

Δ YNAMICS: Language-Based Representation for Inferring Rigid-Body Dynamics From Videos

Anonymous CVPR submission

Paper ID 6659

Abstract

001 *Inferring rigid-body physical states and properties from*
 002 *monocular videos is a fundamental step toward physics-*
 003 *based perception and simulation. Existing approaches as-*
 004 *sume specific underlying physical systems, object types, and*
 005 *camera poses, which are unable to generalize to complex*
 006 *real-world settings. We introduce Δ YNAMICS, a vision-*
 007 *language framework that uses language as a unified rep-*
 008 *resentation of rigid-body dynamics. Instead of directly*
 009 *predicting parameters, Δ YNAMICS generates scene con-*
 010 *figurations in a structured text format for physics simula-*
 011 *tion. We enhance the model’s generalization by integrat-*
 012 *ing natural language motion reasoning and leveraging op-*
 013 *tical flow as a semantic-agnostic input. On the CLEVRER*
 014 *dataset [59], Δ YNAMICS achieves a segmentation IoU of*
 015 *0.30, a $7\times$ improvement over leading VLMs (InternVL3-*
 016 *8B, Qwen2.5-VL-7B and Claude-4-Sonnet). Further, test-*
 017 *time sampling and evolutionary search further boost perfor-*
 018 *mance by 27% and 120% in segmentation IoU, respectively.*
 019 *Finally, we demonstrate strong transfer to a new dataset of*
 020 *235 real-world rigid-body videos, highlighting the poten-*
 021 *tial of language-driven physics inference for bridging per-*
 022 *ception and simulation. Additional results and videos are*
 023 *available in the supplementary material.*

024 1. Introduction

025 Understanding physical dynamics from visual observations
 026 is a foundational capability for intelligent systems operating
 027 in the real world [8, 20, 30, 51]. When perceiving events
 028 such as a ball sliding or bouncing, the system should not
 029 only identify and track objects but also infer their intrinsic
 030 physical attributes, including friction, elasticity, and other
 031 parameters that govern motion. These inferred properties
 032 enable reasoning about cause and effect, anticipating out-
 033 comes under varying conditions, and consequently planning
 034 and control in embodied settings.

035 In this work, we focus on rigid-body motion dynamics.
 036 Given a video, our goal is to infer the underlying physical

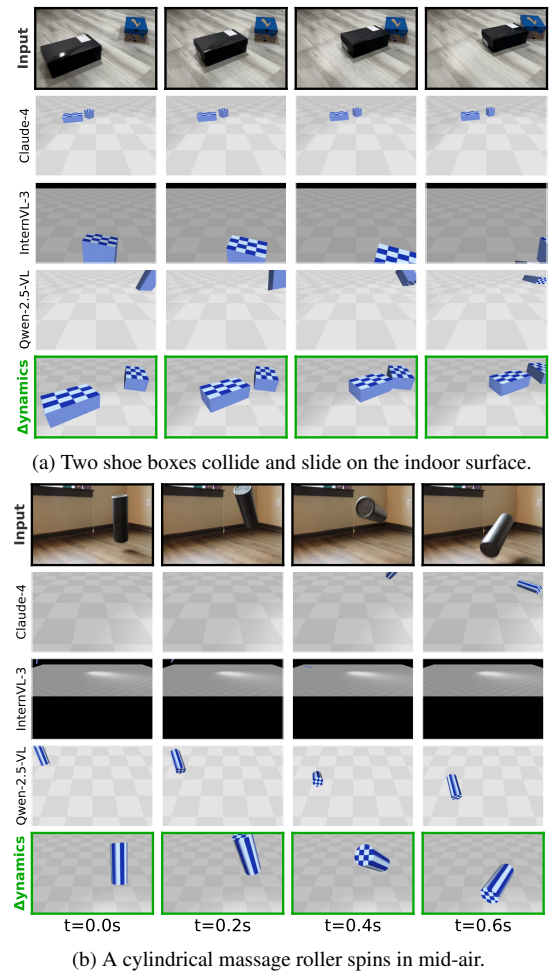


Figure 1. **Motion transfer from real videos to simulation environments.** Δ YNAMICS accurately reproduces the object shapes, initial position and orientation, material properties, and camera pose with respect to the input videos, while competing VLMs (Claude-4-Sonnet, InternVL-3-8B, Qwen-2.5-VL-7B) fail.

model that allows the reproduction of the video trajectories within a simulator. While prior works [4, 14, 19, 23, 26, 53, 54] have made progress in estimating physics parameters from videos in constrained settings, such as slid-

037
 038
 039
 040

041 ing boxes [19, 53], billiards [54], or projectiles [4, 14, 23],
042 they are not yet applicable to complex real-world motion
043 that involves multiple object interactions and motion types.
044 First, they assume a model-specific, fixed-length vector of
045 physics parameters for particular object types (e.g., spheres
046 or boxes) and motion types (e.g., sliding or projectile). This
047 representation is not scalable and does not accommodate the
048 full variety of object interactions. Second, prior works typ-
049 ically assume a known or fixed camera pose, thus failing to
050 generalize across varying object distances and camera view-
051 points. As a result, prior methods solve only a narrow subset
052 of this video-to-simulation problem and fail to generalize to
053 complex, real-world scenes involving multiple moving ob-
054 jects and unconstrained viewpoints.

055 A fundamental challenge lies in how the scene itself
056 is parameterized. In this work, we introduce *a unified,*
057 *language-based representation of rigid body motion* as a
058 bridge between perception and simulation. Instead of re-
059 gressing a fixed-length numeric vector, we reformulate the
060 problem as the generation of symbolic scene configurations
061 specifying object geometry, initial states, material proper-
062 ties, and camera parameters. This language representation
063 is inherently *interpretable* and *scalable* to diverse motion
064 types and object interactions.

065 This language-centric formulation naturally motivates
066 the use of Vision-Language Models (VLMs) [16, 18, 34,
067 35, 52], widely employed for visual reasoning [2, 36] and
068 physics understanding [5, 7, 17, 39, 41]. Taking this direc-
069 tion, we develop Δ YNAMICS, a VLM trained on 400K syn-
070 thetic videos rendered with MuJoCo [50], whose output is a
071 YAML format of the scene configuration. To enhance gen-
072 eralization, we make two key design choices. First, we take
073 optical flow as the input as it is agnostic to visual semantics
074 and background, which provides explicit motion cues and
075 improves full-sequence segmentation IoU by 26% (from
076 0.19 to 0.24) on CLEVRER [59]. Second, we augment su-
077 pervision with *natural-language motion descriptions* that
078 capture trajectories, object visibility, and collision events
079 as an auxiliary textual target. Together, these components
080 make Δ YNAMICS robust to domain shifts.

081 To evaluate cross-domain generalization, we adapt
082 CLEVRER for controlled testing and curate a new dataset
083 of 235 real-world rigid-body motion videos. The reasoning-
084 enhanced model consistently outperforms the vanilla ver-
085 sion on real-world transfer, indicating stronger generaliza-
086 tion capabilities. We also investigate several test-time en-
087 hancement strategies that do not require labeled ground
088 truth in the target domain. We find that best-of-k sampling
089 consistently yields a 10% improvement, and that an addi-
090 tional evolutionary search provides over 50% further gains.
091 For real-world application, we also show the potential of
092 physically plausible video editing using our framework; the
093 corresponding results are deferred to Appendix D.

Our main contributions are summarized below:

- 094 • **Language-based representation for motion dynam-** 095
096 **ics:** We reformulate rigid object motion estimation from 097
098 videos as a *language modeling* problem, where the model 099
100 generates structural textual scene configurations that are 101
102 directly consumable by a physics engine. 103
- 104 • **VLM for rigid-body physics inference:** We present 105
106 Δ YNAMICS, a VLM that directly infers the underlying 107
108 physics parameters of rigid object motion, which enables 109
110 the reconstruction of physically-plausible motion trajec- 111
112 tories from monocular videos. 113
- 114 • **Cross-domain generalization:** We boost the general- 115
116 ization of models for different physics engines and real- 117
118 world videos by introducing two key innovations: using 119
120 optical flow as semantics-agnostic input and training the 121
122 model to predict natural-language motion descriptions. 123
- 124 • **Comprehensive evaluation benchmarks:** We adapt the 125
126 CLEVRER dataset for controlled benchmarking and cu- 127
128 rate a new dataset of 235 real-world rigid-body motion 129
130 videos with corresponding annotations for segmentation 131
132 masks and optical flows. 133

134 2. Related Work 135

136 **Rigid-Body Motion Parameter Estimation.** Estimating 137
138 physics parameters of rigid moving objects and camera ge- 139
140 ometry from videos is a key step towards physics-based 141
142 perception and simulation. Early efforts tackled the prob- 143
144 lem in narrow scenarios, e.g., sliding boxes [19, 53], bil- 145
146 liard games [54], projectile motion [4, 14], articulated rigid 147
148 body [26], or free fall [23]), to keep the parameter esti- 149
150 mation problem tractable. Furthermore, they assume fixed 151
152 camera parameters, preventing their applications from gen- 153
154 eral real-world settings [4, 14, 19, 23, 26, 53, 54]. Our con- 155
156 tribution is a general solution to the physics parameter esti- 157
158 mation problem, which is applicable to a wide spectrum of 159
160 physical motions in unconstrained real-world settings. 161

162 **Structured Representations for Visual Content.** Repre- 163
164 senting image and video content using structured graphics 165
166 programs has been widely adopted for graphics simulation 167
168 engines such as Blender and MuJoCo format [50]. Recent 169
170 works in image simulation and generation make use of pro- 171
172 grammatic formats such as SVG [15, 42, 46, 55–57] and 173
174 TikZ [9, 10, 48] to formulate the problem as conditional 175
176 generation of structured text based on textual and image 177
178 prompts using diffusion models. Further, the use of struc- 179
180 tured graphics programs has also been extended to the do- 181
182 main of inverse rendering and the editing and generation of 183
184 3D scenes. Specifically, the workflow in [11, 31, 32, 60] in- 185
186 volves training VLMs to translate images into a structured 187
188 format (e.g., JSON) and employing graphic engines for ren- 189
190 dering. Other works train VLMs to infer the programmatic 191
192 representation of existing scenes for graphics editing [24] 193

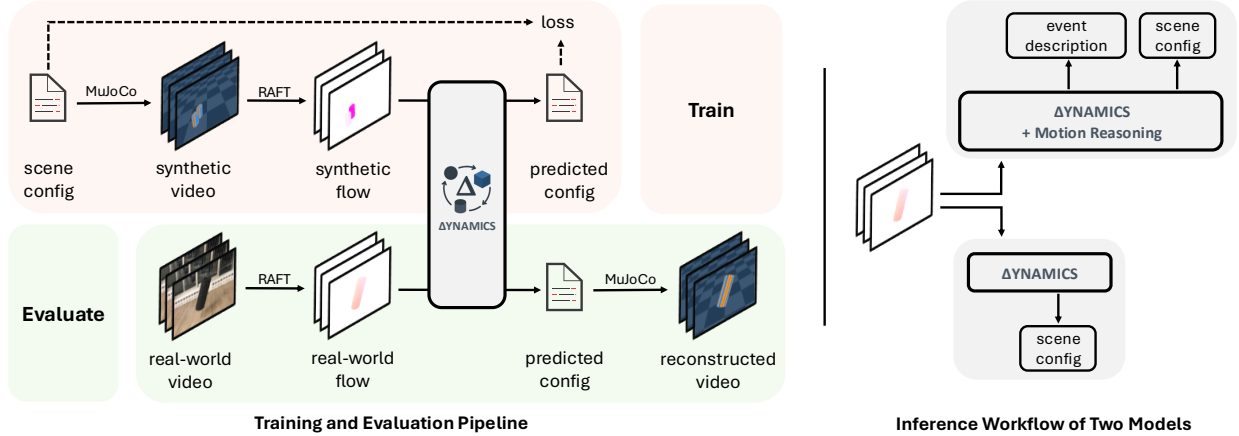


Figure 2. **Training, evaluation and inference workflow for Δ YNAMICS .** **Training (top left):** We sample scene configurations and render corresponding synthetic videos using the MuJoCo physics engine. Next, we compute optical flows using RAFT [49] and train Δ YNAMICS to generate scene configurations in a structured text format given optical flows. **Evaluation (bottom left):** Δ YNAMICS takes input optical flows derived from real-world videos to infer scene configurations. **Inference (right):** The base variant (bottom) directly generates the scene configuration, whereas the motion reasoning variant (top) first generates a motion event description (details illustrated in Figure 3), then predicts the scene configuration.

