

DensePhysNet: Learning Dense Physical Object Representations via Multi-step Dynamic Interactions

Author Names Omitted for Anonymous Review. Paper-ID [149]

Abstract—We study the problem of learning physical object representations for robot manipulation. Understanding object physics is critical for successful object manipulation, but also challenging because physical object properties can rarely be inferred from the object’s static appearance. In this paper, we propose DensePhysNet, a system that actively executes a sequence of dynamic interactions (e.g., sliding and colliding), and uses a deep predictive model over its visual observations to learn dense pixel-wise representations that reflect the physical properties of observed objects. Our experiments in both simulation and real settings demonstrate that the learned representations carry rich physical information, and can directly be used to decode physical object properties such as friction and mass. With knowledge of object physics, these representations also lead to more accurate and efficient manipulation in downstream tasks than the state-of-the-art. The system also generalizes well to novel scenes with more objects than in training.

I. INTRODUCTION

Intelligent manipulation benefits from the ability to distinguish between object materials and infer their physical properties from sight. For example, when rearranging objects on a table, recognizing the physical differences between heavy and lightweight materials enables a robot to better plan its manipulation strategies. Is it possible for robots to self-learn these differences without any explicit supervision?

Although considerable research has been devoted to learning object-centric representations that reflect visual features, they rarely account for latent physical attributes such as mass or friction. Unsupervised learning of physical properties is a less explored problem due to three major challenges:

- Most physical attributes cannot be directly inferred from appearance cues alone in a static environment. For example, while aluminum shares a similar appearance with steel, it is much lighter.
- Most physical attributes are not salient under static or quasi-static interactions: gently pushing a wooden or a metal block results in only subtle differences in their visible motion, despite their different materials and densities.
- Each physical property may only be revealed under specific types of interactions. For example, the sliding distance of an object is determined by both its friction coefficient and mass when given the same initial momentum; while it is only determined by the object’s friction coefficient when given the same initial velocity. Therefore, without an explicit physics model, the system needs not only to explore different types of interactions, but also to infer and decouple physical properties from multiple action outcomes jointly.

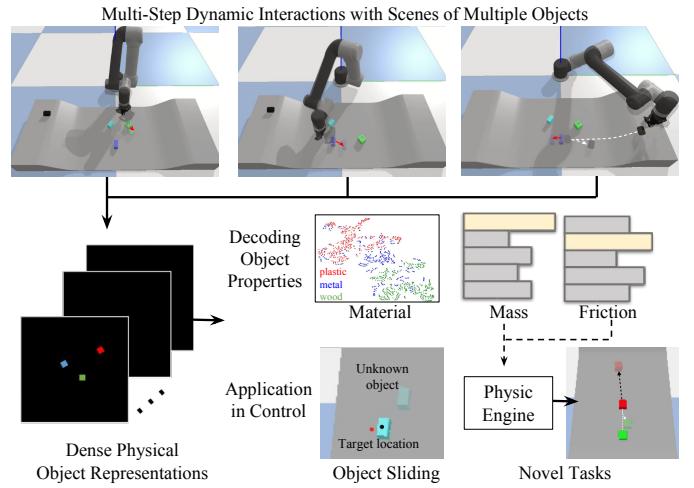


Fig. 1. Our goal is to build a robotic system that learns a dense physical object representation from a few dynamic interactions with objects. The learned representation can then be used to decode object properties such as its material and mass, and be applied in manipulation tasks such as sliding objects with unknown physics or combined with a physic engine to tackle novel tasks.

In this work, we propose to discover and learn the physical properties of objects through visual observations of **multi-step, self-supervised, dynamic interactions**. Our system, DensePhysNet, actively executes a sequence of dynamic interactions (e.g., sliding and colliding), and uses its visual observations to learn physical object properties without any explicit supervision (Figure 1).

DensePhysNet takes in the current state and action as input, and predicts the change in objects’ state after the interaction, which is represented as pixel-wise optical flow. By learning to accurately predict the future states of objects conditioned on different interactions, DensePhysNet acquires an implicit understanding of the objects’ physical properties and how they influence observed motions. We model DensePhysNet with a recurrent structure that aggregates information from multiple interactions, so that it can better infer and encode object’s physical properties over time. We also design DensePhysNet in a modularized way, so that the learned physical representations can be disentangled from the representations which encode information about objects’ visual appearance and actions. This unique design enables it to generalize to new tasks that involve different types of interactions. Because DensePhysNet produces a pixel-wise dense representation, instead of a single feature vector for the entire scene as in many prior methods [18, 15, 1], it also able to generalize to more complex scenes that contain multiple objects.

The main contribution of our paper is to demonstrate that

deep predictive models over visual observations of multiple dynamic interactions can enable an agent to learn the latent physical attributes of manipulated objects. We execute a series of experiments in both simulation and real settings to evaluate our approach qualitatively and quantitatively. The results show that DensePhysNet is more effective at learning physical object properties than other representation learning methods, and that the learned representations can be leveraged to improve the performance of downstream control tasks such as planar sliding. We also show that when combined with a physics engine, the system is capable of leveraging the object physical properties decoded from DensePhysNet to execute novel control tasks unseen in training. Code and data will be released.

II. RELATED WORK

Modeling object physics is a long-standing problem in both robotics and artificial intelligence. Dating back to the 1980s, Atkeson et al. [2] estimates the mass and moments of inertia of a grasped object, using measurements from a wrist force-torque sensor. Yu et al. [27] tackles the same problem by pushing objects using two fingertips equipped with force-torque sensors, and recording the fingertip forces, velocities, and accelerations. Recently, researchers have explored learning to estimate physical properties from their appearance and motion, either in combination with physical simulators [23, 24] or via end-to-end deep learning [10, 13, 26]. These models build upon an explicit physical model, i.e., a model parameterized by physical properties such as mass and force. This enables generalization to new scenarios, but also limits their practical usage: annotations on physical parameters in real-world applications are expensive and challenging to obtain.

