
3D-IntPhys: Towards More Generalized 3D-grounded Visual Intuitive Physics under Challenging Scenes

Haotian Xue¹ Antonio Torralba² Joshua Tenenbaum² Daniel Yamins³
Yunzhu Li^{3, 4*} Hsiao-Yu Tung^{2*}

¹ Georgia Tech ² MIT ³ Stanford University ⁴ UIUC

Abstract

Given a visual scene, humans have strong intuitions about how a scene can evolve over time under given actions. The intuition, often termed visual intuitive physics, is a critical ability that allows us to make effective plans to manipulate the scene to achieve desired outcomes without relying on extensive trial and error. In this paper, we present a framework capable of learning 3D-grounded visual intuitive physics models from videos of complex scenes. Our method is composed of a conditional Neural Radiance Field (NeRF)-style visual frontend and a 3D point-based dynamics prediction backend, using which we can impose strong relational and structural inductive bias to capture the structure of the underlying environment. Unlike existing intuitive point-based dynamics works that rely on the supervision of dense point trajectory from simulators, we relax the requirements and only assume access to multi-view RGB images and (imperfect) instance masks acquired using color prior. This enables the proposed model to handle scenarios where accurate point estimation and tracking are hard or impossible. We generate datasets including three challenging scenarios involving fluid, granular materials, and rigid objects in the simulation. The datasets do not include any dense particle information so most previous 3D-based intuitive physics pipelines can barely deal with that. We show our model can make long-horizon future predictions by learning from raw images and significantly outperforms models that do not employ an explicit 3D representation space. We also show that once trained, our model can achieve strong generalization in complex scenarios under extrapolate settings. The code is released in <https://github.com/xavihart/3D-IntPhys>.

1 Introduction

Humans can achieve a strong intuitive understanding of the 3D physical world around us simply from visual perception [5, 8, 52, 51, 46, 11]. As we constantly make physical interactions with the environment, the intuitive physical understanding applies to objects of a wide variety of materials [6, 56]. For example, after watching videos of water pouring and doing the task ourselves, we can develop a mental model of the interaction process and predict how the water will move when we apply actions like tilting or shaking the cup (Figure 1). The ability to predict the future evolution of the physical environment is extremely useful for humans to plan our behavior and perform everyday manipulation tasks. It is thus desirable to develop computational tools that learn 3D-grounded models of the world purely from visual observations that can generalize to objects with complicated physical properties like fluid and granular materials.

*This work was done when Haotian Xue was a visiting student advised by Joshua Tenenbaum in MIT CSAIL and he was directly mentored by Yunzhu Li and Hsiao-Yu Tung.

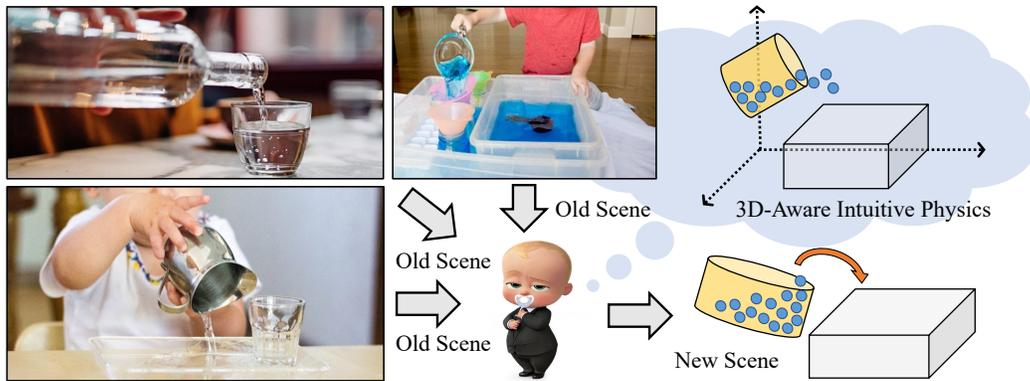


Figure 1: **Visual Intuitive Physics Grounded in 3D Space.** Humans have a strong intuitive understanding of the physical environment. We can predict how the environment would evolve when applying specific actions. This ability roots in our understanding of 3D and applies to objects of diverse materials, which is essential when planning our behavior to achieve specific goals. In this work, we are the first to leverage a combination of implicit neural representation and explicit 3D particle representation to build 3D-grounded visual intuitive physics models of the challenging scenes that applies to objects with complicated physical properties, such as fluids, rigid objects, and granular materials.

There has been a series of works on learning intuitive physics models of the environment from data. However, most existing work either focuses on 2D environments [60, 1, 21, 64, 20, 4, 44, 28, 67, 24, 23, 49, 33, 19, 31, 58, 22, 12, 65] or has to make strong assumptions about the accessible information of the underlying environment [36, 35, 47, 43, 70, 54, 48, 7, 2, 26] (e.g., full-state information of the fluids represented as points). The limitations prevent their use in tasks requiring an explicit 3D understanding of the environments and make it hard to extend to more complicated real-world environments where only visual observations are available. There are works aiming to address this issue by learning 3D-grounded representation of the environment and modeling the dynamics in a latent vector space [34, 32]. However, these models typically encode the entire scene into one single vector. Such design does not capture the structure of the underlying systems, limiting its generalization to compositional systems or systems of different sizes (e.g., unseen container shapes or different numbers of floating ice cubes).

In this work, we propose 3D Visual Intuitive Physics (3D-IntPhys), a framework that learns intuitive physics models of the environment with **explicit 3D** and **compositional** structures with **visual inputs**.

Specifically, the model consists of (1) a perception module based on conditional Neural Radiance Fields (NeRF) [41, 68] that transforms the input images and instance masks into 3D point representations and (2) a dynamics module instantiated as graph neural networks to model the interactions between the points and predict their evolutions over time. Despite advances in graph-based dynamics networks [47, 36], existing methods require strong supervision provided by 3D GT point trajectories, which are hard to obtain in most real setups. To tackle the problem, we train the dynamics model using (1) a distribution-based loss function measuring the difference between the predicted point sets and the actual point distributions at the future timesteps and (2) a spacing loss to avoid degenerated point set predictions. Our perception module learns spatial-equivariant representations of the environment grounded in the 3D space, which then transforms into points as a flexible representation to describe the system’s state. Our dynamics module regards the point set as a graph and exploits the compositional structure of the point systems.

The structures allow the model to capture the compositionality of the underlying environment, handle systems involving objects with complicated physical properties, and perform extrapolated generalization, which we show via experiments greatly outperform various baselines without a structured 3D representation space.

2 Related Work

Visual dynamics learning. Existing works learn to predict object motions from pixels using frame-centric features [1, 20, 4, 24, 23, 53, 30, 59, 69, 10, 25, 62] or object-centric features [21, 61, 28, 44, 27, 58, 14, 22, 45, 66], yet, most works only demonstrate the learning in 2D scenes with objects

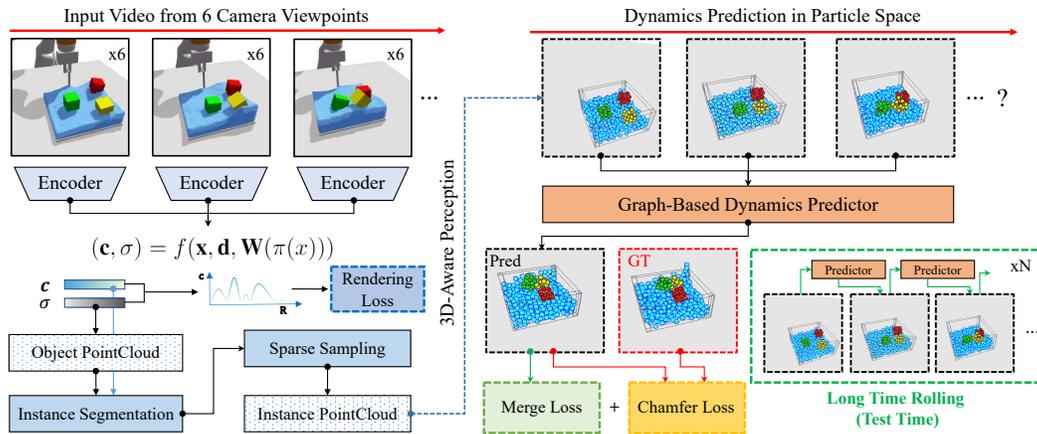


Figure 2: **Overview of 3D Visual Intuitive Physics (3D-IntPhys).** Our model consists of two major components: **Left:** The perception module maps the visual observations into implicit neural representations of the environment. We then subsample from the reconstructed implicit volume to obtain a particle representation of the environment. **Right:** The dynamics module, instantiated as graph neural networks, models the interaction within and between the objects and predicts the evolution of the particle set.

moving only on a 2D plane. We argue that one reason that makes it hard for these existing methods to be applied to general 3D visual scenes is because they often operate on view-dependent features that can change dramatically due to changes in the camera viewpoint, which shouldn't have any effect on the actual motion of the objects. Recent works by [9] have shown that only methods that use 3D view-invariant representations can pave the way toward human-level physics dynamics prediction in diverse scenarios.

