# Learning Manipulation under Physics Constraints with Visual Perception

Wenbin Li[1], Aleš Leonardis[2], Jeannette Bohg[3], and Mario Fritz[1]

[1]Max Planck Institute for Informatics, Saarland Informatics Campus, Germany
[2]School of Computer Science, University of Birmingham, UK
[3]Department of Computer Science, Stanford University, USA

## Abstract

Understanding physical phenomena is a key competence that enables humans and animals to act and interact under uncertain perception in previously unseen environments containing novel objects and their configurations. In this work, we consider the problem of autonomous block stacking and explore solutions to learning manipulation under physics constraints with visual perception inherent to the task. Inspired by the intuitive physics in humans, we first present an end-to-end learning-based approach to predict stability directly from appearance, contrasting a more traditional model-based approach with explicit 3D representations and physical simulation. We study the model's behavior together with an accompanied human subject test. It is then integrated into a real-world robotic system to guide the placement of a single wood block into the scene without collapsing existing tower structure. To further automate the process of consecutive blocks stacking, we present an alternative approach where the model learns the physics constraint through the interaction with the environment, bypassing the dedicated physics learning as in the former part of this work. In particular, we are interested in the type of tasks that require the agent to reach a given goal state that may be different for every new trial. Thereby we propose a deep reinforcement learning framework that learns policies for stacking tasks which are parametrized by a target
structure.

## 1 Introduction

Understanding and predicting physical phenomena in daily life is an important component of human intelligence. This ability enables us to effortlessly manipulate objects in previously unseen conditions. It is an open question how this kind of knowledge can be represented and what kind of models could explain human manipulation behavior (Yildirim et al., 2017). In this work, we explore potential frameworks for the robot to learn manipulation with respect to the corresponding physics constraints underneath the task.

Behind the human's physics reasoning in everyday life, the *intuitive physics* (Smith and Casati, 1994; McCloskey, 1983) plays an important role in the process, representing *raw* knowledge for human to understand the physical environment and interactions. Albeit sometimes erroneous, it works well enough for most situations in daily life. It has been an ongoing research in cognitive science and psychology, among others, to understand computational models (Battaglia et al., 2012) and explain such a mechanism.

It has not yet been shown how to equip machines with a similar set of physics commonsense—thereby bypassing a model-based representation and a physical simulation. In fact, it has been argued that such an approach is unlikely due to e.g., the complexity of the problem (Battaglia et al., 2013). Only recently, several works have revived

this idea and reattempted a fully data driven approach to capturing the essence of physical events via machine learning methods (Mottaghi et al., 2016; Wu et al., 2015; Fragkiadaki et al., 2016; Bhattacharyya et al., 2018).

In the first part of this work, we draw inspiration from the studies in developmental psychology (Baillargeon, 1994, 2008) where the infants acquire the knowledge of physical events gradually through the observation of various event instances. In this context, we revisit the classical setup of Tenenbaum and colleagues (Battaglia et al., 2013) and explore to which extent machines can predict physical stability events directly from appearance cues. We approach this problem by synthetically generating a large set of wood block towers under a range of conditions, including varying number of blocks, varying block sizes, planar vs. multi-layered configurations. We run those configurations through a simulator (*only at training time!*) in order to generate labels whether a tower would fall or not. We show for the first time that the aforementioned stability test can be learned and predicted in a purely data driven way—bypassing traditional model-based simulation approaches. Further, we accompany our experimental study with human judgments on the same stimuli. Then we utilize this framework to guide the robot to stably stack a single block onto the existing block-structure based on the stability prediction. To circumvent the domain shift between the synthesized images and the real world scene images, we extract the foreground masks for both synthesized and captured images. Given a real world block structure, the robot uses the model trained on the synthesized data to predict the stability outcome across possible candidate placements, and performs stacking on the feasible locations afterwards. We evaluate both the prediction and manipulation performance on the very task.

In the second part of this work, we further tackle a more challenging stacking task (target stacking). The goal for the task is to reproduce a shape shown in an image by consecutively stacking multiple blocks while retaining physical stability and avoiding pre-mature collisions with the existing structure. We explore an alternative modeling to learn the block stacking through trial-and-error, bypassing the need to explicitly model the corresponding physics knowledge as in the former part of this work. For this purpose, we build a synthetic environment with physics simulation, where the agent can move and stack blocks and observe the different outcomes of its actions. We apply deep reinforcement learning to directly acquire the block stacking skill in an end-to-end fashion. By introducing the goal-parameterized policies, we learn a single model
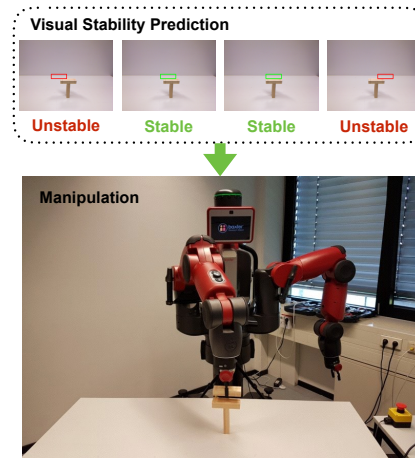


Figure 1: Given a wood block structure, our visual stability classifier predicts the stability for future placements, and the robot then stacks a block among the predicted stable placements.

to guide the agent to build different shapes on request. We first validated this model on a toy example where the agent has to navigate in a grid-world where both the location of the start and end point are randomized for each episode and then the target stacking. On both experiments, we observe generalization across different goals.

We tackle the block stacking task, by focusing the manipulation learning with respect to the corresponding physics constraints underneath the task. The first part of the work emphasizes the visual stability prediction formalism and its application to single block stacking task that was initially presented in Li et al. (2017). The second part of the work features an extension to goal-parametrized deep reinforcement learning framework and its application to target stacking task. In addition, we now include a discussion of the simulation environment implemented for both tasks. While simulation engines and game engines, such as Panda3D (Goslin and Mine, 2004), provide a generic platform to devise environment to different needs, it is still very time-consuming and remains non-trivial to customize one for a specific tasks such as the target stacking in this work. Hence, we briefly recap our design for the environments to provide readers more insights for further practice.

# 2 Related Work

## 2.1 Physics Understanding

Humans possess the amazing ability to perceive and understand ubiquitous physical phenomena occurring in their daily life. This gives rise to the concept of intuitive physics, aiming to describe the knowledge which enables humans to understand physical environment and interact accordingly. In particular, the "intuitive" part emphasizes knowledge which is considered commonsense to ordinary people not reliant on specialized training in physics. Intuitive physics is ubiquitous in guiding humans' actions in daily life, such as where to put a cup stably and how to catch a ball. Along the years, research on intuitive physics has been conducted from many different perspectives across psychology, cognitive science and artificial intelligence.

In developmental psychology, researchers seek to understand how this ability develops. Baillargeon (2002) suggest that infants acquire the knowledge of physical events at a very young age by observing those events, including support events and others. Interestingly, in a recent work Denil et al. (2017), the authors introduce a basic set of tasks that require the learning agent to estimate physical properties (mass and cohesion combinations) of objects in an interactive simulated environment and find that it can learn to perform the experiments strategically to discover such hidden properties in analogy to human's development of physics knowledge. Another interesting question that has been explored in psychology is how knowledge about physical events affects and guides human's actual interaction with objects Yildirim et al. (2017). Yet it is not clear how a machine model trained for physics understanding can directly be applied to real-world interactions with objects and accomplish manipulation tasks.

In cognitive science, Battaglia et al. (2013) proposes an intuitive physics simulation engine as an internal mechanism for such type of ability and found close correlation between its behavior patterns and human subjects' on several psychological tasks.

More recently, there has been renewed interest in physics understanding in computer vision and machine learning communities. For instance, understanding physical events plays an important role in scene understanding in computer vision. By including additional clues from physical constraints into the inference mechanism, mostly from the support event, it has further improved results in segmentation of surfaces Gupta et al. (2010), scenes Sil-

berman et al. (2012) from image data, and object segmentation in 3D point cloud data Zheng et al. (2013). Another research direction is to equip artificial agents with such an ability by letting them learn physical concepts from visual data. Mottaghi et al. (2016) aim at understanding dynamic events governed by laws of Newtonian physics and use proto-typical motion scenarios as exemplars. Fragkiadaki et al. (2016) analyze billiard table scenarios and learn dynamics from observation with explicit object notion. An alternative approach based on boundary extrapolation Bhattacharyya et al. (2018) addresses similar settings without imposing any object notion. Wu et al. (2015) aims to understand physical properties of objects based on explicit physical simulation. Mottaghi et al. (2017) proposes to reason about containers and the behavior of the liquids inside them from a single RGB image.

A most related work is done by Lerer et al. (2016), where the authors propose using a visual model to predict stability and falling trajectories for simple 4 block scenes. In the first part of this work, we investigate if and how the prediction performance of such image-based models changes when trained on block stacking scenes with larger variety and further examine how the human's prediction adapts to the variation in the generated scenes and compare to the learned visual model. Each work requires significant amounts of simulated, physically-realistic data to train the large-capacity, deep models.

## 2.2 Learning from Synthetic Data and Simulation

Learning from synthetic data has a long tradition in computer vision and has recently gained increasing interest (Li and Fritz, 2012; Rematas et al., 2014; Peng et al., 2015; Rematas et al., 2016) due to data hungry deep-learning approaches. In the first part of this work, we use a game engine to render scene images and a built-in physics simulator to simulate the scenes' stability behavior. The data generation procedure is based on the platform used in Battaglia et al. (2013), however as discussed before, their work hypothesized a simulation engine as an internal mechanism for human to understand the physics in the external world while we are interested in finding an image-based model to directly predict the physical behavior from visual channel.

In reinforcement learning domain, it is a common practice to utilize simulation environment to train the reinforcement agent, such as the Atari Games (Bellemare

et al., 2013) and ViZDoom (Kempka et al., 2016). However, there is no off-the-shelf environment supporting trial-and-error interaction for block stacking. Hence, in the second part of this work, we build on existing game engine and create our interactive environment with physics simulation to allow the learning agent to learn stacking skills through interactions.

## 2.3 Blocks-based Manipulation Tasks

To shed more light on the capabilities of our model, in the first part of this work, we explore how the visual stability prediction model can be used in a robotic manipulation task, i.e., stably stacking a wood block given a block structure. In the past, we have seen researchers perform tasks with wood blocks, like playing Jenga from different perspectives. Kröger et al. (2006) demonstrated multi-sensor integration by using a marker-based system with multiple cameras and sensors: a random block is first chosen in the tower, then the robot arm will try to pull the very block, if the force sensor detects large counter force or the CCD cameras detect large motion of tower, then it will stop pulling and try other block.