An alternative line of work is to learn object representations without explicit modeling of physical properties, but in a ‘self-supervised’ way through robot interactions. Pinto et al. [19] proposed to build a ‘curious’ robot that interacts with objects to learn representations for visual recognition. Denil et al. [8] uses reinforcement learning to build object models via physical experiments. Byravan and Fox proposes to use deep networks to approximate rigid object motion [4]. A few follow-ups have extended these models for planning and control, including poking [1] and pushing via transfer learning [18] or visual predictive learning [11]. Such progress is impressive, though the focus of these papers is still on *visual* representation learning. Without considering the underlying physical processes, they fall short to generalizing to novel objects and tasks.

The work closest to ours is Push-Net, proposed by Li et al. [15]. They use a multi-step model to learn physical object properties for planar pushing. While Push-Net focuses on scenes with a single object and encodes the entire image into a single latent representation, we instead learn dense (pixel-wise) representations from interactions – which enables our model to generalize not only to novel objects, but also to novel scenes with multiple objects. Furthermore, while Push-Net only considers planar quasi-static pushing, we consider two types of dynamic interactions: planar sliding and collision. Our

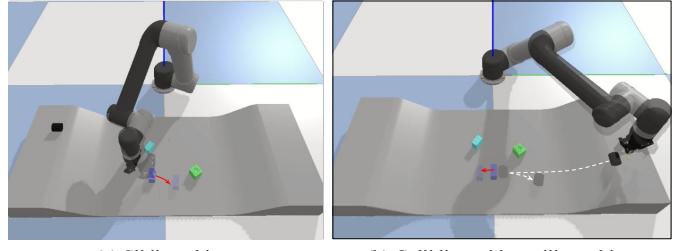


Fig. 2. **Dynamic Interaction.** We design two types of dynamic interactions (sliding and collision) to reveal the object physical property. To slide an object, the robot approaches it and executes a push with a high speed such that the object can slide after the push. For collision, the robot will release the same auxiliary object at the same height to collide it with other objects.

experiments demonstrate that using multiple types of dynamic interactions is key to fully revealing latent physical properties.

We are, of course, not the first to learn dense representations on visual data. Most prior work on this topic revolve around learning correspondences across views in 2D [6, 20] and 3D [28, 22, 21, 3]. Florence et al. [12] proposes dense object nets, learning dense descriptors by multi-view reconstruction and applying the descriptors to manipulation tasks. While their paper primarily focuses on learning object representations which reflect visual appearance, we learn object representations which reflect physical properties and show that they are useful for manipulation tasks.

III. METHOD

The goal of DensePhysNet is to learn latent representations that encode object-centric physical properties (e.g., mass, friction) through self-supervision. To this end, we train a deep predictive model of depth images on a large dataset of observed dynamic robotic interactions. The setting consists of a collection of objects on an inclined ramp laid in front of the robot (Figure 2). When interacting with objects, the robot captures a depth image I_t of the state at time t , executes an action a_t , then captures another depth image I_{t+1} at the next time step. We model DensePhysNet as a neural network that takes as input the visual observation I_t and the executed action a_t , and outputs a prediction of the next state observation I_{t+1} in the form of optical flow $O_{t,t+1}$. The idea is that in order for DensePhysNet to accurately predict the future states of objects conditioned on different interactions, it needs to acquire an implicit understanding of the objects’ physical properties and how they influence observed motions.

A. Dynamic Interactions

The robotic agent executes two types of dynamic interactions to accentuate the physical properties of objects: sliding and collision. The vast majority of prior work in representation learning [15, 11, 18, 1] use static or quasi-static manipulations like pushing, grasping, or poking, where it is challenging to observe latent physical object properties from object motion. This is because in (quasi-)static manipulation, changes in object movements are largely influenced by the actions and dynamics of the manipulator, less so from the object itself. For example, during quasi-static pushing, the object is assumed to move only

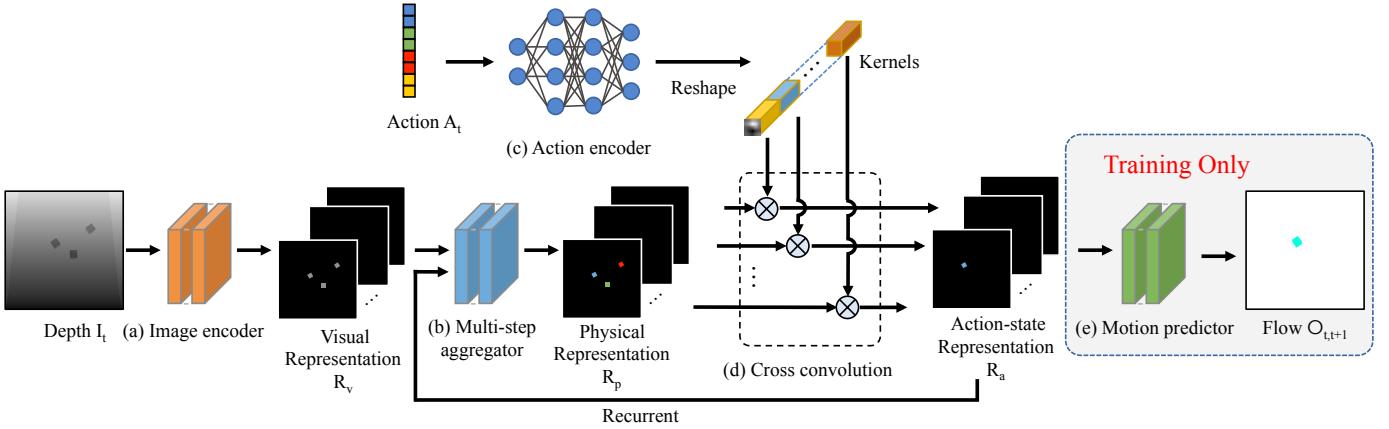


Fig. 3. **DensePhysNet.** DensePhysNet takes in the current state (depth image I_t) and interaction A_t as input and predicts the change of objects’ state after the interaction. The change of objects’ state is represented as pixel-wise optical flow $O_{t,t+1}$. The network consists of five modules: an image encoder, a multi-step information aggregator, an action encoder, a cross convolutional layer, and a motion predictor. These five modules work jointly to learn three object representations: visual representations that encode visual signals of the object, physical representations that encode physical object properties, and action-state representations that encode objects’ states after interaction.