Researchers have attempted to learn object motion in 3D [55, 63, 40, 34], [55] and [63] use object-centric volumetric representations inferred from RGB-D to predict object motion, yet, these volumetric approaches have much higher computation costs than 2D methods due to the 4D representation bottleneck, which hinders them from scaling up to more complex scenes. [40] use self-supervised 3D keypoints and [15, 16] use implicit representations to model multi-object dynamics but cannot handle objects with high degrees of freedom like fluid and granular materials. [34] use neural implicit representation to reduce the potential computational cost, yet the works have not shown how the approach can generalize to unseen scenarios. Our works aim to solve the tasks of learning generalizable object dynamics in 3D by combining the generalization strength of input-feature-conditioned implicit representation and point-based dynamics models.

Point-based dynamics models. Existing works in point- and mesh-based dynamics models [36, 42, 57, 47, 43] have shown impressive results in predicting the dynamics of rigid objects, fluid [36, 42, 57, 47, 3, 13], deformable objects [36, 42, 47], and clothes [43, 38]. Most works require access to full 3D states of the points during training and testing, yet, such information is usually not accessible in a real-world setup. [35] learn a visual frontend to infer 3D point states from images, but still require 3D point states and trajectories during training time. [50] propose to learn point dynamics directly from vision, but they only consider elasto-plastic objects consisting of homogeneous materials. How to learn about 3D point states and their motion from raw pixels remain a question. Our paper tries to build the link from pixels to points using recent advances in unsupervised 3D inference from images using NeRF [41, 68].

3 Methods

We present 3D Visual Intuitive Physics (3D-IntPhys), a model that learns to simulate physical events from unlabeled images (Figure 2). 3D-IntPhys contains a perception module that transforms visual observations into a 3D point cloud that captures the object geometries (Section 3.1) and a point-based simulator that learns to simulate the rollout trajectories of the points (Section 3.2). The design choice of learning physics simulation in a 3D-point representation space enables stronger simulation performance and generalization ability. The performance gain mainly comes from the fact that describing/learning objects' motion and interactions in 3D are easier compared to doing so in 2D since objects live and move persistently in the 3D space. 3D-IntPhys also supports better

generalization ability since its neural architecture explicitly models how local geometries of two objects/parts interact, and these local geometries and interactions can be shared across different and novel object combinations.

Although 3D-IntPhys learns to simulate in a 3D representation space, we show it can learn without any 3D supervision such as dense point trajectories as in previous work [47, 36]. Dense point trajectories are hard and sometimes impossible to obtain in the real world, e.g., capturing the trajectories of each water point. 3D-IntPhys does not require such 3D supervision and can simply learn by observing videos of the scene evolution.

3.1 2D-to-3D Perception Module

Given a static scene, the perception module learns to transform one or a few posed RGB images, $\mathbf{I} = \{(I_i, \pi_i) | i \in \{1, 2, \dots, N_v\}\}$, taken from N_v different views, into a 3D point cloud representation of the scene, \mathbf{X} . We train the model in an unsupervised manner through view reconstruction, using a dataset consisting of N_t videos, where each video has N_f frames, and each frame contains images taken from N_v viewpoints.

Neural Radiance Field (NeRF). NeRF [41] learns to reconstruct a volumetric radiance field of a scene from unlabeled multi-view images. After training, the model learns to predict the RGB color \mathbf{c} and the corresponding density σ of a query 3D point $\mathbf{x} \in \mathbb{R}^3$ from the viewing direction $\mathbf{d} \in \mathbb{R}^3$ with a function $(\mathbf{c}, \sigma) = f(\mathbf{x}, \mathbf{d})$. We can formulate a camera ray as $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$, where $\mathbf{o} \in \mathbb{R}^3$ is the origin of the ray. The volumetric radiance field can then be rendered into a 2D image via $\hat{\mathbf{C}}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(t)\mathbf{c}(t)dt$, where $T(t) = \exp(-\int_{t_n}^t \sigma(s)ds)$ handles occlusion. The rendering range is controlled by the depths of the near and far plane (i.e., t_n and t_f). We can train NeRF through view prediction by:

$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}(\mathbf{p})} \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|, \quad (1)$$

where $\mathcal{R}(\mathbf{p})$ is the set of camera rays sampled from target camera pose \mathbf{p} .

Image-conditioned NeRF. To infer the NeRF function from an image, previous work proposed to encode the input image into a vector, with a CNN encoder, as a conditioning input to the target NeRF function [34]. We found this type of architecture is in general hard to train and does not generalize well. Instead, we adopt pixelNeRF [68], which conditions NeRF rendering with local features, as opposed to global features. Given an image I in a scene, pixelNeRF first extracts a feature volume using a CNN encoder $\mathbf{W} = E(I)$. For a point \mathbf{x} in the world coordinate, we retrieve its feature vector by projecting it onto the image plane, so that we can get the feature vector $\mathbf{W}(\pi(\mathbf{x}))$. PixelNeRF combines the feature vector together with the 3D position of that point and predict the RGB color and density information:

$$\mathbf{V}(\mathbf{x}) = (\mathbf{c}, \sigma) = f(\mathbf{x}, \mathbf{d}, \mathbf{W}(\pi(\mathbf{x}))). \quad (2)$$

In the experiment section, we will show that it is surprisingly effective to train **only one general model** to learn a conditional Neural Radiance Field that can apply to all videos in one type of scene (e.g. FluidPour) with **five different settings** (e.g. extrapolate setting), which provides a better 3D representation for the scene and greatly facilitates the learning of 3D Intuitive Physics.

Explicit 3D representation from pixelNeRF. From a few posed RGB images \mathbf{I} , of a scene s , we infer a set of points for O_s target object (such as fluid, cube) in the scene. We achieve this by first sampling a set of points according to the predicted occupancy measure, then clustering the points into objects using object segmentations. We found that sampling with low resolution will hurt the quality of the rendered point cloud to generate objects with inaccurate shapes, while sampling with high resolution will increase the computation for training the dynamics model since the input size increases. To speed up training while maintaining the quality of the reconstructed point cloud, we first infer the points with higher resolution and do sparse sampling of each point cloud using FPS (Farthest Point Sampling) [17]. Next, we cluster the inferred points into objects according to object segmentation masks. Since solving object segmentation in general is not the main focus of this paper, we resort to using the color information to obtain the masks.

3.2 Point-Based Dynamics Learner

Given the point representation at the current time step, \mathbf{X}_t , the dynamics simulator predicts the points' evolution T steps in the future, $\{\mathbf{X}_{t+1}, \mathbf{X}_{t+2}, \dots, \mathbf{X}_{t+T}\}$, using graph-based networks [47, 36].

We first form a graph (V, E) based on the distance between points. If the distance between two points is smaller than a threshold δ , we include an edge between these two points. Each vertex $v_i = (\hat{x}_i, a_i^v) \in V$ contains the velocity of the point, \hat{x}_i , and point attributes, a_i^v , to indicate the point's type. For each relation, $(i, j) \in E$, we have its associated relation attribute a_{ij}^e , indicating the types of relation and the relative distance between the connected points.

Spatial message passing and propagation. At time step t , we can do message passing to update the points, $v_i \in V$, and relation representations, $(i, j) \in E$, in the graph:

$$g_{ij,t} = Q_e(v_{i,t}, v_{j,t}, a_{ij}^e), (i, j) \in E \quad (3)$$

$$h_{i,t} = Q_v(v_{i,t}, \sum_{k \in \{j | (i,j) \in E\}} g_{ik,t}), v_i \in V \quad (4)$$

where Q_v and Q_e are encoders for vertices and relations respectively. Please refer to [7] for more details. Though this kind of message passing can help with updating representation, it can only share one-hop information in each step, limiting its performance on instantaneous passing of forces. To improve long-range instantaneous effect propagation, we use multi-step message propagation as in [37, 36]. The propagation step is shown in Alg 1:

Algorithm 1: Point-based Dynamics Predictor

Data: Current timestep t , point cloud V_t , vertex encoder Q_v , edge encoder Q_e , vertex propagator P_v , edge propagator P_e , state predictor f_s

Result: V_{t+1}

Form graph $G_t = (V_t, E_t)$

//message passing

$$g_{ij,t} = Q_e(v_{i,t}, v_{j,t}, a_{ij}^e), (i, j) \in E_t$$

$$h_{i,t} = Q_v(v_{i,t}, \sum_{k \in \{j | (i,j) \in E_t\}} g_{ik,t}), v_i \in V_t$$

//message propagation

$$h_{i,t}^0 = h_{i,t}, g_{i,t}^0 = g_{i,t},$$

for $l \in \{1, 2, 3, \dots, L\}$ **do**

$$\quad \left| \begin{array}{l} g_{ij,t}^l = P_e(g_{ij,t}^{l-1}, h_{i,t}^{l-1}, h_{j,t}^{l-1}), (i, j) \in E_t \\ h_{i,t}^l = P_v(h_{i,t}^{l-1}, \sum_{k \in \{j | (i,j) \in E_t\}} g_{ik,t}^l), v_i \in V_t \end{array} \right.$$

end

//state prediction

$$v_{i,t+1} = f_s(h_{i,t}^L)$$

$$V_{t+1} = \{v_{i,t+1}\}$$

where P_e, P_v are propagation functions of nodes and edges, respectively, and $g_{ij,t}^l$ is the effect of relation (i, j) in propagation step l . $h_{i,t}^l$ is the hidden states for each point in the propagation process. Finally, we have the predicted states of points at time step $t + 1$ after L steps of propagation:

$$\hat{v}_{i,t+1} = f_s(h_{i,t}^L). \quad (5)$$

Environments. We assume that the surrounding environment (e.g., the table) is known and the robot/tool/container are of known shape and fully actuated, where the model has access to their complete 3D state information. We convert the full 3D states into points through sampling on the 3D meshes and include these points in the prediction of the graph-based dynamics.