145 and 3D asset creation [61]. Our proposed approach is in-
 146 spired by this school of thought because it offers the capa-
 147 bility to interpret both the underlying physics and control-
 148 lable editing of visual content via scene attributes. However,
 149 our research problem is distinguished from prior works by
 150 the focus on the modeling of motion dynamics in videos.

151 **Physics Simulation.** Another line of research focuses on
 152 end-to-end training with differentiable simulation and ren-
 153 dering pipelines [21, 27, 28], enabling physical scene un-
 154 derstanding via gradient-based optimization. These works [29,
 155 37, 38, 47, 58] jointly optimize object geometry and phys-
 156 ical parameters to directly match the scene dynamics and
 157 image formation. Our approach differs in three ways. First,
 158 we do not require a predefined physics model; instead, our
 159 model learns to infer physical properties and dynamics di-
 160 rectly from video observations. Second, we do not assume
 161 access to differentiable simulators, as most physics engines
 162 are not differentiable in practice. Third, instead of per-
 163 scene optimization, we employ a direct feedforward model
 164 that predicts a complete scene configuration in a single pass.

165 3. Method

166 We present the training, evaluation and inference workflows
 167 of Δ YNAMICS in Figure 2. We now formalize the problem
 168 and describe each component in detail.

169 3.1. Problem Statement

170 We address the problem of recovering physical scene pa-
 171 rameters and dynamics from a monocular video. Given
 172 an input video \mathbf{X} , a model \mathcal{F}_θ predicts a parameter set
 173 $\mathbf{c} = \mathcal{F}_\theta(\mathbf{X})$, which is then provided to a physics engine

174 \mathcal{S} to generate a reconstructed sequence $\hat{\mathbf{X}} = \mathcal{S}(\mathbf{c})$. The ob-
 175 jective is to learn \mathcal{F}_θ such that the simulated dynamics in $\hat{\mathbf{X}}$
 176 faithfully reproduce those observed in the input video.

177 3.2. Unified Scene Representation

178 A fundamental challenge lies in how the scene itself is pa-
 179 rameterized. Prior work typically regresses a fixed-length
 180 parameter vector specific to a particular object set, simula-
 181 tion model, or physics system, which limits generalization
 182 across real-world scenarios.

183 **Language Representation.** To address these limitations,
 184 we shift the core paradigm from numerical regression to
 185 symbolic generation. The key idea is a unified, language-
 186 based representation that acts as a bridge between percep-
 187 tion and simulation. Specifically, we recast physics esti-
 188 mation as a text-generation problem: the model outputs a
 189 YAML-formatted sequence that encodes the entire scene
 190 configuration, including object geometry, initial states, ma-
 191 terial properties, and camera parameters. This provides
 192 three key advantages:

- 193 • **Extensibility and Interpretability.** A textual format
 194 scales naturally to scenes with arbitrary numbers of ob-
 195 jects. It is human-readable, easy to edit, and conducive
 196 to counterfactual analysis. Extended results on physically
 197 plausible video editing are provided in Appendix D.
- 198 • **Natural Integration with VLM.** Casting simulation as
 199 text generation enables end-to-end training of a unified
 200 VLM without engineering multi-stage components. We
 201 simply format the target as `<answer> configuration`
 202 `</answer>`, allowing the model to directly output the
 203 scene description.

Table 1. **Parameter categories.** The complete parameter space spans object properties, initial states, and global parameters.

Category	Parameters
Object Property	<i>Geometry / Inertial:</i> radius, height, width, depth, mass. <i>Material:</i> friction (rolling, sliding) and damping.
Initial State	<i>Kinematics:</i> position, linear and angular velocity. <i>Orientation:</i> quaternion.
Global Parameter	<i>Camera:</i> pose (height, angle, FOV). <i>Environment:</i> gravity.

204 • **Joint Reasoning and Configuration Generation.** Lan-
205 guage models can interleave descriptive reasoning with
206 configuration prediction, enabling richer intermediate
207 representations within the same autoregressive process.

208 **Parameterization Details.** To operationalize this
209 language-based approach, we define a structured schema
210 for the scene configuration. This configuration represents
211 the full set of geometry, physics, and camera parameters re-
212 quired for simulation, as summarized in Table 1. We com-
213 pose our scenes using three primitive shapes (spheres, cylin-
214 ders, and boxes) that can cover common household objects
215 such as tennis balls, soda cans, mugs, books, and crates.
216 These primitives are sufficient to simulate essential rigid-
217 body dynamics, including bouncing, rolling, sliding, and
218 collisions. For the camera, we place it at $(0, -2, h)$, where
219 h denotes its height, and vary the pitch angle while setting
220 roll and yaw to zero. We include gravity as a parameter
221 to account for variations in frame rate or time scale. This
222 scene configuration format supports both single-object and
223 arbitrary multi-object scenes. For example, a scene contain-
224 ing four box-shaped objects includes 20×4 box-specific
225 parameters, along with 3 camera parameters and 1 gravity
226 term, totaling 84 parameters to be estimated.

227 3.3. Motion Reasoning

228 By reasoning about motion and object interactions before
229 predicting scene parameters, the model learn richer repre-
230 sentations of the underlying dynamics, which in turn im-
231 prove the accuracy of the subsequent scene parameter es-
232 timates. To enable motion reasoning, we train a variant of
233 the model that first generates a natural language descrip-
234 tion of the observed dynamics and then produces the scene
235 configuration. As shown in Figure 3, these descriptions
236 are derived from simulation traces and artifacts (i.e., ob-
237 ject state histories, contact logs, segmentation masks) to-
238 gether with ground-truth configurations. We specifically
239 consider events such as visibility (e.g., when the object en-
240 ters or leaves the camera view), motion change (e.g., when
241 it stops rolling or sliding), and collisions (e.g., when it
242 touches the ground or another object), and we design rule-
243 based functions to parse them. These detected events are

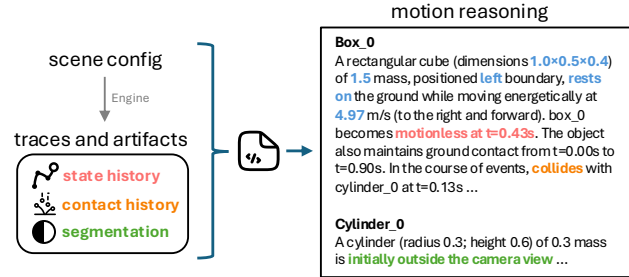


Figure 3. **Synthetic training data generation.** During the data generation process, we create natural language descriptions of motion events. An event-mining script processes the simulation traces and artifacts (left), including state history, contact history, and segmentation maps, to find key dynamic events. The resulting textual descriptions (right) serve as ground-truth targets for the motion reasoning model during training.

then inserted into predefined templates to generate a struc- 244
tured, natural language description of the motion. Finally, 245
we prepend the motion description to the scene configura- 246
tion as `<think> description </think> <answer>` 247
`configuration </answer>`. Table 7 provides an illus- 248
trative example. 249

250 3.4. Motion-Aware Input Representation

251 Raw RGB videos contain visual semantics unrelated to 251
252 motion, which can introduce confounding factors for our 252
253 model. As an alternative, we use optical flow fields to 253
254 represent motion as it is agnostic to visual semantics and 254
255 appearance. Specifically, we compute optical flows using 255
256 RAFT [49] and convert them into a 2D array per color chan- 256
257 nel, which can further be fed into VLM without architec- 257
258 tural changes. We evaluate models trained with both RGB 258
259 video inputs and the RGB-transformed optical flow maps. 259

260 4. Training and Evaluation Method

261 4.1. Training Approach

262 **Synthetic Data Curation.** We generate a dataset of 400K 262
263 unique physical scenes using the MuJoCo simulator [50]. 263
264 Each training example is created by sampling a full YAML 264
265 scene configuration, which is then converted into MuJoCo’s 265
266 XML format to initialize the simulator and assign dynamic 266
267 states (e.g., initial velocities), as detailed in Appendix A.1. 267
268 Each data point includes a rendered RGB video of a scene 268
269 with up to four objects and the corresponding YAML 269
270 file. To specifically test for compositional generalization, 270
271 four distinct object-type combinations in four-object scenes 271
272 (e.g., two boxes and two cylinders) are also held out. 272

273 To ensure physically plausible and visually meaningful 273
274 interactions, we filter out (i) scenes with overlapping objects 274
275 at initialization, identified using MuJoCo’s built-in collision 275
276 detection; (ii) scenes where more than one object remains 276

277 outside the camera’s field of view throughout simulation;
278 and (iii) scenes where any object is too small (with a total
279 area less than 8000 pixels). Simulations are run for 1 second
280 at 30 FPS with eight physics substeps per frame, rendering
281 images at 480×320 resolution.

282 **Learning Objective.** Let $\mathcal{D} = \{(\mathbf{X}_i, \mathbf{c}_i)\}_{i=1}^N$ denote the
283 training set, where \mathbf{X}_i is a video observation and \mathbf{c}_i is the
284 corresponding scene configuration, represented as a tok-
285 enized text sequence. We aim to learn a model \mathcal{F}_θ that max-
286 imizes the conditional likelihood $p_\theta(\mathbf{c} | \mathbf{X})$. This likelihood
287 is modeled autoregressively over the tokens of \mathbf{c} as

$$288 \quad p_\theta(\mathbf{c} | \mathbf{X}) = \prod_{t=1}^{|\mathbf{c}|} p_\theta(c_t | \mathbf{X}, c_{<t}), \quad (1)$$

289 where c_t denotes the t -th token of \mathbf{c} . We train \mathcal{F}_θ by mini-
290 mizing the negative log-likelihood (NLL) over the dataset:

$$291 \quad \mathcal{L}_{\text{VLM}} = - \sum_{(\mathbf{X}, \mathbf{c}) \in \mathcal{D}} \log p_\theta(\mathbf{c} | \mathbf{X}). \quad (2)$$

292 **Model and Training Implementation.** Our architecture
293 is based on Qwen2.5-VL-3B [6]. The inputs consist of 10
294 frames uniformly sampled from 1-second, 30 FPS videos.
295 We train two variants, one predicting only scene configura-
296 tion, and the other generating motion reasoning as an addi-
297 tional output. The format of the target text sequences for
298 these variants has been described in Sections 3.2 and 3.3.
299 We fine-tune the full model for 10 epochs in bfloat16 mixed
300 precision on eight 40 GB A100 GPUs. We employ the
301 AdamW optimizer with a learning rate of 2×10^{-5} , weight
302 decay of 0.01, and a global batch size of 128.

303 4.2. Test-Time Strategy

304 To improve instance-level accuracy at inference time, we
305 explore three complementary test-time optimization strate-
306 gies: best-of-K sampling, preference-based refinement, and
307 evolutionary search.

308 **Best-of-K Sampling.** The greedy decoding strategy in
309 VLMs is not guaranteed to find the parameter set with the
310 best quality, as the optimal parameter set may lie in the long
311 tail of the model’s output distribution. Hence, we adopt a
312 best-of- N evaluation scheme to explore this distribution:
313 for each case, we generate $N = 32$ diverse predictions
314 with a temperature of 0.1 and top-p of 0.9 and report the
315 *Best@32* performance. This reflects the model’s ability to
316 recover accurate physical dynamics with multiple attempts.

317 **Preference Optimization.** While ground-truth configura-
318 tions are typically unavailable for novel environments, the
319 similarities between the forward rendering and the input
320 videos, such as object mask IoU, can serve as *implicit re-*
321 *ward signals* for configuration selection without explicit su-
322 pervision. We conduct experiments with preference rank

323 optimization [45], where details and relevant results are pro-
324 vided in Appendix A.4.