with the end effector (and to stop when the end effector stops). In these scenarios, it is naturally more difficult to observe the subtleties in motion due to the physical properties of the object.

Dynamic manipulations, in contrast, can lead to object motions (beyond contact with the manipulator) that are more likely to reveal its physical properties. For example, when objects are pushed with high speeds, the forces of acceleration can cause objects to slide by themselves for a distance. The differences in their motion after contact (e.g. distance traveled) can serve as a visual cue for differences in surface friction. These cues are relatively more salient than what can be observed in quasi-static manipulations. However, the distance traveled by an object after sliding alone (assuming known initial velocity) can only help derive the object’s friction. Therefore, we need other types of dynamic manipulations (e.g. collision) to reveal and distinguish other physical properties.

Sliding. The sliding action is parameterized by a direction θ and velocity v with respect to the robot. To slide an object, the robot approaches the it from θ and executes a push with a high speed such that the object can slide after the push. To achieve high pushing velocities without exceeding the physical force-torque safety limits of the real-world robot, we constrain the motion planning so that joints closer to the base of the arm (e.g. elbow, shoulder) move slower than joints closer to the end effector (e.g. wrist). The direction θ and velocity v are quantized as discrete variables and we use one-hot vectors to encode them. There are sixteen possible directions and four possible speeds. The sixteen directions are uniformly distributed on $[0, 2\pi]$. Four robot’s speeds are $\{0.96, 1.28, 1.44, 1.6\}$ rps, representing the robot’s joint rotation speed.

Collision. For collisions, we set two inclined ramps, one on each side of the workspace. The robot grasps an auxiliary cylinder and places it on the top of one of the ramps. The cylinder then rolls down the ramp and collides with an object in the middle. The only parameters of the collision action are the cylinder’s starting position $[d \times X_s, y, H_s]$, where $d \in \{-1, 1\}$ indicates whether the cylinder is rolling from the left

or right ramp, X_s is the distance from the ramp to the center of workspace, y is the same as the y -coordinate of the target object, and H_s is the height of the ramp. The auxiliary cylinder is fixed, with radius 2cm, height 4cm, and weight 0.4kg.

Interaction policy. We use a balanced-random interaction policy to ensure interaction diversity and to encourage a full exploration of each action. Specifically, for each step, we choose one type of action to execute via $P(x_{i,j}) = \frac{(1/2)^{a_{i,j}}}{\sum_{n=1}^N \sum_{m=1}^M (1/2)^{a_{n,m}}}$, where $P(x_{i,j})$ is the probability to apply action i to object j , N is the number of objects in the scene, M is the size of action space, and $a_{i,j}$ is the number of action i applied to object j up to now. We randomly sample the parameters for the chosen action, while enforcing the constraint that the object needs to stay inside the workspace.

B. DensePhysNet

We design DensePhysNet based on three key insights. First, it needs to be modularized in a way that the learned physical representations are disentangled from the representations that encode information about object visual appearance and actions. It is critical when applying the learned physical representations to new tasks that involve different environment setting and action types. Second, DensePhysNet has a recurrent structure, so that it can aggregate information from multiple interactions to better infer physical properties. Third, to handle complex scenes with multiple objects, our model produces a dense pixel-wise representation, instead of encoding the entire scene into one single latent representation as in previous papers [15].

DensePhysNet (Figure 3) consists of five modules: an image encoder, a multi-step information aggregator, an action encoder, a cross convolutional layer, and a motion predictor. These five modules work jointly to learn three object representations: visual representations R_v that encode visual signals of the object, physical representations R_p that encode physical object properties, and action-state representations R_a that encode objects’ states after interaction.

Each iteration of learning works as follows. First, given a

depth image I_t , the image encoder (Figure 3a) extracts visual signals from the image and outputs the visual representation R_v . Then, the information aggregator (Figure 3b) learns to integrate the visual representation R_v with the object representation after the last interaction R_a to extract the physical representation R_p . Intuitively, after interactions, objects with different physical properties will end up in different positions and poses, and such signals are now available from the visual input and therefore should lie within the visual representation R_v . The goal of the information aggregator is thus to distill physical knowledge R_p by analyzing R_v and R_a jointly.

In parallel, the action encoder (Figure 3c) encodes action a_t into convolution kernels. The cross convolutional layer (Figure 3d) then applies the encoded action kernels on the physical representation R_p to produce the effect of the action on the objects. Here, the cross convolutional layer can be seen as a learned, latent physical simulator, learning physics to approximate the effect of actions. As shown in Xue et al. [25], this cross convolution layer better combines the action and the physical representations compared to tensor concatenation. It outputs the action-state representation R_a .

Finally, the action-state representation R_a is fed into the motion predictor (Figure 3e) to predict the optical flow $O_{t,t+1}$ across two images I_t and I_{t+1} . This is the only supervisory signal we have during training; during testing, the motion predictor is no longer needed.

Network architecture. Here we provides more details on each of the network module. The image encoder takes as input the depth image with a size of 160×160 , and outputs a 32-channel feature map. It applies one 11×11 , one 5×5 , and two 3×3 convolution layers. Between convolutional layers, there are batch normalization and ReLU layers.

