (ICML)*, 2020. 2
- [25] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision (ECCV)*, 2024. 12

677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734

- 735 [26] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao,  
736 Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jian-  
737 feng Gao, et al. Sora: A review on background, technology,  
738 limitations, and opportunities of large vision models. *arXiv*  
739 *preprint arXiv:2402.17177*, 2024. 1
- 740 [27] Haoyu Lu, Guoxing Yang, Nanyi Fei, Yuqi Huo, Zhiwu  
741 Lu, Ping Luo, and Mingyu Ding. Vdt: An empirical  
742 study on video diffusion with transformers. *arXiv preprint*  
743 *arXiv:2305.13311*, 2023. 2
- 744 [28] J Norman, Elizabeth Wiesemann, Hideko Norman, M Tay-  
745 lor, and Warren Craft. The visual discrimination of bending.  
746 *Perception*, 2007. 2
- 747 [29] Vivian Paulun, Takahiro Kawabe, Shin’ya Nishida, and  
748 Roland Fleming. Seeing liquids from static snapshots. *Vi-*  
749 *sion research*, 2015.
- 750 [30] Vivian Paulun, Philipp Schmidt, Jan Assen, and Roland Flem-  
751 ing. Shape, motion, and optical cues to stiffness of elastic  
752 objects. *Journal of Vision*, 2017.
- 753 [31] Vivian C. Paulun and Roland W. Fleming. Visually infer-  
754 ring elasticity from the motion trajectory of bouncing cubes.  
755 *Journal of Vision*, 2020. 2
- 756 [32] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang  
757 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman  
758 Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junt-  
759 ing Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-  
760 Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feicht-  
761 enhofer. Sam 2: Segment anything in images and videos.  
762 In *International Conference on Learning Representations*  
763 *(ICLR)*, 2025. 12
- 764 [33] Tianhe Ren, Qing Jiang, Shilong Liu, Zhaoyang Zeng, Wen-  
765 long Liu, Han Gao, Hongjie Huang, Zhengyu Ma, Xiaoke  
766 Jiang, Yihao Chen, Yuda Xiong, Hao Zhang, Feng Li, Peijun  
767 Tang, Kent Yu, and Lei Zhang. Grounding dino 1.5: Ad-  
768 vance the ”edge” of open-set object detection. *arXiv preprint*  
769 *arXiv:2405.10300*, 2024.
- 770 [34] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-  
771 chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen,  
772 Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang,  
773 Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam:  
774 Assembling open-world models for diverse visual tasks.  
775 *arXiv preprint arXiv:2401.14159*, 2024. 12
- 776 [35] Ronan Riochet, Mario Ynocente Castro, Mathieu Bernard,  
777 Adam Lerer, Rob Fergus, Véronique Izard, and Emmanuel  
778 Dupoux. Intphys: A framework and benchmark for visual in-  
779 tuitive physics reasoning. *arXiv preprint arXiv:1803.07616*,  
780 2018. 2
- 781 [36] Prafull Sharma, Julien Philip, Michaël Gharbi, Bill Freeman,  
782 Fredo Durand, and Valentin Deschaintre. Materialistic: Se-  
783 lecting similar materials in images. *ACM Transactions on*  
784 *Graphics (TOG)*, 2023. 2
- 785 [37] Hui Shen, Taiqiang Wu, Qi Han, Yunta Hsieh, Jizhou Wang,  
786 Yuyue Zhang, Yuxin Cheng, Zijian Hao, Yuansheng Ni, Xin  
787 Wang, et al. Phyx: Does your model have the ”wits” for  
788 physical reasoning? *arXiv preprint arXiv:2505.15929*, 2025.  
789 2
- 790 [38] Aleksandar Shtedritski, Christian Rupprecht, and Andrea  
791 Vedaldi. What does clip know about a red circle? vi-  
792 sual prompt engineering for vlms. In *Proceedings of the*  
*IEEE/CVF International Conference on Computer Vision*  
*(ICCV)*, 2023. 6
- [39] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob  
Fergus. Indoor segmentation and support inference from  
rgbd images. In *European Conference on Computer Vision*  
*(ECCV)*, 2012. 2
- [40] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng  
Phoo, and Bharath Hariharan. Emergent correspondence  
from image diffusion. *Advances in Neural Information Pro-*  
*cessing Systems (NeurIPS)*, 2023. 5
- [41] Hsiao-Yu Tung, Mingyu Ding, Zhenfang Chen, Daniel Bear,  
Chuang Gan, Josh Tenenbaum, Dan Yamins, Judith Fan, and  
Kevin Smith. Physion++: Evaluating physical scene under-  
standing that requires online inference of different physical  
properties. *Advances in Neural Information Processing Sys-*  
*tems (NeurIPS)*, 2023. 2
- [42] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher  
Pal. Masked conditional video diffusion for prediction, gen-  
eration, and interpolation. *arXiv preprint arXiv:2205.09853*,  
2022. 2
- [43] Bin Wang, Paul Kry, Yuanmin Deng, Uri Ascher, Hui Huang,  
and Baoquan Chen. Neural material: Learning elastic consti-  
tutive material and damping models from sparse data. *arXiv*  
*preprint arXiv:1808.04931*, 2018. 2
- [44] Kun Wang, Mridul Aanjaneya, and Kostas Bekris. A first  
principles approach for data-efficient system identification of  
spring-rod systems via differentiable physics engines. In  
*Learning for Dynamics and Control*, 2020. 2
- [45] Tianyu Wang, Xiaowei Hu, Qiong Wang, Pheng-Ann Heng,  
and Chi-Wing Fu. Instance shadow detection. In *Proceed-*  
*ings of the IEEE Conference on Computer Vision and Pattern*  
*Recognition (CVPR)*, 2020. 2
- [46] Tianyu Wang, Xiaowei Hu, Chi-Wing Fu, and Pheng-Ann  
Heng. Single-stage instance shadow detection with bidirec-  
tional relation learning. In *Proceedings of the IEEE Confer-*  
*ence on Computer Vision and Pattern Recognition (CVPR)*,  
2021. 2
- [47] Wikipedia contributors. Coefficient of restitution —  
wikipedia, the free encyclopedia, 2025. 17
- [48] Wikipedia contributors. Viscosity — wikipedia, the free en-  
cyclopedia, 2025. 17
- [49] Wikipedia contributors. Wetting — wikipedia, the free en-  
cyclopedia, 2025. 17
- [50] Jiajun Wu, Ilker Yildirim, Joseph J Lim, Bill Freeman, and  
Josh Tenenbaum. Galileo: Perceiving physical object prop-  
erties by integrating a physics engine with deep learning. *Ad-*  
*vances in neural information processing systems (NeurIPS)*,  
2015. 2
- [51] Jiajun Wu, Joseph J Lim, Hongyi Zhang, Joshua B Tenen-  
baum, and William T Freeman. Physics 101: Learning phys-  
ical object properties from unlabeled videos. In *British Ma-*  
*chine Vision Conference (BMVC)*, 2016. 2
- [52] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen,  
Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying  
Shan, and Tien-Tsin Wong. Dynamicrafter: Animating  
open-domain images with video diffusion priors. In *Euro-*  
*pean Conference on Computer Vision (ECCV)*, 2024. 1, 4