Wang et al. (2009) improved on Kröger et al. (2006) by further incorporating a physics engine to initialize the candidates for pulling test. In comparison, we do not fit 3D models or run physics simulation at test time for the given scene but instead use the scene image as input to directly predict the physics of the structure.

A different line of research is Kimura et al. (2010) where physical force is explicitly formulated with respect to the tower structure for planning. In our work, we do not do explicit formation of contact force as in Kimura et al. (2010), nor do we perform trials on-site for evaluating the robot's operation. We only use physics engine to acquire synthesized data for training the visual-physics model. At test time, the planning system for our robot mainly exploits the knowledge encoded in the visual-physics model to evaluate the feasibility of individual candidates and performs operations accordingly.

## 2.4 Reinforcement Learning

In the second part of this work, reinforcement learning is used to learn an end-to-end model directly from the experience collected during interaction with a physically-realistic environment. The majority of work in reinforcement learning focuses on solving task with a single goal. However, there are also tasks where the goal may change for every trial. It is not obvious how to directly apply the model learned towards a specific goal to a different one. An early idea has been proposed by Kaelbling (1993) for a maze navigation problem in which the goal changes. The author introduces an analogous formulation to the Q-learning by using shortest path in replacement of the value functions. Yet there are two major limitations for the framework: 1) it is only formulated in tabular form which is not practical for application with complex states 2) the introduced shortest path is very specific to the maze navigation setting and hence cannot be easily adapt to handle task like target stacking. In contrast, we propose a goal-parameterized model to integrate goal information into a general learning-based framework that facilitates generalization across different goals. The model has been shown to work on both a navigation task and target stacking.

A few other works also explored the idea of integrating goal information into learning. Schaul et al. (2015) propose the universal value function approximators (UVFAs) to integrate goal information into learning. Oh et al. (2017) proposes zero-shot task generation with explicit modeling of skills and subtasks where our approach is end-to-end and thereby bypasses this kind of modeling. Dosovitskiy and Koltun (2017) use the terms "goal" and "target". However, their notion is different from ours. Their framework is using auxiliary information (intermediate measurements as goal module) to improve learning to maximize gaming score in visually-richer task like vizDoom. There is no specific goal to enforce the agent to reach certain state which is exactly what our approach facilitates. The metacontroller (Hamrick et al., 2017), the imagination strategy (Pascanu et al., 2017) and the I2A model (Weber et al., 2017) are all appealing generic frameworks to boost overall performance over existing RL methods. While these approaches – similar to some of the Atari games – have a goal encoded in the observation/state, this goal remains implicit. This is applicable in many robotic scenarios, where goals are specified externally. Therefore, we explore an approach for target stacking that has an explicit notion and model of goals. Our work is the first to bring this concept to bear towards manipulation and planing under physical constraints – breaking with more conventional simulation and planning approaches (e.g. Yildirim et al. (2017)).

## 3 Simulation Environment

Modern machine learning techniques build on data. With the rise of deep learning models, the need of data be-

comes even more prominent. One of the key elements for the recent success of deep learning model in domains like generic image and speech recognition is the abundance of data in related fields. However, when it comes to a specific domain, it often runs short of data and in reality, collecting data is still an expensive process. One remedy for this issue is to exploit the domain knowledge to synthesize data and utilize the generated data for learning. This is also the case for our study as there is no obvious source of clean and sufficient data to learn physics from. Hence, we make use of the Panda3D [1] to create both simulation environments in our work.

Panda3D is an open source game engine for Python and C++ programs. It was originally developed by the Disney's VR studio to be "flexible enough to support everything from realtime graphics applications to the development of high-end virtual reality theme park attractions or video games" (Goslin and Mine, 2004) and it has evolved significantly along the years.

As a game engine, the Panda3D provides additional capabilities besides 3D rendering including the physics system with integration of different physics engines. Asides from its own built-in basic physics engine, it also supports more advanced ones including the Open Dynamics Engine (ODE) [2] (Smith, 2005) and the Bullet [3] (Coumans, 2010) for physics simulation back-end. We used the Bullet in both of our simulation environments.

The workflow of Panda3D builds on the concept of scene graph. The **Scene graph** is a general data structure to represent and organize a graphical scene. In Panda3D, the scene graph is maintained as a tree of objects to be rendered. The tree consists of objects of class `PandaNode`. As shown in Figure 2, the root node is called the `render` and the rest define different perspectives of the scene with various attributes. For example, the `LensNode` controls the camera such as the perspective, the `LightNode` manages the lighting in the scene such as the color and type of the lighting and the `ModelNode` encodes the 3D object in the frame.

Another important concept is the **task**. Tasks are functions called by Panda3D at every frame or for every specified amount of time. Together with **event handler** which is called upon special conditions (**events**) occur, update can be made to the scene in Panda3D between rendering steps as shown in Figure 3. For instance, the task can be the simulation subroutine that updates the states of
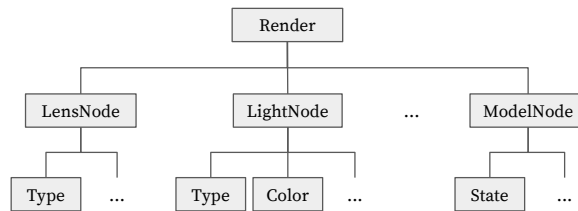


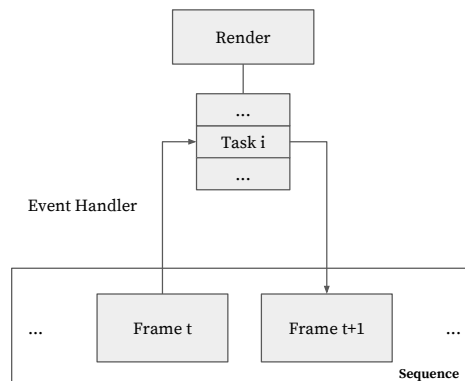Figure 2: An example of scene graph in Panda3D.



Figure 3: The render process in Panda3D.

objects in the scene caused by physics.

# 4 Part I: From Visual Stability Prediction to Single Block Stacking

## 4.1 Stability Prediction from Still Images

**Inspiration from Human Studies** Research in Hamrick et al. (2011); Battaglia et al. (2013) suggests the combinations of the most salient features in the scenes are insufficient to capture people's judgments, however, contemporary study reveals human's perception of visual information, in particular some geometric feature, like critical angle Cholewiak et al. (2013, 2015) plays an important role in the process. Regardless of the actual inner mechanism for humans to parse the visual input, it is clear there is a mapping $f$ involving visual input $I$ to the stability
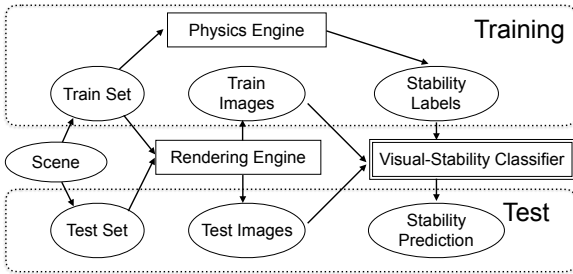
---

Figure 4: An overview of our approach for learning visual stability. Note that physics engine is only used during training time to get the ground truth to train the deep neural network while at test time, only rendered scene images are given to the learned model to predict the physical stability of the scenes.

prediction $P$.

$$f : I, * \rightarrow P$$

Here, $*$ denotes other possible information, i.e., the mapping can be inclusive, as in Hamrick et al. (2011) using it along with other aspects, like physical constraint to make judgment or the mapping is exclusive, as in Cholewiak et al. (2013) using visual cues alone to decide.

**Image Classifier for Stability Prediction** In our work, we are interested in the mapping $f$ exclusive to visual input and directly predicts the physical stability (visual stability test). To this end, we use deep convolutional neural networks as it has shown great success on image classification tasks (Krizhevsky et al., 2012). Such networks have been shown to be able to adapt to a wide range of classification and prediction task (Razavian et al., 2014) through re-training or adaptation by fine-tuning. Therefore, these approaches seem to be adequate methods to study visual prediction on this challenging task with the motivation that by changing conventional image classes labels to stability labels the network can learn "physical stability salient" features. By setting up a data generation process that allows us to control various degrees of freedom induced by the problem as well as generation of large quantities of data in a repeatable setting, we can then validate our approach. An overview of our approach is shown in Figure 4.
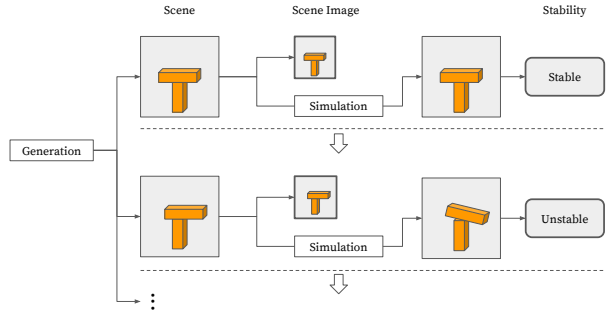


Figure 5: Overview of the data generator for visual stability prediction. The scene images and their corresponding stability labels are collected.

## 4.2 Data Generator for Visual Stability Prediction

Based on the scene simulation framework used in (Hamrick et al., 2011; Battaglia et al., 2013), we build a data generator to synthesize data for visual stability tests. Figure 5 gives an overview of the data generator. The first key component is the **tower generation**. It automatically generate a large number of different blocks structures (tower) under the *scene parameters*, including the number of blocks, stacking depth and block size (We will elaborate on this later). These towers are recorded as scene files which only track all the locations of blocks in the scene.

The next component is the **stability simulation**. It loads the stored scene files and simulates their stability with physics engine. The images for the towers before running the physics engine are captured as the scene images, the stability labels from the simulation are automatically determined and recorded for the corresponding towers. Both the scene images and the obtained stability labels are later put into the deep convolutional neural network to learn the visual stability classifier.