The action encoder uses seven fully connected layers. The numbers of hidden units are 64, 128, 256, 256, 512, and 800. It takes in the action vector (a 37-dim one-hot vector) as input and outputs 32 action kernels, each with a size of 5×5 .

The multi-step aggregator takes the visual representation R_v and the action-state representation R_a in the last step as input (both have a size of $32 \times 160 \times 160$), and combines them with a gate C via

$$C = S(g(R_v))f(R_v) + [1 - S(g(R_v))]R_a, \quad (1)$$

where $S(\cdot)$ is the sigmoid function and $f(\cdot)$ and $g(\cdot)$ are a 1×1 convolution layer. The gate acts as a filter to block or pass on information based on the calculated weight. Then fifty residual blocks are applied to the combination C to get the physical representation.

The cross convolutional layer applies the convolutional kernels learned by the action encoder to the physical representation R_p learned by the multi-step aggregator. Here, the convolution operations are carried out in a channel-wise manner. In between, there are batch normalization and ReLU layers.

The motion predictor takes the action-state representation as input and outputs pixel-wise optical flow map. This module consists of four convolution layers with a number of channels

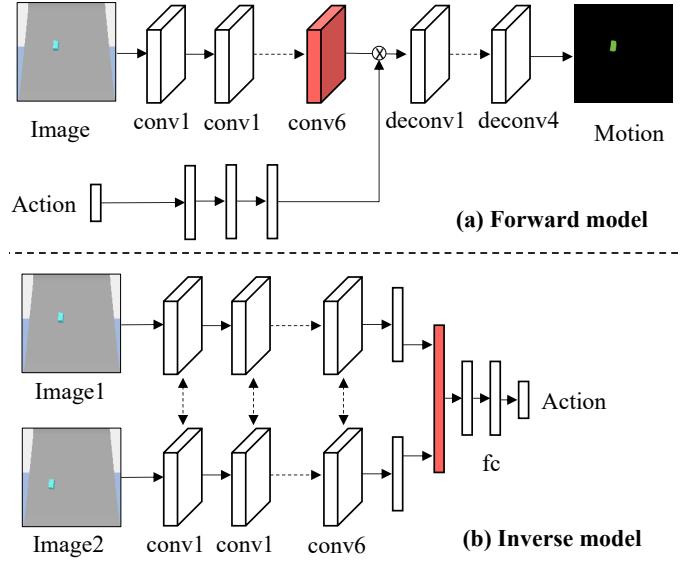


Fig. 4. **Baselines models.** We compare our approach with the the model proposed by Agrawal et al. [1]. Their model consist of a forward and an inverse model, both designed to handle a single-step interaction. The forward model takes the current frame and the action as input, and predict the motion of the object. The inverse model takes the frames before and after the action as input, and predict the action parameters.

[32, 32, 32, 2] and a kernel size of 3×3 .

Self-supervised training. We use Mean Square Error (MSE) of the optical flow between our predictions and ground truth as the training loss. Since the optical flow can be automatically computed based on visual observations (i.e. images taken before and after the interaction), the whole training process of DensePhysNet can be fully self-supervised without any human annotations. We implement our model in PyTorch [17]. Optimization is carried out using ADAM [14] with $\beta_1 = 0.9$ and $\beta_2 = 0.95$. We use an initial learning rate of 10^{-3} and a learning rate decay of 0.9 after each epoch. The model is trained for 40 epochs with a mini-batch size of 32.

IV. EXPERIMENTS

We execute a series of experiments to evaluate our learned physical representation R_p . The goal of the experiments is to understand, through both qualitative analysis and quantitative evaluation, whether the learned physical representations encode information about physical object properties; if so, how accurate the encoding is; and further, how useful the representation is for object manipulation.

A. Decoding Object Material

We first analyze whether the learned physical representation can be used to distinguish objects of different materials.

Data. We use three visually indistinguishable objects, made of one of the three materials:

- Plastic with mass $m \in [0.11, 0.14]\text{kg}$ and friction coefficient $\mu \in [0.4, 0.6]$.
- Wood with mass $m \in [0.11, 0.14]\text{kg}$ and friction coefficient $\mu \in [0.8, 1.0]$.

	(a)	(b)	(c)
Action space	Sliding and collision	Sliding	Sliding and collision
Representation type	Visual representation	Physical representation	Physical representation
t-SNE visualization • Red: plastic • Green: metal • Blue: wood			
SVM classification	0.34	0.69	0.88

Fig. 5. **t-SNE embeddings of learned features.** OurDensePhysNet (c) designed to learn physical object representations, learns to cluster objects based on their material. The representation learned from sliding alone (b) is only informative on distinguishing the difference in friction (wood vs. others), but has trouble telling the difference in mass (plastic vs. metal). The visual representation is not informative on distinguishing object material.

- Metal with mass $m \in [0.17, 0.20]\text{kg}$ and friction coefficient $\mu \in [0.4, 0.6]$.

Each object’s physical properties are uniformly sampled according to its material. Hence, objects of different materials have distinct physical properties, and those made of the same material have similar, but not identical, properties. There are 8,000 sequences for training and 2,000 sequences for testing. We use PyBullet for all our experiments in simulation [7].

For each trial, the robot interacts with the objects 20 times (in test mode without optical flow supervision). We use background subtraction to compute the silhouette of each object, which we then use to extract pixel-wise features for each object from the representation learned after the interactions.

Baselines. We compare our model for learning physical representations R_p to following baselines to understand the role of active interaction.

- Visual representations: we first compare the physical representations R_p with the learned visual representations extracted from our model R_v in Figure 3. The learned visual representations are computed directly from static images, without the knowledge of object motion or actions.
- Our model with sliding: We also compare with our model trained with only sliding, but not collisions, to validate the importance of having multiple types of interactions.