325 **Evolutionary Search.** Since preference optimization gener-
326 ally requires training data, it is generally impractical in
327 real-world scenarios. Thus, we explore Covariance Mat-
328 rix Adaptation Evolution Strategy (CMA-ES) [25], which
329 is an evolutionary algorithm suited for non-convex black-
330 box optimization. When using CMA-ES, we initialize
331 the search with the *Best@32* sample, and we optimize
332 scene configurations, including object sizes, initial states,
333 physics parameters, and camera poses, while keeping ob-
334 ject types fixed. We employ a heuristic fitness function that
335 maximizes the segmentation Intersection-over-Union (IoU)
336 while minimizing the optical flow end-point error (EPE),
337 formulated as (IoU – EPE). We use a population size of
338 128 and optimize for 100 iterations.

339 4.3. Evaluation Method

340 **Evaluation via Simulation.** We evaluate the quality of the
341 predicted configuration $\hat{\mathbf{c}}$ through re-simulation. The gener-
342 ated text is passed to the physics engine \mathcal{S} to produce a sim-
343 ulated RGB video $\hat{\mathbf{X}} = \mathcal{S}(\hat{\mathbf{c}})$, along with auxiliary outputs
344 such as object segmentation masks and optical flow fields.
345 We then compare the simulated outputs against the ground-
346 truth counterparts from \mathbf{X} using segmentation Intersection-
347 over-Union (IoU) and optical flow end-point error (EPE).
348 For synthetic data, ground-truth masks are available from
349 the renderer, while for real-world videos, we use pretrained
350 models [40, 49] for pseudo-annotation.

351 **Metrics.** We evaluate across three dimensions:

- 352 • Object Composition: Accuracy of object composition.
- 353 • Motion Reconstruction Quality: Similarity between the
354 re-simulated and reference videos, measured by segmen-
355 tation IoU and flow EPE.
- 356 • Physics Parameter: L_1 distance between estimated and
357 ground-truth parameters, computed only when the object
358 combination is correct.

359 **Baselines.** Since our work is the first to estimate a com-
360 plete scene configuration for diverse physical systems from
361 a single monocular video, it is not directly comparable
362 to prior methods on physics parameter estimation. For
363 evaluation, we establish baselines using both proprietary
364 and open-source vision–language models (VLMs), includ-
365 ing InternVL3-8B [63], Qwen2.5-VL-7B [6], and Claude-
366 4-Sonnet [1]. Each model is evaluated using *three-shot*
367 *in-context learning (ICL)*. We adopt ICL over zero-shot
368 prompting because (1) zero-shot generation of a MuJoCo
369 XML file is insufficient for running a simulation, since dy-
370 namic states can only be set after engine initialization; and
371 (2) Generating both the XML file and the dynamic-state ini-
372 tialization code is difficult. Details about a few-shot exam-
373 ples are deferred to Appendix A.3.

374 We also establish non-VLM baselines that directly pre-
375 dict scene configuration parameters from videos. We con-
376 catenate these parameters into a fixed-length vector, repre-
377 sent the object type with one-hot encoding, and perform
378 zero-padding for missing objects. Objects are ordered by
379 their x - and then y -positions. We adopt the pretrained
380 ViViT [3] model, designed for video classification, and fully
381 fine-tune it using an ℓ_2 loss on the regression targets.

382 5. Results on Synthetic Dataset

383 We divide the synthetic data into two subsets and evaluate
384 the model in the following settings

- 385 1. Comparative evaluation: Scenes containing 1–3 objects,
386 with 100 samples for each object count.
- 387 2. Complex scene dynamics: 400 four-object scenes con-
388 structed from the four specific object-type combinations
389 that were excluded from the training set. 400 scenes with
390 five objects and 400 scenes with six objects, which ex-
391 ceed the object counts seen during training.

392 5.1. Comparison with Baselines

393 As shown in Table 2, Claude-4 is the best model among the
394 baselines. While sufficient to roughly identify object com-
395 position in a scene, they perform poorly in reconstructing
396 motion trajectories and estimating the physics parameters,
397 with low segmentation IoU (≤ 0.09), and high optical flow
398 errors (> 11) observed.

399 Meanwhile, the RGB-based Δ YNAMICS model outper-
400 forms all baselines in object composition accuracy, segmen-
401 tation IoU, and parameter estimation, but underperforms in
402 optical flow end-point error (EPE). The higher EPE is pri-
403 marily due to occasional interpenetration in the predicted
404 initial states, which causes MuJoCo to apply large correc-
405 tive contact forces to separate overlapping objects, resulting
406 in abrupt motions that increase flow error.

407 When explicitly conditioned on optical flow,
408 Δ YNAMICS achieves 97% object composition accu-
409 racy, improves segmentation IoU, and substantially reduces
410 EPE. One exception is the damping estimation, where the
411 raw-RGB model performs slightly better, likely because
412 the checkerboard ground plane provides additional visual
413 cues helpful to estimating damping parameters. Finally,
414 adding motion reasoning, as shown in the last row, further
415 improves overall performance.

416 5.2. Evaluation on Complex Scenes

417 Next, we evaluate the model’s generalization ability on un-
418 seen scene configurations. In particular, we focus on the be-
419 havior of Δ YNAMICS variant with motion reasoning (which
420 is the best one based on the previous section). In Table 3, the
421 model shows only a marginal degradation of segmentation
422 map IoU (compared to the last row of Table 2) for scenes
423 with 4 or 5 objects. Even for scenes with 6 objects, the

degradation is gradual and slow. These results show that in-
corporating motion reasoning adds robustness to more com-
plex, unseen multi-object dynamics.

6. Cross-Engine and Real-World Results

6.1. Cross-Engine Generalization

We assess our model’s ability to generalize across a fun-
damentally different simulation and rendering engine. In
particular, we perform this evaluation on the CLEVRER
dataset [59]. CLEVRER is a video question-answering
benchmark rendered by Blender, which covers sliding mo-
tion dynamics for three object types: cubes, spheres, and
cylinders. We sampled 100 test videos from the CLEVRER
for evaluation. To evaluate against the baseline VLMs, we
adopt a similar few-shot, in-context prompting approach
(details in Appendix A.3).

Comparison with Baselines In Table 4, Δ YNAMICS con-
sistently demonstrates superior performance compared to
baseline models. Furthermore, the previous observations
on our synthetic dataset hold true: (i) the model using op-
tical flow inputs outperforms that using RGB videos, and
(ii) incorporating reasoning further increases motion recon-
struction accuracy, e.g., full-sequence segmentation map
IoU increases from 0.24 to 0.29, a 21% relative improve-
ment. Complementary to numerical results, Figure 4 shows
that Δ YNAMICS accurately captures multi-object motion
trajectories, even in an unseen domain such as CLEVRER.

Test-Time Enhancement. We evaluate the model’s per-
formance with different testing time strategies. As shown
in Table 5, the vanilla Δ YNAMICS model gains a marginal
improvement by sampling more scene configurations. For
example, the full-sequence IoU increases from 0.24 (with
the greedy sampling approach) to 0.28 (the best out of 32
sampled configurations), a 14% increase.

Consistent with earlier findings, the motion-reasoning
variant significantly outperforms the base model, and using
best-of-32 sampling further boosts performance: the first-
frame IoU increases from 0.30 to 0.38 (+27%), and the
full-sequence optical flow EPE decreases by 13%. These
relative gains are larger than those achieved by the vanilla
model under the same sampling strategy. We hypothesize
that the intermediate motion-reasoning step provides a more
structured and physically meaningful representation, which
enables sampling to explore a broader and effective set of
plausible solutions rather than drifting into implausible re-
gions of the parameter space. This broader yet more guided
search helps resolve long-tailed errors and yields higher-
quality reconstructions.

Lastly, we perform evolutionary search with an initial-
ization from the best-of-32 sample. This method yields the
highest accuracy for the full sequence, partly thanks to the

Table 2. **Evaluation metrics for the in-distribution setting on the synthetic evaluation data.** When Δ YNAMICS takes optical flows as the input, it consistently outperforms baseline methods across most evaluation dimensions. Note that parameter estimation metrics are unavailable for non-VLM baselines since they do not correctly predict the object composition. **Best** and **runner-up** results are highlighted.

Input	Obj. Comp. Acc. (\uparrow)	Segmentation Map IoU (\uparrow)		Optical Flow EPE (\downarrow)		Physics Parameter MAE (\downarrow)			
		First-Frame	Full Sequence	First-Frame	Full Sequence	Damping	Roll Friction	Slide Friction	
Non-VLM Models									
ViViT [3]	RGB	0.00	0.08	0.07	18.52	9.38	–	–	–
ViViT [3]	Opt. Flow	0.00	0.07	0.06	8.54	8.90	–	–	–
VLM Models									
InternVL3-8B [63]	RGB	0.02	0.05	0.05	25.13	15.77	2.94	0.35	0.81
Qwen2.5-VL-7B [6]	RGB	0.27	0.03	0.03	39.98	16.33	1.97	0.32	0.66
Claude-4-Sonnet [1]	RGB	0.45	0.09	0.07	13.79	11.07	1.71	0.24	0.43
Ours									
Δ YNAMICS	RGB	0.60	0.52	0.32	27.58	19.66	1.52	0.16	0.16
Δ YNAMICS	Opt. Flow	0.97	0.88	0.49	5.75	9.24	1.72	0.15	0.16
+ Motion Reasoning	Opt. Flow	0.99	0.91	0.54	4.88	8.52	1.60	0.16	0.15

Table 3. **Robustness to complex scene dynamics.** Δ YNAMICS models, trained on up to four objects, generalize effectively to more complex scenes with up to six interacting objects. Structured motion reasoning enhances robustness and consistency under increasing scene complexity. Note that four four-object configurations are held out during training and evaluated here to assess true out-of-distribution generalization.

# Objects	Model	Segmentation Map IoU (\uparrow)	
		First Frame	Full Sequence
4	Δ YNAMICS	0.88	0.53
	+ Motion Reasoning	0.89	0.54
5	Δ YNAMICS	0.87	0.51
	+ Motion Reasoning	0.88	0.54
6	Δ YNAMICS	0.85	0.50
	+ Motion Reasoning	0.81	0.52

Table 4. **Cross-engine generalization.** Evaluating transfer from MuJoCo (training) to Blender (CLEVRER [59]) demonstrates that Δ YNAMICS maintains its performance in a zero-shot setting despite domain shifts. Incorporating structured motion description consistently improves segmentation map IoU.

Modality		Segmentation Map IoU (\uparrow)	
		First Frame	Full Sequence
VLM Models			
InternVL3-8B	RGB	0.01	0.02
Qwen2.5-VL-7B	RGB	0.01	0.01
Claude-4-Sonnet	RGB	0.03	0.04
Ours			
Δ YNAMICS	RGB	0.43	0.19
Δ YNAMICS	Opt. Flow	0.63	0.24
+ Motion Reasoning	Opt. Flow	0.67	0.30

474 quality initialization. This result shows that CMA-ES is the
475 method of choice for optimal accuracy during test-time.

476 6.2. Real-World Applications

477 In this section, we evaluate Δ YNAMICS in the real world.
478 Additional results on physically plausible video editing are
479 provided in Section D.