### 4.2.1 Tower Generation

An example of the scene setup is shown Figure 6. The tower is placed on a plane, and a camera is positioned at the front-facing location of which the elevation is adjusted for the towers of different heights. For different scenes, the towers are generated differently, and the scene images are captured through the camera. To make the scene images more realistic, wood texture is added to all the blocks.
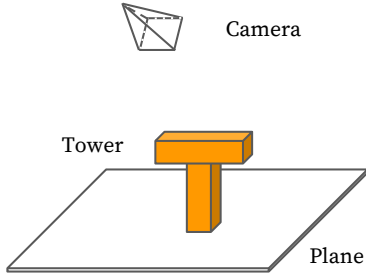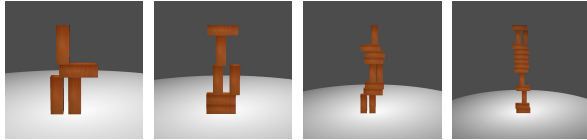
Figure 6: Example set-up of the scene for the data generator.



(a) 4 Blocks  (b) 6 Blocks  (c) 10 Blocks  (d) 14 Blocks

Figure 7: Example of generated scene images with different total number of blocks.

The basic tower generation system is based on the framework by Battaglia et al. (2013). Given a specified number of total blocks in the tower, the blocks are sequentially added into the scene under the geometrical constraints, such as no collision between blocks. Some examples of obtained scene images are shown in Figure 7.

### 4.2.2 Stability Simulation

During the stability simulation, we set the simulation time universally to 2 seconds at $1000Hz$ for all the scenes. Surface friction and gravity are enabled in the simulation. The system records the configuration of a scene of $N$ blocks at time $t$ as $(p_1, p_2, ..., p_N)_t$, where $p_i$ is the location for block $i$. The stability is then automatically decided as a Boolean variable:

$$S = \bigvee_{i=1}^{N} (\Delta((p_i)_{t=T} - (p_i)_{t=0}) > \tau)$$

where $T$ is the end time of simulation, $\delta$ measures the displacement for the blocks between the starting point and end time, $\tau$ is the displacement threshold, $\bigvee$ denotes the logical Or operator, that is to say it counts as unstable $S = $ True if any block in the scene moved in simulation, otherwise as stable $S = $ False. This is necessary as

the blocks in stable towers can generate small displacement during the simulation process, simply using zero displacement to determine the stability can lead to lots of erroneous cases where the stable towers are labeled as unstable. In practice, we picked the threshold based on the evaluation of a small set of towers.

## 4.3 Synthetic Data

The number of blocks, blocks' size and stacking depth are varied in the generated scenes, to which we will refer as *scene parameters*.

**Numbers of Blocks**  We expect that varying the size of the towers will influence the difficulty and challenge the competence of "eye-balling" the stability of a tower in humans and machine. While evidently the appearance becomes more complex with the increasing number of blocks, the number of contact surfaces and interactions equally make the problem richer. Therefore, we include scenes with four different number of blocks, i.e., 4 blocks, 6 blocks, 10 blocks and 14 blocks as $\{4B, 6B, 10B, 14B\}$.

**Stacking Depth**  As we focus our investigations on judging stability from a monocular input, we vary the depth of the tower from a one layer setting which we call $2D$ to a multi-layer setting which we call $3D$. The first one only allows a single block along the image plane at all height levels while the other does not enforce such constraint and can expand in the image plane. Visually, the former results in a single-layer stacking similar to Tetris while the latter ends in a multiple-layer structure as shown in Table 1. The latter most likely requires the observer to pick up on more subtle visual cues, as many of its layers are heavily occluded.

**Block Size**  We include two groups of block size settings. In the first one, the towers are constructed of blocks that have all the same size of $1 \times 1 \times 3$ as in the Battaglia et al. (2013). The second one introduces varying block sizes where two of the three dimensions are randomly scaled with respect to a truncated Normal distribution $N(1, \sigma^2)$ around $[1 - \delta, 1 + \delta]$, $\sigma$ and $\delta$ are small values. These two settings are referred to as $\{Uni, NonUni\}$. The setting with non-uniform blocks introduces small visual cues where stability hinges on small gaps between differently sized blocks that are challenging even for human observers.

7

**Scenes** Combining these three scene parameters, we define 16 different scene groups. For example, group 10B-2D-Uni is for scenes stacked with 10 Blocks of same size, stacked within a single layer. For each group, 1000 candidate scenes are generated where each scene is constructed with non-overlapping geometrical constraint in a bottom-up manner. There are 16K scenes in total. For prediction experiments, half of the images in each group are for training and the other half for test, the split is fixed across the experiments.

**Rendering** While we keep the rendering basic, we like to point out that we deliberately decided against colored bricks as in Battaglia et al. (2013) in order to challenge perception and make identifying brick outlines and configurations more challenging. The lighting is fixed across scenes and the camera is automatically adjusted so that the whole tower is centered in the captured image. Images are rendered at resolution of $800 \times 800$ in color.

## 4.4 Prediction Performance

In this part of the experiments, the images are captured before the physics engine is enabled, and the stability labels are recorded from the simulation engine as described before. At the training time, the model has access to the images and the stability labels. At test time, the learned model predicts the stability results against the results generated by the simulator.

In a pilot study, we tested on a subset of the generated data with LeNet (LeCun et al., 1995), a relatively small network designed for digit recognition, AlexNet (Krizhevsky et al., 2012), a large network and VGG Net (Simonyan and Zisserman, 2014), an even larger network than AlexNet. We trained from scratch for the LeNet and fine-tuned for the large network pre-trained on ImageNet Deng et al. (2009). VGG Net consistently outperforms the other two, hence we use it across our experiment. We use the Caffe framework (Jia et al., 2014) in all our experiments.

We divide the experiment design into 3 sets: the intragroup, cross-group and generalization. The first set investigates influence on the model's performance from an individual scene parameter, the other two sets explore generalization properties under different settings.

### 4.4.1 Intra-Group Experiment

In this set of experiments, we train and test on the scenes with the same scene parameters in order to assess the

feasibility of our task.

**Number of Blocks (4B, 6B, 10B, 14B)** In this group of experiment, we fix the stacking depth and keep the all blocks in the same size but vary the number of blocks in the scene to observe how it affects the prediction rates from the image trained model, which approximates the relative recognition difficulty from this scene parameter alone. The results are shown in Table 2. A consistent drop of performance can be observed with increasing number of blocks in the scene under various block sizes and stacking depth conditions. More blocks in the scene generally leads to higher scene structure and hence higher difficulty in perception.

**Block Size (Uni. vs. NonUni.)** In this group of experiment, we aim to explore how same size and varied blocks sizes affect the prediction rates from the image trained model. We compare the results at different number of blocks to the previous group, in the most obvious case, scenes happened to have similar stacking patterns and same number of blocks can result in changes visual appearance. To further eliminate the influence from the stacking depth, we fix all the scenes in this group to be 2D stacking only. As can be seen from Table 2, the performance decreases when moving from 2D stacking to 3D. The additional variety introduced by the block size indeed makes the task more challenging.

**Stacking Depth (2D vs. 3D)** In this group of experiment, we investigate how stacking depth affects the prediction rates. With increasing stacking depth, it naturally introduces ambiguity in the perception of the scene structure, namely some parts of the scene can be occluded or partially occluded by other parts. Similar to the experiments in previous groups, we want to minimize the influences from other scene parameters, we fix the block size to be the same and only observe the performance across different number of blocks. The results in Table 2 show a little inconsistent behaviors between relative simple scenes (4 blocks and 6 blocks) and difficult scenes (10 blocks and 14 blocks). For simple scenes, prediction accuracy increases when moving from $2D$ stacking to $3D$ while it is the other way around for the complex scene. Naturally relaxing the constraint in stacking depth can introduce additional challenge for perception of depth information, yet given a fixed number of blocks in the scene, the condition change is also more likely to make the scene structure lower which reduces the difficulty in perception.

| Block Numbers | | | | Stacking Depth | | Block Size | |
|---|---|---|---|---|---|---|---|



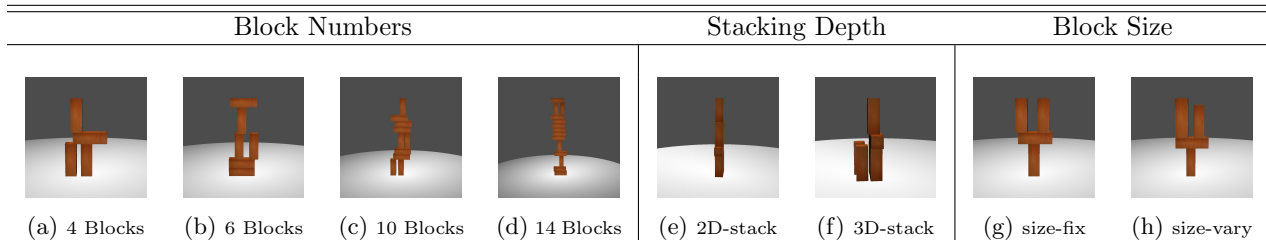| (a) 4 Blocks | (b) 6 Blocks | (c) 10 Blocks | (d) 14 Blocks | (e) 2D-stack | (f) 3D-stack | (g) size-fix | (h) size-vary |
|---|---|---|---|---|---|---|---|

Table 1: Overview of the scene parameters in our rendered scenes. There are 3 groups of scene parameters across number of blocks, stacking depth and block size.

| # Blks | Uni. | | NonUni. |
|---|---|---|---|
| | 2D | 3D | 2D |
| 4B | 93.0 | 99.2 | 93.2 |
| 6B | 88.8 | 91.6 | 88.0 |
| 10B | 76.4 | 68.4 | 69.8 |
| 14B | 71.2 | 57.0 | 74.8 |

Table 2: Intra-group experiment by varying scene parameters.



| Setting | Simple → Complex | Complex → Simple |
|---|---|---|
| Accuracy (%) | 69.9 | 86.9 |

Table 3: The upper figure shows the experiment settings for Cross-group classification where we train on simpler scenes and test on more complex scenes. The lower table shows the results.

A combination of these two factors decides the final difficulty of the task, for simple scenes, the height factor has stronger influence and hence exhibits better prediction accuracy for $3D$ over $2D$ stacking while for complex scenes, the stacking depth dominates the influence as the significant higher number of blocks can retain a reasonable height of the structure, hence receives decreased performance when moving from $2D$ stacking to $3D$.