Fluids, rigid bodies, and granular materials. We distinguish different materials by using different point attributes a_i^v . We also set different relation attributes a_{ij}^e in Equation 3 to distinguish different interaction (e.g., Rigid-Fluids, Fluids-Fluids, Granular-Pusher). For rigid objects, to ensure the object shapes remain consistent throughout the rollout predictions, we add a differentiable rigid constraint in the prediction head following [36].

Training dynamics model without point-level correspondence. Since our perception model parses each RGB image into object-centric point clouds independently, there does not exist an explicit one-to-one correspondence for points across frames. To handle this, we measure the Chamfer distance between the prediction $\hat{\mathbf{X}}_t = (\hat{V}_t, \hat{E}_t)$ from the dynamics network and the inferred point state $\mathbf{X}_t = (V_t, E_t)$ from the perception module and treat it as the objective function. The Chamfer

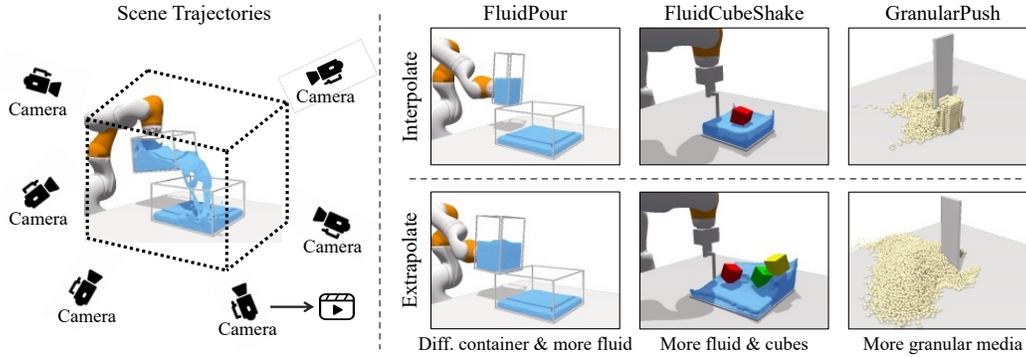


Figure 3: **Data Collection and Evaluation Setups.** **Left:** We collect multi-view videos of the environment from six cameras. **Right:** We consider a diverse set of evaluating environments involving fluids, rigid objects, granular materials, and their interactions with the fully-actuated container and the environment. We evaluate the learned visual intuitive physics model on both the interpolated settings (i.e., seen environment but with different action sequences) and extrapolated settings (i.e., unseen environment with different amounts of fluids, cubes, granular pieces, and containers of different sizes).

distance between two point cloud \hat{V} and V is defined as:

$$L_c(\hat{V}, V) = \frac{1}{\|\hat{V}\|} \sum_{x \in \hat{V}} \min_{y \in V} \|x - y\|_2^2 + \frac{1}{\|V\|} \sum_{x \in V} \min_{y \in \hat{V}} \|x - y\|_2^2. \quad (6)$$

We found that training the model with Chamfer distance in dense scenes with granular materials will often lead to predictions with unevenly distributed points where some points stick too close to each other. To alleviate this issue, we further introduce a spacing loss L_s , which penalizes the gated distance (gated by d_{\min}) of nearest neighbor of each point to ensure enough space between points:

$$L_s(\hat{V}) = \sum_{v \in \hat{V}} (\text{ReLU}(d_{\min} - \min_{v' \in \{\hat{V} \setminus v\}} \|v' - v\|_2^2))^2. \quad (7)$$

The one-step prediction loss L_{dy} for training the dynamics model is $L_c(\hat{V}, V) + \sigma L_s(\hat{V})$ where σ reweights the second loss. To improve long-term rollout accuracy, we train the model with two-step predictions using the first predicted state as input and feed it back into the model to generate the second predicted state. With the two-step loss, the model becomes more robust to errors generated from its own prediction. Finally, the L_{dy} losses for all rolling steps are summed up to get the final loss for this trajectory. More implementation details are included in the supplementary material.

4 Experiments

The experiment section aims to answer the following three questions. (1) How well can the visual inference module capture the content of the environment (i.e., can we use the learned representations to reconstruct the scene)? (2) How well does the proposed framework perform in scenes with objects of complicated physical properties (e.g., fluids, rigid and granular objects) compared to baselines without explicit 3D representations? (3) How well do the models generalize in extrapolate scenarios?

Datasets. We generated three simulated datasets using the physics simulator Nvidia FleX [39]. Each of the datasets represents one specific kind of manipulation scenario, where a robot arm interacts with rigid, fluid, and granular objects (Figure 3). For each of the three scenarios, we apply randomized input actions and change some properties of objects in the scene, e.g., the shape of the container, the amount of water, and the color/number of cubes, to make it diverse. To test the generalization capability of the trained model, we design extrapolated datasets where the data is generated from an extrapolated set of parameters outside the training distribution.

a) FluidPour. This scenario contains a fully-actuated cup pouring fluid into a container. We design the extrapolate dataset to have a larger container, more quantity of fluid, and different pouring actions.

b) FluidCubeShake. This scenario contains a fully-actuated container that moves on top of a table. Inside the container are fluids and cubes with diverse colors. We design the extrapolate dataset to have different container shapes, number of cubes, cube colors, and different shaking actions.

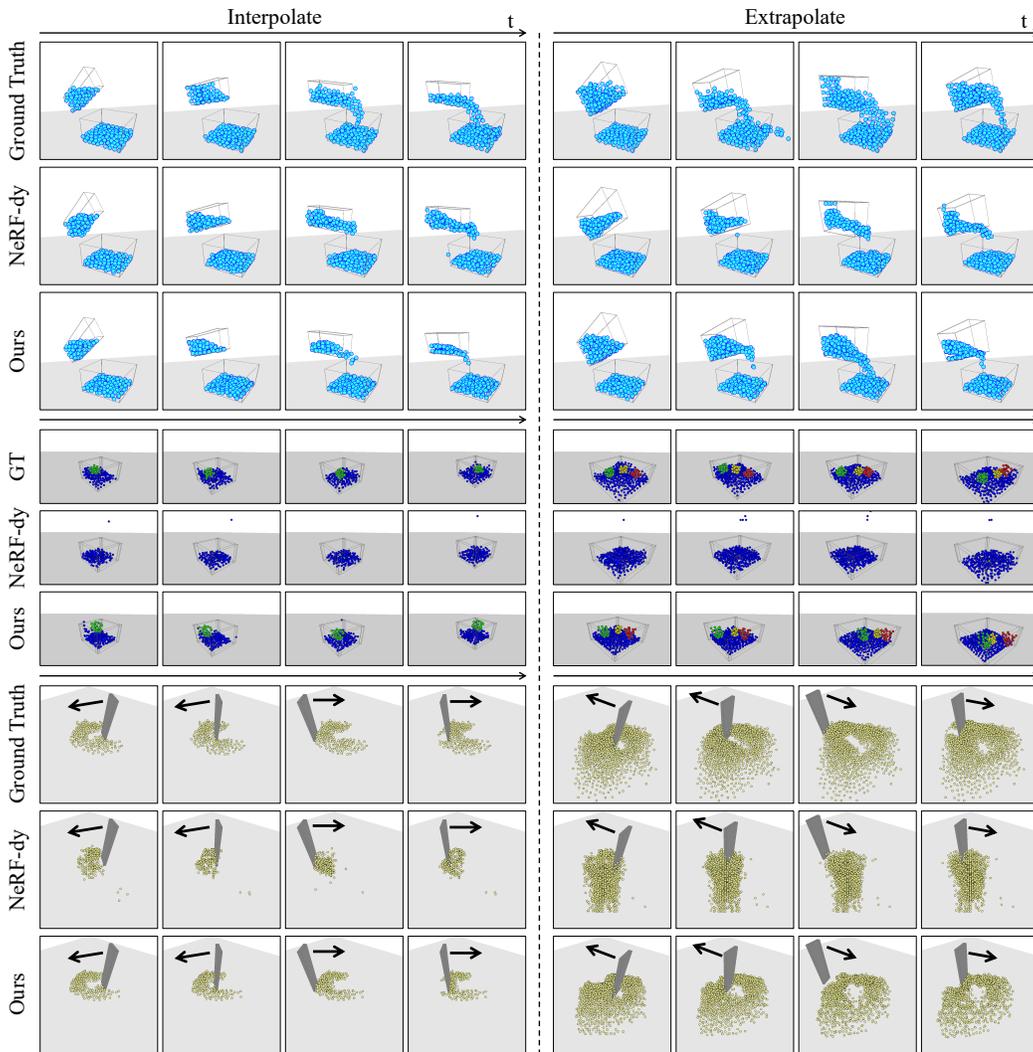


Figure 4: **Qualitative Results of the Dynamics Module on Future Prediction.** Here we visualize our model’s predicted future evolution of the particle set as compared with the NeRF-dy [34] baseline in both interpolate and extrapolate settings. Our method correctly identifies the shape/distribution of the fluids, rigid objects, and granular pieces with much better accuracy than NeRF-dy. The future evolution predicted by our method also matches the ground truth much better and produces reasonable results even in extrapolate settings.

c) GranularPush. This environment contains a fully-actuated board pushing a pile of granular pieces. We design the extrapolate dataset to have a larger quantity of granular objects in the scene than the model has ever seen during training.