480 **Dataset.** We collected real-world videos using an iPhone 13

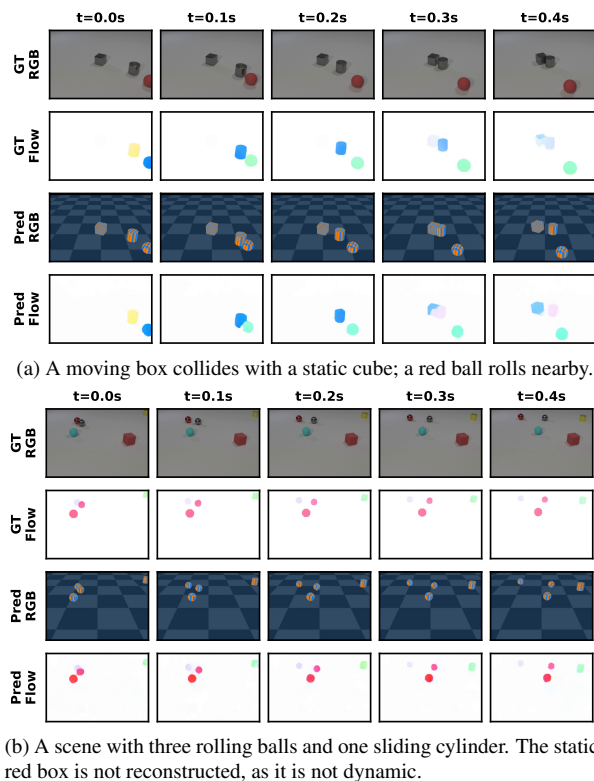


Figure 4. **Zero-shot generalization between engines, from MuJoCo to Blender.** We train Δ YNAMICS on MuJoCo data and evaluate it on CLEVRER [59]. For each example, we show (from top to bottom) (1) the original RGB video, (2) the ground truth optical flow, (3) our model’s reconstructed video, and (4) the optical flow of our reconstruction.

and Canon cameras in landscape orientation. To assess ro-
bustness to diverse surface conditions, we collected data
for multiple environments, including indoor floors, outdoor
running tracks, and basketball courts. Test objects include
everyday items such as shoe boxes, balls, massage roll,
cookie containers, and some irregular-shaped objects such

481
482
483
484
485
486

Table 5. **Evaluation of test-time optimization strategies on CLEVRER.** We compare the base and motion reasoning variants of Δ YNAMICS with greedy decoding, best-of-32 sampling under temperature of 0.1, and evolutionary search (CMA-ES). *Best@1* denotes the average of the 32 samples, while *Best@32* reports the best.

	Segmentation Map IoU (\uparrow)						Optical Flow EPE (\downarrow)					
	Greedy	First Frame		Greedy	Full Sequence		Greedy	First Frame		Greedy	Full Sequence	
		Best@1	Best@32		Best@1	Best@32		Best@1	Best@32		Best@1	Best@32
Δ YNAMICS	0.63	0.63	0.67	0.24	0.24	0.28	3.66	3.65	2.92	6.91	6.86	6.21
+ Motion Reasoning	0.67	<u>0.68</u>	0.76	0.30	0.30	<u>0.38</u>	2.92	2.93	<u>2.22</u>	5.94	5.95	<u>5.17</u>
+ + CMA-ES	0.62	-	-	0.66	-	-	0.13	-	-	0.11	-	-

Table 6. **Performance of real-world rigid-body motion reconstruction.** We evaluate Δ YNAMICS on real-world video dataset. Δ YNAMICS successfully generalizes from synthetic training to real scenes. Incorporating motion reasoning improves segmentation and flow alignment, while Best-of-32 sampling further refines accuracy. CMA-ES optimization provides the best full sequence alignment results.

	Segmentation Map IoU (\uparrow)		Optical Flow EPE (\downarrow)	
	First Frame	Full Seq.	First Frame	Full Seq.
Δ YNAMICS	0.57	0.26	1.62	0.67
+ Motion Reasoning	0.54	0.29	1.39	0.58
+ + Best@32	0.72	<u>0.41</u>	1.06	<u>0.46</u>
+ + CMA-ES	<u>0.57</u>	0.65	<u>1.26</u>	0.36

487 as apples. The dataset details are provided in Appendix C.1.

488 **Results.** As shown in Table 6, the motion reasoning variant
 489 improves segmentation IoU by 12% and flow EPE by
 490 13%, while best-of-32 sampling further enhances motion
 491 accuracy. Evolutionary search provides the largest gains in
 492 the metrics. Qualitatively, Figure 5 shows that our model
 493 captures the trajectories and locations of two-object motion
 494 precisely, implying a high level of accuracy in the estimated
 495 initial states and physics parameters. We provide more real-
 496 world examples and failure analysis in Appendix C.2.

497 7. Conclusion

498 We have presented a novel viewpoint on the problem of pre-
 499 dicting physics configurations for rigid-body motion from
 500 monocular videos. Our main contribution is a general struc-
 501 tured textual representation of the physics states and pa-
 502 rameters for a wide range of motion dynamics and ob-
 503 ject interaction. In addition, we trained Δ YNAMICS, a
 504 vision-language model to generate physics configurations
 505 and camera geometry in a structured textual format. We
 506 also incorporate motion reasoning and test-time optimiza-
 507 tion techniques to enhance our model’s accuracy. Being
 508 trained on 400K synthetically generated scenes in MuJoCo,
 509 our model shows robust generalization across rendering en-
 510 gines and to real-world data, and consistently outperforms
 511 off-the-shelf vision-language models. Our results demon-
 512 strate a promising line of research on using language model-
 513 ing to provide a common physics representation for physics
 514 perception and physics simulation.

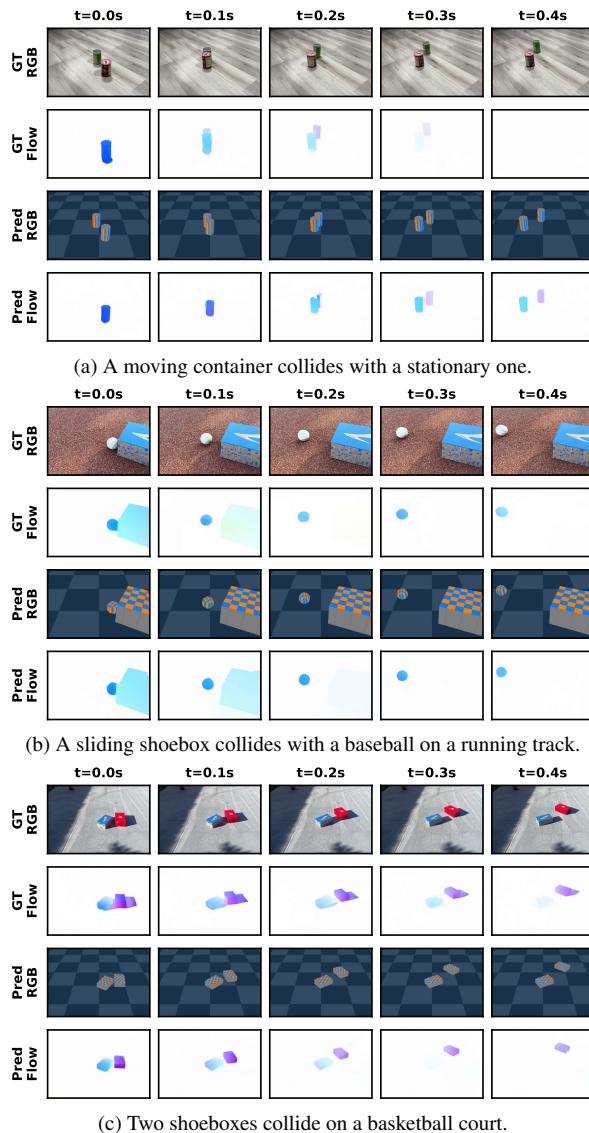


Figure 5. **Motion capture for real-world videos.** Δ YNAMICS is able to reproduce motion trajectory and object location on real-world surfaces and complex lighting. It can also capture multi-body collision dynamics despite the domain gap between synthetic and real data.

515

References

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

- [1] Anthropic. The claude 3 model family: Opus, sonnet, haiku. <https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf>, 2024. 5, 7
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 2
- [3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, pages 6836–6846, 2021. 6, 7
- [4] Martin Asenov, Michael Burke, Daniel Angelov, Todor Davchev, Kartic Subr, and Subramanian Ramamoorthy. Vid2param: Modeling of dynamics parameters from video. *IEEE Robotics and Automation Letters*, 5(2):414–421, 2019. 1, 2
- [5] Tayfun Ates, M Samil Atesoglu, Cagatay Yigit, Ilker Kesen, Mert Kobas, Erkut Erdem, Aykut Erdem, Tilbe Goksun, and Deniz Yuret. Craft: A benchmark for causal reasoning about forces and interactions. *arXiv preprint arXiv:2012.04293*, 2020. 2
- [6] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 5, 7
- [7] Fabien Baradel, Natalia Neverova, Julien Mille, Greg Mori, and Christian Wolf. Cophy: Counterfactual learning of physical dynamics. *arXiv preprint arXiv:1909.12000*, 2019. 2
- [8] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. *Advances in neural information processing systems*, 29, 2016. 1
- [9] Jonas Belouadi, Simone Ponzetto, and Steffen Eger. Detikzify: Synthesizing graphics programs for scientific figures and sketches with tikz. *Advances in Neural Information Processing Systems*, 37:85074–85108, 2024. 2
- [10] Jonas Belouadi, Eddy Ilg, Margret Keuper, Hideki Tanaka, Masao Utiyama, Raj Dabre, Steffen Eger, and Simone Paolo Ponzetto. Tikzero: Zero-shot text-guided graphics program synthesis. *arXiv preprint arXiv:2503.11509*, 2025. 2
- [11] Siyuan Bian, Chenghao Xu, Yuliang Xiu, Artur Grigorev, Zhen Liu, Cewu Lu, Michael J Black, and Yao Feng. Chatgarment: Garment estimation, generation and editing via large language models. In *CVPR*, pages 2924–2934, 2025. 2
- [12] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. 1
- [13] Ryan Burgert, Yuancheng Xu, Wenqi Xian, Oliver Pilarski, Pascal Clausen, Mingming He, Li Ma, Yitong Deng, Lingxiao Li, Mohsen Mousavi, et al. Go-with-the-flow: Motion-controllable video diffusion models using real-time warped noise. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13–23, 2025. 11, 12
- [14] Pradyumna Chari, Chinmay Talegaonkar, Yunhao Ba, and Achuta Kadambi. Visual physics: Discovering physical laws from videos. *arXiv preprint arXiv:1911.11893*, 2019. 1, 2
- [15] Yamei Chen, Haoquan Zhang, Yangyi Huang, Zeju Qiu, Kaipeng Zhang, Yandong Wen, and Weiyang Liu. Symbolic graphics programming with large language models. *arXiv preprint arXiv:2509.05208*, 2025. 2
- [16] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024. 2
- [17] Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Benchmarking and enhancing vision-language models for physical world understanding. *arXiv preprint arXiv:2501.16411*, 2025. 2
- [18] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv e-prints*, pages arXiv–2409, 2024. 2
- [19] 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*, 34:887–899, 2021. 1, 2
- [20] Sebastien Ehrhardt, Aron Monszpart, Niloy Mitra, and Andrea Vedaldi. Unsupervised intuitive physics from visual observations. In *Asian Conference on Computer Vision*, pages 700–716. Springer, 2018. 1
- [21] C Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax—a differentiable physics engine for large scale rigid body simulation. *arXiv preprint arXiv:2106.13281*, 2021. 3
- [22] Hiroki Furuta, Kuang-Huei Lee, Shixiang Shane Gu, Yutaka Matsuo, Aleksandra Faust, Heiga Zen, and Izzeddin Gur. Geometric-averaged preference optimization for soft preference labels. *Advances in Neural Information Processing Systems*, 37:57076–57114, 2024. 2
- [23] 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 *CVPR*, pages 27924–27933, 2025. 1, 2
- [24] Yunqi Gu, Ian Huang, Jihyeon Je, Guandao Yang, and Leonidas Guibas. Blendergym: Benchmarking foundational model systems for graphics editing. In *CVPR*, pages 18574–18583, 2025. 2
- [25] Nikolaus Hansen and Andreas Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. In *Proceedings of IEEE international conference on evolutionary computation*, pages 312–317. IEEE, 1996. 5
- [26] Eric Heiden, Ziang Liu, Vibhav Vineet, Erwin Coumans, and Gaurav Sukhatme. Learning articulated rigid body dynamics simulations from video. In *ICLR Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022. 1, 2