### 4.4.2 Cross-Group Experiment

In this set of experiment, we want to see how the learned model transfers across scenes with different complexity, so we further divide the scene groups into two large groups by the number of blocks, where a *simple scene* group for all the scenes with 4 and 6 blocks and a *complex scene* for the rest of scenes with 10 and 14 blocks. We investigate in two-direction classification, shown in the figure in Table 3:

1. Train on simple scenes and predict on complex scenes: Train on 4 and 6 blocks and test on 10 and 14 blocks

2. Train on complex scenes and predict on simple scenes: Train on 10 and 14 blocks and test on 4 and 6 blocks

As shown in Table 3, when trained on simple scenes and predicting on complex scenes, it gets 69.9%, which is significantly better than random guess at 50%. This is understandable as the learned visual feature can transfer across different scene. Further we observe significant performance boost when trained on complex scenes and tested on simple scene. This can be explained by the richer feature learned from the complex scenes with better generalization.

### 4.4.3 Generalization Experiment

In this set of experiment, we want to explore if we can train a general model to predict stability for scenes with any scene parameters, which is very similar to human's prediction in the task. We use training images from all

| # Blks | Uni. | | NonUni. | |
|---|---|---|---|---|
| | 2D | 3D | 2D | 3D |
| 4B | 93.2 | 99.0 | 95.4 | 99.8 |
| 6B | 89.0 | 94.8 | 87.8 | 93.0 |
| 10B | 83.4 | 76.0 | 77.2 | 74.8 |
| 14B | 82.4 | 67.2 | 78.4 | 66.2 |

Table 4: Results for generalization experiments.

| # Blks | Uni. | | NonUni. | |
|---|---|---|---|---|
| | 2D | 3D | 2D | 3D |
| 4B | 79.1/**91.7** | 93.8/**100.0** | 72.9/**93.8** | 92.7/**100.0** |
| 6B | 78.1/**91.7** | 83.3/**93.8** | 71.9/**87.5** | 89.6/**93.8** |
| 10B | 67.7/**87.5** | 72.9/72.9 | 66.7/**72.9** | 71.9/68.8 |
| 14B | 71.9/**79.2** | 68.8/66.7 | 71.9/**81.3** | 59.3/**60.4** |

Table 5: Results from human subject test $a$ and corresponded accuracies from image-based model $b$ in format $a/b$ for the sampled data.

different scene groups and test on any groups. The Result is shown in Table 4. While the performance exhibits similar trend to the one in the intra-group with respect to the complexity of the scenes, namely increasing recognition rate for simpler settings and decreasing rate for more complex settings, there is a consistent improvement over the intra-group experiment for individual groups. Together with the result in the cross-group experiment, it suggests a strong generalization capability of the image trained model.

### 4.4.4 Discussion

Overall, we can conclude that direct stability prediction is possible and in fact fairly accurate at recognition rates over 80% for moderate difficulty levels. As expected, the 3D setting adds difficulties to the prediction from appearance due to significant occlusion for towers of more than 10 blocks. Surprisingly, little effect was observed for small tower sizes switching from uniform to non-uniform blocks

- although the appearance difference can be quite small. To better understand our results, we further discuss the following two questions:

**How does the model performs compared to human?** To answer this, we conduct a human subject test. We recruit human subjects to predict stability for give scene images. Due to large number of test data, we sample images from different scene groups for human subject test. 8 subjects are recruited for the test. Each subject is presented with a set of captured images from the test split. Each set includes 96 images where images cover all 16 scene groups with 6 scene instances per group. For each scene image, subject is required to rate the stability on a scale from $1 - 5$ without any constraint for response time:

1. Definitely unstable: definitely at least one block will move/fall
2. Probably unstable: probably at least one block will move/fall
3. Cannot tell: the subject is not sure about the stability
4. Probably stable: probably no block will move/fall
5. Definitely stable: definitely no block will move/fall

The predictions are binarized, namely 1) and 2) are treated as unstable prediction, 4) and 5) as stable prediction, "Cannot tell" will be counted as 0.5 correct prediction.

The results are shown in Table 5. For simple scenes with few blocks, human can reach close to perfect performance while for complex scenes, the performance drops significantly to around 60%. Compared to human prediction in the same test data, the image-based model outperforms human in most scene groups. While showing similar trends in performance with respect to different scene parameters, the image-based model is less affected by a more difficult scene parameter setting, for example, given the same block size and stacking depth condition, the prediction accuracy decreases more slowly than the counter part in human prediction. We interpret this as image-based model possesses better generalization capability than human in the very task.

**Does the model learn something explicitly interpretable?** Here we apply the technique from Zhou et al. (2016) to visualize the learned discriminative image regions from CNN for individual category. The approach is illustrated in Figure 8a. With Global Average Pooling (GAP), the resulted spatial average of the feature maps from previous convolutional layers forms
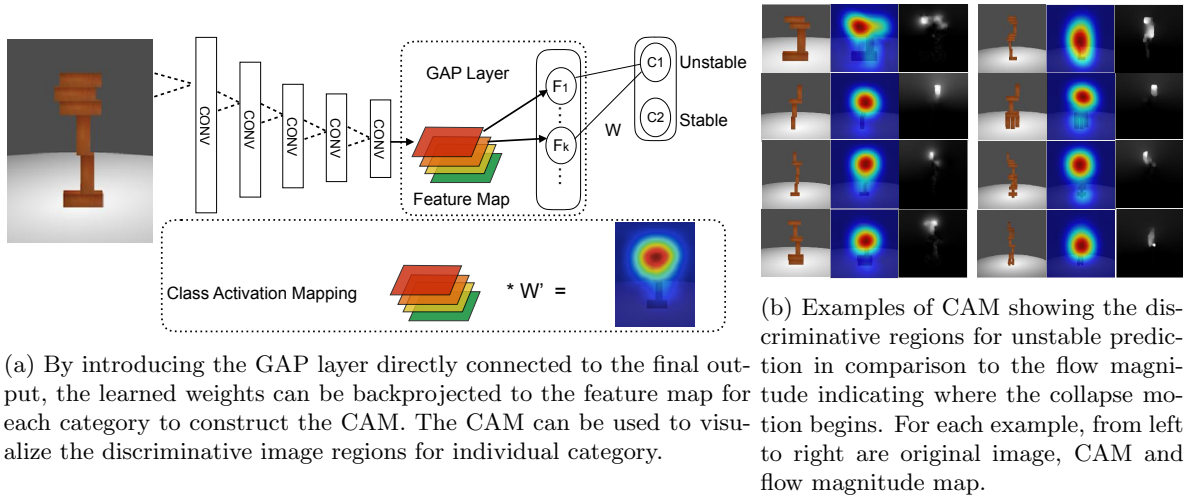
(a) By introducing the GAP layer directly connected to the final output, the learned weights can be backprojected to the feature map for each category to construct the CAM. The CAM can be used to visualize the discriminative image regions for individual category.

(b) Examples of CAM showing the discriminative regions for unstable prediction in comparison to the flow magnitude indicating where the collapse motion begins. For each example, from left to right are original image, CAM and flow magnitude map.

Figure 8: We use CAM to visualize the results for model interpretation.

fully-connected layer to directly decides the final output. By back-projecting the weights from the fully-connected layer from each category, we can hence obtain Class Activation Map (CAM) to visualize the discriminative image regions. In our case, we investigate discriminative regions for unstable predictions to see if the model can spot the weakness in the structure. We use deep flowWeinzaepfel et al. (2013) to compute the optical flow magnitude between the frame before the physics engine is enabled and the one afterwards to serve as a coarse ground truth for the structural weakness where we assume the collapse motion starts from such weakness in the structure. Though not universal among the unstable cases, we do find significant positive cases showing high correlation between the activation regions in CAM for unstable output and the regions where the collapse motion begins. Some examples are shown in Figure 8b.

In the previous section, we have shown that an appereance-based model can predict physical stability relatively well on the synthetic data. Now we want to further explore if and how the synthetic data trained model can be utilized for a real world application, especially for robotic manipulation. Hence, we decide to set up a testbed where a Baxter robot's task is to stack one wood block on a given block structure without breaking the structure's stability as shown in Figure 1. The overview of our system is illustrated in Figure 9. In our experiment, we use Kapla blocks as basic unit, and tape 6 blocks into a bigger one as shown in Figure 10a. To simplify the task,

adjustments were made to the free-style stacking:

- The given block structure is restricted to be single layer as the $2D$ case in the previous section. For the final test, we report results on the 6 scenes as shown in Table 6.

- The block to be put on top of the given structure is limited two canonical configurations {vertical, horizontal} as shown in Figure 10b. and assumed to be held in hand of robot before the placement.

- The block is constrained to be placed on the topmost horizontal surface (stacking surface) in the given structure.

- The depth of the structure (perpendicular distance to the robot) is calibrated so that we only need to decide the horizontal and vertical displacements with respect to the stacking surface.

## 4.5   Prediction on Real World Data

Considering there are significant difference between the synthesized data and real world captured data, including factors (not limited to) as texture, illumination condition, size of blocks and accuracy of the render, we performed a pilot study to directly apply the model trained on the RGB images to predict stability on the real data, but only got results on par with random guessing. Hence we decided to train the visual-stability model on the binary-valued foreground mask on the synthesized data and deal
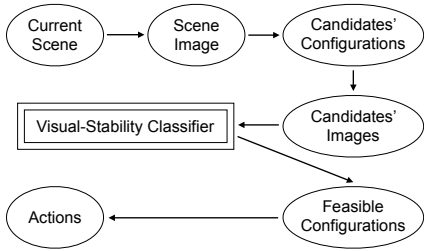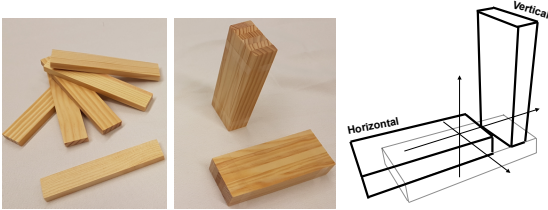
Figure 9: An overview of our manipulation system. Our visual-stability classifier is integrated to recognize feasible candidate placement to guide manipulation.



(a) Kapla block (left), block in test (right).

(b) Allowed configurations in test.

Figure 10: Blocks used in our experiment.

with the masks at test time also for the real scenes. In this way, we significantly reduce the effect from the aforementioned factors. Observing comparable results when using the RGB images, we continue to the approach on real world data.