Baselines. We compare our method with two baselines, NeRF-dy [34] and autoencoder (AE) (similar to GQN [18] augmented with a latent-space dynamics model). NeRF-dy is a 3D-aware framework that also learns intuitive physics from multi-view videos. Yet, instead of learning the object dynamics with explicit and compositional 3D representations, the model learns dynamics models with implicit 3D representations in the form of a single latent vector. We also compare our method with an autoencoder-based reconstruction model (AE) [18] that can perform novel-view synthesis but is worse at handling 3D transformations than neural implicit representations. AE first learns scene representations through per-frame image reconstruction, and then it learns a dynamics model on top of the learned latent representations. All methods take RGB images and camera parameters as inputs. To incorporate object-level information, we perform color-based segmentation to obtain object masks as additional inputs to the baselines. The implementation details and parameter settings of our method can be found in the supplementary materials.

| Metrics | Model | FluidPour | | FluidCubeShake | | GranularPush | |
|---------|---------|---------------|---------------|----------------|--------------|---------------|---------------|
| | | InD | OoD | InD | OoD | InD | OoD |
| MSE(↓) | AE | 451.03 | 542.86 | 869.3 | 1727.55 | 562.06 | 1537.2 |
| | NeRF-dy | 202.95 | 317.27 | 527.46 | 1585.97 | 481.95 | 1020.0 |
| | Ours | 111.66 | 124.33 | 66.52 | 81.38 | 147.97 | 646.85 |
| SSIM(↑) | AE | 0.86 | 0.84 | 0.71 | 0.86 | 0.81 | 0.62 |
| | NeRF-dy | 0.89 | 0.86 | 0.73 | 0.65 | 0.81 | 0.61 |
| | Ours | 0.90 | 0.89 | 0.94 | 0.93 | 0.89 | 0.69 |

Table 1: **Quantitative Results of the Perception Module.** We compare our method with autoencoder (AE) and NeRF-dy [34] with additional instance masks based on color. We measure the quality of rendered images by computing the Mean Squared Error (MSE) and Structural Similarity Index Measure (SSIM) compared to the ground truth. InD stands for in-distribution tests, and OoD stands for out-of-distribution tests.

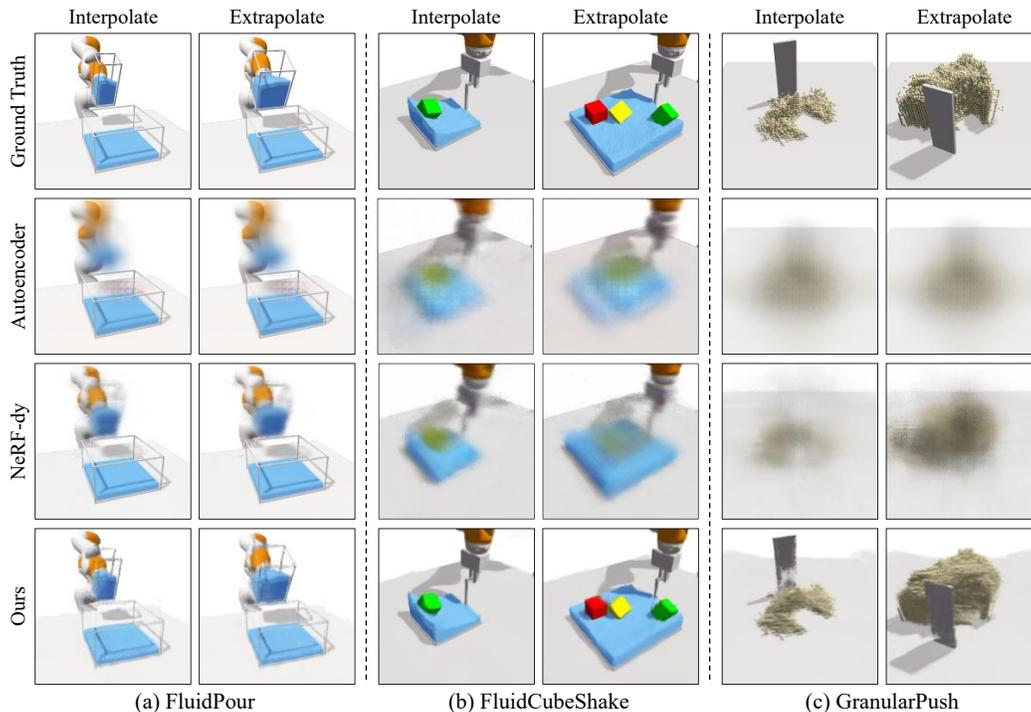


Figure 5: **Qualitative Reconstruction Results of the Perception Module.** The images generated by our method contain more visual details and are much better aligned with the ground truth. Our model is much better at handling large scene variations than NeRF-dy, especially in extrapolate settings.

4.1 Image Reconstruction From Learned Scene Representations

We test how well the perception modules capture scene information by evaluating the visual front-end of all models on their ability to reconstruct the observed scene from the inferred representations. We measure the difference between the reconstructed and ground truth images with Mean Squared Error (MSE) and Structural Similarity (SSIM) in pixel level (Table 1). Our perception module outperforms all baselines in all three environments. The performance gap is exaggerated in extrapolate settings, especially in scenarios that involve complex interactions between rigid and deformable materials (Figure 5 qualitative comparisons).

4.2 Learned Visual Dynamics On In-Distribution Held-Out Scenes

Next, we compare long-term rollouts in the 3D space. We evaluate the models using the Chamfer distance between the predicted point cloud and the ground truth. For NeRF-dy, we decode the predicted rollouts latent vectors into the point cloud with the learned NeRF decoder. We exclude the comparison with AE since it is unclear how to decode the learned representations into point clouds. We show quantitative comparison in Figure 6 and qualitative results in Figure 4. 3D-IntPhys can

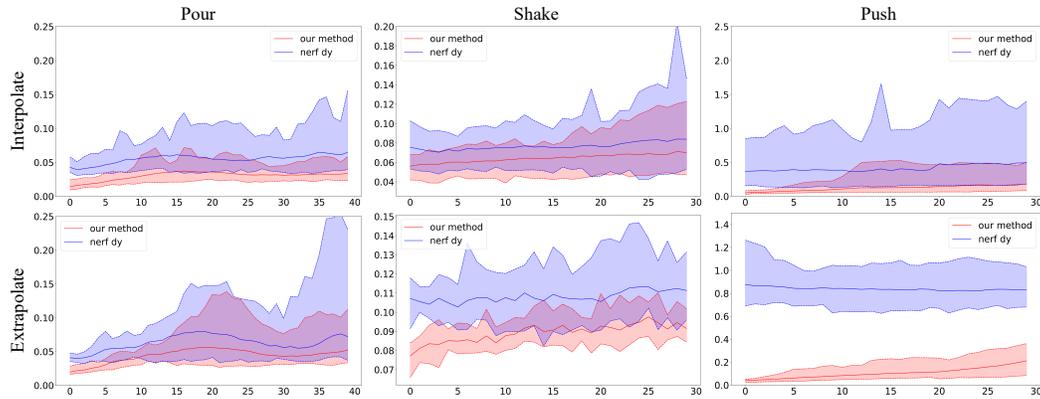


Figure 6: **Quantitative Results of the Dynamics Module.** This figure compares our method and NeRF-dy [34] on their long-horizon open-loop future prediction loss. The loss is measured as the Chamfer distance between the predicted particle set evolution and the actual future. Our method outperforms the baseline in both interpolate and extrapolate settings, showing the benefits of explicit 3D modeling.

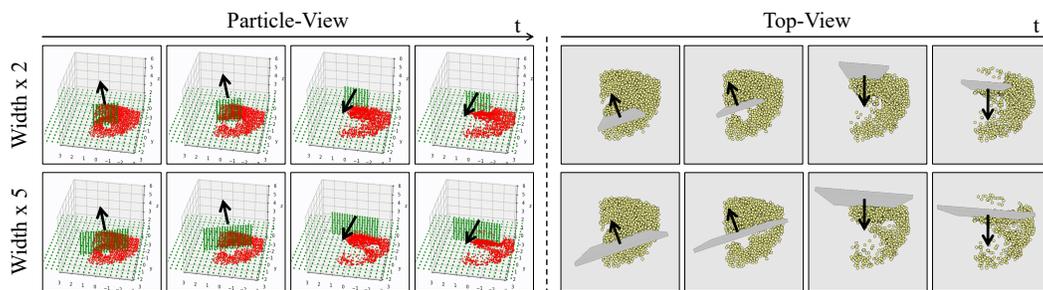


Figure 7: **Strong Generalization Ability of the Dynamics Module to Wider Pushers.** We evaluate our dynamics model on unseen width of pushers in **GranularPush** environment. The left part shows in 3D space where red indicates granular materials, green shows the table and pusher, and the arrow shows about to move. The right part shows from the top view of the rendering results.

learn reasonable scene dynamics in all scenarios and significantly outperforms NeRF-dy. While NeRF-dy can learn relatively reasonable movements of fluids, it fails to learn complex dynamics such as the floating cube and the morphing of the granular materials. The results suggest that the proposed explicit 3D point-based representations are critical to learning complex multi-material dynamics.