Results. Figure 5 shows 2D t-SNE embeddings [16] of 1,000 learned feature vectors from 100 test sequences, where colors represent different materials. Our model, designed to learn physical object representations, learns to cluster objects based on their material. The representations learned from sliding alone is only informative on distinguishing the difference in friction (wood vs. others), but has trouble differentiating different mass (plastic vs. metal). The visual representation is not informative on distinguishing object material.

We further verify each representation’s discriminative power by training a 3-way linear-SVM on one sequence and then evaluate its material classification accuracy on test data. The representation learned by DensePhysNet achieves a 88%

accuracy on material classification, while the model trained only on sliding achieves an accuracy of 69%. The visual representation achieves an accuracy of 34%, close to random guess (33.3%).

B. Decoding Physical Object Properties

In this experiment, we examine how accurately we can decode physical properties from the learned latent representations.

Setup. To this end, we train a linear classifier to decode physical properties from the latent representations on a small annotated dataset, and test it on a set of novel objects that have different physical properties and shapes. During testing, the robot interacts with the objects to update the latent representations but no optical supervision is used. We conduct this experiment in simulation, where ground truth physical properties are available for training. We only evaluate on scenes of a single object, as both the forward and the inverse model in Agrawal et al. [1] only handle single-object scenarios and predict one feature vector for the whole scene. The shape of each object is randomly sampled from ShapeNet [5], where training and testing objects are of different shapes. Each object’s physical properties (friction and mass) are uniformly randomly sampled from 30 discrete values: $\mu \in \{0.4, 0.42, 0.44, \dots, 0.7\}$ and $m \in \{0.11, 0.113, 0.116, \dots, 0.2\}\text{kg}$. There are 8,000 sequences for training and 2,000 for testing. The robot has 10 interactions with the object using three types of policies: sliding only, collision only, and sliding and collision. We extract the pixel-wise features of the object from the last step using its bounding box (automatically via background subtraction), and flatten it into one feature vector. For the forward and the inverse model, we directly take its hidden layer features (marked red in Figure 4) as the features of the object.

For evaluation, we train both of the friction and mass decoder as a 30-way classifier. We then calculate the average distance for each piece of data as the evaluation metric: $D = \sum_{i=1}^{30} p_i \times |i - y|$, where p_i is the probability of category i predicted by our model and y is the ground truth.

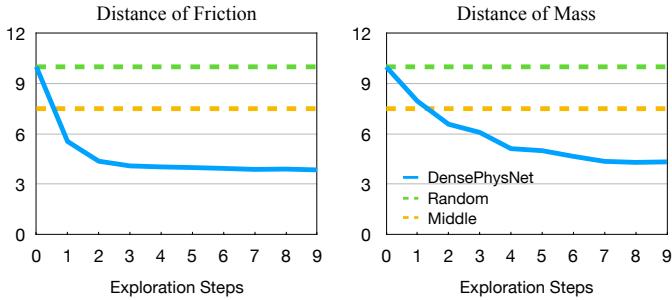


Fig. 6. **Decoding physical object properties.** The plot shows the error in distance on decoding object’s friction coefficient and mass at each step. The plot shows that the prediction accuracy at the beginning is, as expected, the same as random guess; over the course of interactions, the average distance decreases quickly, which means our model gradually accumulates knowledge of object physics.

TABLE I
QUANTITATIVE RESULTS OF DECODING PHYSICAL PROPERTY.

	Sliding		Collision		Sliding & Collision	
	Friction	Mass	Friction	Mass	Friction	Mass
Forward [1]	10.00	10.00	10.00	10.00	10.00	10.00
Inverse [1]	6.41	9.88	8.59	7.96	6.36	6.87
Ours	4.07	9.23	6.58	7.03	3.81	4.24

Baselines. We compare with the following baselines.

- Our model with sliding or collisions: the performance of our model trained with only sliding or collisions, but not both.
- Single-step forward model [1]: We also compare our approach with the model proposed by Agrawal et al. [1]. Their model consists of a forward and an inverse model, both designed to handle a single-step interaction instead of aggregating the information across multiple interactions. The forward model takes the current frame and the action as input, and predict the motion of the object for one step (see Figure 4a).
- Single-step inverse model [1]: The inverse model from Agrawal et al. takes the frames before and after the action as input, and predict the action parameters (see Figure 4b).

Results. Table I shows that DensePhysNet outperforms the baselines, demonstrating the importance of information aggregation in multi-step interactions in modeling physical object properties. The forward model only takes the current image and the action as input. Its representation only contains visual information, which does not help to predict physical properties. Hence its result is the same as random guessing. The inverse model takes images before and after the interaction and learns some information about physical properties from object motion. However, it cannot handle long-range data; its performance is therefore limited.

Moreover, the comparison between different action types demonstrate that sliding helps a lot for the learning of friction, while the mass can only be inferred from the combination of sliding and collisions. Hence the diversity of action space is of vital importance.

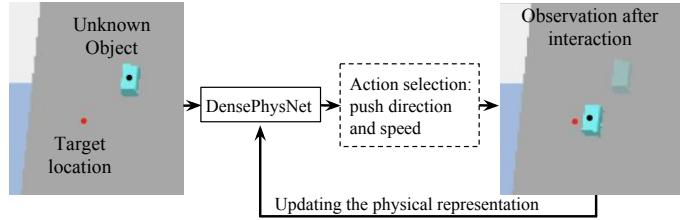


Fig. 7. **Application in object sliding.** The goal for this task is to push an object with selected push direction and speed, so that the object will slide to the target position. At each step, DensePhysNet enumerates all possible actions and predicts the motion of the object, then selects the action whose predicted position is closest to the target. After each interaction the new action-state representation is fed back to the multi-step aggregator to improve the object’s physical representation, and consequentially to improve the action predictions for the following interactions.