At test time, a background image is first captured for the empty scene. Then for each test scene (shown in Table 6), an image is captured and converted to foreground mask via background subtraction. The top-most horizontal boundary is detected as the stacking surface and then used to generate candidate placements: the surface is divided evenly into 9 horizontal candidates and 5 vertical candidates, resulting in 84 candidates. The process is shown in Figure 11. Afterwards, these candidates are put to the visual-stability model for stability prediction. Each generated candidate's actual stability is manually tested and recorded as ground truth. The final recognition result is shown in Table 6. The model trained with synthetic data is able to predict with overall accuracy of 78.6% across different candidates in real world.
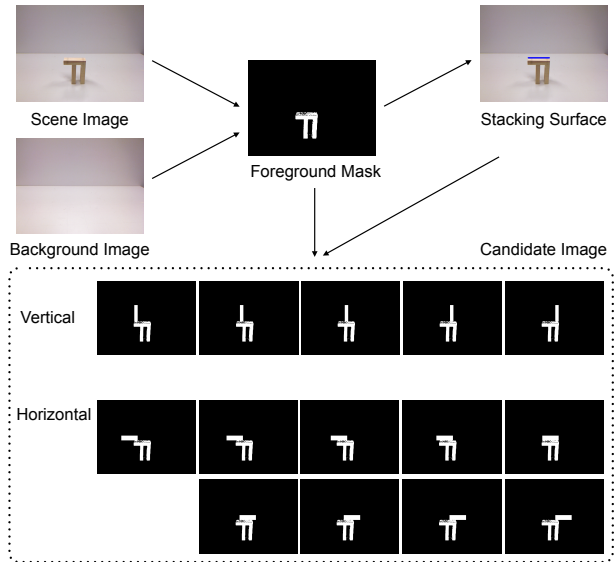


Figure 11: The procedure to generate candidates placement images for a give scene in our experiment.

## 4.6 Manipulation Test

At test time, when the model predicts a give candidate placement as stable, the robot will execute routine to place the block with 3 attempts. We count the execution as a success if any of the attempt works. The manipulation success rate is defined as:

$$\frac{\#\{\text{successful placements}\}}{\#\{\text{all stable placements}\}}$$

where #{successful placements} is the number of successful placements made by the robot, and #{all stable placements} is the number of all ground truth stable placements.

As shown in Table 6, the manipulation performance is good across most of the scenes for both horizontal and vertical placements except for the 6-th scene where the classifier predicts all candidates as unstable hence no attempts have been made by the robot.

## 4.7 Discussion

Simply putting the block along the center of mass (COM) of the given structure may often be a feasible option, yet, there are two limitations to this approach: first, it is nontrivial to compute the COM of a given structure; second,

12

| Id. | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Scene | | | | | | | | | | | | |
| Pred.(%) | 66.7 | 100.0 | 66.7 | 60.0 | 88.9 | 100.0 | 77.8 | 80.0 | 100.0 | 40.0 | 66.7 | 60.0 |
| Mani.(%) | $80.0_{(4/5)}$ | $100.0_{(5/5)}$ | $66.7_{(2/3)}$ | $100.0_{(3/3)}$ | $66.7_{(2/3)}$ | $100.0_{(1/1)}$ | $66.7_{(2/2)}$ | $66.7_{(2/3)}$ | $100.0_{(3/3)}$ | $25.0_{(1/4)}$ | $0.0_{(0/3)}$ | $0.0_{(0/1)}$ |
| Placement | H | V | H | V | H | V | H | V | H | V | H | V |

Table 6: Results for real world test. "Pred." is the prediction accuracy. "Mani." is the manipulation success rate with counts for successful placements/all possible stable placements for each scene. "H/V" refer to horizontal/vertical placement.

it only gives one possible stable solution (assuming it actually stay stable). In comparison, our method does not rely the COM of the structure and provide a search over multiple possible solutions.

So far, we have successfully applied our visual stability prediction model to guide the stable placement of a single block onto the existing structure. However, this approach has its own limitations. Firstly, it only deals with one manipulation step and does not take into account the interaction process which is also under physics constraint. For instance, the block can collide with existing structure before its placement. Secondly, it is also not obvious how to extend the approach to plan a consecutive stacking of multiple blocks. In part II, we will explore an alternative approach to overcome these limitations.

# 5 Part II: Beyond Single Block Stacking — Target Stacking

In this part, we look at a manipulation task beyond the single block placement as discussed in part I. We now seek solutions to plan the placement of multiple blocks and consider more complex interactions in physical environment. In particular, just like children stack blocks towards a certain blueprint, we further introduce the "target" into the stacking task so that the sequence of manipulations are planned to a achieve the target — similar to planning. Yet we want to stay fully in an end-to-end learning approach – concerning the perception, physics as well as the plan, where intuitive physics as well as plan execution is learned jointly.

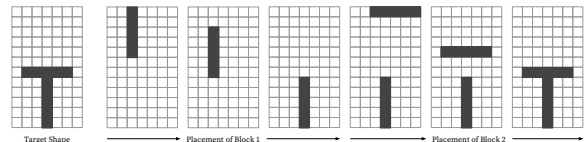The extended block stacking task is called *target stack-*



Figure 12: Target stacking: Given a target shape image, the agent is required to move and stack blocks to reproduce it.

*ing*. In this task, an image of a target structure made of stacked blocks is provided. Given the same number of blocks as in the target structure, the goal is to reproduce the structure shown in the image. The manipulation primitives in this task include moving and placing blocks. This is inspired by the scenario where young children learn to stack blocks to different shapes given an example structure. We want to explore how an artificial agent can acquire such a skill through trial and error.

## 5.1 Task Description

For each task instance, a target structure is generated and its image is provided to the agent along with the number of blocks. Each of these blocks has a fixed orientation. The sequence of block orientations is such that reproducing the target is feasible. The agent attempts to construct the target structure by placing the blocks in the given sequence. The spawning location for each block is randomized along the top boundary of the environment. A illustrative sample for the task is shown in Figure 12.

## 5.2 Task Distinction

The following characteristics distinguish this task from other tasks commonly used in the literature.

**Goal-Specific** A widely-used benchmark for deep reinforcement learning algorithm are the Atari games (Bellemare et al., 2013) that were made popular by Mnih et al. (2013). While this game collection has a large variety, the games are defined by a single goal or no specific goal is enforced at a particular point in time. For example in Breakout, the player tries to bounce off as many bricks as possible. In Enduro, the player tries to pass as many cars as possible while simultaneously avoiding cars.

In the target stacking task, each task instance differs in the specific goal (the target structure), and all the moves are planned towards this goal. Given the same state, moves that were optimal in one task instance are unlikely to be optimal in another task instance with a different target structure. This is in contrast to games where one type of move will most likely work in similar scenes. This argument also applies to AI research platforms with richer visuals like VizDoom (Kempka et al., 2016).

**Longer sequences** Target stacking requires looking ahead over a longer time horizon to simultaneously ensure stability and similarity to the target structure. This is different from learning to poke (Agrawal et al., 2016) where the objective is to select a motion primitive that is the optimal next action. It is also different from the work by Li et al. (2017) that reasons about the placement of one block.

**Rich Physics Bounded** Besides stacking to the assigned target shape the agent needs to learn to move the block without colliding with the environment and existing structure and to choose the block's placement wisely not to collapse the current structure. The agent has no prior knowledge of this. It needs to learn everything from scratch by observing the consequence (collision, collapse) of its actions.

## 5.3 Target Stacking Environment

Compared to the environment for visual stability prediction as a data generator, we now need more a more capable task environment supports the learning process with which the agent can interact. This requires us to significantly extend the implementation.

The environment now implements a more dynamic system where different moves from the agent can lead to different consequences: when a block (not the last one for the episode) is placed stably in the scene, a new block should be spawned into the environment whereas if the block collapses the existing structure, the current episode is terminated so that a new one can be started. The key to react correctly under these different conditions is to detect them effectively.

While we keep the essential parts of the task, at its current stage the simulated environment remains an abstraction of a real-world robotics scenario. This generally requires an integration of multiple modules for a full-fledged working system, such as Toussaint et al. (2010), which is out of scope of this paper.

An overview of the stacking environment is shown in Figure 13. The environment is again implemented in Panda3D as the data generator described in the previous part. The block size follows a ratio of $l : w : h = 5 : 2 : 1$, where $l,w,h$ denote length, width and height respectively.

For each episode, a target structure is randomly picked from a pool of structures and then the blocks are spawned according to the target. The agent moves the spawned block in the scene until the placement is done.

During this process, if the agent's move causes the block to (1) move out over the boundary of the scene or (2) collide with existing structure (3) collapse the structure with the placement, the episode is terminated and a new episode can start.

A new block is spawned upon the previous one is placed successfully till reaching the total number of blocks in the target. When all the blocks are placed in the scene, the environment determines if the agent achieves the assigned target structure.

We ignore the impact during block placement and focus on the resulting stability of the entire structure. Once the block makes contact with the existing structure, it is treated as releasing the block for a placement. The physics simulation runs at $60Hz$. However considering the cost of simulation we only use it when there is contact between the moving block and the boundary or the existing structure. Otherwise, the current block is moving without actual physics simulation.

### 5.3.1 State Representation

To further reduce the appearance difference caused by varying perspective, the state in the stacking environment is rendered using orthographic projection. It is implemented by the `OrthographicLens` in Panda3D where
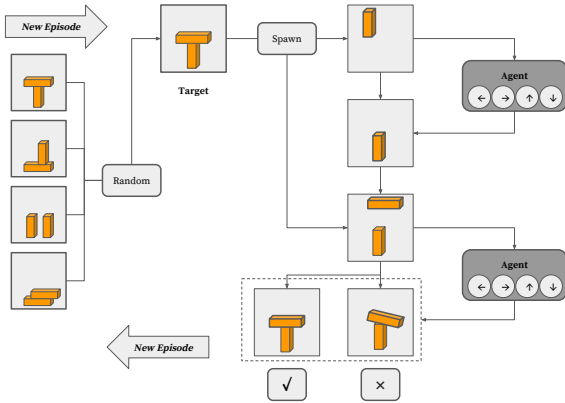
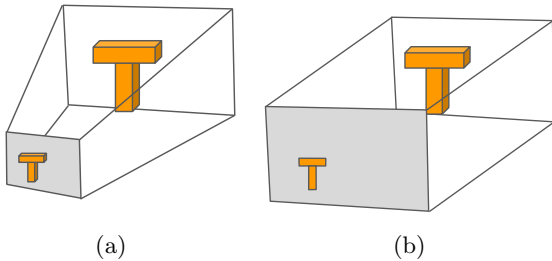Figure 13: Overview of design for the target stacking environment.