4.3 Generalization on Out-of-Distribution Scenes

To test the generalization ability of the models, we introduce extrapolate settings of all of the three scenarios. See “Extrapolate” results in Table 1, Figure 5, 6, and 4. The proposed 3D-IntPhys generalizes well to extrapolate settings both at the visual perception stage and the dynamics prediction stage, whereas NeRF-dy and autoencoder both fail at generalizing under extrapolate settings. For example, in **FluidShake**, both baselines cannot capture the number and the color of the rigid cubes (Figure 5). And in **GranularPush**, both baselines fail to capture the distributions of the granular materials. NeRF-dy performs much worse on extrapolation scenes compared to in-distribution scenes, suggesting that incorporating 3D information in an explicit way, as opposed to implicit, is much better at capturing the structure of the underlying environment, thus leading to better generalization. We further test our model on completely unseen changes to the environment – in the **GranularPush** environment, we extend the width of the pusher by a factor of 2 and 5. Though the stretched pusher has never shown in the training data, our model can make reasonable pushing predictions (see Fig 7).

5 Conclusions

In this work, we propose a 3D-aware and compositional framework, 3D-IntPhys, to learn intuitive physics from unlabeled visual inputs. Our framework can work on complex scenes involving fluid,

rigid objects, and granular materials, and generalize to unseen scenes with containers of different sizes, more objects, or larger quantities of fluids and granular pieces. We show the proposed model outperforms baselines by a large margin, highlighting the importance of learning dynamics models in an explicit 3D representations space. The major limitation of our work is the assumption of access to object masks. However, with the progress on segmentation in the wild [71, 29], we believe that it will be possible to get such kinds of masks in real-world 3D environments. Our work serves as a pioneer in visual intuitive physics learning of complex scenes, and it is an exciting future direction to learn more complex intuitive physics from real-world data with the help of these large models.

References

- [1] P. Agrawal, A. V. Nair, P. Abbeel, J. Malik, and S. Levine. Learning to poke by poking: Experiential learning of intuitive physics. *Advances in neural information processing systems*, 29, 2016.
- [2] A. Ajay, M. Bauzá, J. Wu, N. Fazeli, J. B. Tenenbaum, A. Rodriguez, and L. P. Kaelbling. Combining physical simulators and object-based networks for control. *CoRR*, abs/1904.06580, 2019.
- [3] K. R. Allen, T. Lopez-Guevara, K. L. Stachenfeld, A. Sanchez-Gonzalez, P. W. Battaglia, J. B. Hamrick, and T. Pfaff. Physical design using differentiable learned simulators. *CoRR*, abs/2202.00728, 2022.
- [4] M. Babaeizadeh, M. T. Saffar, S. Nair, S. Levine, C. Finn, and D. Erhan. Fitvid: Overfitting in pixel-level video prediction. *CoRR*, abs/2106.13195, 2021.
- [5] R. Baillargeon, E. S. Spelke, and S. Wasserman. Object permanence in five-month-old infants. *Cognition*, 20:191–208, 1985.
- [6] C. Bates, I. Yildirim, J. B. Tenenbaum, and P. W. Battaglia. Modeling human intuitions about liquid flow with particle-based simulation. *CoRR*, abs/1809.01524, 2018.
- [7] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. F. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, Ç. Gülçehre, H. F. Song, A. J. Ballard, J. Gilmer, G. E. Dahl, A. Vaswani, K. R. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.
- [8] P. W. Battaglia, J. B. Hamrick, and J. B. Tenenbaum. Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110:18327 – 18332, 2013.
- [9] D. M. Bear, E. Wang, D. Mrowca, F. J. Binder, H.-Y. F. Tung, R. Pramod, C. Holdaway, S. Tao, K. Smith, L. Fei-Fei, et al. Physion: Evaluating physical prediction from vision in humans and machines. *arXiv preprint arXiv:2106.08261*, 2021.
- [10] Y. Burda, H. Edwards, D. Pathak, A. Storkey, T. Darrell, and A. A. Efros. Large-scale study of curiosity-driven learning. In *ICLR*, 2019.
- [11] S. Carey and F. Xu. Infants’ knowledge of objects: beyond object files and object tracking. *Cognition*, 80(1):179–213, 2001. Objects and Attention.
- [12] M. B. Chang, T. D. Ullman, A. Torralba, and J. B. Tenenbaum. A compositional object-based approach to learning physical dynamics. *CoRR*, abs/1612.00341, 2016.
- [13] F. de Avila Belbute-Peres, T. D. Economou, and J. Z. Kolter. Combining differentiable PDE solvers and graph neural networks for fluid flow prediction. *CoRR*, abs/2007.04439, 2020.
- [14] D. Ding, F. Hill, A. Santoro, and M. M. Botvinick. Object-based attention for spatio-temporal reasoning: Outperforming neuro-symbolic models with flexible distributed architectures. *CoRR*, abs/2012.08508, 2020.
- [15] D. Driess, J.-S. Ha, M. Toussaint, and R. Tedrake. Learning models as functionals of signed-distance fields for manipulation planning. In *Conference on Robot Learning*, pages 245–255. PMLR, 2022.
- [16] D. Driess, Z. Huang, Y. Li, R. Tedrake, and M. Toussaint. Learning multi-object dynamics with compositional neural radiance fields. *arXiv preprint arXiv:2202.11855*, 2022.
- [17] Y. Eldar, M. Lindenbaum, M. Porat, and Y. Zeevi. The farthest point strategy for progressive image sampling. *IEEE Transactions on Image Processing*, 6(9):1305–1315, 1997.
- [18] S. A. Eslami, D. Jimenez Rezende, F. Besse, F. Viola, A. S. Morcos, M. Garnelo, A. Ruderman, A. A. Rusu, I. Danihelka, K. Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018.

- [19] C. Finn, I. J. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. *CoRR*, abs/1605.07157, 2016.
- [20] C. Finn and S. Levine. Deep visual foresight for planning robot motion. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2786–2793. IEEE, 2017.
- [21] K. Fragkiadaki, P. Agrawal, S. Levine, and J. Malik. Learning visual predictive models of physics for playing billiards. In Y. Bengio and Y. LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [22] R. Girdhar, L. Gustafson, A. Adcock, and L. van der Maaten. Forward prediction for physical reasoning. *CoRR*, abs/2006.10734, 2020.
- [23] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- [24] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR, 2019.
- [25] D. Hafner, T. P. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. *CoRR*, abs/1912.01603, 2019.
- [26] M. Janner, J. Fu, M. Zhang, and S. Levine. When to trust your model: Model-based policy optimization. *CoRR*, abs/1906.08253, 2019.
- [27] M. Janner, S. Levine, W. T. Freeman, J. B. Tenenbaum, C. Finn, and J. Wu. Reasoning about physical interactions with object-oriented prediction and planning. In *International Conference on Learning Representations*, 2019.
- [28] T. Kipf, E. van der Pol, and M. Welling. Contrastive learning of structured world models. *arXiv preprint arXiv:1911.12247*, 2019.
- [29] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- [30] A. X. Lee, R. Zhang, F. Ebert, P. Abbeel, C. Finn, and S. Levine. Stochastic adversarial video prediction. *CoRR*, abs/1804.01523, 2018.
- [31] A. Lerer, S. Gross, and R. Fergus. Learning physical intuition of block towers by example. *CoRR*, abs/1603.01312, 2016.
- [32] T. Li, M. Slavcheva, M. Zollhoefer, S. Green, C. Lassner, C. Kim, T. Schmidt, S. Lovegrove, M. Goesele, and Z. Lv. Neural 3d video synthesis. *arXiv preprint arXiv:2103.02597*, 2021.
- [33] W. Li, S. Azimi, A. Leonardis, and M. Fritz. To fall or not to fall: A visual approach to physical stability prediction. *CoRR*, abs/1604.00066, 2016.
- [34] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba. 3d neural scene representations for visuomotor control. *arXiv preprint arXiv:2107.04004*, 2021.
- [35] Y. Li, T. Lin, K. Yi, D. Bear, D. L. Yamins, J. Wu, J. B. Tenenbaum, and A. Torralba. Visual grounding of learned physical models. In *International Conference on Machine Learning*, 2020.
- [36] Y. Li, J. Wu, R. Tedrake, J. B. Tenenbaum, and A. Torralba. Learning particle dynamics for manipulating rigid bodies, deformable objects, and fluids. In *ICLR*, 2019.
- [37] Y. Li, J. Wu, J.-Y. Zhu, J. B. Tenenbaum, A. Torralba, and R. Tedrake. Propagation networks for model-based control under partial observation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1205–1211. IEEE, 2019.
- [38] X. Lin, Y. Wang, Z. Huang, and D. Held. Learning visible connectivity dynamics for cloth smoothing. In *Conference on Robot Learning*, 2021.