We also conduct an ablation study to understand how the number of interaction steps affects the results. Figure 6 shows DensePhysNet’s performance in each step. For calibration, ‘Random’ shows the performance of random guessing and ‘Middle’ shows the performance of predicting the average mass and friction. At the beginning, our model’s prediction is similar to random guessing. Over the course of interactions, the error decreases quickly, which means our model gradually accumulates knowledge of object physics.

C. Application in Sliding Objects with Unknown Physics

The ability to understand physical properties is important for flexible manipulation. In this experiment, we test whether our model can efficiently infer physical properties of unknown objects through multi-step interactions, and further, make use of its knowledge of object physics to suggest more accurate policies in object manipulation.

Setup. Our task is to push an unknown object with a selected push direction and speed, so that the object will slide to the target position. The target position is uniformly randomly chosen in a circle centered on the object and with a radius of 0.25m. The robot needs to propose the parameters of the action of sliding, including its direction and speed. The sixteen possible directions are uniformly distributed on $[0, 2\pi]$ and the speed can be chosen from $\{0.96, 1.28, 1.44, 1.6\}$ rps. After each interaction, the new action-state representation is fed back to the multi-step aggregator to improve the physical representation. Intuitively, if the target object is heavy, the system should be able to infer its physical property through interactions and then push the object with a higher speed.

The physical properties of the objects are randomly chosen from two distribution families:

- common objects: $m \in [0.15, 0.16]$ kg and $\mu \in [0.6, 0.8]$,
- uncommon objects: $m \in [0.11, 0.13]$ kg or $[0.18, 0.2]$ kg and $\mu \in [0.4, 0.5]$ or $[0.9, 1.0]$.

We compare our model in Section IV-B with the forward and inverse models as introduced in Section IV-A. Because our model needs a sequence of interactions to learn physical properties, we also evaluate our model after 0, 3, and 7 interaction steps to understand how interaction helps control. Here, the interaction steps are just for updating the latent

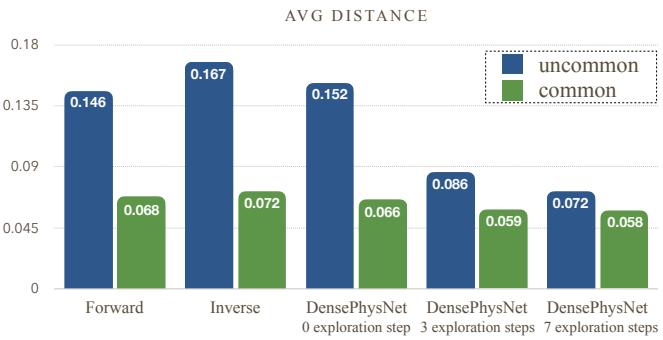


Fig. 8. **Results on object sliding in simulation.** For objects with common physical properties, both our model and the baseline work well in this control problem. However, for objects with uncommon physical properties, only our model (after a few exploration steps) performs well.

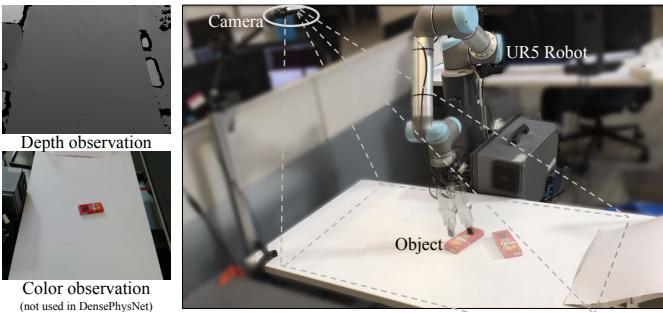


Fig. 9. **Real-world experiment Setup.**

representations; no optical supervision or finetuning is needed. The policy used by different models are as follows:

- Our model and the forward model: enumerate all possible actions and predict the motion of the object, and then choose the action whose predicted position is closest to the target.
- The inverse model: use the action predicted by the model from the current and the target image.

We test each model 100 times and calculate the mean distance between the target position and the object’s position after interactions.

Results. Figure 8 shows the results. For objects with common physical properties, all three models have similar performance with an error distance of 0.06m. However, if the physical properties are uncommon, our model outperforms the other two baselines after a few explorations. This suggests that interactions help to refine our model’s estimate of physical properties and lead to better performance.

Real-world experiments. We also evaluate DensePhysNet in real-world settings, where we use a UR5 robot arm with an RG2 gripper to push a collection of objects on a flat table. Figure 9 shows the setup. We capture RGB-D images using a calibrated Intel RealSense D415, mounted on a fixed tripod overlooking the table of objects from the side. The camera is localized using an automatic calibration procedure from Zeng et al. [29]. The pushing primitives are open-loop, with robot arm motion planning executed using stable, collision-free IK solvers [9]. In real-word experiments, we use the model pre-trained on 8,000 interaction sequences in simulation. In

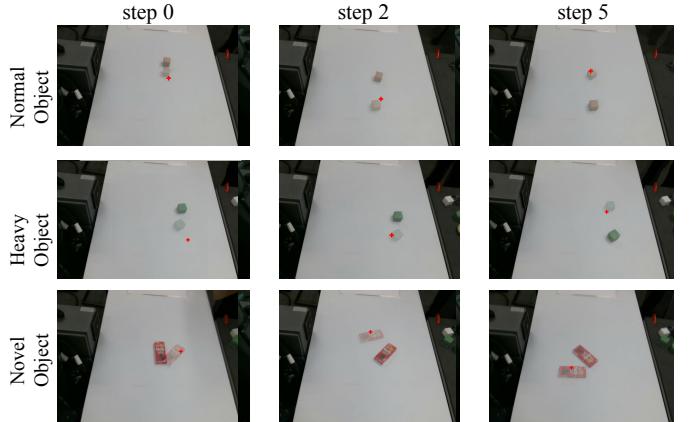


Fig. 10. **Qualitative result of the real-world object sliding experiment.** Each row shows the object sliding result for one object after 0, 2, and 5 exploration steps, where the transparent object indicates its position after applying the suggested action. The red cross indicates the target location. The ‘normal object’ is a wooden block that has similar physical property and shape as the training objects in simulation. The ‘heavy object’ is a plastic box filled with heavy metal balls inside. The ‘novel object’ is a snack box that has not been seen by the algorithm in training. We can see that our DensePhysNet is not only able to generalize to the real-world setting, but also make use of the exploration steps to improve its action prediction (especially for the heavy object in the second row).