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

- 629 [27] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan
630 Carr, Jonathan Ragan-Kelley, and Frédo Durand. DiffTaichi:
631 Differentiable programming for physical simulation. *arXiv
632 preprint arXiv:1910.00935*, 2019. 3
- 633 [28] Zhiao Huang, Yuanming Hu, Tao Du, Siyuan Zhou, Hao
634 Su, Joshua B Tenenbaum, and Chuang Gan. Plasticinellab:
635 A soft-body manipulation benchmark with differentiable
636 physics. *arXiv preprint arXiv:2104.03311*, 2021. 3
- 637 [29] Krishna Murthy Jatavallabhula, Miles Macklin, Florian
638 Golemo, Vikram Voleti, Linda Petrini, Martin Weiss, Brean-
639 dan Considine, Jérôme Parent-Lévesque, Kevin Xie, Kenny
640 Erleben, et al. gradsim: Differentiable simulation for sys-
641 tem identification and visuomotor control. *arXiv preprint
642 arXiv:2104.02646*, 2021. 3
- 643 [30] James R Kubricht, Keith J Holyoak, and Hongjing Lu. Intu-
644 itive physics: Current research and controversies. *Trends in
645 cognitive sciences*, 21(10):749–759, 2017. 1
- 646 [31] Peter Kulits, Haiwen Feng, Weiyang Liu, Victoria Abrevaya,
647 and Michael J Black. Re-thinking inverse graphics with large
648 language models. *arXiv preprint arXiv:2404.15228*, 2024. 2
- 649 [32] Peter Kulits, Michael J Black, and Silvia Zuffi. Reconstruct-
650 ing animals and the wild. In *CVPR*, pages 16565–16577,
651 2025. 2
- 652 [33] Long Le, Jason Xie, William Liang, Hung-Ju Wang, Yue
653 Yang, Yecheng Jason Ma, Kyle Vedder, Arjun Krishna, Di-
654 nesh Jayaraman, and Eric Eaton. Articulate-anything: Auto-
655 matic modeling of articulated objects via a vision-language
656 foundation model. *arXiv preprint arXiv:2410.13882*, 2024.
657 7
- 658 [34] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
659 Visual instruction tuning. *Advances in neural information
660 processing systems*, 36:34892–34916, 2023. 2
- 661 [35] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
662 Improved baselines with visual instruction tuning. In *Pro-
663 ceedings of the IEEE/CVF conference on computer vision
664 and pattern recognition*, pages 26296–26306, 2024. 2
- 665 [36] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar.
666 Docvqa: A dataset for vqa on document images. In *Proceed-
667 ings of the IEEE/CVF winter conference on applications of
668 computer vision*, pages 2200–2209, 2021. 2
- 669 [37] Himangi Mittal, Peiye Zhuang, Hsin-Ying Lee, and Shub-
670 ham Tulsiani. Uniphy: Learning a unified constitutive model
671 for inverse physics simulation. In *Proceedings of the Com-
672 puter Vision and Pattern Recognition Conference*, pages
673 16208–16218, 2025. 3
- 674 [38] Yi-Ling Qiao, Alexander Gao, and Ming Lin. Neu-
675 physics: Editable neural geometry and physics from monoc-
676 ular videos. *Advances in Neural Information Processing Sys-
677 tems*, 35:12841–12854, 2022. 3
- 678 [39] Nazneen Fatema Rajani, Rui Zhang, Yi Chern Tan, Stephan
679 Zheng, Jeremy Weiss, Aadit Vyas, Abhijit Gupta, Caiming
680 Xiong, Richard Socher, and Dragomir Radev. Esprit: Ex-
681 plaining solutions to physical reasoning tasks. *arXiv preprint
682 arXiv:2005.00730*, 2020. 2
- 683 [40] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang
684 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
685 Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2:
Segment anything in images and videos. *arXiv preprint
arXiv:2408.00714*, 2024. 5
- [41] Ronan Riochet, Mario Ynocente Castro, Mathieu Bernard,
Adam Lerer, Rob Fergus, Véronique Izard, and Emmanuel
Dupoux. Intphys 2019: A benchmark for visual intuitive
physics understanding. *IEEE Transactions on Pattern Anal-
ysis and Machine Intelligence*, 44(9):5016–5025, 2021. 2
- [42] Juan A Rodriguez, Abhay Puri, Shubham Agarwal, Issam H
Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez,
Christopher Pal, and Marco Pedersoli. Starvector: Gener-
ating scalable vector graphics code from images and text. In
*Proceedings of the Computer Vision and Pattern Recognition
Conference*, pages 16175–16186, 2025. 2
- [43] Aadarsh Sahoo, Vansh Tibrewal, and Georgia Gkioxari.
Aligning text, images, and 3d structure token-by-token.
arXiv preprint arXiv:2506.08002, 2025. 7
- [44] Arsalan Sharifnassab, Saber Salehkaleybar, Sina Ghiassian,
Surya Kanoria, and Dale Schuurmans. Soft preference opti-
mization: Aligning language models to expert distributions.
arXiv preprint arXiv:2405.00747, 2024. 2
- [45] Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei
Huang, Yongbin Li, and Houfeng Wang. Preference rank-
ing optimization for human alignment. In *Proceedings of
the AAAI Conference on Artificial Intelligence*, pages 18990–
18998, 2024. 5, 1
- [46] Yiren Song, Danze Chen, and Mike Zheng Shou. Layer-
tracer: Cognitive-aligned layered svg synthesis via diffusion
transformer. *arXiv preprint arXiv:2502.01105*, 2025. 2
- [47] Priya Sundareshan, Rika Antonova, and Jeannette Bohgl. Dif-
fcloud: Real-to-sim from point clouds with differentiable
simulation and rendering of deformable objects. In *2022
IEEE/RSJ International Conference on Intelligent Robots
and Systems (IROS)*, pages 10828–10835. IEEE, 2022. 3
- [48] Cheng Tan, Qi Chen, Jingxuan Wei, Gaowei Wu, Zhangyang
Gao, Siyuan Li, Bihui Yu, Ruifeng Guo, and Stan Z Li.
Sketchagent: Generating structured diagrams from hand-
drawn sketches. *arXiv preprint arXiv:2508.01237*, 2025. 2
- [49] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
transforms for optical flow. In *ECCV*, pages 402–419.
Springer, 2020. 3, 4, 5, 11
- [50] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A
physics engine for model-based control. In *2012 IEEE/RSJ
international conference on intelligent robots and systems*,
pages 5026–5033. IEEE, 2012. 2, 4
- [51] Nicholas Watters, Daniel Zoran, Theophane Weber, Peter
Battaglia, Razvan Pascanu, and Andrea Tacchetti. Visual
interaction networks: Learning a physics simulator from
video. *Advances in neural information processing systems*,
30, 2017. 1
- [52] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan
Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei
Chen, et al. Qwen-image technical report. *arXiv preprint
arXiv:2508.02324*, 2025. 2
- [53] 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.
NeurIPS, 28, 2015. 1, 2

- 743 [54] Jiajun Wu, Erika Lu, Pushmeet Kohli, Bill Freeman, and
744 Josh Tenenbaum. Learning to see physics via visual de-
745 animation. *NeurIPS*, 30, 2017. 1, 2
- 746 [55] Ronghuan Wu, Wanchao Su, and Jing Liao. Chat2svg: Vec-
747 tor graphics generation with large language models and im-
748 age diffusion models. In *Proceedings of the Computer Vision
749 and Pattern Recognition Conference*, pages 23690–23700,
750 2025. 2
- 751 [56] Ximing Xing, Haitao Zhou, Chuang Wang, Jing Zhang,
752 Dong Xu, and Qian Yu. Svgdreamer: Text guided svg gener-
753 ation with diffusion model. In *Proceedings of the IEEE/CVF
754 Conference on Computer Vision and Pattern Recognition*,
755 pages 4546–4555, 2024.
- 756 [57] Ximing Xing, Qian Yu, Chuang Wang, Haitao Zhou, Jing
757 Zhang, and Dong Xu. Svgdreamer++: Advancing editability
758 and diversity in text-guided svg generation. *IEEE Transac-
759 tions on Pattern Analysis and Machine Intelligence*, 2025. 2
- 760 [58] Gengshan Yang, Shuo Yang, John Z Zhang, Zachary Manch-
761 ester, and Deva Ramanan. Ppr: Physically plausible re-
762 construction from monocular videos. In *Proceedings of the
763 IEEE/CVF International Conference on Computer Vision*,
764 pages 3914–3924, 2023. 3
- 765 [59] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun
766 Wu, Antonio Torralba, and Joshua B Tenenbaum. Clevrer:
767 Collision events for video representation and reasoning.
768 *arXiv preprint arXiv:1910.01442*, 2019. 1, 2, 6, 7, 4
- 769 [60] Yunzhi Zhang, Zizhang Li, Matt Zhou, Shangzhe Wu, and
770 Jiajun Wu. The scene language: Representing scenes with
771 programs, words, and embeddings. In *Proceedings of the
772 Computer Vision and Pattern Recognition Conference*, pages
773 24625–24634, 2025. 2
- 774 [61] Wang Zhao, Yan-Pei Cao, Jiale Xu, Yuejiang Dong, and
775 Ying Shan. Di-pcg: Diffusion-based efficient inverse pro-
776 cedural content generation for high-quality 3d asset creation.
777 In *CVPR*, pages 11061–11072, 2025. 3
- 778 [62] Xian Zhou, Yiling Qiao, Zhenjia Xu, TH Wang, Z Chen, J
779 Zheng, Z Xiong, Y Wang, M Zhang, P Ma, et al. Genesis:
780 A generative and universal physics engine for robotics and
781 beyond. *arXiv preprint arXiv:2401.01454*, 2024. 7
- 782 [63] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shen-
783 glong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su,
784 Jie Shao, et al. Internvl3: Exploring advanced training and
785 test-time recipes for open-source multimodal models. *arXiv
786 preprint arXiv:2504.10479*, 2025. 5, 7

ΔYNAMICS: Language-Based Representation for Inferring Rigid-Body Dynamics From Videos

Supplementary Material

787	A. Method Details		832
788	A.1. YAML-to-XML Conversion and Initialization		833
789	Our dataset uses YAML as the canonical representation of		834
790	the full scene configuration. The YAML file contains both		835
791	static attributes (e.g., object geometries, masses, friction co-		836
792	efficients, camera pose) and dynamic initial states (linear		837
793	and angular velocities). MuJoCo, however, accepts model		838
794	definitions only in its XML-based scene description format,		839
795	which encodes static model properties but does not support		840
796	specifying initial velocities.		841
797	To run each simulation, we therefore proceed in two		842
798	steps. First, we convert the static components of the YAML		843
799	configuration into a MuJoCo XML file, defining the bod-		844
800	ies, inertial properties, joints, geoms, and camera. Second,		
801	we load this XML into MuJoCo, initialize the engine and		
802	then set all remaining dynamic quantities (e.g., initial linear		
803	and angular velocities) directly in the simulator state before		
804	simulation rollout.		
805	A.2. Dataset Preparation		
806	MuJoCo Rendering Details. Each simulation is rendered		
807	for one second at 30 FPS, with eight physics steps per		
808	frame. The output resolution is 480×320 pixels. We intro-		
809	duce randomized lighting conditions with varying shadow		
810	configurations and apply diverse object textures, including		
811	checkerboard and gradient patterns, to improve robustness		
812	to visual appearance. During simulation, we record RGB		
813	frames, segmentation masks, contact information, and full		
814	state histories, which are later used to generate structured		
815	visual reasoning annotations.		
816	Optical flow estimation is often unreliable on smooth		
817	or textureless surfaces. To mitigate this, we fill the back-		
818	ground and floor with natural scene images, which intro-		
819	duce sufficient texture and yield substantially higher-quality		
820	flow fields. Optical flow is computed between consecutive		
821	frames (30 FPS), producing 29 flow maps per video; we		
822	then sample every third frame to obtain the inputs used for		
823	model training.		
824	Structured Visual Reasoning. To ensure consistent and in-		
825	terpretable spatial reasoning, we discretize continuous ob-		
826	ject positions and motions into categorical linguistic de-		
827	scriptors along three spatial axes. For example, positions		
828	along the x -axis are quantized into seven categories: far left		
829	($x < -2$), moderately left ($-2 \leq x < -1$), slightly left		
830	($-1 \leq x < -0.5$), near center ($-0.5 \leq x < 0.5$), slightly		
831	right ($0.5 \leq x < 1$), moderately right ($1 \leq x < 2$), and far		
		right ($x \geq 2$).	832
		Segmentation maps are used to determine whether each	833
		object is initially visible and to record when it leaves the	834
		camera’s field of view. During simulation, we log colli-	835
		sion and contact histories from the physics engine to cap-	836
		ture fine-grained interaction events such as ground contact	837
		and inter-object collisions. We also analyze object state	838
		trajectories to identify when each object comes to rest, en-	839
		abling precise annotation of temporal dynamics such as mo-	840
		tion duration and stopping time. To increase linguistic di-	841
		versity, the transformation script converts these auxiliary	842
		signals into textual descriptions using multiple paraphrased	843
		sentence templates while preserving structural consistency.	844
		Example of Structured Visual Reasoning. An example	845
		of synthetic data is shown in Table 7. The model analyzes	846
		object properties, motion patterns, and physical interactions	847
		in natural language, which serves as an intermediate step to	848
		improve parameter estimation accuracy.	849
		A.3. Details on In-Context Example Preparation	850
		To guide the models effectively, we construct three in-	851
		context examples that (i) include all primitive shapes	852
		present in our dataset and (ii) cover the full range of physi-	853
		cal parameters by selecting configurations at the minimum	854
		and maximum values of the training distribution. Each ex-	855
		ample consists of ten video raw RGB frames paired with its	856
		YAML-based scene configuration.	857
		For CLEVRER [59], we first select three target exam-	858
		ple videos and manually annotate three corresponding scene	859
		configurations. The examples are chosen to (i) include all	860
		three primitive shapes and (ii) span the minimum and maxi-	861
		mum values of our physical-parameter ranges, ensuring that	862
		the model does not extrapolate beyond the provided exam-	863
		ples. To maintain fidelity, we iteratively refine each con-	864
		figuration annotation so that the resulting simulated mo-	865
		tions and object trajectories closely match those in the target	866
		videos.	867
		A.4. Details on Preference Optimization	868
		Preliminary: Preference Rank Optimization. Following	869
		the Bradley–Terry formulation [12], a reward model (RM)	870
		estimates pairwise preferences by contrasting two responses	871
		y^1 and y^2 for a given input x . Preference Ranking Op-	872
		timization (PRO) [45] extends this idea by directly fine-	873
		tuning the policy π_θ , treating it as both the RM and the	874