(a)                          (b)

Figure 14: Different perspective lens used in this thesis.(a): Perspective projection. (b): Orthographic projection.

parallel lines stay parallel and don't converge as shown in Figure 14b. This is in contrast with regular perspective camera used in the visual stability prediction environment as shown in Figure 14a.

Additionally, as we consider the action in discrete space, namely {left,right,down} for the task, the state representation can also be formalized correspondingly as shown in Figure 15. The individual block size follows a ratio of $l : w : h = 5 : 2 : 1$, where $l,w,h$ denote length, width and height respectively. With the aforementioned orthographic camera, only length and height are preserved on the projected image with a ratio of $l : h = 5 : 1$. Each action moves the block by a displacement of the block's height along its direction. Hence the state can be represented in a grid map of unit square cells with length of the block's height ($u$ in the image space). In actual im-

plementation, the blocks have to be set with coordinates in continuous values, so a separate binary matrix is used to maintain the discrete representation of the state where 1 denotes the block presence in the cell and 0 otherwise.
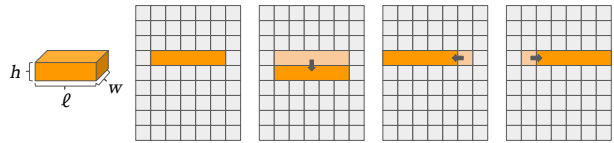


Figure 15: State representation in the stacking environment

The environment can query both the rendered image from the camera and the grid map for the state's representation, though in our current experiments, we use mostly the grid map for a simplified representation. Another benefit from the grid map representation is that it can be used as a occupancy map which is crucial for our implemented collision detection mechanism discussed in the next subsection.

### 5.3.2   Collision and Stability Detection

As the block is maneuvered in the environment, it will eventually make contact with the scene. The occurrence of contact marks an important transition of the state for the environment. As shown in Figure 16:
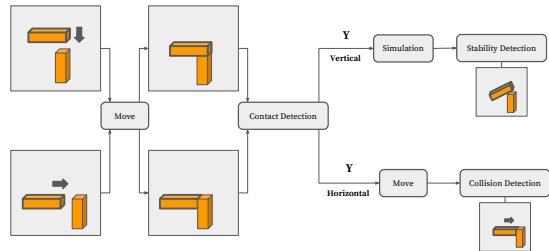


Figure 16: Process of collision and stability detection

If the contact is made by a block from above (vertical), it suggests a placement [4], then physics simulation is called to obtain the stability. The stability detection is implemented the same way as the one used in the data

---

[4]In our environment, we simplify the release movement from the real world block placement, i.e. the block is deemed to be released once it makes vertical contact with the scene, either the ground or the structure in the scene.

generator for visual stability prediction. Here the simulation only activates upon vertical contact is detected to save runtime. If the block remains stable in the scene and is not the last block for the episode, a new block is then spawned at the top of the scene with random horizontal location. If the block is the final one and remains stable or it collapses the structure in the scene, then the current episode is terminated and a new episode can be started.

If the contact is made by a horizontal movement, it will not directly affect the current episode. If an additional move causes the collision, then the episode will be terminated, otherwise the block can be moved further until either vertical contact happens, it goes into the simulation and stability-detection branch or collision and terminated.

Note if the block is about to move out of the view of the camera, it is also considered a sub-type of collision. This is implemented straightforwardly by comparing the block's incoming location with the boundary of the scene.

The contact and collision detection are implemented with the grid map representation described in the previous subsection as shown in Figure 17. We keep two grid map for each episode, one for the existing structure in the scene (**background map**) and the other for the moving block (**foreground map**). A state (grid map) is internally represented as the summation of the background map and foreground map. Given a state and a specified action, the system will compute the incoming foreground map after the action, if the incoming foreground map overlaps with the background map, then collision is detected as shown in Figure 17b; if any of the adjacent cells below the block in the incoming foreground map is non-empty in the incoming grid map, then a vertical contact is detected as shown in Figure 17a.

### 5.3.3   Interface Design

With recent progress in reinforcement learning, there is an increasing need for the standardized benchmarks for different learning agents. To this end, OpenAI recently introduced Gym (Brockman et al., 2016) to provide access to a standardized set of environments, such as Atari games and board games. By encapsulating the environment's internal mechanism from the agent, different agents can be evaluated on the environment with minimal modification.

We also adopt a similar design in our environment. The diagram in Figure 18 shows part of the design of the environment with interactions with an agent. The agent's `act` method takes the current state from the en-
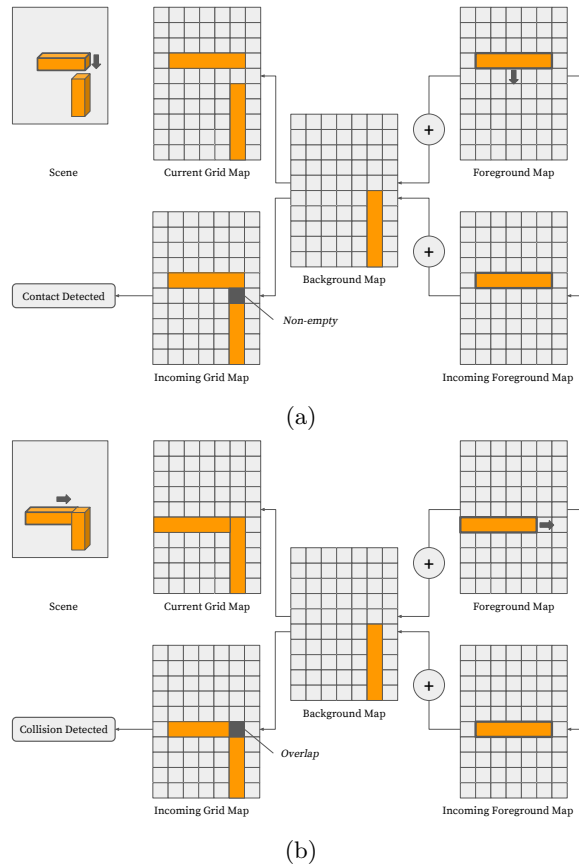


(a)



(b)

Figure 17: Internal representation to detect contact and collision conditions: (a) Detect vertical contact. (b) Detect collision.
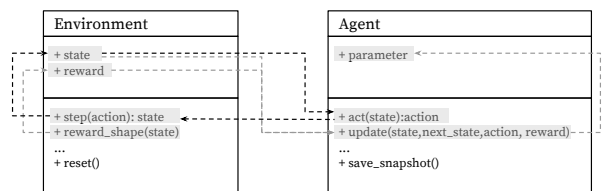


Figure 18: Interface design for the environment with interaction of the agent.

vironment, decides its action and pass the decision to the environment. The `step` method receives the action from the agent and updates the environment's state. The `reward_shape` method evaluates the current state and transforms the reward with specific designs. Then the pair of states before and after the action, together with the reward, are input to the `update` method of the agent to update its parameter depending on different learning algorithms.

## 5.4 Goal-Parameterized Deep Q Networks (GDQN)

As one major characteristic of this task is that it requires goal-specific planning: given the same or similar states under different objectives, the optimal move can be different. To this end, we extend the typical reinforcement learning formulation to incorporate additional goal information.

### 5.4.1 Learning Framework

In a typical reinforcement learning setting, the agent interacts with the environment at time $t$, observes the state $s_t$, takes action $a_t$, receives reward $r_t$ and transits to a new state $s_{t+1}$. A common goal for a reinforcement learning agent is to maximize the cumulative reward. This is commonly formalized in form of a value function as the expected sum of rewards from a state $s$, $\mathbf{E}[\sum\limits_{i=0}^{T} \gamma^i r_{t+i+1}|s_t = s, \pi]$ when actions are taken with respect to a policy $\pi(a|s)$, with $0 \leq \gamma \leq 1$ being the discount factor, T for the final time step. The alternative formulation to this is the action-value function $Q^\pi(s,a) = \mathbf{E}[\sum\limits_{i=0}^{T} \gamma^i r_{t+i+1}|s_t = s, a_t = a]$.

Value-based reinforcement learning algorithms, such as Q-learning (Watkins and Dayan, 1992) directly search for optimal Q-value function. Recently by incorporating deep neural network as a function approximator for $Q$-function, the DQN (Mnih et al., 2015) has shown impressive results across a variety of Atari games.

**DQN** For our task, we apply a *Deep Q Network* (DQN) which uses a deep neural network for approximating the action-value function $Q(s, a; \theta)$, mapping from an input state $s$ and action $a$ to Q values. In particular, two important improvements have been proposed by Mnih et al. (2015) for the learning process, including (1) experience

replay, the agent stores observed transitions in a memory buffer for some time, and uniformly samples from the memory to update the network (2) the target network, agent maintains two networks for the loss function — one for the current estimator of Q function and one for the surrogate of the true Q function. For the current estimator, the parameters are constantly updated. For the surrogate, the parameters are only updated for every certain number of steps from the current estimator network otherwise kept fixed.

**Learning Goal-Parameterized Policies** To plan with respect to the specific goal, we can parametrize the agent's policy $\pi$ by the goal $g$:

$$\pi(s, g, a) \qquad (1)$$

Since in this work, we applies DQN as value-based method, this corresponds to the update to original Q function with the additional goal information. The new Q-value function is hence defined as:

$$Q^\pi(s, g, a) = \mathbf{E}[\sum_{i=0}^{T} \gamma^i r_{t+i+1}|s_t = s, g, a_t = a] \qquad (2)$$

As shown in Figure 19, in contrast to the original DQN model, where state and action are used to estimate Q-value, the new model further include the current goal into the network to produce the estimate. We call this model as Goal-Parametrized Q Network (GDQN).

The resulted loss function is as:

$$L_Q = \mathbf{E}[(R + \gamma \max'_a Q^\pi(s', g, a'; \theta^-) - Q(s, g, a; \theta))^2] \qquad (3)$$

where $\theta^-$ are the previous parameters and the optimization is with respect to $\theta$.

### 5.4.2 Implementation Details

The DQN agent is implemented in Theano and Keras to adapt to the settings in our experiment, while we use a 2 hidden layer (each with 64 hidden units and rectified linear activation) multilayer perceptron (MLP) for most cases, we additionally swap the MLP with the CNN and follow the reported parameter settings as in the original paper (Mnih et al., 2015) to ensure our implementation can reach similar performance.