- [39] M. Macklin, M. Müller, N. Chentanez, and T.-Y. Kim. Unified particle physics for real-time applications. *ACM Transactions on Graphics (TOG)*, 33(4):1–12, 2014.
- [40] L. Manuelli, Y. Li, P. Florence, and R. Tedrake. Keypoints into the future: Self-supervised correspondence in model-based reinforcement learning. *arXiv preprint arXiv:2009.05085*, 2020.
- [41] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020.
- [42] D. Mrowca, C. Zhuang, E. Wang, N. Haber, L. Fei-Fei, J. B. Tenenbaum, and D. L. K. Yamins. Flexible neural representation for physics prediction. *CoRR*, abs/1806.08047, 2018.
- [43] T. Pfaff, M. Fortunato, A. Sanchez-Gonzalez, and P. W. Battaglia. Learning mesh-based simulation with graph networks. In *International Conference on Learning Representations*, 2021.
- [44] H. Qi, X. Wang, D. Pathak, Y. Ma, and J. Malik. Learning long-term visual dynamics with region proposal interaction networks. In *ICLR*, 2021.
- [45] R. Riochet, J. Sivic, I. Laptev, and E. Dupoux. Occlusion resistant learning of intuitive physics from videos. *CoRR*, abs/2005.00069, 2020.
- [46] A. N. Sanborn, V. K. Mansinghka, and T. L. Griffiths. Reconciling intuitive physics and newtonian mechanics for colliding objects. *Psychological review*, 120 2:411–37, 2013.
- [47] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning to simulate complex physics with graph networks. In *International Conference on Machine Learning*, pages 8459–8468. PMLR, 2020.
- [48] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. A. Riedmiller, R. Hadsell, and P. W. Battaglia. Graph networks as learnable physics engines for inference and control. *CoRR*, abs/1806.01242, 2018.
- [49] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- [50] H. Shi, H. Xu, Z. Huang, Y. Li, and J. Wu. Robocraft: Learning to see, simulate, and shape elasto-plastic objects with graph networks. *arXiv preprint arXiv:2205.02909*, 2022.
- [51] K. Smith, L. Mei, S. Yao, J. Wu, E. Spelke, J. Tenenbaum, and T. Ullman. Modeling expectation violation in intuitive physics with coarse probabilistic object representations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [52] E. S. Spelke. Principles of object perception. *Cognitive Science*, 14(1):29–56, 1990.
- [53] H. Suh and R. Tedrake. The surprising effectiveness of linear models for visual foresight in object pile manipulation. In *International Workshop on the Algorithmic Foundations of Robotics*, pages 347–363. Springer, 2020.
- [54] A. Tacchetti, H. F. Song, P. A. M. Mediano, V. F. Zambaldi, N. C. Rabinowitz, T. Graepel, M. M. Botvinick, and P. W. Battaglia. Relational forward models for multi-agent learning. *CoRR*, abs/1809.11044, 2018.
- [55] H.-Y. F. Tung, Z. Xian, M. Prabhudesai, S. Lal, and K. Fragkiadaki. 3d-oes: Viewpoint-invariant object-factorized environment simulators. *arXiv preprint arXiv:2011.06464*, 2020.
- [56] T. Ullman, E. Kosoy, I. Yildirim, A. A. Soltani, M. H. Siegel, J. Tenenbaum, and E. S. Spelke. Draping an elephant: Uncovering children’s reasoning about cloth-covered objects. In *Proceedings of the 41st Annual Conference of the Cognitive Science Society*, pages 3008–3014, 2019.

- [57] B. Ummenhofer, L. Prantl, N. Thuerey, and V. Koltun. Lagrangian fluid simulation with continuous convolutions. In *International Conference on Learning Representations*, 2019.
- [58] R. Veerapaneni, J. D. Co-Reyes, M. Chang, M. Janner, C. Finn, J. Wu, J. B. Tenenbaum, and S. Levine. Entity abstraction in visual model-based reinforcement learning. *CoRR*, abs/1910.12827, 2019.
- [59] C. Vondrick, H. Pirsiavash, and A. Torralba. Anticipating the future by watching unlabeled video. *CoRR*, abs/1504.08023, 2015.
- [60] M. Watter, J. Springenberg, J. Boedecker, and M. Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. *Advances in neural information processing systems*, 28, 2015.
- [61] N. Watters, D. Zoran, T. Weber, P. Battaglia, R. Pascanu, and A. Tacchetti. Visual interaction networks: Learning a physics simulator from video. *Advances in neural information processing systems*, 30, 2017.
- [62] B. Wu, S. Nair, R. Martín-Martín, L. Fei-Fei, and C. Finn. Greedy hierarchical variational autoencoders for large-scale video prediction. *CoRR*, abs/2103.04174, 2021.
- [63] Z. Xu, Z. He, J. Wu, and S. Song. Learning 3d dynamic scene representations for robot manipulation. In *Conference on Robotic Learning (CoRL)*, 2020.
- [64] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song. Densephysnet: Learning dense physical object representations via multi-step dynamic interactions. In *Robotics: Science and Systems (RSS)*, 2019.
- [65] T. Xue, J. Wu, K. L. Bouman, and W. T. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances In Neural Information Processing Systems*, 2016.
- [66] Y. Ye, D. Gandhi, A. Gupta, and S. Tulsiani. Object-centric forward modeling for model predictive control. In *CoRL*, 2019.
- [67] Y. Ye, M. Singh, A. Gupta, and S. Tulsiani. Compositional video prediction. In *International Conference on Computer Vision (ICCV)*, 2019.
- [68] A. Yu, V. Ye, M. Tancik, and A. Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021.
- [69] M. Zhang, S. Vikram, L. Smith, P. Abbeel, M. J. Johnson, and S. Levine. SOLAR: deep structured latent representations for model-based reinforcement learning. *CoRR*, abs/1808.09105, 2018.
- [70] R. Zhang, J. Wu, C. Zhang, W. T. Freeman, and J. B. Tenenbaum. A comparative evaluation of approximate probabilistic simulation and deep neural networks as accounts of human physical scene understanding. *CoRR*, abs/1605.01138, 2016.
- [71] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Gao, and Y. J. Lee. Segment everything everywhere all at once. *arXiv preprint arXiv:2304.06718*, 2023.

A Additional Results

To better understand the performance of our framework visually, we prepare test time rollouts of our framework as well as those of various baselines in the [supplementary video](#). The video is published anonymously and can be accessed in <https://sites.google.com/view/3d-intphys>

A.1 Ablation Study

We find that training the model with Chamfer distance in dense scenes with granular materials will often lead to predictions with unevenly distributed points where some points stick too close to each other. To alleviate the issue, we introduce the spacing loss to penalize the distance between these points. We set the threshold of penalty d_{min} to be 0.08 and the loss weight σ to be 10. We find that spacing loss can help improve the performance of the dynamics learner especially under extrapolate settings, as shown in Figure 8. We provide qualitative results in the [supplementary video](#).

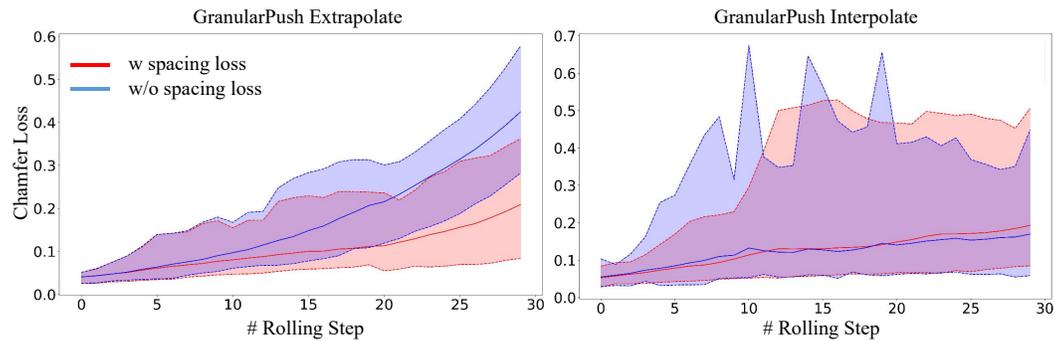


Figure 8: **Ablation Study on the Spacing Loss.** Training dynamics models in the GranularPush scenario with spacing loss results in better rolling prediction. The performance gap is even more substantial in the extrapolate setting.

B Implementation Details

B.1 Dataset Generation

Our datasets are generated by the NVIDIA Flex simulator. Each of the three scenarios (Pour, Shake and Push) has 500 videos of trajectories taken from 6 views, with each trajectory consisting of 300 frames. We manually select the 6 views with reasonable coverage of the tabletop space to minimize the occlusion. The 500 trials are generated from five different sets of environmental parameters, detailed in Table 3. We take one set of parameters that are outside the training distribution as the **extrapolate** dataset for evaluating model generalization. For the rest of the four settings, we randomly split them into train and test sets with a ratio of 0.8.

Next, we provide more details for each scenario:

- In the FluidPour environment, we randomly initialize the position of the upper container and then generate random back-and-forth actions by tilting the container. The action space is then the position and tilting angle of the upper container.
- In FluidCubeShake, we also randomly initialize the position of the container and the cubes inside the container. We then generate random but smooth action sequences moving the container in the 2D plane. The action space is then the x-y location of the container.
- In GranularPush, we randomly initialize the position of the granular pile. Then, for each push, we randomly generate the starting and ending positions of the pusher and move the pusher along the straight line with an angle perpendicular to the pushing direction. The action space is a four-number tuple stating the starting and ending position on the 2D plane.

The following table shows the moving range of the robot arms in the FluidPour and FluidCubeShake environments after normalizing the robot into a size that is the same as in the real world (unit:

| | X-Range | Y-Range | Z-Range |
|----------------|------------------|----------------|-----------------|
| FluidPour | [-29.11, -12.66] | [42.00, 60.00] | [-7.78, 7.78] |
| FluidCubeShake | [-3.25, 42.25] | [19.25, 19.25] | [-24.50, 24.00] |

Table 2: **Robot Action Space(centimeters):** we show the range the robot arms can move in the FluidPour and FluidCubeShake environments.