875 policy network. The PRO loss is defined as:

$$876 \quad \mathcal{L}_{\text{PRO}} = -\log \frac{e^{r_{\pi}(x, y^1)}}{e^{r_{\pi}(x, y^1)} + e^{r_{\pi}(x, y^2)}}, \quad (3)$$

877 where $r_{\pi}(x, y)$ denotes the implicit reward $r_{\pi_{\text{PRO}}}$ for a
878 given input x and candidate response y^k . It is defined as
879 the average token-level log-likelihood:

$$880 \quad r_{\pi_{\text{PRO}}}(x, y^k) = \frac{1}{|y^k|} \sum_{t=1}^{|y^k|} \log P_{\pi}(y_t^k | x, y_{<t}^k). \quad (4)$$

881 Intuitively, $r_{\pi_{\text{PRO}}}(x, y^k)$ measures the normalized sequence
882 log-likelihood (i.e., the mean per-token log-probability) and
883 serves as a scalar proxy for how confidently the model as-
884 signs probability mass to the response y^k given x .

885 **Soft Preference Weighting.** However, Eq. 3 assumes a
886 *one-hot* preference—one response is strictly preferred ($y^1 \succ$
887 y^2) while the other is not. In our setting, this binary assump-
888 tion is overly rigid: two simulated rollouts may each excel
889 in distinct aspects (e.g., one accurately reproduces geome-
890 try while the other better matches damping or velocity). To
891 capture such nuanced trade-offs, we introduce a *soft pref-*
892 *erence weighting* that transforms the simulator-derived re-
893 wards into continuous targets:

$$894 \quad \tilde{r}(y^i) = \frac{e^{s(y^i)/\tau}}{\sum_{j \in \{1,2\}} e^{s(y^j)/\tau}},$$

895 where τ is a temperature and $s(\cdot)$ denotes the simulation-
896 derived score function, i.e., segmentation IoU. Thus, the
897 optimization objective then becomes:

$$898 \quad \mathcal{L}_{\text{soft-PRO}} = -\sum_{i \in \{1,2\}} \tilde{r}(y^i) \log \frac{e^{r_{\pi}(x, y^i)}}{e^{r_{\pi}(x, y^1)} + e^{r_{\pi}(x, y^2)}} \quad (5)$$

899 which can be interpreted as a *reward-weighted cross-*
900 *entropy*—analogous to replacing a binary BCE loss with a
901 soft-label CE loss. This formulation better accommodates
902 partially correct rollouts and encourages the policy to allo-
903 cate probability mass in proportion to their normalized sim-
904 ulator rewards, leading to smoother and more stable test-
905 time adaptation. This approach bears similarity with soft-
906 preference concept in previous work [22, 44]

Table 7. **Example Synthetic Training Data Instance.** We illustrate the target output format for both the vanilla and reasoning-augmented variants of our model. The blue box shows the structured YAML configuration. The red box shows the corresponding natural-language reasoning describing object motions and interactions. In the vanilla setting, the model is trained to generate only the configuration text wrapped within the `<answer>` tag. In the reasoning-enhanced setting, the model first outputs the reasoning text enclosed by `<think>` tags, followed by the configuration text within `<answer>` tags.

`configuration` =

```
- type: box
  name: box_0
  size: [1.0, 0.5, 0.4]
  state:
    angular_velocity: [0, 0, 0]
    linear_velocity: [4.3, 2.5, 0.0]
    orientation: [0.97, 0.0, 0.0, 0.23]
    position: [-5.0, -0.3, 0.4]
  physics:
    friction: [1.1, 0.3]
    mass: 1.0
    damping: -4
- type: cylinder
  name: cylinder_0
  radius: 0.3
  height: 0.5
  state:
    angular_velocity: [0, 0, 0]
    linear_velocity: [-0.1, -0.5, 0.0]
    orientation: [0.69, 0.69, 0.15, 0.15]
    position: [-0.5, 0.8, 0.3]
  physics:
    friction: [0.5, 0.3]
    mass: 1.0
    damping: -4
- type: camera
  fovy: 45
  orientation: 45
  position: [0, -2, 3.5]
- type: gravity
  gravity: [0, 0, -7.0]
```

`reasoning` =

```
This physics simulation showcases objects interacting under realistic physics.

- Box_0: A cuboid with dimensions 1.0 x 0.5 x 0.4 m, mass 1.0 kg, positioned leftmost and in the foreground, close to the surface. It moves at 4.97 m/s (rightward and forward), stops at t=0.4s, and stays grounded from 0.00-0.90s. Collides with cylinder_0 at t=0.1s.

- Cylinder_0: A solid column (radius=0.3m, height=0.5m, mass=1.0kg) located at the X-axis origin, near the camera and base plane. Moves at 0.51 m/s, momentum 2.79 m/s until end, collides with box_0 at t=0.1s.

- Observation Data: Visible entities: box_0, cylinder_0. Both visible in 10/10 frames.

- Dynamic Interactions: Contact event between cylinder_0 and box_0 at t=0.1s.
```

Target Sequence (Vanilla Δ YNAMICS):

`<answer>` `configuration` `</answer>`

Target Sequence (Δ YNAMICS + Motion Reasoning):

`<think>` `reasoning` `</think>` `\n\n` `<answer>` `configuration` `</answer>`

907 **B. Cross-Engine Generalization**

908 **B.1. Dataset Preparation**

909 We subsample 400 one-second clips from the validation
910 split of CLEVRER [59]. Object segmentation masks are
911 obtained from the official release. We first clean the masks
912 by checking whether each object’s segmentation changes
913 within the subsampled clip—objects with static masks are
914 treated as non-moving and removed from motion evalua-
915 tion. This information is then used to refine the correspond-
916 ing optical flow maps, ensuring that static regions are not
917 misinterpreted as motion.

918 **B.2. Baseline Evaluation Outcomes**

919 The full quantitative evaluations, including optical flows,
920 are shown in Table 8. Qualitative results are shown in Fig-
921 ure 6. These are the results without using CMA-ES.

922 **Test-Time Optimization.** We therefore explore preference
923 optimization to learn from unlabeled videos as detailed pre-
924 viously in Section 4.2. Specifically, we use 1000 cases from
925 CLEVRER training data, generate 32 sampled predictions
926 per case, render rollouts with MuJoCo, and compute seg-
927 mentation IoU to construct pairwise preferences. We then
928 draw three paired data per case to construct the prefer-
929 ence learning dataset and finetune our reasoning model.
930 As shown in Table 9, preference optimization yields con-
931 sistent improvements—modest IoU gains (+1%) and a no-
932 table reduction in first-frame EPE (2.22 to 1.85) under the
933 Best@32 metric. This demonstrates that preference opti-
934 mization provides a practical route to label-free adaptation,
935 enabling models to refine directly through simulation feed-
936 back. However, the marginal gain in performance is less
937 than CMA-ES.

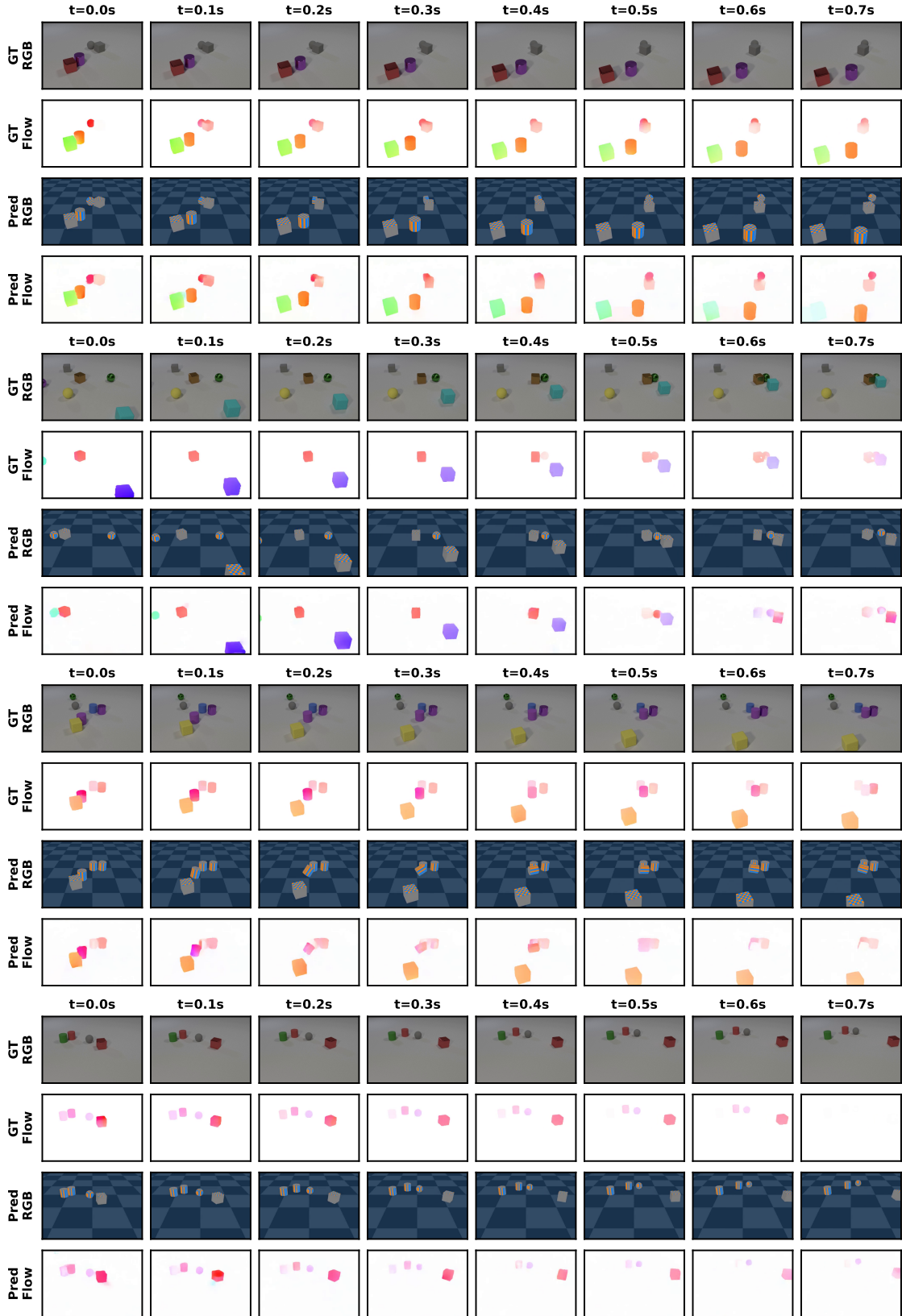


Figure 6. CLEVRER Dataset Results.