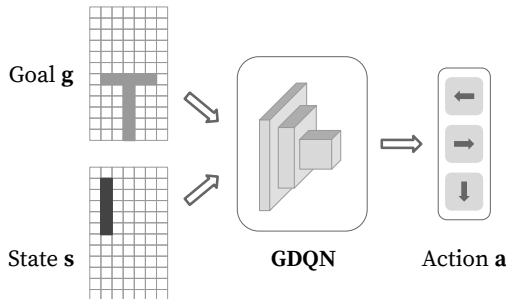Note we don't apply the frame-skipping technique (Bellemare et al., 2012) used for Atari games (Mnih et al.,

Figure 19: Our proposed model GDQN which extends the $Q$-function approximator to integrate goal information.

2015) allowing the agent sees and selects actions on every $k$th frame where its last action is repeated on skipped frames. It does not suit our task, in particular when the moving block is getting close to the existing structure, simply repeating action decided from previous frame can cause unintended collision or collapse.

**Reward**  In the target stacking task, the agent gets reward $+1$ when the episode ends with complete reproduction of the target structure, otherwise 0 reward.

Further, we explore reward shaping (Ng et al., 1999) in the task providing more prompt intermediate reward. Two types of reward shaping are included: overlap ratio and distance transform.

For the overlap ratio, for each state $s_t$ under the same target $g_i$, an overlap ratio is counted as the ratio between the intersected foreground region (of the current state and the target state) and the target foreground region (shown in Figure 20a):

$$o(s_t, g_i) = \frac{s_t \cap g_i}{g_i} \qquad (4)$$

For each transition $(s_t, a_t, s_{t+1})$, the reward is defined by the change of overlap ratio before and after the action:

$$r_t = \begin{cases} 1, & \text{if } \Delta o_{t \to t+1} = o(s_{t+1}) - o(s_t) > 0 \\ -1, & \text{if } \Delta o_{t \to t+1} = o(s_{t+1}) - o(s_t) < 0 \\ 0, & \text{otherwise} \end{cases} \qquad (5)$$

The intuition is that actions increasing the current state to become more overlapped with the target scene should be encouraged.
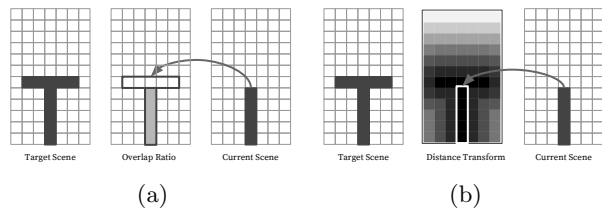


Figure 20: Reward shaping used in target stacking. (a): overlap ratio to the target. The gray area in the middle figure denotes the intersected foreground region between current and target scene, and the overlap ratio is the ratio between the areas of the two. (b): distance under the distance transform of the target. The middle figure denotes the distance transform under the target shown in the left. The distance from current scene to the target is the sum of distances masked by the current scene in the distance transform.

For the distance transform (Fabbri et al., 2008), it generates a map $D$ whose value in each pixel $p$ is the smallest distance from it to a target object $O$:

$$D(p) = \min\{\text{dist}(p, q) | q \in O\} \qquad (6)$$

where dist can be any valid distance metric, like Euclidean or Manhattan distance.

For each state $s_t$ under the same target $g_i$, a distance to the goal is the sum of all the element-wise distance in $s_t$ to $g_i$ under $D_{g_i}$ (shown in Figure 20b) as:

$$d(s_t, g_i) = \sum_j D_{g_i}(s_t^j), s_t^j \in s_t \qquad (7)$$

For each transition $(s_t, a_t, s_{t+1})$, the reward is defined as:

$$r_t = \begin{cases} 1, & \text{if } \Delta d_{t \to t+1} = d(s_{t+1}) - d(s_t) < 0 \\ -1, & \text{if } \Delta d_{t \to t+1} = d(s_{t+1}) - d(s_t) > 0 \\ 0, & \text{otherwise} \end{cases} \qquad (8)$$

The intuition behind this is that action decreasing the distance between the current state and the target scene should be encouraged.

### 5.4.3  Toy Example with Goal Integration

At first, we introduce a type of navigation task in the classic gridworld environment. As shown in Figure 21a, for
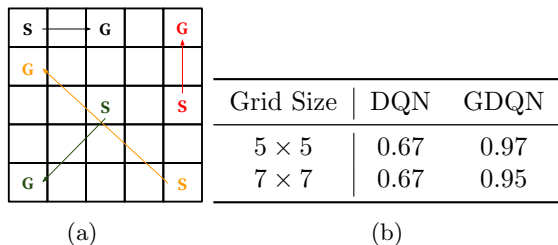
| Grid Size | DQN | GDQN |
|---|---|---|
| $5 \times 5$ | 0.67 | 0.97 |
| $7 \times 7$ | 0.67 | 0.95 |

(a) (b)

Figure 21: a: Navigation task, each color denotes a different episode, per episode, a random pair of starting and goal location are generated, the agent needs to reach the goal. b: Results from navigation task.

each episode, we generate a random start and goal point in the grid world whereas the start and the goal have to be different from each other, the agent has to find the way to this randomly generated goal. The agent is rewarded only when it reaches it. The agent needs to reach the goal with four possible actions {left, right, up, down}. Action that will make the agent go off the grid will leave it stay in the same location. The episode only terminates once the agent reaches the goal. The agent only receive reward +1 when reaching the current goal. Two different sizes of gridworld are tested at $5 \times 5$ and $7 \times 7$. We learn a single model for each size of the mazes.

The training epoch size is 1000 in steps for the smaller gridworld and 3000 for the larger one, the test sizes are the same for both at 100. All the agents run for 100 epochs and the $\epsilon$ for $\epsilon$-greedy anneals linearly from 1.0 to 0.1 over the first 20 epochs, and fixed at 0.1 thereafter. The memory buffer size is set the same to the annealing length, i.e. for the smaller gridworld, the buffer size equals to the length 20 epochs in training with 20000 steps whereas for the larger one, the buffer size is 30000 steps. We measure the proportion of episodes in the test epoch that reaches the goal in shortest distance as the success ratio. The results are shown in Table 21b for the best agents throughout the training process.

As in this simple task with relative small state space, DQN gets some performance due to running an average policy across all the the goals, but this is not addressing the task we set out to do. In contrast, GDQN parametrized specifically to include goal information achieves significant better results on both sizes of the environment.
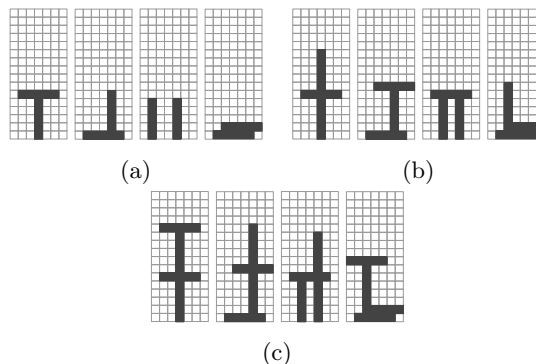


Figure 22: a: Targets for 2 blocks. b: Targets for 3 blocks. c: Targets for 4 blocks.

### 5.4.4 Target Stacking

We set up 3 groups of target structures consisted of different number of blocks $\{2, 3, 4\}$ in the scene as shown in Figure 22. Within each group of target shapes, a random target (with the accompanied orientation order) is first picked at the very beginning for individual episode and then each block is spawned at randomized location as described earlier in the target stacking environment. Each training epoch consists of 10000 steps and each test epoch with 1000 steps. Similar to the setting in the toy example, all the agents run for 100 epochs and the $\epsilon$ anneals for the first 20 epochs, and the memory buffer size is set as long as the annealing steps at $200K$ steps.

We computed both average overlap ratio (OR) and success rate (SR) for the finished stacking episodes in each test epoch. Here overlap ratio is the same as defined in the reward shaping in Equation 4, but simply measures the end scene over the assigned target scene. This tells the relative completion of the stacked structure in comparison to the assigned target structure, the higher the value is, the better completion it is to the target. At the maximum of 1, it suggests completely reproduction of the target. The success rate counts the ratio how many episodes complete the exact same shape as assigned over the total number of episodes finished in the test epoch. This is the absolute metric counting overall successful stacking. The results are shown in Table 7 for the best agents throughout the training process.

Over all groups on both metrics, we observe GDQN outperforms DQN, showing the importance of integrating goal information. In general, the more blocks in the task, the more difficult it becomes. When there are only

| Num. of Blks. | DQN | | GDQN | | GDQN + OR | | GDQN + DT | |
|---|---|---|---|---|---|---|---|---|
| | OR | SR | OR | SR | OR | SR | OR | SR |
| 2 | 0.70 | 0.70 | 0.82 | 0.72 | 0.84 | 0.77 | **0.88** | **0.78** |
| 3 | 0.43 | 0.43 | 0.76 | **0.67** | **0.86** | 0.63 | 0.83 | 0.65 |
| 4 | 0.03 | 0.0 | 0.41 | 0.17 | 0.73 | 0.55 | **0.79** | **0.56** |

Table 7: Results for target stacking. For "GDQN + X", X denotes different ways for reward shaping as described in previous section, OR for overlap ratio, DT for distance transform. For metrics, OR stands for average overlap ratio, SR for average success rate.



Figure 23: Example scenes constructed by the agent.

small number of blocks (2 blocks and 3 blocks) in the scene, the single policy learned by DQN averages over the few target shapes can still work to some extent. However when introducing more blocks into the scene, it becomes more and more difficult for this averaged model to handle. As we can see from the result, there is already a significant decrease of performance (success rate drops from 0.70 to 0.43) when increasing the blocks number from 2 to 3, whereas GDQN's performance only decreases slightly from 0.72 to 0.67. In 4 blocks scene, the DQN can no longer reproduce any single target (0.0 for success rate, 0.03 for overlap ratio) while GDQN parametrized specifically to include goal information can still do. Though the success rate (absolute completion to the target) for the basic GDQN is relatively low at 0.17 but the average overlap ratio (relative completion to the target) still holds up pretty well at 0.41. Also we see reward shaping can further improves GDQN model, in particular distance transform can boost the performance more than overlap ratio.

# 6    Conclusion

In this work, we answer the question if and how well we can build up a machinery to learn a manipulation task under physics constraints with visual perception. In particular, We focus on a wood block stacking problem and explored frameworks from two inter-related but different perspectives.