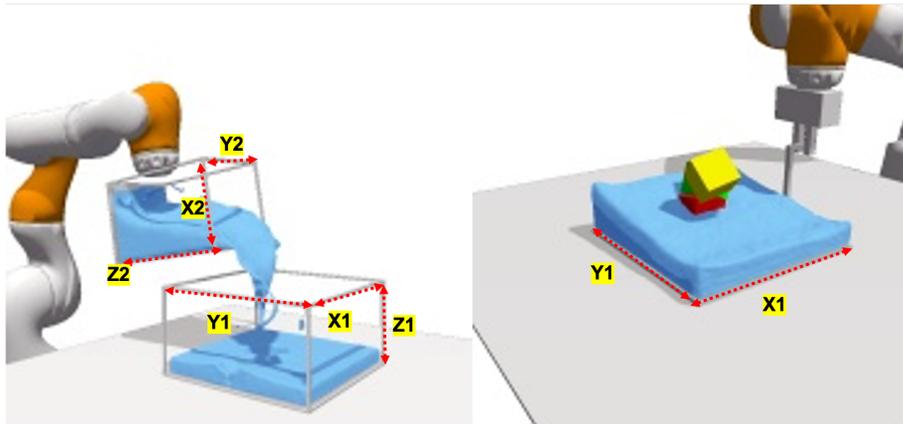


Figure 9: **Illustration of the Environment Settings.** In the FluidPour scenario, a robot arm holds a container and tries to pour some fluid into another container. In the FluidShake scenario, a robot moves a container with some fluid and cubes. We show the parameters for the container shape referred in Table 3.

| SceneName | Params | Env1 | Env2 | Env3 | Env4 | Extrapolate |
|----------------|----------------|------|------|------|------|-------------|
| FluidPour | X2 | 0.53 | 0.53 | 0.81 | 0.81 | 0.81 |
| | Y2 | 0.53 | 0.81 | 0.53 | 0.81 | 0.81 |
| | Z2 | 1.24 | 1.24 | 1.24 | 1.24 | 1.24 |
| | X1 | 1.35 | 1.35 | 1.35 | 1.35 | 1.35 |
| | Y1 | 1.35 | 1.35 | 1.35 | 1.35 | 1.35 |
| | Z1 | 0.74 | 0.74 | 0.74 | 0.74 | 0.74 |
| | AmountofWater | 5125 | 5125 | 6125 | 5375 | 7625 |
| FluidCubeShake | X1 | 0.88 | 0.88 | 1.32 | 1.32 | 1.32 |
| | Y1 | 0.88 | 1.32 | 0.88 | 1.32 | 1.32 |
| | CubeNumber | 1 | 1 | 2 | 2 | 3 |
| | Water | 2173 | 3322 | 3322 | 4858 | 4983 |
| GranularPush | GranularNumber | 2197 | 4032 | 5832 | 9261 | 12167 |

Table 3: **Scene Parameters for Generating the Interpolate and Extrapolate Datasets.** We generate the datasets by varying the shape of container, amount of water, number of cubes, and quantity of the granular material. Z_i, X_i, Y_i are the height, width, and depth for a container i . Please refer to Figure 9 for more details.

centimeters). For GranularPush, the pusher is moving over the entire table; we ignore the specific number in this environment as we do not have robot arms as a reference.

Additional dataset samples. We show samples from the FluidPour, FluidCubeShake and GranularPush dataset in Figure 10, 11 and 12, respectively. Note that all trajectories for the extrapolate settings are used only for testing and will not show up during the training process. We include more samples from the dataset in the video format in the [supplementary video](#).

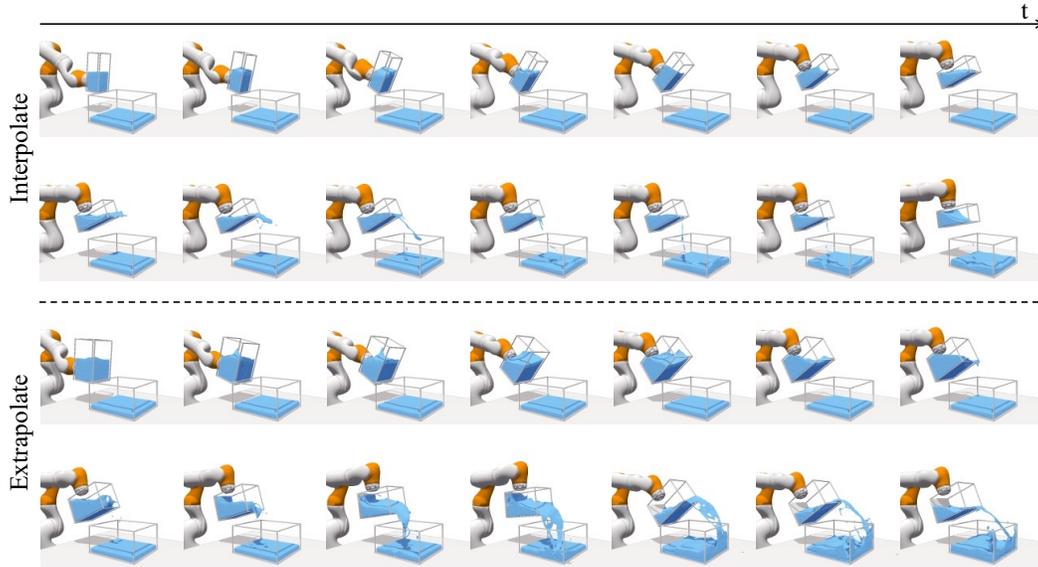


Figure 10: **Samples from FluidPour Dataset.** We show sequences of frames over time with an interval of 20 frames. The sequences above the dashed line are for **interpolate** data, and the bottom images illustrate the **extrapolate** data.

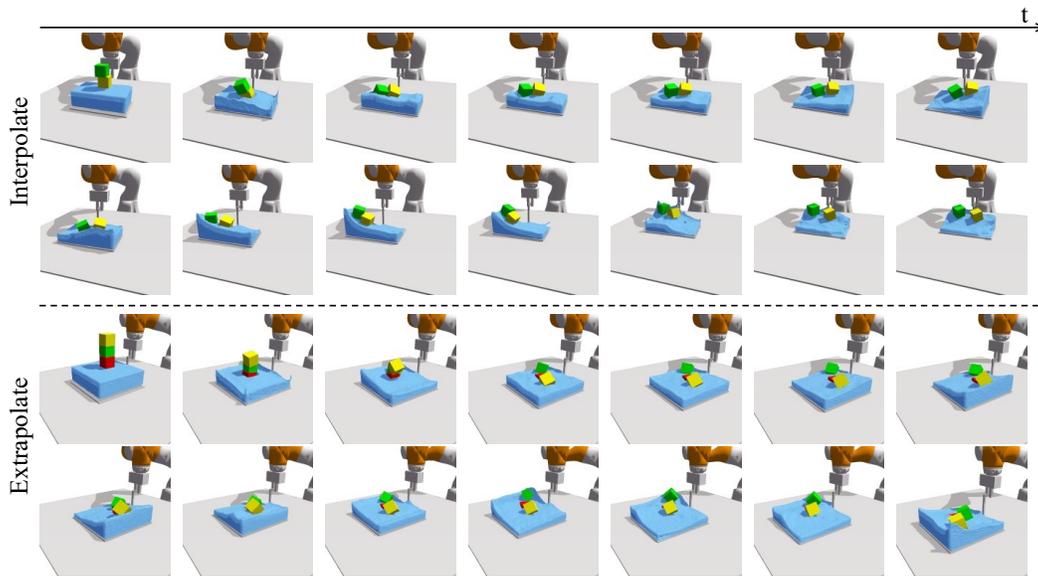


Figure 11: **Samples from FluidCubeShake Dataset.** We show sequences of frames over time with an interval of 20 frames. The sequences above the dashed line are for **interpolate** data, and the bottom images illustrate the **extrapolate** data.

B.2 Model Architecture

Image-conditional NeRF. We follow the architectural design by [68]. For the feature encoder, we employ a ResNet-34 backbone to extract features. We use the output layers prior to the first four pooling layers, upsampling them using bilinear interpolation to the same size, and then concatenating these four feature maps. We initialize the weight of the feature extractor of the scene using ImageNet pre-trained weight. For the NeRF function f , We use fully-connected ResNet architecture with 5 ResNet blocks with a width of 512.

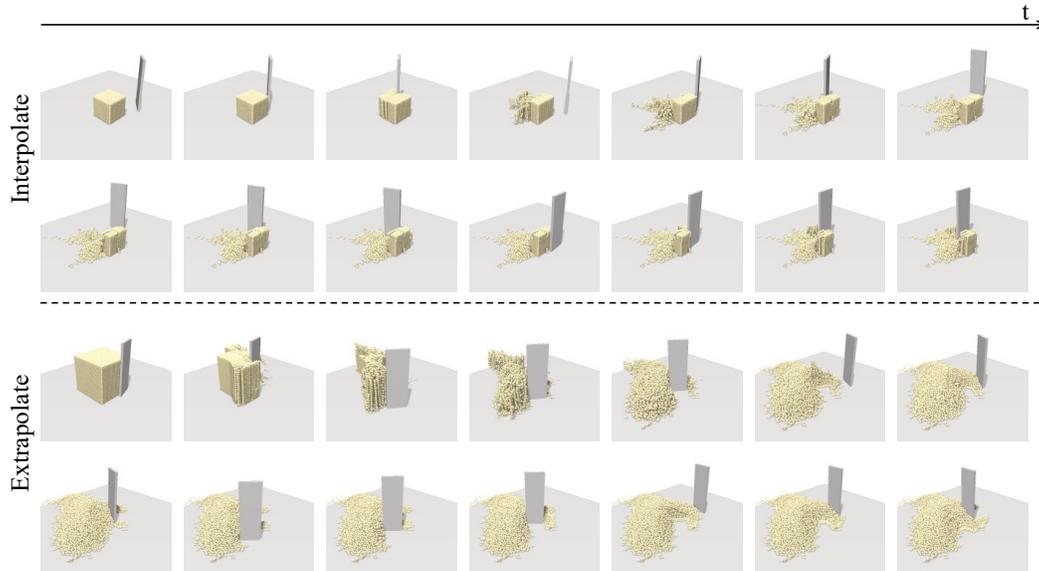


Figure 12: **Samples from GranularPush Dataset.** We show sequences of frames over time with an interval of 20 frames. The sequences above the dashed line are for **interpolate** data, and the bottom images illustrate the **extrapolate** data.