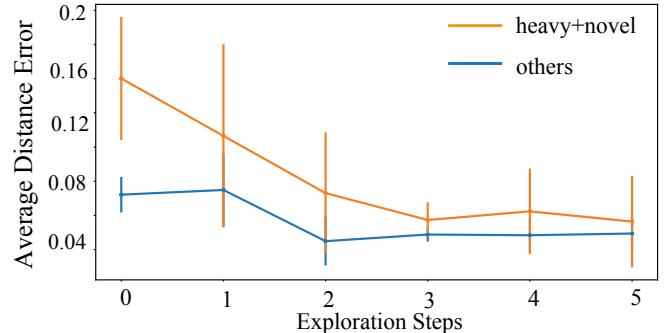


Fig. 11. **Object sliding result in real world.** The algorithm is able to predict more accurate actions and achieves lower error after a few explorations steps.

total, we have tested 18 sequences for six different objects with distinct physical properties, appearance, and shape. Each sequence contains six interaction steps. The target location is manually selected for each step.

Figure 10 shows examples of quantitative results for different objects. DensePhysNet not only generalizes from simulation to real-world settings, but also makes use of the exploration steps to improve its action predictions. In particular, for the objects with larger mass (e.g., the ‘heavy object’ in Figure 10), the algorithm in the first step predicts a push speed that is too small for the objects to slide to the target position, due to the unexpected mass. However, after three steps of interactions, the model learns to adjust its predictions according to the estimated physical properties. Figure 11 shows the mean distance error for different objects and steps. On average, the algorithm predicts more accurate actions and achieves lower errors after a few explorations steps, consistent with the experiments in simulation.

D. Generalization

Practical robotic systems need to generalize to complex real-life scenarios. We evaluate on two cases: generalizing to scenes with more objects and generalizing to novel tasks.

Generalizing to scenes with more objects. Pixel-wise dense representations work for scenes with multiple objects, and further, generalizes to scenes with more objects than those in training scenes. In this experiment, we train the model using two objects and test it with three objects. Each scene consists of objects with different shapes and physical properties. The robot is allowed to have 20 interactions for both training and testing. Other setups of this experiment is the same as in Section IV-B.

There are 5,000 sequences with two objects for training and 1,000 sequences for testing: 500 with two objects and 500 with three objects. We calculate the average distance between the predicted physical properties and ground truth for each object in each piece of data. For our model, the average distances of friction are 4.02 (two objects) and 4.48 (three objects); as to the mass, the distances are 4.57 (two objects) and 5.04 (three objects). The small gap between scenes with two or three objects suggests that our DensePhysNet generalizes to new scenes with more objects; in contrast, baselines [15, 1] do not have an object-wise representation and cannot directly work on scenes with multiple objects.

Generalizing to novel tasks. Once we have decoded the physical properties of objects, we can integrate them into a physics engine for planning and control in alternative tasks. As a demonstration, we study a new task, where the goal is to slide an auxiliary cube so that it hits the object to a target position. Here, we need to select the mass and speed of the auxiliary cube. This task is different from object collisions in the training phase: during training, the auxiliary object is a cylinder with a fixed size, mass, initial position and speed.

The robot first interacts with the target object for 8 steps using the pre-trained model in Section IV-B, without finetuning; it then uses the decoder, also trained in Section IV-B, to decode physical properties from the physical representation. We set the target object’s initial position to always be at center of the workspace, and its target position to the north of the initial position with a distance uniformly sampled from [0.1, 0.3]m. We set the position of the auxiliary cube to be 0.05m south to the center. The mass of auxiliary objects can be chosen from {0.4, 0.43, 0.46 … 0.7}kg and the speed can be chosen from {0.5, 0.53, 0.56, … 0.8}m/s, and its friction is fixed at 0.2.

For each model, we simulate the collision with the decoded physical properties and each possible pair of mass and speed for the auxiliary object; we choose the pair that gives the best prediction. Figure 12 demonstrates the pipeline of this experiment. We test each model 100 times and calculate the mean distance between the target position and the objects position final position after collision.

Figure 13 shows the results. Our model outperforms other baselines significantly. Further, the performance of using only sliding during exploration is not as good as the one using both sliding and collisions. Again, it demonstrates that a diverse set

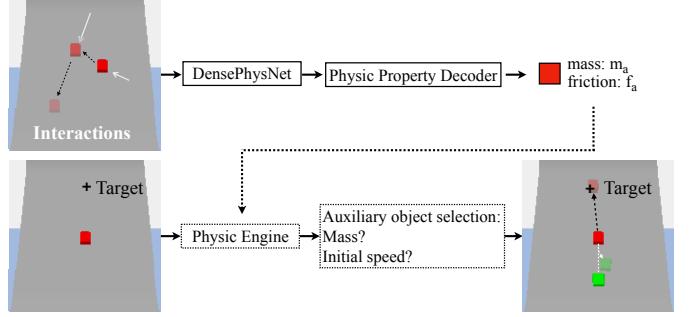


Fig. 12. **Experiments on task generalization.** Here, the robot first interacts with the objects. The learned representation is then used to decode physical properties, which are later used in a physics engine for planning in other tasks.