In the first part of the work, we aim to guide the robot to stably place a single block in the scene. To achieve this, we envision a mechanism of explicit physics understanding to predict physical stability directly from visual input, bypassing explicit 3D representations and physical simulation. We initially evaluate our model on a synthetic dataset, covering a range of conditions including variations in number of blocks, size of blocks and 3D structure of the overall tower. The results reflect the challenges of inference with growing complexity of the structure. To further understand the results, we conduct a human subject study on a subset of our synthetic data and show that our model achieves comparable or even better results than humans in the same setting. Moreover, we investigate the discriminative image regions found by the model and spot correlation between such regions and initial collapse area in the structure. Finally, We apply our approach to a block stacking setting and show that our model successfully guide a robot for placements of new blocks by predicting the stability of future states.

In the following part of the work, we explore further on a more challenging target stacking task where the agent stacks blocks to reproduces a tower shown in an image. It presents a distinct type of challenge requiring the agent to reach a given goal state that may be different for every new trial. To this end, we propose an alternative framework where the manipulation is learned end-to-end through pure trial and error, bypassing the explicitly modeling for the corresponding physics knowledge as we did in the previous task. We create a synthetic block stacking environment with physics simulation in which the agent can interactively stack the block. In particular, we propose a goal-parametrized GDQN model to plan with respect to the specific goal, allowing better generalization across different goals. We validate the model on both a navigation task in a classic gridworld environment with different start and goal positions and the block stacking task itself with different target structures.

# References

Agrawal P, Nair AV, Abbeel P, Malik J and Levine S (2016) Learning to poke by poking: Experiential learning of intuitive physics. In: *Advances in Neural Information Processing Systems.* pp. 5074–5082.

Baillargeon R (1994) How do infants learn about the physical world? *Current Directions in Psychological Science* .

Baillargeon R (2002) The acquisition of physical knowledge in infancy: A summary in eight lessons. *Blackwell handbook of childhood cognitive development* .

Baillargeon R (2008) Innate ideas revisited: For a principle of persistence in infants' physical reasoning. *Perspectives on Psychological Science* .

Battaglia P, Ullman T, Tenenbaum J, Sanborn A, Forbus K, Gerstenberg T and Lagnado D (2012) Computational models of intuitive physics. In: *Proceedings of the Cognitive Science Society*, volume 34.

Battaglia PW, Hamrick JB and Tenenbaum JB (2013) Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences* .

Bellemare MG, Naddaf Y, Veness J and Bowling M (2013) The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research* 47: 253–279.

Bellemare MG, Veness J and Bowling M (2012) Investigating contingency awareness using atari 2600 games. In: *AAAI*.

Bhattacharyya A, Malinowski M, Schiele B and Fritz M (2018) Long-term image boundary extrapolation. In: *AAAI*.

Brockman G, Cheung V, Pettersson L, Schneider J, Schulman J, Tang J and Zaremba W (2016) Openai gym. *arXiv preprint arXiv:1606.01540* .

Cholewiak SA, Fleming RW and Singh M (2013) Visual perception of the physical stability of asymmetric three-dimensional objects. *Journal of vision* .

Cholewiak SA, Fleming RW and Singh M (2015) Perception of physical stability and center of mass of 3-d objects. *Journal of vision* .

Coumans E (2010) Bullet physics engine. *Open Source Software: http://bulletphysics. org* 1.

Deng J, Dong W, Socher R, Li LJ, Li K and Fei-Fei L (2009) Imagenet: A large-scale hierarchical image database. In: *CVPR*.

Denil M, Agrawal P, Kulkarni TD, Erez T, Battaglia P and de Freitas N (2017) Learning to perform physics experiments via deep reinforcement learning. In: *ICLR*.

Dosovitskiy A and Koltun V (2017) Learning to act by predicting the future. In: *ICLR*.

Fabbri R, Costa LDF, Torelli JC and Bruno OM (2008) 2d euclidean distance transform algorithms: A comparative survey. *ACM Computing Surveys (CSUR)* 40(1): 2.

Fragkiadaki K, Agrawal P, Levine S and Malik J (2016) Learning visual predictive models of physics for playing billiards. In: *ICLR*.

Goslin M and Mine MR (2004) The panda3d graphics engine. *Computer* 37(10): 112–114.

Gupta A, Efros AA and Hebert M (2010) Blocks world revisited: Image understanding using qualitative geometry and mechanics. In: *ECCV*.

Hamrick J, Battaglia P and Tenenbaum JB (2011) Internal physics models guide probabilistic judgments about object dynamics. In: *Proceedings of the 33rd annual conference of the cognitive science society*. Cognitive Science Society Austin, TX.

Hamrick JB, Ballard AJ, Pascanu R, Vinyals O, Heess N and Battaglia PW (2017) Metacontrol for adaptive imagination-based optimization. In: *ICLR*.

Jia Y, Shelhamer E, Donahue J, Karayev S, Long J, Girshick R, Guadarrama S and Darrell T (2014) Caffe: Convolutional architecture for fast feature embedding. In: *Proceedings of the ACM International Conference on Multimedia*. ACM.

Kaelbling LP (1993) Learning to achieve goals. In: *IJCAI*. pp. 1094–1099.

Kempka M, Wydmuch M, Runc G, Toczek J and Jaśkowski W (2016) Vizdoom: A doom-based ai research platform for visual reinforcement learning. In: *Computational Intelligence and Games (CIG)*. IEEE, pp. 1–8.

Kimura S, Watanabe T and Aiyama Y (2010) Force based manipulation of jenga blocks. In: *IROS*.

Krizhevsky A, Sutskever I and Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: *NIPS*.

Kröger T, Finkemeyer B, Winkelbach S, Molkenstruck S, Eble L and Wahl FM (2006) Demonstration of multisensor integration in industrial manipulation (poster). In: *ICRA*.

LeCun Y, Jackel L, Bottou L, Brunot A, Cortes C, Denker J, Drucker H, Guyon I, Muller U, Sackinger E, Simard P and Vapnik V (1995) Comparison of learning algorithms for handwritten digit recognition. In: *International conference on artificial neural networks*.

Lerer A, Gross S and Fergus R (2016) Learning physical intuition of block towers by example. In: *ICML*.

Li W and Fritz M (2012) Recognizing materials from virtual examples. In: *ECCV*.

Li W, Leonardis A and Fritz M (2017) Visual stability prediction for robotic manipulation. In: *ICRA*. IEEE.

McCloskey M (1983) Intuitive physics. *Scientific american* .

Mnih V, Kavukcuoglu K, Silver D, Graves A, Antonoglou I, Wierstra D and Riedmiller M (2013) Playing atari with deep reinforcement learning. In: *NIPS Deep Learning Workshop*.

Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller MA, Fidjeland A, Ostrovski G, Petersen S, Beattie C, Sadik A, Antonoglou I, King H, Kumaran D, Wierstra D, Legg S and Hassabis D (2015) Human-level control through deep reinforcement learning. *Nature* .

Mottaghi R, Bagherinezhad H, Rastegari M and Farhadi A (2016) Newtonian image understanding: Unfolding the dynamics of objects in static images. In: *CVPR*.

Mottaghi R, Schenck C, Fox D and Farhadi A (2017) See the glass half full: Reasoning about liquid containers, their volume and content. In: *ICCV*.

Ng AY, Harada D and Russell S (1999) Policy invariance under reward transformations: Theory and application to reward shaping. In: *ICML*, volume 99. pp. 278–287.

Oh J, Singh S, Lee H and Kohli P (2017) Zero-shot task generalization with multi-task deep reinforcement learning. In: *ICML*.

Pascanu R, Li Y, Vinyals O, Heess N, Buesing L, Racanière S, Reichert D, Weber T, Wierstra D and Battaglia P (2017) Learning model-based planning from scratch. *arXiv preprint arXiv:1707.06170* .

Peng X, Sun B, Ali K and Saenko K (2015) Learning deep object detectors from 3d models. In: *ICCV*.

Razavian A, Azizpour H, Sullivan J and Carlsson S (2014) Cnn features off-the-shelf: an astounding baseline for recognition. In: *CVPR workshops*.

Rematas K, Ritschel T, Fritz M, Gavves E and Tuytelaars T (2016) Deep reflectance maps. In: *CVPR*.

Rematas K, Ritschel T, Fritz M and Tuytelaars T (2014) Image-based synthesis and re-synthesis of viewpoints guided by 3d models. In: *CVPR*.

Schaul T, Horgan D, Gregor K and Silver D (2015) Universal value function approximators. In: *ICML*.

Silberman N, Hoiem D, Kohli P and Fergus R (2012) Indoor segmentation and support inference from rgbd images. In: *ECCV*.

Simonyan K and Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* .

Smith B and Casati R (1994) *Naive Physics: An Essay in Ontology*. Philosophical Psychology.

Smith R (2005) Open dynamics engine .

Toussaint M, Plath N, Lang T and Jetchev N (2010) Integrated motor control, planning, grasping and high-level reasoning in a blocks world using probabilistic inference. In: *ICRA*. IEEE, pp. 385–391.

Wang J, Rogers P, Parker L, Brooks D and Stilman M (2009) Robot jenga: Autonomous and strategic block extraction. In: *IROS*.

Watkins CJ and Dayan P (1992) Q-learning. *Machine learning* 8(3-4): 279–292.

Weber T, Racanière S, Reichert DP, Buesing L, Guez A, Rezende DJ, Badia AP, Vinyals O, Heess N, Li Y, Pascanu R, Battaglia P, Silver D and Wierstra D (2017) Imagination-augmented agents for deep reinforcement learning. *arXiv preprint arXiv:1707.06203* .

Weinzaepfel P, Revaud J, Harchaoui Z and Schmid C (2013) Deepflow: Large displacement optical flow with deep matching. In: *ICCV*.

Wu J, Yildirim I, Lim JJ, Freeman B and Tenenbaum J (2015) Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In: *NIPS*.

Yildirim I, Gerstenberg T, Saeed B, Toussaint M and Tenenbaum J (2017) Physical problem solving: Joint planning with symbolic, geometric, and dynamic constraints. In: *Proceedings of the Thirty-Night Annual Conference of the Cognitive Science Society.*

Zheng B, Zhao Y, Yu J, Ikeuchi K and Zhu SC (2013) Beyond point clouds: Scene understanding by reasoning geometry and physics. In: *CVPR*.

Zhou B, Khosla A, A L, Oliva A and Torralba A (2016) Learning Deep Features for Discriminative Localization. *CVPR* .