Dynamics predictor. For the edge and vertex encoders, Q_e and Q_v , we use 3-layer fully-connected networks activated by the ReLU function with 150 hidden units. For the propagators, P_e and P_v , we use a 1-layer fully-connected network followed by ReLU activation. The output dimension of the linear layer is 150.

Sampling 3D points from the trained visual perception module. We sample points on a $40 \times 40 \times 40$ grid from an area of $55\text{cm} \times 55\text{cm} \times 55\text{cm}$ and $63\text{cm} \times 63\text{cm} \times 63\text{cm}$ at the center of the table for FluidPour and FluidCubeShake respectively, and on a $70 \times 70 \times 70$ grid from an area of $6\text{cm} \times 6\text{cm} \times 6\text{cm}$ for GranularPush. We evaluate and include points with a density (measured by the occupancy in the predicted neural radiance fields) larger than 0.99. To reduce the total number of points, we subsample the inferred points with FPS with a ratio of 5% for FluidPour and 10% for FluidCubeShake and GranularPush.

Graph building. We set the neighbour distance threshold δ to be 0.2, 0.15, 0.15 for FluidPour, FluidCubeShake and GranularPush respectively. We select the threshold so that each point will have on average 20 30 neighbors. Since, in FluidPour, we sample the points with lower density $2000\text{points}/\text{m}^2$, we use a larger threshold for this scenario. For FluidShape and GranularPush, since the density is around $3000\text{points}/\text{m}^2$, we cut down the number by 25%.

We found that if the threshold is too small, the performance will degrade significantly since each particle will only receive messages from a few neighbors (and miss out on the larger context). On the other hand, setting the threshold too large will cause the training time to increase since the graph will have more edges. We found that setting the threshold around the right scale generally leads to more effective training of a reasonable dynamics network.

B.3 Training Details

The models are implemented in PyTorch. We train the perception module using Adam optimizer with a learning rate of $1e-4$, and we reduce the learning rate by 80% when the performance on the validation set has stopped improving for 3 epochs. To compute the rendering loss when training the perception module, we sample 64 points through each ray in the scene and set the ray-batch size of the NeRF query function f to be 1024×32 . Training the perception module on a single scenario takes around 5 hours on one RTX-3090.

We train the dynamics simulator using Adam optimizer with a learning rate of $1e-4$, and we reduce the learning rate by 80% when the performance on the validation set has stopped improving for 3

epochs. The batch size is set to 4. We train the model for 20, 30, and 40 epochs for FluidPour, FluidCubeShake, and GranularPush, respectively. It takes around 10 ~ 15 hours to train the dynamics model in one environment on one single RTX-3090.

B.4 Graph-Based Dynamics Model without Particle-level Correspondence

The velocity of an object provides critical information on how the object will move in the future, yet, we do not have access to such information when tracking the object is impossible. As described in Section 3.2, the attributes a_i^v of a vertex v_i in the built graph consists of (1) velocity of this point in the past frames and (2) attributes of the point (rigid, fluid, granular). To get the velocity of a vertex v , we should have the history position of this vertex. However, since the point clouds are inferred from each frame independently, we do not know how each point moves over time since we do not have point correspondence between frames.

To address the problem, we leverage the fact that some objects in the scene are easier to track, and we try to use the motion of these trackable objects to infer motion for the untrackable units. We assume that we know the dense-labeled states of some known fully-actuated shapes like desks and cups connected to the robot arms. Here we will list one specific scenario where a cup of water is poured into another cup. In this case, we have two different types of points: points for fluid and points for cups, we name the states of them in time step t as $V_P^t = \{v_{P,i}^t\}$ and $V_S^t = \{v_{S,i}^t\}$ respectively. For the particle encoder Q_v , if the particle belongs to the cups, then the input of particle encoder contains n_s history states before t_0 : $\{V_S^{(t_0-n_s):t_0}\}$. If the particle belongs to the water, then we have no history states, so the input of Q_v is all-zero.

By adding the relative position between receiver and sender points, we can pass the momentum of V_P to V_S . Compared with human intuition, we can get an intuitive prediction of the movement of water by simply knowing the past movement of the cup without knowing the past movement of water.

Following [47], we use the velocity of points and their relative position as inputs to the dynamics module instead of using the absolute positions of the points. This ensures the model is translation-invariant so the learned dynamics model can be shared across different spatial locations.

B.5 Inference Speed of Our Model

The prediction speed of the dynamics module depends on the number of input particles, and it takes around 0.1s for graphs with around 300 nodes in FluidShake and FluidPour, and around 0.2s for scenes with 700+ nodes in GranularPush.

For our visual module, the main time consumption comes from NeRF sampling, it takes 0.2s to sample from a grid space introduced in the experiment section of our paper, this was run in blocks, with block-size=1000, made up 4G of a V100 GPU. And it can be even faster with larger blocks. The sub-sampling process (FPS, segmentation) is fast since they are all written in parallel versions, which takes less than 5ms.

C Potential Society Impact

Our work shows the possibility of learning dynamics models from raw sensory inputs, opening up opportunities to automate the design of differentiable physics engines through data-driven learning algorithms. The resulting system can potentially benefit many downstream tasks, including general scene understanding, robotics manipulation, the construction of 3D generative models, and inverse tasks like planning/control and inverse design. Furthermore, predictions from our model are highly interpretable, which makes it straightforward to explain model behaviors and re-purpose the outputs for other downstream applications.

Though data-driven approaches are potentially more scalable with enough data, concerns still exist that it might be hard to ensure the robustness of the model under sensor noise and adversarial attacks. It also becomes less clear how to fully mitigate data biases. Therefore, bringing in advanced techniques from ML robustness will be one critical future avenue to pursue.



Figure 13: **SAM Working on FluidCube Shake:** Recent large segmentation models can well generate masks for different objects in the scene.

D Important Clarifications and Discussions

Q: *What is the input of 3D-IntPhys?*

Video inputs and object instance masks based on color assumption.

Q: *What is the difference of our work compared with [34, 35, 16]?*

1. Compared with nerf-dy [34]: our method use explicit 3D representation instead of implicit representation, which we show in our paper that our method can generalize better than [34].
2. Compared with VPGL [35]: (1)[35] requires GT 3D particle set as supervision to train the visual frontend, while our 3D-IntPhys does not need particle-level supervision (2) In [35], the 3D representation is a particle set, which is an ordered list. As a result, it can only generate a fixed number of points. In contrast, our method produces dense representations that are not limited to ordered sets, making it more flexible and adaptable to systems of varying sizes. (3) [35] assumes that we have a learned dynamics model from the simulation as a dynamics prior, while we learn the dynamics model from the data.
3. Compared with Comp nerfdyn [16]: [16]only works on rigid objects and a rope that exhibits only slight deformation, most of the objects do not have topological changes, and they are all constrained to move in a 2D plane. So while the object-centric dynamics model used in [16] can solve the tasks in their paper, the object-centric representation is not suitable for learning the complex dynamics of fluids or granular materials, as in our paper. Our settings contain much more diverse 3D dynamics of challenging materials, which can not be solved using [16].

Q: *Is the color segmentation of the fluid objects a reasonable assumption?*

It should be noted that the color-based segmentation will not degrade the challenging problem of learning 3D Intuitive Physics, since the task focuses more on learning complex visual dynamics from images.

We want to emphasize that the work focuses more on learning complex visual dynamics from images, as opposed to solving object segmentation in general. Learning fluids dynamics from videos is a challenging task, and there are only a few existing works. NeRF-dy is the closest to us, yet the model's generalization ability is limited. We have shown in the proposed work that we can significantly improve the generalization ability by operating with a hybrid of implicit and explicit, as opposed to pure implicit, 3D representations. We agree object segmentation is a critical visual understanding

problem, and solving it is an important next step to getting a more general visual dynamics learning framework.

With recent advancements such as SAM [29] and SEER [71], which focus on segmentation in real-world scenarios, the possibility of video segmentation without the need for annotations has emerged (as is shown in Figure 13). This development paves the way for leveraging existing large-scale models to enhance the segmentation pipeline, offering great promise for future applications.

Q: *Since the fluid has zero velocity, how to predict the intuitive dynamics?*

The intuition is that we can infer the water movement from the container's movement. We also assume that the initial velocity of water is **nearly zero**, which is also used in [50], so the momentum can be gradually passed from the container to the water.

We propose this assumption so that the intuitive physics model can be learned from (1) particles sampled from the neural radiance field, which is not stable (2) point clouds without one-to-one correspondence. The results show that we can learn reasonable dynamics (water poured out from a cup, water falling in the container, cubes moving in water, and granular materials pushed away by a pusher). It also shows the potential of distribution-based loss in learning visual dynamics.