Fig. 13. **Results on task generalization.** DensePhysNet outperforms baselines significantly. In particular, using both types of interactions is important to achieve good generalization results.

of action types is important to accurate prediction and better performance in control.

V. DISCUSSION AND FUTURE WORK

We have proposed DensePhysNet, a model that learns dense, physical object representations from self-supervised interactions. We have demonstrated that DensePhysNet learns about object materials and physics through both qualitative and quantitative analyses. Further, we have shown that the learned representations can be used in downstream control tasks such as planar sliding to suggest more accurate action policies, both in simulation and on a real robot.

The design of action space is important for the robot to learn object physics. We have shown that, with only both planar sliding and collisions, the model learns to infer both object mass and friction. Infants interact with the world in many possible ways – pushing, poking, sliding, dropping, grasping, and even tasting. Extending DensePhysNet to accommodate a much richer set of interaction types would fully demonstrate its potentials.

DensePhysNet learns physical representations from object motions observed through geometric cues in depth images; our model does not take color images as input. Color signals can however, be very informative of both object geometry and potentially physics, as demonstrated by previous papers on visual representation learning [19, 12]; they can also help to scale our model to a richer set of objects. Therefore, a promising research direction is to build models that learn from both color and geometric cues for better object representations.

REFERENCES

- [1] Pulkit Agrawal, Ashvin Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to poke by poking: Experiential learning of intuitive physics. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 1, 2, 4, 5, 6, 8
- [2] Christopher G Atkeson, Chae H An, and John M Hollerbach. Estimation of inertial parameters of manipulator loads and links. *The International Journal of Robotics Research (IJRR)*, 5(3):101–119, 1986. 2
- [3] Eric Brachmann, Alexander Krull, Frank Michel, Stefan Gumhold, Jamie Shotton, and Carsten Rother. Learning 6d object pose estimation using 3d object coordinates. In *European Conference on Computer Vision (ECCV)*, 2014. 2
- [4] Arunkumar Byravan and Dieter Fox. Se3-nets: Learning rigid body motion using deep neural networks. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017. 2
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *arXiv:1512.03012*, 2015. 5
- [6] Christopher B Choy, JunYoung Gwak, Silvio Savarese, and Manmohan Chandraker. Universal correspondence network. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 2
- [7] Erwin Coumans. Bullet physics engine. *Open Source Software*: <http://bulletphysics.org>, 2010. 5
- [8] Misha Denil, Pulkit Agrawal, Tejas D Kulkarni, Tom Erez, Peter Battaglia, and Nando de Freitas. Learning to perform physics experiments via deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017. 2
- [9] Rosen Diankov. *Automated Construction of Robotic Manipulation Programs*. PhD thesis, Carnegie Mellon University, Robotics Institute, 2010. 7
- [10] Sebastien Ehrhardt, Aron Monszpart, Niloy Mitra, and Andrea Vedaldi. Taking visual motion prediction to new heightfields. *arXiv:1712.09448*, 2017. 2
- [11] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 2786–2793. IEEE, 2017. 2
- [12] Peter R Florence, Lucas Manuelli, and Russ Tedrake. Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. In *Conference on Robot Learning (CoRL)*, 2018. 2, 8
- [13] Katerina Fragkiadaki, Pulkit Agrawal, Sergey Levine, and Jitendra Malik. Learning visual predictive models of physics for playing billiards. In *International Conference on Learning Representations (ICLR)*, 2016. 2
- [14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 4
- [15] Jue Kun Li, David Hsu, and Wee Sun Lee. Push-net : Deep planar pushing for objects with unknown physical properties. In *Robotics: Science and Systems (RSS)*, 2018. 1, 2, 3, 8
- [16] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research (JMLR)*, 9(Nov):2579–2605, 2008. 5
- [17] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *Advances in Neural Information Processing Systems (NeurIPS) Workshops*, 2017. 4
- [18] Lerrel Pinto and Abhinav Gupta. Learning to push by grasping: Using multiple tasks for effective learning. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 1, 2
- [19] Lerrel Pinto, Dhiraj Gandhi, Yuanfeng Han, Yong-Lae Park, and Abhinav Gupta. The curious robot: Learning visual representations via physical interactions. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 8
- [20] Tanner Schmidt, Richard Newcombe, and Dieter Fox. Self-supervised visual descriptor learning for dense correspondence. *IEEE Robotics and Automation Letters (RA-L)*, 2(2):420–427, 2017. 2
- [21] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene coordinate regression forests for camera relocalization in rgbd images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 2
- [22] James Thewlis, Hakan Bilen, and Andrea Vedaldi. Unsupervised learning of object frames by dense equivariant image labelling. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 2
- [23] Jiajun Wu, Ilker Yildirim, Joseph J Lim, William T Freeman, and Joshua B Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. 2
- [24] Jiajun Wu, Erika Lu, Pushmeet Kohli, William T Freeman, and Joshua B Tenenbaum. Learning to see physics via visual de-animation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 2
- [25] Tianfan Xue, Jiajun Wu, Katherine Bouman, and William T Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 4
- [26] Tian Ye, Xiaolong Wang, James Davidson, and Abhinav Gupta. Interpretable intuitive physics model. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [27] Yong Yu, Tetsu Arima, and Showzow Tsujio. Estimation of object inertia parameters on robot pushing operation. In *IEEE International Conference on Robotics and*

Automation (ICRA), 2005. 2

- [28] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, and Jianxiong Xiao. 3dmatch: Learning the matching of local 3d geometry in range scans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [29] Andy Zeng, Shuran Song, Stefan Welker, Johnny Lee, Alberto Rodriguez, and Thomas Funkhouser. Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. In *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2018. 7