Table 8. **Generalization Across Simulation Engines.** Evaluating transfer from MuJoCo (training) to Blender (CLEVRER [59]) demonstrates that Δ YNAMICS maintains its performance in a zero-shot setting. Despite domain shifts in rendering and dynamics, incorporating structured motion description consistently improves segmentation IoU and optical flow EPE. Note that some off-the-shelf models achieve low optical flow EPE by generating fewer objects than present in the scene, which artificially reduces motion and lowers EPE compared to models attempting more complete and realistic physics simulation. **Best** and runner-up results are highlighted.

	Input Modality	Segmentation Map IoU (\uparrow)		Optical Flow EPE (\downarrow)	
		First Frame	Full Sequence	First Frame	Full Sequence
VLM Models					
InternVL3-8B	RGB	0.01	0.02	7.12	6.10
Qwen2.5-VL-7B	RGB	0.01	0.01	9.22	7.41
Claude-4-Sonnet	RGB	0.03	0.04	6.34	5.43*
Ours					
Δ YNAMICS	RGB	0.43	0.19	3.68	7.13
Δ YNAMICS	Opt. Flow	0.63	0.24	3.51	6.85
+ Motion Reasoning	Opt. Flow	0.67	0.30	2.79	5.94

Table 9. **Evaluation of Test-Time Sampling and Preference Optimization on CLEVRER.** We compare the base, reasoning-enhanced, and preference-optimized Δ YNAMICS under greedy decoding and best-of- N sampling. For each case, 32 samples are generated with a temperature of 0.1; *Best@1* denotes the average, while *Best@32* reports the best. **Best** and runner-up results are highlighted.

	Segmentation Map IoU (\uparrow)						Optical Flow EPE (\downarrow)					
	First Frame			Full Sequence			First Frame			Full Sequence		
	Greedy	Best@1	Best@32	Greedy	Best@1	Best@32	Greedy	Best@1	Best@32	Greedy	Best@1	Best@32
Δ YNAMICS	0.63	0.63	0.67	0.24	0.24	0.28	3.66	3.65	2.92	6.91	6.86	6.21
+ Motion Reasoning (MR)	0.67	0.68	0.76	0.30	0.30	0.38	2.92	2.93	2.22	5.94	5.95	5.17
+ MR + PRO	0.68	<u>0.69</u>	0.77	0.31	0.31	<u>0.39</u>	2.90	2.94	<u>1.85</u>	5.78	5.81	<u>4.78</u>
+ MR + CMA-ES	0.62	-	-	0.66	-	-	0.13	-	-	0.11	-	-

938 **C. Real-World Benchmark and Qualitative** 939 **Results**

940 **C.1. Real-World Dataset Statistics.**

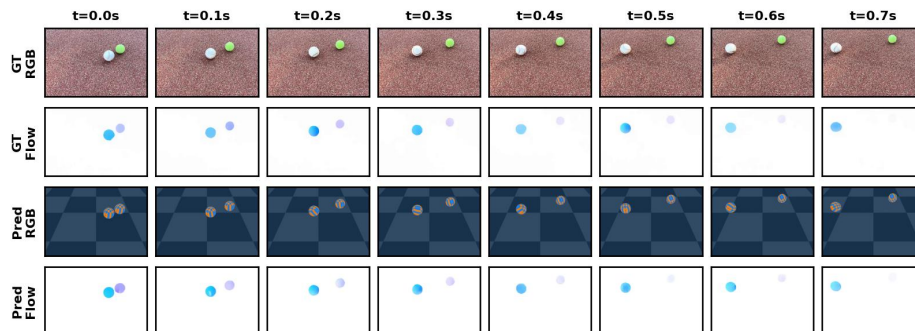
941 Our real-world evaluation set comprises a total of 235
942 videos capturing diverse rigid-body motion scenarios.
943 Among them, 155 videos were recorded with an iPhone 13
944 and 80 videos with a Canon camera. The iPhone videos
945 were captured at 210 FPS or 240 FPS, while the Canon
946 videos were recorded at 50 FPS. To create variations in tem-
947 poral resolution, we downsampled the original recordings
948 by uniformly sampling frames to obtain videos at 25, 30,
949 50, 60, and 70 FPS. The resulting frame-rate distribution
950 is: 30 FPS (76 videos), 70 FPS (46), 25 FPS (40), 50 FPS
951 (40), and 60 FPS (33). Each scene contains between one
952 and five objects, with 139 single-object, 69 two-object, 17
953 three-object, 6 four-object, and 4 five-object configurations.
954 The objects span a wide variety of everyday items, includ-
955 ing shoebox, drug spray, Pringles can, baseball, tennis ball,
956 tissue box, soccer ball, basketball, pool ball, massage roller,
957 gel container, aerosol can, handcrafted vehicle with wheels,
958 apple, tumbler, soda can, whiteboard eraser, napkin roll, in-
959 sulation pad, box, banana, and plum. Among all videos, 86
960 exhibit object collisions.

961 **C.2. Qualitative Results.**

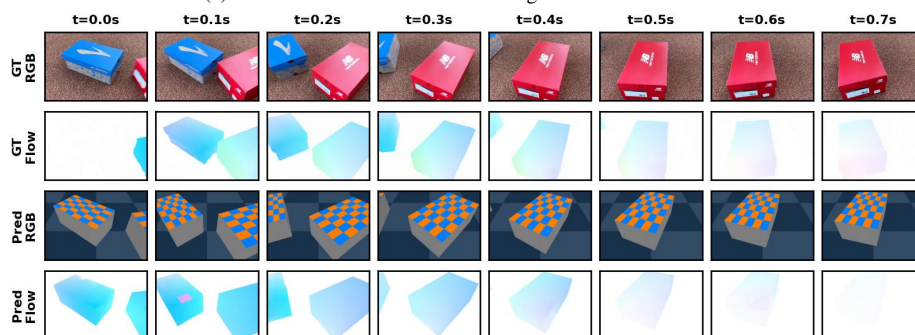
962 In Figure 7 we show more real-world examples; in Figure 8
963 we show the results of irregular-shaped objects such as ap-
964 ples, wooden aircraft, and soda cans. We also present fail-
965 ure cases in Figure 9, where the failure modes comes from
966 irregular shapes that are not seen during training or that
967 results in unexpected bouncing trajectory. In some other
968 cases, the model can predict wrong primitive shapes.

969 **C.3. Future Work**

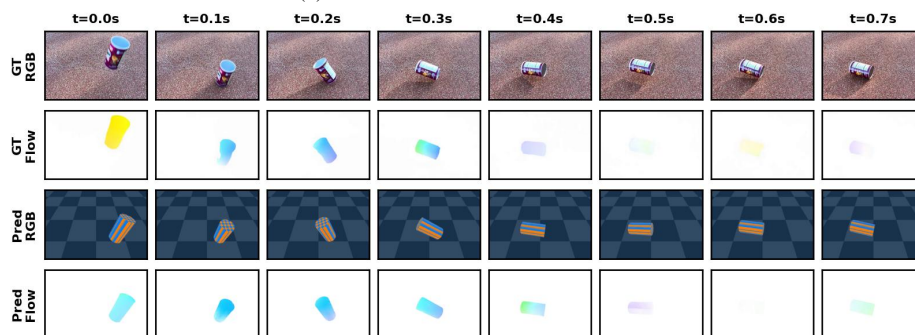
970 While our results are promising in capturing rigid-body mo-
971 tions using a language-centric, VLM-based approach, we
972 identify several directions for future work: (1) incorporat-
973 ing 3D shape tokens [43] to move beyond primitive shapes,
974 (2) extending to articulated objects [33] and sloped envi-
975 ronments to cover more types of rigid-body motion, and
976 (3) adopting more powerful engines such as Genesis [62]
977 to model deformable objects and motions.



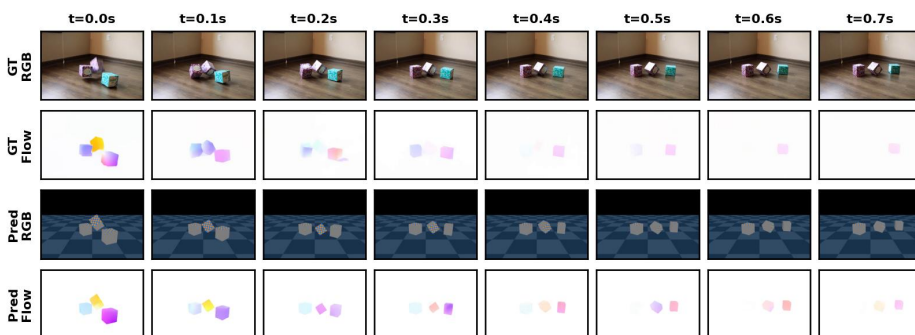
(a) A baseball and a tennis ball moving in different directions.



(b) A red shoe box hits a blue shoe box.

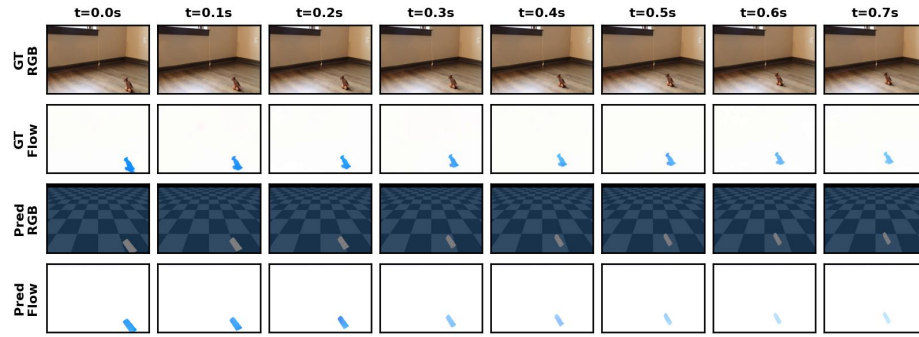


(c) A cylindrical container bouncing on the ground.

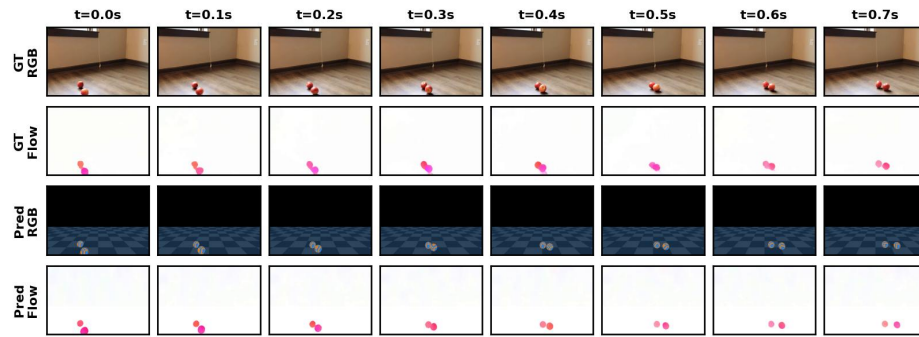


(d) Three tissue boxes dropping on the floor.