- 850 [53] Guanqi Zhan, Weidi Xie, and Andrew Zisserman. A tri-layer  
851 plugin to improve occluded detection. *British Machine Vi-*  
852 *sion Conference (BMVC)*, 2022. 2
- 853 [54] Guanqi Zhan, Chuanxia Zheng, Weidi Xie, and Andrew Zis-  
854 serman. Amodal ground truth and completion in the wild.  
855 In *Proceedings of the IEEE/CVF Conference on Computer*  
856 *Vision and Pattern Recognition (CVPR)*, 2024. 2
- 857 [55] Guanqi Zhan, Chuanxia Zheng, Weidi Xie, and Andrew Zis-  
858 serman. A general protocol to probe large vision models for  
859 3d physical understanding. *Advances in Neural Information*  
860 *Processing Systems (NeurIPS)*, 2024. 5
- 861 [56] Chenyu Zhang, Daniil Cherniavskii, Antonios Tragoudaras,  
862 Antonios Vozikis, Thijmen Nijdam, Derck WE Prinzhorn,  
863 Mark Bodraczka, Nicu Sebe, Andrii Zadaianchuk, and Efs-  
864 tratios Gavves. Morpheus: Benchmarking physical reason-  
865 ing of video generative models with real physical experi-  
866 ments. *arXiv preprint arXiv:2504.02918*, 2025. 2
- 867 [57] Tianyuan Zhang, Hong-Xing Yu, Rundi Wu, Brandon Y  
868 Feng, Changxi Zheng, Noah Snavely, Jiajun Wu, and  
869 William T Freeman. Physdreamer: Physics-based interac-  
870 tion with 3d objects via video generation. In *European Con-*  
871 *ference on Computer Vision (ECCV)*, 2024. 8
- 872 [58] Xian Zhou, Yiling Qiao, Zhenjia Xu, Tsun-Hsuan Wang,  
873 Zhehuan Chen, Juntian Zheng, Ziyan Xiong, Yian Wang,  
874 Mingrui Zhang, Pingchuan Ma, Yufei Wang, Zhiyang Dou,  
875 Byungchul Kim, Yunsheng Tian, Yipu Chen, Xiaowen Qiu,  
876 Chunru Lin, Tairan He, Zilin Si, Yunchu Zhang, Zhanlue  
877 Yang, Tiantian Liu, Tianyu Li, Kashu Yamazaki, Hongxin  
878 Zhang, Huy Ha, Yu Zhang, Michael Liu, Shaokun Zheng,  
879 Zipeng Fu, Qi Wu, Yiran Geng, Feng Chen, Milky, Yuan-  
880 ming Hu, Guanya Shi, Lingjie Liu, Taku Komura, Zack-  
881 ory Erickson, David Held, Minchen Li, Linxi "Jim" Fan,  
882 Yuke Zhu, Wojciech Matusik, Dan Gutfreund, Shuran Song,  
883 Daniela Rus, Ming Lin, Bo Zhu, Katerina Fragkiadaki, and  
884 Chuang Gan. Genesis: A universal and generative physics  
885 engine for robotics and beyond, 2024. 2, 3