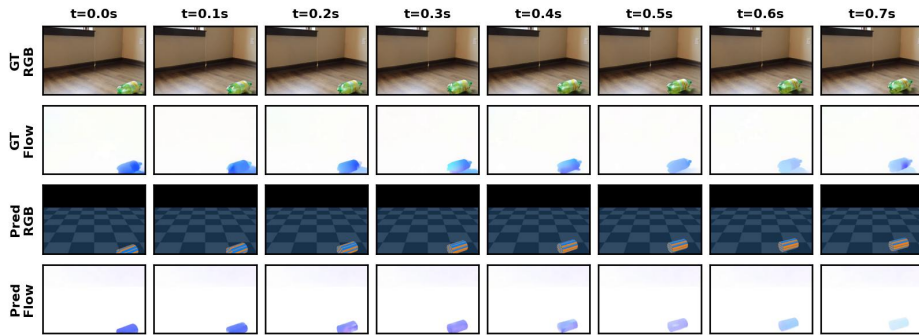
Figure 7. **Rigid-Body Motion Estimation on Our Real-World Dataset.** Δ YNAMICS reconstructs physically plausible trajectories from real-world videos of rigid-body motion, capturing object interactions, material properties, and dynamics across diverse conditions.



(a) A handcrafted wooden eagle with wheels.

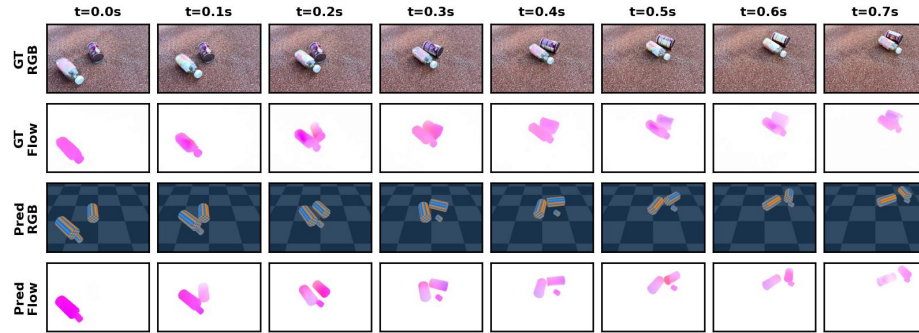


(b) Two moving (rolling) apples collide.

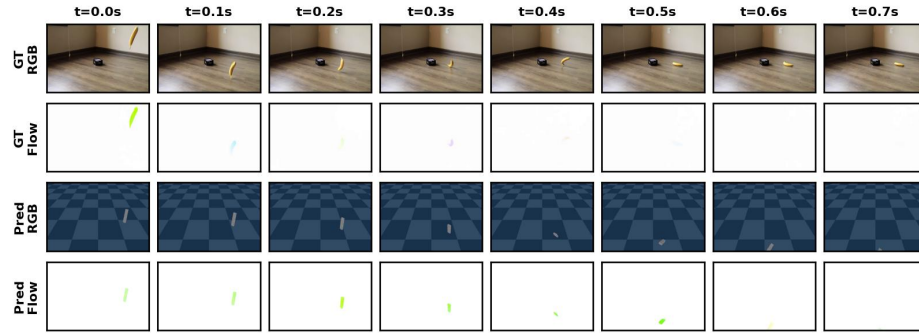


(c) A near-cylindrical soda can with rounded top surfaces rolls slowly .

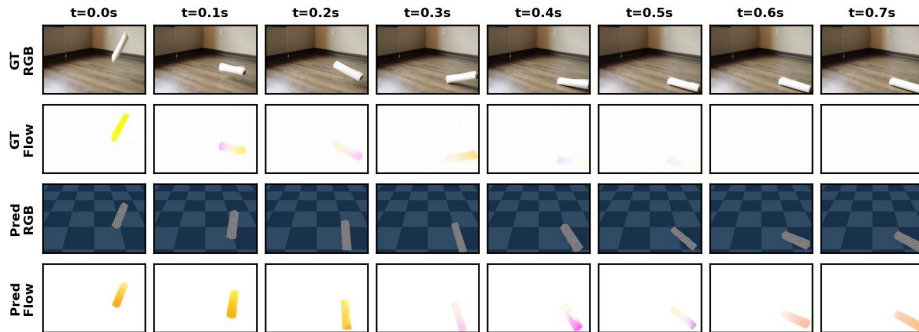
Figure 8. Rigid-Body Motion Estimation on Our Real-World Dataset, Focusing on Irregularly Shaped Objects.



(a) Due to the bottle's irregular shape, our model approximates it with two cylinders and focuses on faithfully reconstructing the motion in the first frame ($t = 0.0s$).



(b) A banana falls to the floor, bounces once, and then comes to rest. Its irregular shape produces an unexpected trajectory, making it difficult to capture object positions and camera poses accurately.



(c) In this case, the model predicts an incorrect primitive shape and an imprecise camera angle.

Figure 9. Rigid-Body Motion Estimation on Our Real-World Dataset, Focusing on Failure Cases.

978 D. Physically Plausible Editing

979 Our framework naturally supports physically consistent
980 editing of object dynamics and scene parameters. Because
981 Δ YNAMICS represents each scene using an explicit and inter-
982 pretable YAML configuration, user-provided editing in-
983 structions can be translated directly into modified physi-
984 cal parameters. This enables a closed-loop editing system
985 that integrates a physics engine, a language model, and
986 a video synthesis model to generate physically plausible
987 edited videos.

988 D.1. Editing Pipeline Overview

989 Figure 10 illustrates our four-stage editing pipeline:

- 990 • **Configuration Extraction with Δ YNAMICS.** Given an
991 input video, we first infer its underlying physical config-
992 uration using Δ YNAMICS. The model outputs a complete
993 YAML file specifying object geometries, initial poses, ve-
994 locities, masses, gravity, friction, and other physical at-
995 tributes. This YAML file serves as an editable and fully
996 interpretable *source code* for the scene.
- 997 • **Language-Guided Configuration Editing.** To incorpo-
998 rate a user instruction (e.g., “reduce the x-velocity by
999 80%” or “decrease gravity by 50%”), we prompt Claude-
1000 3-Haiku with (i) the full YAML configuration predicted
1001 by Δ YNAMICS, and (ii) the editing instruction. Claude
1002 outputs a revised YAML file with localized and seman-
1003 tically appropriate modifications (e.g., updating only the
1004 fields for initial velocity, gravity, or angular velocity). Be-
1005 cause YAML is structured, line-addressable, and seman-
1006 tically meaningful, the language model reliably edits only
1007 the intended parameters while leaving the rest of the con-
1008 figuration intact. This enables precise and controllable
1009 manipulation of physical properties that would otherwise
1010 be entangled in a latent space.
- 1011 • **Simulation and Flow Generation.** The edited configu-
1012 ration is executed in MuJoCo, producing a modified mo-
1013 tion trajectory consistent with the user’s edit. We then
1014 compute dense optical flow using RAFT [49], yielding a
1015 physically grounded flow field that encodes the new dy-
1016 namics.
- 1017 • **Video Synthesis via Go-With-The-Flow.** Following
1018 Burgert et al. [13], the Go-With-The-Flow model synthe-
1019 sizes the edited video by warping noise according to the
1020 edited optical flow, conditioned on the first RGB frame of
1021 the original video. This preserves the appearance of the
1022 scene while enforcing the motion cues determined by the
1023 edited flow.

1024 Empirical results are shown in Fig. 11, which demon-
1025 strates edits to both box dynamics (e.g., reducing x- or y-
1026 velocities) and ball dynamics (e.g., modifying velocity di-
1027 rection or reducing gravity). The pipeline produces physi-
1028 cally correct motion and high-quality visual results in most

cases.

A primary limitation arises from **appearance preser-
vation under complex motion**. Although Go-With-The-
Flow accurately follows the edited optical flow, it some-
times struggles with fine-grained dynamic appearance de-
tails, e.g., maintaining the texture fidelity of a spinning or
rolling basketball. As this remains an open challenge, fu-
ture work could explore fine-tuning the generative model
conditioning on edited optical flow fields to achieve better
physically plausible video editing results.

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

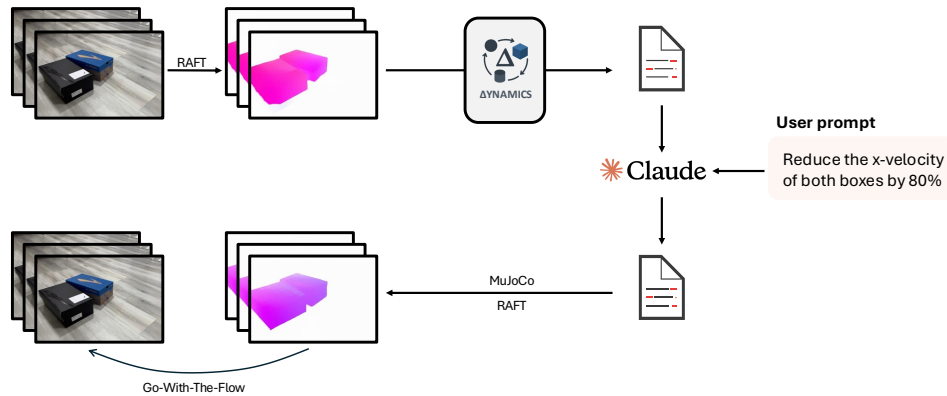
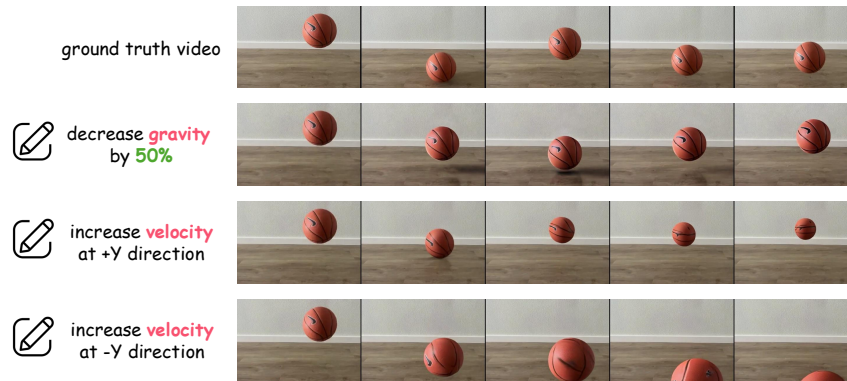
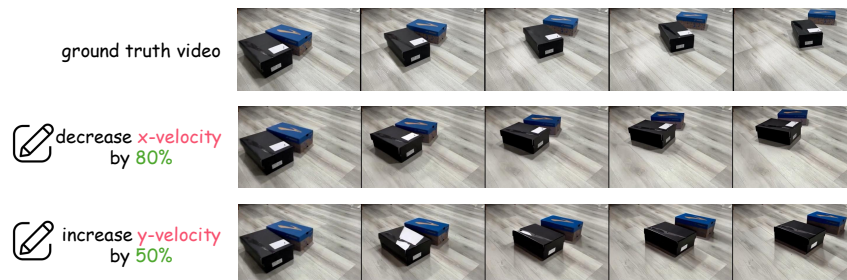


Figure 10. **Physics Editing Pipeline.** Given a user-provided editing instruction (e.g., “reduce the x-velocity by 80%”), we first infer the original scene configuration using Δ YNAMICS. Next, we prompt a large language model (Claude) with both the predicted configuration and the user instruction to generate a revised, physically consistent configuration. The edited configuration is then executed in MuJoCo to produce a modified motion trajectory, from which we compute RAFT optical flow. Finally, we feed the edited flow fields to Go-With-The-Flow [13] to synthesize the edited video.



(a) **Ball Motion Editing.** Example edits include decreasing gravity by 50% and modifying velocity magnitude or direction. These edits are reflected in the regenerated simulated trajectories and corresponding edited videos. Here, positive y motion indicates movement *away from the camera* (deeper into the scene), while negative y motion brings the object *closer toward the viewer*.



(b) **Box Collision Editing.** Here we apply velocity reductions along specified axes (e.g., reducing x-velocity by 80% or y-velocity by 50%). The resulting interactions follow physically plausible adjustments in motion and contact behavior.

Figure 11. CLEVRER Dataset Results.